

# Gait Library Synthesis for Quadruped Robots via Augmented Random Search

Sashank Tirumala, Aditya Sagi, Kartik Paigwar, Ashish Joglekar,  
Shalabh Bhatnagar, Ashitava Ghosal, Bharadwaj Amrutur, Shishir Kolathaya

**Abstract**—In this paper, with a view toward fast deployment of learned locomotion gaits in low-cost hardware, we generate a library of walking trajectories, namely, forward trot, backward trot, side-step, and turn in our custom built quadruped robot, Stoch 2, using reinforcement learning. There are existing approaches that determine optimal policies for each time step, whereas we determine an optimal policy, in the form of end-foot trajectories, for each half walking step i.e., swing phase and stance phase. The way-points for the foot trajectories are obtained from a linear policy, i.e., a linear function of the states of the robot, and cubic splines are used to interpolate between these points. Augmented Random Search, a model-free and gradient-free learning algorithm, is used to learn the policy in simulation. This learned policy is then deployed on hardware, yielding a trajectory in every half walking step. Different locomotion patterns are learned in simulation by enforcing a preconfigured phase shift between the trajectories of different legs. Transition from one gait to another is achieved by using a low-pass filter for the phase, and the sim-to-real transfer is improved by a linear transformation of the states obtained through regression.

**Keywords:** *Quadrupedal walking, Reinforcement Learning, Random Search*

## I. INTRODUCTION

Quadrupedal locomotion with multiple types of gaits has been successfully demonstrated via numerous techniques— inverted pendulum models [1], zero-moment point [2], hybrid zero dynamics [3], and central pattern generators (CPG) [4]—to name a few. However, due to the increase in computational resources, data driven approaches like reinforcement learning (RL) [5] are gaining popularity today. From a user’s point of view, RL only requires the determination of a scalar reward, like distance travelled, energy consumed etc, and then the algorithm identifies the best controller for walking by itself. The success of this methodology was enabled by the use of deep neural networks in RL, popularly known as deep reinforcement learning (D-RL) [6], [7].

The considerable success shown by D-RL are mostly for simulated robotic tasks [5] even today. Translating these results in real hardware are faced with numerous challenges.

\*This work is supported by DST INSPIRE Fellowship IFA17-ENG212, and Robert Bosch Center for Cyber-Physical Systems, IISc, Bengaluru.

S. Tirumala is with the Department of Engineering Design at the Indian Institute of Technology-Madras, Chennai, A. Sagi, K. Paigwar and A. Joglekar are with the Robert Bosch Center for Cyber-Physical Systems at the Indian Institute of Science, Bangalore

S. Kolathaya is an INSPIRE Faculty Fellow and S. Bhatnagar, A. Ghosal, B. Amrutur are with the faculty of the Robert Bosch Center for Cyber-Physical Systems, Indian Institute of Science, Bengaluru, India. email: {shishirk, shalabh, asitava, amrutur} at iisc.ac.in

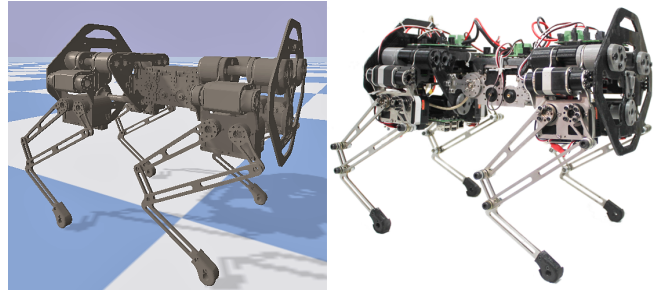


Fig. 1: Figure showing the custom built quadruped robot, Stoch 2. Simulated version is shown in the left, and the actual hardware is shown in the right.

The use of deep neural networks (DNN) entails a large computational overhead to obtain an inference. This not only requires expensive hardware, but also results in higher power consumption, which may not be economical for a large commercial deployment of robots. In addition, D-RL based techniques require a significant number of training iterations (to the order of millions), which are detrimental to hardware. A possible alternative would be to train in a simulated environment, and then deploy on hardware. This was demonstrated by [6], [7], [8] for quadrupeds and bipeds, where the neural networks used sensory feedback to control the motor joints in real-time. However, this type of sim-to-real transfer presents with itself several challenges. Some of these challenges are mainly attributed to the modeling uncertainty, sensory noise, control saturations etc.

It is worth noting that some of the solutions to the above presented challenges were addressed in part via the traditional RL techniques [9], [10], [11], [12]. In fact, RL was used as an optimization tool (gait generating tool) by most of the teams participating in the quadruped RoboCup soccer league in 2004 as stated by [13]. These techniques have been used to obtain the parameterized foot trajectories that can be deployed on the robot. This can be implemented in low cost hardware and requires no extra computing power. However, these works implemented open-loop controllers to playback the trajectories on the robots, thereby, limiting their capability to handle external disturbances due to the lack of feedback in the system.

Similar to the methodologies followed in the past in RL, we would like to realize locomotion on a low-cost embedded platform and, at the same time, incorporate some form of feedback that will help improve robustness in real-time.

By limiting our controller to a linear policy, i.e., a single matrix with no activation functions, our goal is to learn this policy with the modern tools available for learning. We chose to use a model-free and gradient-free algorithm, Augmented Random Search (ARS), which was shown to be at least as sample efficient as DNN based policy gradient algorithms [14] existing in literature. We present a control architecture that combines a trajectory-based framework with RL to obtain a controller that is capable of handling external disturbances and exhibiting multiple complex quadruped locomotion behaviours while requiring little manual tuning.

### A. Related work on ARS

ARS was successfully demonstrated on a variety of robotic control tasks (including quadrupedal locomotion) with a linear policy [14]. These types of policies have also been deployed successfully on the quadruped robot Minitaur [15]. In particular, [15] used a policy modulating trajectory generator (PMTG) that modulated a trajectory and simultaneously modified the outputs to obtain the final commands for the motors. This was successfully used with linear policies to achieve walking and bounding gaits. [16] extended the PMTG framework to include turning gaits in Minitaur. However, the turning demonstrated by Minitaur was limited due to the absence of abduction motors. In addition, the PMTG framework updates the angle commands after every time step. Therefore, with a view toward lowering the computational overhead in hardware, and including abduction based gaits like side-stepping, in-place turning, we make a slight deviation from this approach, by directly obtaining the end-foot trajectories from a trained policy. The policy is updated every half-step, and the motors are commanded to track the resulting trajectories. We demonstrate a library of gaits forward trot, backward trot, side step and turn, in the custom built quadruped robot Stoch 2. Towards the end we demonstrate robustness by showing that the linear policy is able to reject disturbances such as pushing and pressing.

### B. Organization

The paper is structured as follows: Section II will provide a brief background on reinforcement learning (RL) along with a description of the robot hardware and the associated control framework, Section III will describe the methodology used for training followed by the experimental results in Section IV.

## II. ROBOT MODEL AND CONTROL

In this section, we will discuss the reinforcement learning framework used for our custom built quadruped robot Stoch 2. Specifically, we will provide details about the hardware, the associated model, and the trajectory tracking control framework used for the various gaits of the robot.

### A. Hardware description of Stoch 2

*Stoch 2* is a quadruped robot designed and developed in-house at the Indian Institute of Science (IISc), Bengaluru, India. It is a second generation robot in the *Stoch* series [17],

[18]. It is similar to *Stoch* [18] in form factor, and weighs approximately 4kg. Each leg contains three joints—hip flexion/extension, hip abduction and knee flexion/extension. Servo motors from Kondo Kogaku (model: KRS6003) are used to actuate these joints. The main controller on the robot is an ARM processor (the Tiva TM4C123GH6PM) running at 80 MHz. This controller is placed on the central module that connects the front and back parts of *Stoch 2*. The URDF model used in the simulator was created directly from the SolidWorks assembly (see Fig. 1). Overall, the robot simulation model consists of 6 base degrees-of-freedom and 12 actuated degrees-of-freedom.

### B. Kinematic description

Each leg comprises of a parallel five-bar linkage mechanism where two of the links are actuated as shown in Fig. 2. This enables the end foot point to follow the given trajectory in a plane. The two actuators which control the motion of upper hip and knee linkages are mounted on a fixed link. These actuated linkages in turn connect to the lower linkages via revolute joints.

In this paper, we focus on realizing trajectories of the feet in polar coordinates (see Section II-F ahead). As seen from Fig 2, the five-bar linkage is divided into two serial 2-R linkages and solved for each branch. The details of the equations for a serial 2-R linkage can be found in [18].

### C. Reinforcement learning for walking

We formulate the problem of locomotion as a Markov Decision Process (MDP). Here an MDP is a 5-tuple  $\{S, A, r, P, \gamma\}$  where  $S \subset \mathbb{R}^n$  refers to set of states of the robot,  $A \subset \mathbb{R}^m$  refers to the set of actions,  $r : S \times A \rightarrow \mathbb{R}$  refers to the reward received for every  $(S, A)$  pair,  $P : S \times A \times S \rightarrow [0, 1]$  refers to the transition probabilities between two states for a given action, and  $\gamma \in (0, 1)$  is the discount factor of the MDP. Description of the states, actions for *Stoch 2* are explained later on in this section. Rewards are chosen depending upon the gait, which are explained in Section III-B. Given a policy  $\pi : S \rightarrow A$ , we evaluate it for an episode. In this formulation, the optimal policy is the policy that maximizes the return ( $R$ ):

$$R = \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots], \quad (1)$$

where the subscript for  $r_t$  denotes the step index. Since walking is periodic in nature, each step here corresponds to one cycle/loop of the foot trajectory. Therefore, our policy is evaluated two times for every gait cycle/loop of the robot. Henceforth, we will refer to this half-cycle by “gait step”. We will describe the trajectories that yield these gait steps next.

### D. Cubic splines

A cubic Hermite spline is a piece-wise third degree polynomial, which not only interpolates the values between two points, but also their derivatives. The equation of a Hermite

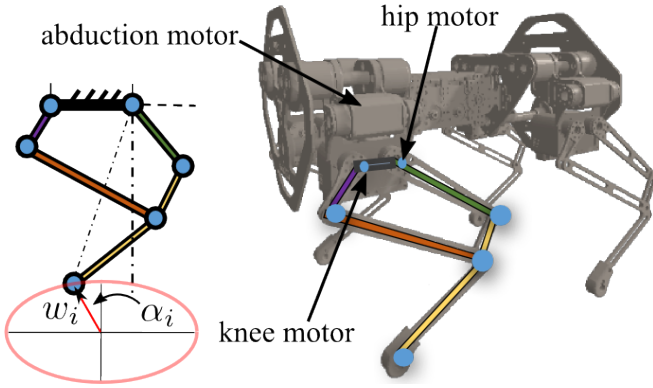


Fig. 2: A five-bar mechanism is used as the legs of the quadruped robot. This mechanism is actuated by the motors located at the main torso of the robot.

spline specified by its endpoints  $\{w_0, w_1\}$  and tangents  $\{w'_0, w'_1\}$  is

$$w(t) = \begin{bmatrix} 2t^3 - 3t^2 + 1 \\ -2t^3 + 3t^2 \\ t^3 - 2t^2 + t \\ t^3 - t^2 \end{bmatrix}^T \times \begin{bmatrix} w_0 \\ w_1 \\ w'_0 \\ w'_1 \end{bmatrix}, \quad (2)$$

where  $t$  is the phasing variable. Therefore, given  $n$  points  $(w_0, t_0), (w_1, t_1), \dots, (w_{n-1}, t_{n-1})$ , we can connect the adjacent pair of points resulting in a loop. Since the tangents are user definable, we choose them to be:

$$w'_i = \frac{w_{i+1} - w_{i-1}}{t_{i+1} - t_{i-1}}, \quad (3)$$

and if  $i = 0$ , we choose  $w_{-1} = w_{n-1}$ . It is worth noting that depending upon the gait, the plane containing these points is oriented appropriately w.r.t. the foot, and the corresponding joints angle trajectories are obtained via inverse kinematics. See Fig. 3 and Section III-B for more details.

### E. State space

The state is represented by motor angles of the robot. Hence our state space is a 12 dimensional vector space:  $S \subset \mathbb{R}^{12}$ . Note that other sensory values such as angular velocities, accelerations and currents were ignored due to the fact that the policy was updated after every gait step.

### F. Action space

In reference [19], the actuator joint angles were selected for the action space. In reference [17], the legs' end-point positions in polar coordinates were selected for the action space. The polar coordinates for the four legs collectively provided an 8-dimensional action space. In this paper, since we are using splines to determine the end-foot trajectories, we require the actions to yield the control points for these splines (see Fig. 3). We express the splines in polar coordinates with the center located directly under the hip for each leg. Let  $w_i$  denote the radius, and  $\alpha_i$  (used in place of  $t$ ) denote the phase angle (see Figs. 2, 3). Accordingly, we have  $(w_0, \alpha_0), \dots, (w_{n-1}, \alpha_{n-1})$ , where the  $\alpha_i$ 's are obtained by

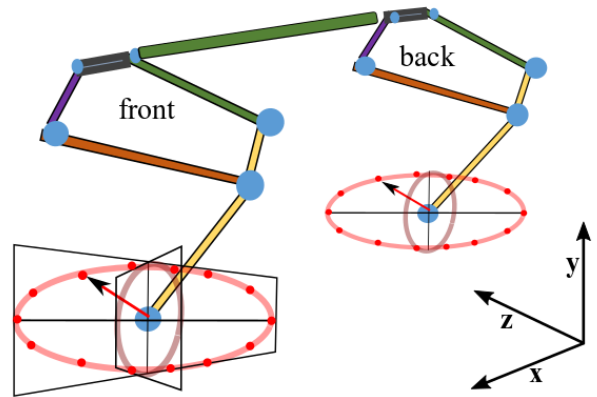


Fig. 3: Figure showing the construction of the splines for the front and back legs of the robot. Two types of ellipses are shown each on the  $x-y$  and  $y-z$  plane respectively. These ellipses are drawn to indicate the plane on which the splines are constructed. For illustration, the control points are shown in the form of red dots on one of the ellipses. For trotting, the  $x-y$  plane is used, and for side-stepping, the  $y-z$  plane is used. For turning, the planes are oriented at  $\pm 45^\circ$  yaw w.r.t. the  $x$  axis were used.

evenly spacing from 0 to  $2\pi$  rad, and  $w_i$ 's are obtained from a policy. Note that  $w_i$ 's are updated as a function of  $s$  after every gait step i.e., whenever the phase crosses 0,  $\pi$ .

Having obtained the end-foot trajectory from the control points, the joint angles are obtained via an inverse kinematics solver. More details about the inverse kinematics for parallel 5-bar linkages are provided in [18]. Inclusion of more points allow the spline to take more complex shapes. It was empirically observed that 6 control points were incapable of generating stable gaits, and after 18 points, there was not much improvement in the gaits obtained. It is worth noting that the plane on which the control points are placed is predetermined depending upon the gait (more details are given in Section III-B). In addition, with different phase shifts and signs, we run the same trajectory on all the four legs. Therefore, we choose the action space to be an 18 dimensional vector that is constrained to remain within the leg's work-space using a bounding box.

## III. TRAINING AND SIMULATION

Having defined the model and the control methodology, we are now ready to discuss the policy and training algorithm used for Stoch 2.

### A. Algorithm

Since the goal is to realize a library of gaits in low-cost hardware, we need a simple representation of the policy, one that is capable of running on a single embedded Tiva-M Series Microcontroller that contains only one floating point unit. We chose the simplest representation, a linear policy, thus our policy is a matrix that multiplies with the state at every step to output the action. Augmented Random Search (ARS) [14], is a learning algorithm which

is designed for finding linear deterministic policies, and is known to be on par with other model-free Reinforcement Learning algorithms for robot control tasks. In literature, this algorithm was successfully demonstrated on the simulated MuJoCo robot control tasks to obtain rewards that were on par with the rewards obtained by deep nets trained using traditional Reinforcement Learning Algorithms like TRPO and PPO [14]. Our experiments on the Stoch2 environment also demonstrate the same. Therefore we choose the policy to be  $\pi(s) := Ms$ , where  $M \in \mathbb{R}^{18 \times 12}$  is the matrix that maps the 12 motor angles to the 18 control points. Let the parameters of  $M$  be denoted by  $\theta$ . In this case as we are not adding any constraints to the matrix  $M$ , the parameters  $\theta$  of  $M$  are simply each element of  $M$ . Then the goal of ARS is to determine the  $\theta$ , of the matrix  $M$ , that yields the best rewards which in turn leads to the best locomotion gaits for Stoch 2.

ARS is a policy gradient algorithm, i.e., it optimizes a parameterized policy by moving along the gradient of the function that maps the parameters of the policy to the expected reward. Since the function is stochastic in nature, different algorithms use different estimates of the gradient. ARS estimates the gradient through a finite difference method as opposed to likelihood ratio methods used in other algorithms like PPO [20] or TRPO [21]. We use Version **V-1t** of ARS from [14]. **V-1t** performs the gradient descent step without normalization of state or action space as in **V-2**, and averages a subset of the top performing directions to determine the final direction of the gradient along which the policy should move. We do not require normalization of state space since all the motor angles vary between the same limits. In addition, we move along the average of a subset of top performing directions **V-1t**, unlike in **V-1**, which uses the average of all directions. This speeds up the training process by ignoring poor directions. Since a properly tuned **V-1t** includes the possibility of choosing the set of all directions, **V-1t** can not perform worse than **V-1**. In particular, we pick  $N$  i.i.d. directions  $\{\delta_i\}_{i=1, \dots, N}$  from a normal distribution, where  $N$  is the dimension of the policy parameters, i.e.,  $\theta \in \mathbb{R}^N$ . In our case,  $N = 18 \times 12$ , since we have 18 actions and 12 states. With a scaling factor of  $\nu > 0$ , we perturb the policy parameters  $\theta$  across each of these directions with the scaling  $\nu$ . The policy is perturbed both along the direction  $\theta + \nu\delta_{(i)}$  and away from the direction  $\theta - \nu\delta_{(i)}$ . Executing this perturbed policy for one episode yields the returns for each direction. Therefore, for  $N$  directions, we collect  $2N$  returns. We choose the best  $N/2$  directions corresponding to the maximum returns. Let  $\delta_{(i)}$ ,  $i = 1, 2, \dots, N/2$  correspond to these best directions in decreasing order of returns. We update  $\theta$  as follows:

$$\theta = \theta + \frac{\beta}{\frac{N}{2} * \sigma_R} \sum_{i=1}^{N/2} (R(\theta + \nu\delta_{(i)}) - R(\theta - \nu\delta_{(i)})) \delta_{(i)}, \quad (4)$$

where  $\beta > 0$  is the step size, and  $\sigma_R$  is the standard deviation of the  $N/2$  returns obtained.  $\theta$  is updated in the matrix  $M$ ,

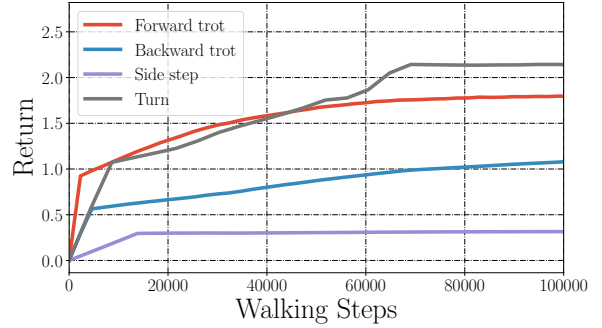


Fig. 4: Figure showing the training curve for all the four gaits. The returns saturated after about 60000 gait steps.

|               | Step-Size ( $\beta$ ) | Noise ( $\nu$ ) |
|---------------|-----------------------|-----------------|
| Forward trot  | 0.09                  | 0.03            |
| Backward trot | 0.1                   | 0.03            |
| Side step     | 0.1                   | 0.03            |
| Turn          | 0.1                   | 0.05            |

TABLE I: Table showing the hyperparameters used for the ARS algorithm.

and the gait is tested. This is repeated in every iteration.

### B. Training results

An online open source physics engine, PyBullet, was used to simulate and train our agent. Accurate measurements of link-lengths, moment of inertia's and masses were stored in the URDF file. An asynchronous version of Augmented Random Search with 20 parallel agents was used to speed up the training on the robot. The reward function ( $r$ ) chosen was

$$r = W_{vel} \cdot \Delta x - W_E \cdot \Delta E. \quad (5)$$

Here  $\Delta x$  is the difference between the current and the previous base positions/orientations. This difference changes depending upon the gait. For example, for trotting we seek for positive increments along the  $x$ -axis, for side-stepping we seek for increments along the  $z$ -axis, and for turning we seek for increments in yaw.  $\Delta E$  is the energy expended to execute the motion in each step.  $\Delta E$  is calculated by the formula  $E = \tau * v * \Delta t$  where  $\tau$  is the joint torque,  $v$  is the joint velocity and  $\Delta t$  is length of the simulation time step. The energy  $E$  is summed over one entire half-step of the robot.  $W_{vel}$ ,  $W_E$  are weights corresponding to each of the terms in (5).

We trained for the following gaits in the pybullet simulator: a) forward trot, b) backward trot, c) side-step and d) turn. These gaits were realized by 1) appropriate orientation of the plane of control points (see Fig. 3) and 2) by adding a phase offset between the  $\alpha$ 's of different legs. For forward and backward trotting, the plane was aligned with the  $x - y$  axes, while for side-stepping the plane was aligned with the  $y - z$  axes. For turning, the plane trajectory was oriented at  $+/- 45^\circ$  yaw angle w.r.t. the  $x$  axis. The planes for the front two legs were facing inwards, and the planes for the

back two legs were facing outwards. The phase offsets for all the four behaviors were  $\{0, \pi, \pi, 0\}$ , which correspond to front left, front right, back left, and back right respectively. Fig. 4 shows the training results for all the four gaits of the robot. The hyper-parameters used for training are shown in Tab. I. We used 20 parallel agents for training every gait. The forward and backward trot took 20 iterations to learn, which corresponds to about 400,000 simulation time-steps. Similarly sidestep took 20 iterations, and turn took 35 iterations to learn. The training times ranged between 2-6 hours, which are, in fact, on par with the training times obtained with D-RL for Stoch [17].

### C. Training with multiple gaits

We also have preliminary results on training with multiple gaits simultaneously. Multiple agents were spawned in a single environment with different phase differences and were run with the same policy. The new reward was the combined reward of all the agents. To the best of our knowledge ours is the only network capable of exhibiting multiple different quadruped gaits with a single linear policy. This also stands to the testament of the surprising capability of fully linear policies for robotic control tasks first observed in simulation in [14] and experimentally in [15]. Our training in all took about 20 iterations to saturate with about 20 parallel agents at a time. So in all it took about 400,000 simulation time-steps.

## IV. EXPERIMENTAL RESULTS

In this section, we will describe the experimental results. We will first describe the methodology used to address the sim-to-real gap.

### A. Bridging the sim-to-real gap

The desired trajectories obtained from the splines are tracked via a joint level PD control law. We observed that, despite having sufficient tracking performances at the joint level, the control points obtained in experiments from the linear policy were not matching with the control points obtained in simulation. This is due to the mismatch between the joint angle trajectories obtained from the simulation and experiment. In simulation, the tracking performance significantly differs from that in the real world. In particular, contact forces, joint and gearbox friction, inaccurate motor models cause the deviation in tracking performance between the simulated and real world data. The policy outputs the control points assuming implicitly a tracking performance observed in simulation and thus the same controllers do not work on the real world robot.

To address this sim-to-real problem [6] used domain randomization and a custom motor model. However domain randomization makes the learning unstable and causes early saturation of rewards. Since we have a linear policy, we observed marginal improvements with domain randomization. [7] aimed to bridge this gap by building an accurate motor model via supervised learning. In a similar vein, we addressed the sim-to-real gap by transforming the observed

motor angles in the real world in such a way that it matches with the observed motor angles in simulation. We first extracted the control points that our policy outputs when it reaches steady state ( 5 steps) in simulation. These control points will be used to compare the tracking performance in simulation and experiment. In particular, these control points are used to obtain the desired trajectories, which are tracked both in experiment and in simulation. Motor angle data was collected for 5 steps at a rate of 3 kHz. Next, a simple weighted linear regression between the experimental and the simulated motor angles collected was used. More weights were given for the stance phase of the legs where the inaccuracies are high due to contact with the ground. Fig. 5 shows the comparison between the two data sets obtained. We denote the linear transformation obtained from the regression as  $\hat{M}$ . We have the actions obtained from the new policy  $\hat{\pi}$  as

$$\hat{\pi}(s) := M\hat{M}s + M\bar{b}, \quad (6)$$

where  $\bar{b}$  is the offset obtained from the data. The new policy is used to obtain the desirable action values (control points) that match with the control points used in simulation. Thus the new policy outputs the same desired trajectory, for the tracking performance observed in the real world and since the tracking performance of the PD Control Law is good, now the policy works well on the robot, bridging the "sim-to-real" gap. There was a net error of 0.01 rad between the simulated and experimental state values after performing the weighted linear regression (see Fig. 5).

### B. Results

We tested all the four gaits on Stoch 2 experimentally. As mentioned previously, the states observed at the end of every gait step were transformed to yield more desirable state values, which were then used to obtain the desired trajectories for the next step. Fig. 7 shows the comparison between the desired and actual end-foot trajectories of the robot, and Fig. 8 shows the corresponding gait tiles for all the four gaits. While each of these gaits were obtained from training separately, the transition from one type of gait to another was achieved by using a low pass filter (see [18] for more details). Video results showing all the four gaits are provided in the submission.

### C. Empirical analysis of limit cycle behaviour

Fig. 6 shows the robustness of the trotting gait to external disturbances. The data recorded was for fifty gait steps. The video submission shows the disturbances applied on the robot, and the desired trajectories converging to a limit cycle. Experiments show that joint angles achieve self-correcting behaviors, which implies that a linear policy is able to calculate the control points that reject the disturbances.

## V. CONCLUSION

We successfully demonstrated multiple types of locomotion gaits learned via Augmented Random Search (ARS) in the custom built quadruped robot Stoch 2. Gaits were first

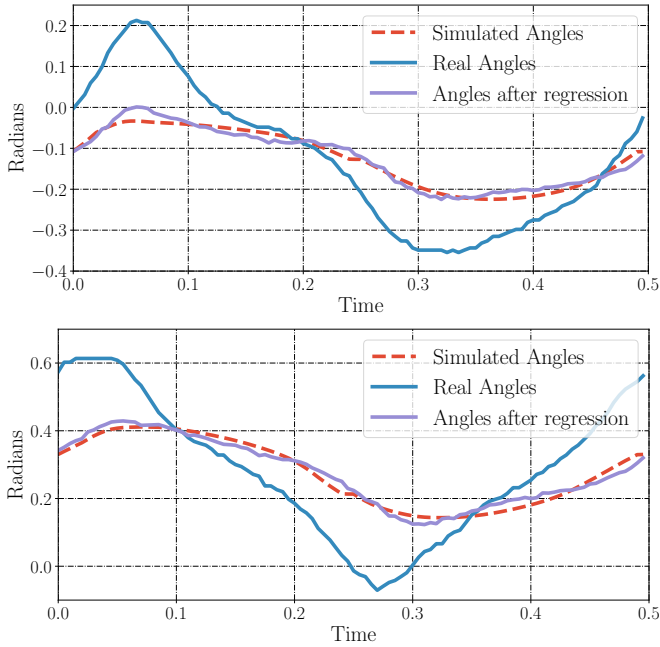


Fig. 5: Linear regression maps the real world data to the corresponding simulation data. The top figure is for the hip angle, and the bottom is for the knee angle. The dashed red line is for the knee angle. The dashed red line is the simulated trajectory, and the solid blue line is the experimental trajectory obtained for the same action values (control points) used. The purple line is the experimental trajectory obtained after the transformation.

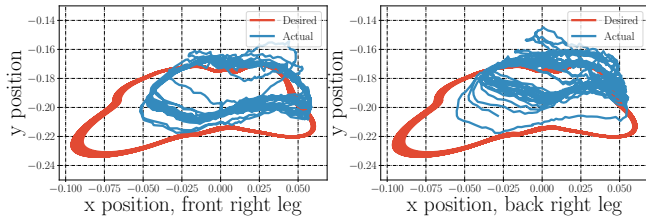


Fig. 6: Figure showing the closed operation of the trotting gait on Stoch 2 for 50 gait steps. It can be observed that the desired trajectories are stable despite the external disturbances.

tested in a simulation environment, and then deployed in Stoch 2 experimentally. Linear policy was used to determine the trajectory for each gait step. These types of policies are easy to deploy on hardware, and, at the same time, do not require large compute capability, unlike the Deep Neural Network (DNN) based control algorithms existing in literature. It is worth noting that training times with ARS are comparable with D-RL with the same compute platform. Future work will involve learning and testing more gaits on a diverse set of terrains like slopes and stairs.

#### ACKNOWLEDGMENT

The authors would like to thank Alok Rawat for assisting with the design and manufacturing of Stoch 2.

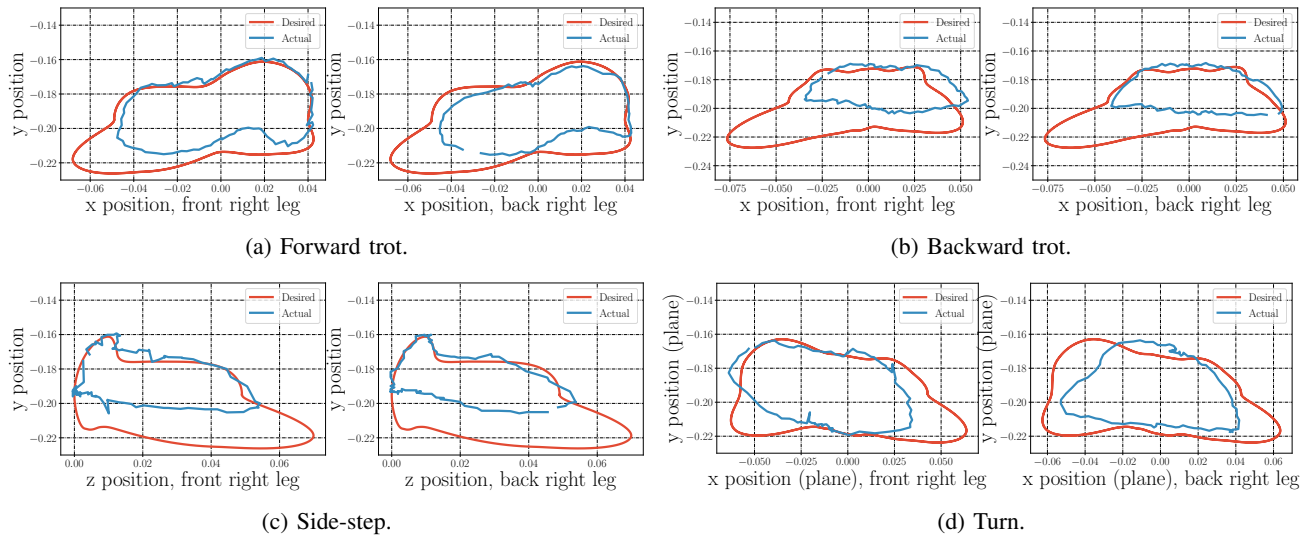


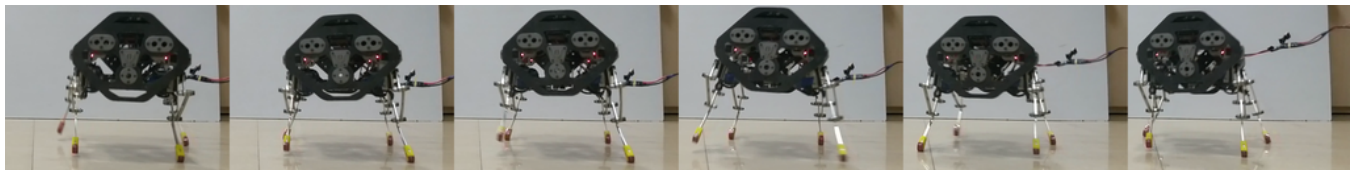
Fig. 7: Figure showing the comparison between the desired (red) and the actual end-foot trajectories for the various gaits tested on the robot. The plots are for one full cycle, i.e., two gait steps. The plots for turning are shown in the plane that is at  $45^\circ$  w.r.t. the  $x$  axis. Note that all the values are specified in meters.



(a) Forward trot.



(b) Backward trot.



(c) Side-step.



(d) Turn.

Fig. 8: Figure showing the sequence of tiles for all the four gaits tested on Stoch 2. The tiles for forward-backward trot, and side-step are for one gait step, while the tiles for turn are for approximately eight gait steps.

## REFERENCES

- [1] M. H. Raibert, "Legged robots," *Communications of the ACM*, vol. 29, no. 6, pp. 499–514, 1986.
- [2] M. Vukobratović and B. Borovac, "Zero-moment point thirty five years of its life," *International journal of humanoid robotics*, vol. 1, no. 01, pp. 157–173, 2004.
- [3] E. R. Westervelt, J. W. Grizzle, and D. E. Koditschek, "Hybrid zero dynamics of planar biped walkers," *IEEE Transactions on Automatic Control*, vol. 48, no. 1, pp. 42–56, Jan 2003.
- [4] A. J. Ijspeert, "Central pattern generators for locomotion control in animals and robots: a review," *Neural networks*, vol. 21, no. 4, pp. 642–653, 2008.
- [5] J. Kober, J. A. Bagnell, and J. Peters, "Reinforcement learning in robotics: A survey," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1238–1274, 2013.
- [6] J. Tan, T. Zhang, E. Coumans, A. Iscen, Y. Bai, D. Hafner, S. Bohez, and V. Vanhoucke, "Sim-to-real: Learning agile locomotion for quadruped robots," *CoRR*, vol. abs/1804.10332, 2018. [Online]. Available: <http://arxiv.org/abs/1804.10332>
- [7] J. Hwangbo, J. Lee, A. Dosovitskiy, D. Bellicoso, V. Tsounis, V. Koltun, and M. Hutter, "Learning agile and dynamic motor skills for legged robots," *Science Robotics*, vol. 4, no. 26, 2019. [Online]. Available: <https://robotics.sciencemag.org/content/4/26/eaau5872>
- [8] Z. Xie, P. Clary, J. Dao, P. Morais, J. Hurst, and M. van de Panne, "Iterative reinforcement learning based design of dynamic locomotion skills for cassie," *arXiv preprint arXiv:1903.09537*, 2019.
- [9] G. Hornby, M. Fujita, S. Takamura, T. Yamamoto, and O. Hanagata, "Autonomous evolution of gaits with the sony quadruped robot," 1999.
- [10] M. J. Quinlan, S. K. Chalup, R. H. Middleton *et al.*, "Techniques for improving vision and locomotion on the sony aibo robot."
- [11] N. Kohl and P. Stone, "Policy gradient reinforcement learning for fast quadrupedal locomotion," in *Robotics and Automation, 2004. Proceedings. ICRA'04. 2004 IEEE International Conference on*, vol. 3. IEEE, pp. 2619–2624.
- [12] R. Tedrake, T. W. Zhang, and H. S. Seung, "Stochastic policy gradient reinforcement learning on a simple 3d biped," in *Intelligent Robots and Systems, 2004.(IROS 2004). Proceedings. 2004 IEEE/RSJ International Conference on*, vol. 3. IEEE, pp. 2849–2854.
- [13] S. K. Chalup, C. L. Murch, and M. J. Quinlan, "Machine learning with aibo robots in the four-legged league of robocup," *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, vol. 37, no. 3, pp. 297–310, 2007.
- [14] H. Mania, A. Guy, and B. Recht, "Simple random search of static linear policies is competitive for reinforcement learning," in *Advances in Neural Information Processing Systems*, 2018, pp. 1800–1809.
- [15] A. Iscen, K. Caluwaerts, J. Tan, T. Zhang, E. Coumans, V. Sindhwani, and V. Vanhoucke, "Policies modulating trajectory generators," in *Conference on Robot Learning*, 2018, pp. 916–926.
- [16] D. Jain, A. Iscen, and K. Caluwaerts, "Hierarchical reinforcement learning for quadruped locomotion," *arXiv preprint arXiv:1905.08926*, 2019.
- [17] A. Singla, S. Bhattacharya, D. Dholakiya, S. Bhatnagar, A. Ghosal, B. Amrutur, and S. Kolathaya, "Realizing learned quadruped locomotion behaviors through kinematic motion primitives," *arXiv preprint arXiv:1810.03842*, 2018.
- [18] D. Dholakiya, S. Bhattacharya, A. Gunalan, A. Singla, S. Bhatnagar, B. Amrutur, A. Ghosal, and S. Kolathaya, "Design, development and experimental realization of a quadrupedal research platform: Stoch," *CoRR*, vol. abs/1901.00697, 2019. [Online]. Available: <http://arxiv.org/abs/1901.00697>
- [19] X. B. Peng and M. van de Panne, "Learning locomotion skills using deeprl: Does the choice of action space matter?" *CoRR*, vol. abs/1611.01055, 2016. [Online]. Available: <http://arxiv.org/abs/1611.01055>
- [20] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *CoRR*, vol. abs/1707.06347, 2017. [Online]. Available: <http://arxiv.org/abs/1707.06347>
- [21] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel, "Trust region policy optimization," *CoRR*, vol. abs/1502.05477, 2015. [Online]. Available: <http://arxiv.org/abs/1502.05477>