

Imitation and Adaptation Based on Consistency: A Quadruped Robot Imitates Animals from Videos Using Deep Reinforcement Learning

Qingfeng Yao^{1,2,3*}, Jilong Wang^{4*}, Shuyu Yang⁴, Cong Wang^{1,2,3}, Hongyin Zhang⁴,
Qifeng Zhang^{1,3}, Donglin Wang⁴⁺

Abstract—The essence of quadrupeds’ movements is the movement of the center of gravity, which has a pattern in the action of quadrupeds. However, the gait motion planning of the quadruped robot is time-consuming. Animals in nature can provide a large amount of gait information for robots to learn and imitate. Common methods learn animal posture with a motion capture system or numerous motion data points. In this paper, we propose a video imitation adaptation network (VIAN) that can imitate the action of animals and adapt it to the robot from a few seconds of video. The deep learning model extracts key points during animal motion from videos. The VIAN eliminates noise and extracts key information of motion with a motion adaptor, and then applies the extracted movements function as the motion pattern into deep reinforcement learning (DRL). To ensure similarity between the learning result and the animal motion in the video, we introduce rewards that are based on the consistency of the motion. DRL explores and learns to maintain balance from movement patterns from videos, imitates the action of animals, and eventually, allows the model to learn the gait or skills from short motion videos of different animals and to transfer the motion pattern to the real robot.

I. INTRODUCTION

Currently, three major types of robots are well developed in locomotion: wheeled, tracked, and legged robots. Compared to robots in other categories, legged robots have shown significantly higher performance when facing rough terrains due to their unique ability to have discrete ground contact points through their gait [1]. They can adjust their foothold position to alter their gaits according to the characteristics of terrain and tasks. However, legged robots must consider a massive number of factors to search for an optimum in high dimensional action space [2].

When straddling, it is critical to that the robot’s center of gravity projection and zero moment point remain within its geometry to maintain movement stability [3]. Traditional algorithms often plan the center of gravity trajectory in the form of discrete cycles and reduce the speed or increase the frequency to ensure stability. Nature is the teacher of humans, and many designs of robots are inspired by animals. To allow the robot to learn the action of animals more efficiently, it is reasonable to apply imitation learning to robots. Imitation

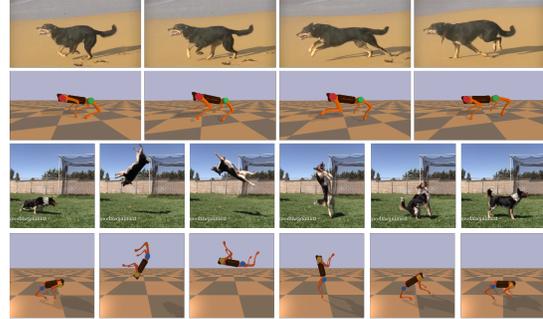


Fig. 1. Quadrupedal robots can learn different skills from few-second videos (the first and third rows), such as galloping (the second row) and backflip (the fourth row).

learning helps robots imitate the action of animals from existing animal data and transfer this knowledge to identical tasks, and the demonstration of animals is regarded as fundamental knowledge.

At the same time, deep reinforcement learning (DRL) has emerged rapidly from machine learning in recent years, and has gradually shown its potential in decision-making problems. In this paper, we propose a network that combines deep learning and DRL for quadruped robots to learn cyclical and noncyclical movements through animal videos. The imitation learning framework includes recognition, imitation and adaptation. Our framework combines the benefits of DRL and deep learning, allowing quadruped robots to learn various behavior patterns using only a few seconds of animal videos. DRL control networks can investigate and learn how to adapt information from various animals to complete tasks by adjusting their posture and maintaining motion stability. The framework proposed in this paper is trained in a simulation environment and transferred directly to the real quadruped robot. In summary, this work made the following contributions:

- The key animal nodes are extracted from the video through the deep neural network as a reference trajectory and a motion adaptor is used to remove the offset and error from the key nodes.
- Features extracted by the motion adaptor are prior information of the DRL network. DRL explores how to maintain balance based on prior information from animal videos.
- Simulation and real experiment with different type of motion demonstrate the ability of our algorithm.

* Contributed equally

+ Corresponding author. Email: wangdonglin@westlake.edu.cn

The main work was done at MiLAB, Westlake University.

¹State Key Laboratory of Robotics, Shenyang Institute of Automation, Chinese Academy of Sciences, Shenyang 110016, China

²Institutes for Robotics and Intelligent Manufacturing, Chinese Academy of Sciences, Shenyang 110169, China

³University of Chinese Academy of Sciences, Beijing 100049, China

⁴School of Engineering, Westlake University, Hangzhou 310024, China

II. RELATED WORK

Traditional legged robot control methods focus on modeling the robot and the surrounding environment, followed by positional estimation and trajectory planning using forward and inverse kinematics [4]. Carlo et al. [5] proposed a model-based method for the calculation of forces, which can be summarized as a convex optimization problem improving the effectiveness of MPC in quadruped robots. Bellicoso et al. [6] proposed an online motion planner based on the zero moment point. The position of the center of mass was optimized for the execution of different gaits.

In recent years, machine learning has played an increasingly important role in quadruped robot control. Villarreal et al. [7] combined a convolutional neural network with MPC. Yang et al. [8] used a gating neural network combined with pretrained neural networks to adapt to changes in the environment and tested the combined method on a real robot.

DRL is a learning process that enables agents to learn from scratch by feedback from exploring the environment [9]. Bellegarda et al. [10] proposed a framework combining DRL with nonlinear trajectory optimization to assist quadruped robots in jumping. Lee et al. [1] proposed a robust DRL-based controller resulted from curriculum training to guide a quadruped robot to locomote on complex terrains without any visual aid. Escontrela et al. [11] proposed a method using simulated 3D lidar sensors to complete navigation tasks on different terrains with a unified policy.

Imitation learning can swiftly master new tasks by observing expert demonstrations and achieving a feasible reproducibility of the expert strategy [12]. A strategy that directly maps trajectory actions by state features with a method called behavior cloning (BC) [13] was developed to directly learn the state-to-action mapping. BC transforms imitation learning into a supervised learning problem, which is easily affected by changes in state distribution, therefore, a regularized behavioral cloning imitation learning method was proposed [14]. In this method, the soft-Q-learning algorithm was introduced to overcome the disturbance from the imbalance of state distributions. Another task is to convert the imitation learning problem into a DRL problem by inferring the reward function from the expert demonstration [15].

Imitation learning already has applications in robots. Xie et al. [16] proposed a deterministic action stochastic state method that uses the policy gradient to help the biped robot Cassie generate gaits that are similar to those of humans. A method that learns basketball dribbling control from motion capture data was proposed by Liu [17]. The system developed a strategy based on trajectory optimization and DRL to learn the skill of dribbling. A controller was developed by Yamane [18] to allow the robot to maintain balance while tracking the movement of a given reference. Peng et al. [19] proposed a method that could handle key-frame motion. By combining imitation targets with task goals, the agents were trained to react intelligently to the environment and perform a variety of skills.

3D pose estimation is a way to take full advantage

of animal demonstrations with multiple cameras [20]. A motion-sensing camera is used to extract the demonstration information, and the Levenberg-Marquardt method is used to optimize the solution of inverse kinematics to improve the stability and similarity of the robot movement from human posture [21].

A framework based on pose estimation and DRL was used to extract the whole body 3D reference motion from public video clips [22]. Kearney et al. [23] studied the 3D pose estimation problem of dogs based on RGBD images by using a single motion capture system to obtain the skeleton. Ou et al. [24] addressed this issue by planning the zero moment point and driving angle through human body demonstration data so that the robot could imitate human actions and gait patterns.

The pose skills of animals can be used by quadruped robots because of their similar morphologies. An imitation learning system was proposed by Peng [25] that enables quadruped robots to learn agile motion skills by imitating animals, where mocap 3D data is used to match and retarget various behaviors for imitation learning. Capturing motion data with a 3D camera were a requirement for the type of equipment and data [26].

Monocular camera videos are fairly easy to obtain; therefore Vondrak et al. [27] took a more economical approach by estimating human motion patterns from monocular videos. A state-space biped controller with a balanced feedback mechanism executed the required movement by directly establishing it from the monocular video. Da et al. [28] proposed a method to imitate different sports, where the retargeting method was applied to the motion capture data, and the inverse kinematics algorithm was used for imitation. The generation of a motion trajectory was transformed into a nonlinear optimization problem to minimize the trajectory error between the animal and the robot.

Learning from animal video technology provides a simple and effective method for animal behavior imitation, but manually extracting key nodes is time-consuming. To solve this problem, an unlabeled pose estimation method [29] based on deep neural network transfer learning is proposed to track different key body nodes of different species with minimal training data.

III. METHOD

We propose a novel *video imitation adaptive network* (VIAN) framework based on animal videos to imitate learning and adapt to the robot by mapping the states to the corresponding actions with neural networks.

As shown in Fig. 2, first, we analyzed the motion of the animal and extracted key nodes with the deep neural network. To compensate for the structural difference between the animal and the quadruped robot, we selected key points on the body of the animal that could be mapped to the corresponding body structure of the quadruped robot. For example, the paw of a dog matched the foot of the robot, and the neck of a dog matched the front of the body of the robot. We designed a motion adaptor that filtered out

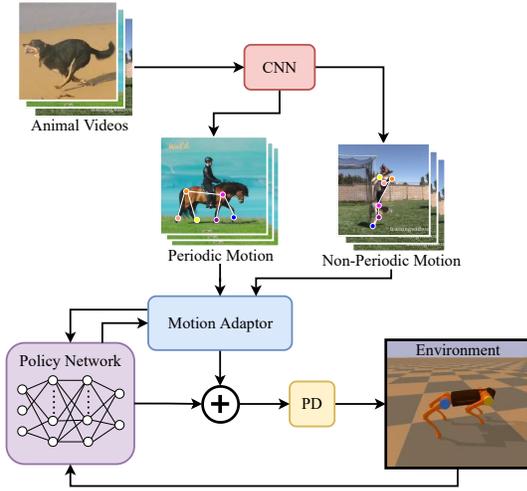


Fig. 2. The overview of VIAN. The imitation learning framework includes the part of recognition and imitation based on deep learning and DRL.

noise and offsets from key nodes to formulate a reference trajectory for the robot to imitate. Finally, we used DRL to train the robot to imitate the motion and maintain balance with a consistency reward.

Here we employ a key frame extraction method based on *DeepLabCut* [29] to obtain the demonstration features from the videos. *DeepLabCut* is a toolbox used for unmarked pose estimation of animals performing various tasks. We tested animals such as a running dog or a pacing horse. In the movement of a quadruped robot, the position of feet and legs is the control variable. Therefore, we extracted the positions of the head, buttocks and feet from moving animals as the initial information.

A. Motion Adaptor

To allow the robot to complete different tasks by imitating animals in the videos, we separate the animal motion in the videos into two types: periodic motion and aperiodic motion. To map the body of the animals onto the structure of the robot, we selected a set of key points that were specified on the animal and that were paired with corresponding target key points on the robot.

Moving limbs tend to deflect and cause additional noise due to angular variations in periodic motion video. To avoid video noise and angle deviation in the moving process, we separated the movements in the animal video into three sections for periodic motion such as walking and running. The time series data method X11 [30] was used to separate the data into long-term trends, seasonal trends and random components. The X11 decomposition method is based on the classic time series decomposition method, and it overcomes the shortcomings of the classic time series decomposition method. Specifically, it can estimate the trend-period term of each period, including the endpoints, and it allows seasonal slow changes. The model uses iterative process fitting to estimate three X11 factors: trend cycle, seasonality, and noise. The X11 method assumes that the main components

of the time series follow the additive model:

$$y_t = TC_t + S_t + I_t, \quad (1)$$

where y_t represents the original series, TC_t represents the trend, S_t represents the periodicity, and I_t represents the residual. The components were estimated using a smoothing linear filter. The filters were applied in order, and finally, a periodic function was obtained as a function of the seasonality of motion. The X11 method was used to process the movement of the x-axis and z-axis of key nodes obtained in video over time. Fig. 3 shows the trajectory of motion and the effect after processing with the X11 method. The X11 method removes the noise and trends in motion and extracts the periodic function S_t as the imitation information. The period is the X11 input parameter, which is calculated from the average peak distance of animal videos.

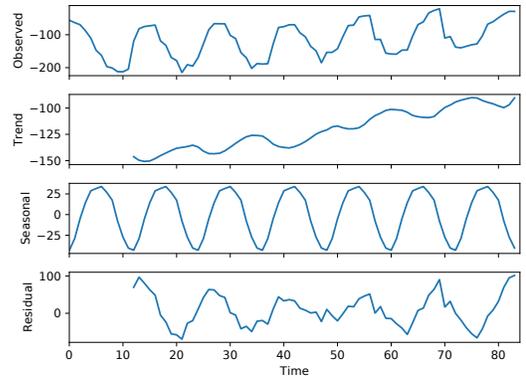


Fig. 3. Original message divided into the trend, seasonal and residual factors by X11.

Due to the different body proportions of the animals in the videos, we preliminarily process key points of the video according to the body proportions of the quadruped robot and the extracted key node positions. The scale of the zoom is:

$$f_i(t) = \frac{X_i^{lim}}{Max(S_i)} * S_i(t). \quad (2)$$

The maximum range that the i -th leg can move is X_i^{lim} and $Max(S_i)$ is the range of the i -th leg periodic function $S_i(t)$. To avoid the inaccuracy and unsteady motion behavior induced by applying $f_i(t)$ directly to a quadruped robot

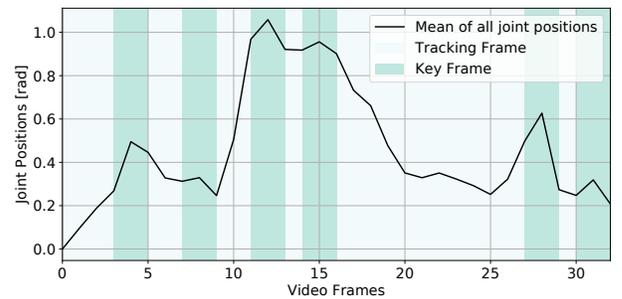


Fig. 4. The key frame is extracted according to foot movement

and performing the movement, the quadruped robot needs to adapt video motion information to its own movement process and learn to adjust $f_i(t)$ according to its current state.

By combining posture imitation and DRL, the robot learns to adapt to the periodic function according to its posture and completes different movements, after which the DRL network adjusts the periodic function $f_i(t)$. DRL outputs are determined by the internal state, which includes three historic motor positions, the IMU, and the foot position. We calculate the final quadruped position as:

$$g_i(t) = A_i(t) (f_i(t)) + b_i(t), \quad (3)$$

where $A_i(t)$ is the prior feature scaling and $b_i(t)$ is the adjustment offset output by the neural network. Finally, the angles of the leg joints are obtained by solving the inverse kinematics with $g_i(t)$. At the same time, the network outputs ΔT to control the motion frequency of the robot.

For aperiodic motion, since there are no seasonal data for the robot to learn the repeated motion cycle, we marked the key frames of actions that had the most impact on finishing the complete movement. For example, in the backflip task, actions such as jumping and pitching back are essential for the rotation of the body. To determine the key frames, we determined that the average angle of the video joint yields the key frame. It becomes a key frame when the average angle $\theta(t)$ reaches the local maximum value with threshold δ_{angle} :

$$\theta_{mean}(t) - \theta_{mean}(t-1) > \theta_{mean}(t+1) - \theta_{mean}(t) + \delta_{angle} \quad (4)$$

The robot is rewarded for walking forward and maintaining a consistent body direction. The episode terminates when the robot loses its balance or walks out of a designated area. The reward function for this task is written as follow:

$$r_{lin} := \begin{cases} \exp(-5.0(v_x - 0.3)^2) & v_x < 0.3 \\ 1 & v_x \geq 0.3 \end{cases}, \quad (5)$$

$$r_{ang} := \exp(-3.0(\omega_y)^2), \quad (6)$$

and

$$r_{body} := \exp(-3v_y^2) + \exp(-3.0(\theta_r + \theta_y)^2), \quad (7)$$

where r_{lin} indicates that the robot moves forward and r_{ang} and r_{body} limit the linear velocity tangent to the target direction and the attitude angle, and ensure that the robot does not deviate from the forward direction. v_x and v_y represent linear velocities along the x-axis and y-axis, respectively and θ_r , θ_y and ω_y represent roll, yaw and the yaw velocity, respectively.

B. Consistency Reward

When analyzing the animal walking data, we only focus on one set of data that contains time steps and the Cartesian coordinates of the corresponding key points. To better adapt to the motion, we introduced gait rewards to strengthen the ability of the robot to imitate the capture video actions. They are divided into two parts. One reward type is the stamp

consistency reward that keeps the robot landing in sync with the video animal. One key parameter in quadruped walking is the ratio of the leg swinging phase to the entire cycle. According to the range of landing, the robot learns how to synchronize its feet to keep track of the animal. Another reward is the motion consistency reward. It is difficult to directly capture the speed and angular velocity in the 2D video, so we use the direction of the relative speed for their determination. This reward helps the robot move in the same direction as the animal, and to have the same orientation of momentum.

We calculate the changes in foot height by analyzing the motion trajectory of the foot in the video. The trajectory threshold line is determined from $(h_i^{max} - h_i^{min}) * \gamma$, where γ is the threshold coefficient and h_i^{max} and h_i^{min} are the maximum and minimum of the foot height for different feet i in $f_i(t)$. When the foot motion trajectory is lower than the threshold line, the animal foot position is close to the ground; hence, the robot foot is on the ground in the simulation. Punishment is applied when the gait trajectories appear dissimilar, as shown in Fig. 5.

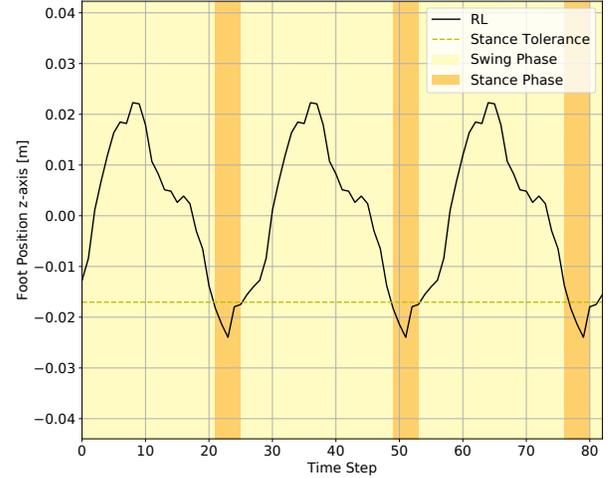


Fig. 5. Foothold location and foot lifting area of the animal in the video obtained by calculating the swing leg thresholds. The consistency of the gaits is proven if the foot of the robot is on the ground.

The curve of the foot swing z-axis value of the animal changing with time is shown in Fig. 5, where we set threshold γ to 0.1, as shown by the dotted line. When the curve is above the threshold, the i -th foot is considered to be in the air and an indicator of the video foot position, with $d_i = 0$; when the curve is beneath the threshold, the foot is considered to be on the ground and $d_i = 1$.

p_i is the foot position indicator for the robot. The i -th foot p_i is on the ground when $p_i = 1$; otherwise, $p_i = 0$. At the beginning of training, the robot did not learn to lift its leg. Hence, to encourage its robot to learn to lift the leg in the same manner as the animal in the video, we introduce the following stamp reward:

$$r_s^i = \begin{cases} -1 & \text{if } d_i = 0 \cup p_i = 1 \\ 0 & \text{other} \end{cases} \quad (8)$$

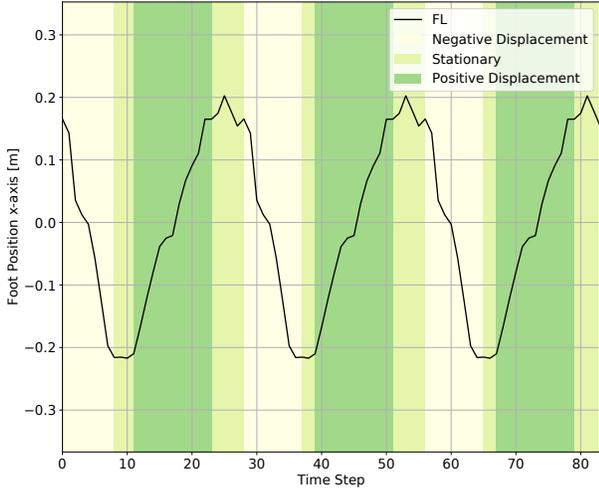


Fig. 6. The front left foot position trending. Different colors present different statuses of motion: low speed is light green, moving forward is dark green, and backward is white.

We further propose motion consistency to enhance the imitation effect. We estimate the accuracy of imitation by calculating the positions of the feet relative to the rear and front. To calculate the positions of feet relative to the body, the position angle of the head(x^1, y^1) and rear(x^2, y^2) is computed:

$$\beta = \arctan \left(\frac{y^1 - y^2}{x^1 - x^2} \right). \quad (9)$$

The feet positions are processed based on the relative angle of the front and rear. We obtain the i -th foot's motion state(x_i, y_i) related to location (x_{body}^i, y_{body}^i) based on the animal body frame in the video as follows:

$$\begin{cases} x_{body}^i = x_i \cos(\beta) + y_i \sin(\beta) \\ y_{body}^i = y_i \cos(\beta) - x_i \sin(\beta) \end{cases} \quad (10)$$

The consistency of the motion is evaluated by calculating the foot's motion direction relative to the body frame. To encourage the foot of the robot to move in the same direction as the reference action, the motion consistency reward is calculated by comparing the computed moving speed of feet and the reference moving direction as shown below:

$$r_m^i = \begin{cases} 0 & \text{if } \|v_{body}^i\| < \delta \\ \frac{v_{body}^i \cdot v_i^{real}}{\|v_{body}^i\| \cdot \|v_i^{real}\|} & \text{if } \|v_{body}^i\| \geq \delta \end{cases} \quad (11)$$

where δ is the speed threshold, v_{body}^i is the i -th foot speed of movement in the video calculated from x_{body}^i and y_{body}^i , and v_i^{real} is the speed of the i -th foot of the robot.

The foot motion of the animal along the x-axis is shown in Fig. 6. We partition the moving process into 3 stages based on velocity, which is presented in different colors based on their values.

The above rewards are summed into the total reward in periodic tasks:

$$r_{total} = 2 * r_{lin} + 2 * r_{ang} + r_{body} + r_f + r_g \quad (12)$$

We further extended the angle consistency index to fix this situation in the backflip task. We refer to the angle of animal β_{ref} in the video and return the corresponding reward according to the angle β_{real} of the robot in the simulation as follows:

$$r_a = \exp(-5(\beta_{ref} - \beta_{real})^2) \quad (13)$$

Rewards are composed of angle consistency, stamp consistency and motion consistency in backflip tasks:

$$r_{total} = 2 * r_a + r_m + r_s \quad (14)$$

Reinforcement learning directly outputs the joint angle and time interval in this scenario.

IV. EXPERIMENT

We first trained our policies in an open-source simulation, PyBullet, and then evaluated them on a *Unitree A1* robot, which is a quadruped robot with 12 degrees of freedom (3 per leg). Meanwhile, to implement the RL algorithm on real robots, sim-to-real problems [31] still need to be solved to bridge the gap between simulation and reality. Common solutions include domain adaptation [32] and randomization [33]. We randomize the parameters in the simulation so that the algorithm can be more robust to measure inaccuracies and uncertainty in the real world. Details of parameter randomization are listed in TABLE I.

TABLE I
RANGE OF PHYSICAL PARAMETERS. AT THE BEGINNING OF EACH TRAINING, THE PHYSICAL PARAMETERS ARE UNIFORMLY SAMPLED WITHIN THESE RANGES.

Physical parameters	Range
Mass	[0.8, 1.2]*defaults
Inertia	[0.5, 1.5]*defaults
Motor strength	[0.8, 1.2]*defaults
Motor friction	[0, 0.05]Nm/s/rad
Latency	[0.0, 0.04]s
Lateral friction	[0.5, 1.25]Ns/m
Battery	[14.0, 16.8]V
Joint friction	[0, 0.05]Nm
CoM position noise	[-5,5]cm
External force	[-3,3]N
Step height	[0.02, 0.06] m
Step width	[0.18, 0.23] m

We extract video clips of walking and trotting from different animals, such as horses and dogs, with each clip being between only 3 seconds and 8 seconds long and containing several motion cycles. Then, the periodic function is obtained by the X11 process. We extract the x and z coordinates of the foot positions using X11 methods to extract seasonal motion, as shown in Fig. 7. The result illustrates how the X11 process marks the positions of the dog's feet during running and the increase in periodicity and stability.

To maintain the stability and consistency of motion, we scale seasonal motion to the range of the quadruped robot and obtain the basic motion trajectory, DRL is used to complete the movement with seasonal motion. Adam is used

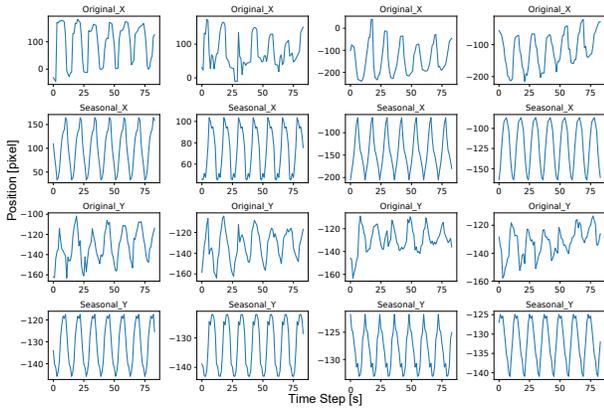


Fig. 7. Effect of the motion trajectory of the dog’s feet extracted by X11. There is an observable improvement in periodicity and the elimination of noise.

TABLE II
EFFECTS OF DIFFERENT TASKS.

Type	Success Rate	Speed	Stamp	Movement
RL(dog)	0.55	0.21	0.16	0.37
X11(dog)	0.0	0.05	0.2	0.52
VIAN(dog)	0.8	0.32	0.21	0.36
RL(horse)	0.55	0.42	0.06	0.41
X11(horse)	0.0	-0.02	0.12	0.29
VIAN(horse)	0.75	0.48	0.08	0.42

as the training optimizer to train the dataset for the network with a learning rate of 0.001. The strategy network completes the convergence for each animal after $3e7$ iterations of training. Robot training is performed with NVIDIA-2080 for 5 hours using Ubuntu 18.04.

Videos of moving horses and running dogs were used to test the effect of VIAN. We compared the effect of using DRL to train the data without X11 processing and applying X11 data directly to the robot. The task was to move steadily forward for 10 meters. We conducted 20 experiments and recorded the success rate, corresponding speed and consistency reward. The results are shown in Table II. The quadruped robot cannot complete the task by directly using the extracted data. Meanwhile, due to data noise, the result of training with data not processed by X11 has a lower effect.

To illustrate the effect of the model, we visualize the experimental results. The scaling curve of the limbs from the network when learning from a horse walking is shown in Fig. 8. The result shows the result of the robot imitating the horse walking. By visualizing the transformation of the motion trajectory, we observe that the robot can learn to maintain its body balance by adjusting the Y-axis even without information about the motion in the Y direction.

Meanwhile, RL provides a suitable scaling ratio and residual to transfer the original curve to a new curve that is suitable for the robot structure. The robot learns to scale up the moving distance along the x-axis, outputs a periodical residual to adjust the position along the y-axis, and smooths the original curve to help maintain the balance of its moving

status on the z-axis.

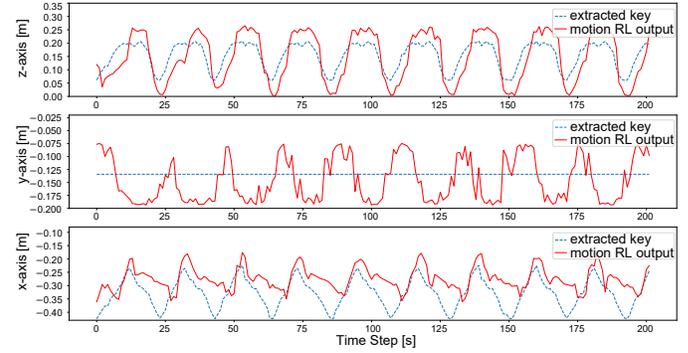


Fig. 8. Result of DRL training. A periodic function was learned by DRL to help the robot maintain its gait according to the baseline knowledge.

To test the adaptability of our method, we attempted to implement the backflip motion on the quadruped robot. In this task, the key frames of the motion are extracted, and the DRL adjusts the joint angle of the corresponding key frame and the time interval of different key moments according to the reward. The final motion trajectory is shown in Fig. 9. Combined with our framework, we used data from a few key frames of the video to complete the robot’s backflip task. To the best of our knowledge, this is the first time an agile robot has completed the backflip task with DRL.

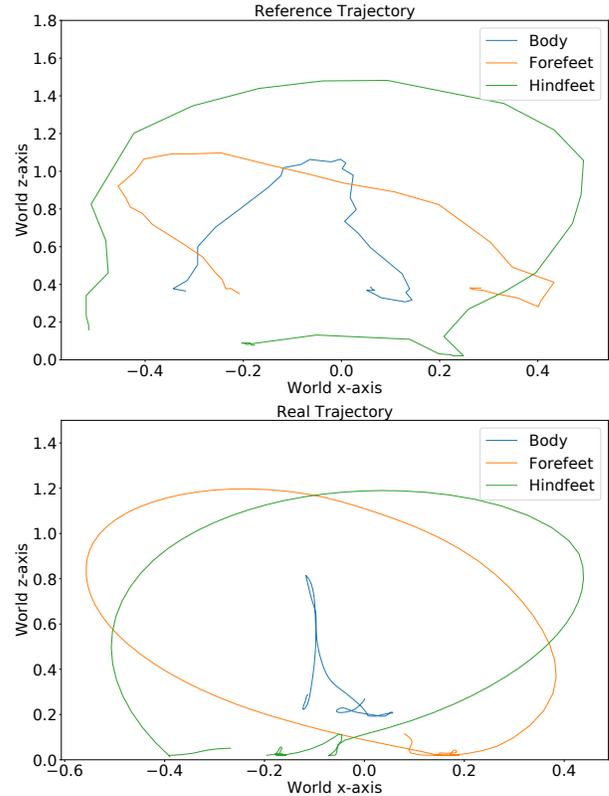


Fig. 9. Result of DRL training Obtained by referring to the video and the final result of the x- and z-axis coordinates of the body center and front and back feet. Through key frames, the backflip skill is learned by DRL.

To solve this sim-to-real problem, we added domain

randomization with randomized variables. We successfully transferred the skills learned from a walking horse video and backflip to the real robot as shown in Fig. 10, which illustrates the stability of our model.

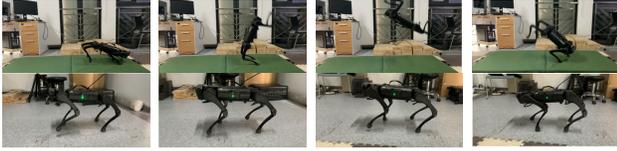


Fig. 10. Real robot experiment of periodic and aperiodic motion.

V. CONCLUSIONS

In this paper, we proposed an adaptive imitation framework that learns skills from a few seconds of animal movement video based on consistency and transfers the periodical gait and nonperiodic motion from different animals to the quadruped robot without any additional information.

In the future, we will establish a high-level network that dynamically allows robots to learn new skills with few sets of data based on the existing skills and choose an optimal skill according to the tasks and learned skills.

REFERENCES

- [1] Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning quadrupedal locomotion over challenging terrain. *Science robotics*, 5(47), 2020.
- [2] Ruben Grandia, Andrew J Taylor, Aaron D Ames, and Marco Hutter. Multi-layered safety for legged robots via control barrier functions and model predictive control. *arXiv preprint arXiv:2011.00032*, 2020.
- [3] Alexander W Winkler, Farbod Farshidian, Diego Pardo, Michael Neunert, and Jonas Buchli. Fast trajectory optimization for legged robots using vertex-based zmp constraints. *IEEE Robotics and Automation Letters*, 2(4):2201–2208, 2017.
- [4] Marc H Raibert. *Legged robots that balance*. MIT press, 1986.
- [5] Jared Di Carlo, Patrick M Wensing, Benjamin Katz, Gerardo Bleddt, and Sangbae Kim. Dynamic locomotion in the mit cheetah 3 through convex model-predictive control. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–9. IEEE, 2018.
- [6] C Dario Bellicoso, Fabian Jenelten, Christian Gehring, and Marco Hutter. Dynamic locomotion through online nonlinear motion optimization for quadrupedal robots. *IEEE Robotics and Automation Letters*, 3(3):2261–2268, 2018.
- [7] Octavio Villarreal, Victor Barasuol, Patrick Wensing, and Claudio Semini. Mpc-based controller with terrain insight for dynamic legged locomotion. *arXiv preprint arXiv:1909.13842*, 2019.
- [8] Chuanyu Yang, Kai Yuan, Qiuguo Zhu, Wanming Yu, and Zhibin Li. Multi-expert learning of adaptive legged locomotion. *Science Robotics*, 5(49), 2020.
- [9] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [10] Guillaume Bellegarda and Quan Nguyen. Robust quadruped jumping via deep reinforcement learning. *arXiv preprint arXiv:2011.07089*, 2020.
- [11] Alejandro Escontrela, George Yu, Peng Xu, Atil Iscen, and Jie Tan. Zero-shot terrain generalization for visual locomotion policies. *arXiv preprint arXiv:2011.05513*, 2020.
- [12] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2):1–35, 2017.
- [13] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991.
- [14] Siddharth Reddy, Anca D Dragan, and Sergey Levine. Sqil: Imitation learning via reinforcement learning with sparse rewards. *arXiv preprint arXiv:1905.11108*, 2019.
- [15] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *arXiv preprint arXiv:1606.03476*, 2016.
- [16] Zhaoming Xie, Patrick Clary, Jeremy Dao, Pedro Morais, Jonathan Hurst, and Michiel van de Panne. Iterative reinforcement learning based design of dynamic locomotion skills for cassie. *arXiv preprint arXiv:1903.09537*, 2019.
- [17] Libin Liu and Jessica Hodgins. Learning basketball dribbling skills using trajectory optimization and deep reinforcement learning. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [18] Katsu Yamane, Stuart O Anderson, and Jessica K Hodgins. Controlling humanoid robots with human motion data: Experimental validation. In *2010 10th IEEE-RAS International Conference on Humanoid Robots*, pages 504–510. IEEE, 2010.
- [19] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018.
- [20] Tanmay Nath, Alexander Mathis, An Chi Chen, Amir Patel, Matthias Bethge, and Mackenzie Weygandt Mathis. Using deeplabcut for 3d markerless pose estimation across species and behaviors. *Nature protocols*, 14(7):2152–2176, 2019.
- [21] Kai Hu, Christian Ott, and Dongheui Lee. Online human walking imitation in task and joint space based on quadratic programming. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages 3458–3464. IEEE, 2014.
- [22] Xue Bin Peng, Angjoo Kanazawa, Jitendra Malik, Pieter Abbeel, and Sergey Levine. Sfvr: Reinforcement learning of physical skills from videos. *ACM Transactions On Graphics (TOG)*, 37(6):1–14, 2018.
- [23] Sinéad Kearney, Wenbin Li, Martin Parsons, Kwang In Kim, and Darren Cosker. Rgbd-dog: Predicting canine pose from rgbd sensors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8336–8345, 2020.
- [24] Yongsheng Ou, Jianbing Hu, Zhiyang Wang, Yiqun Fu, Xinyu Wu, and Xiaoyun Li. A real-time human imitation system using kinect. *International Journal of Social Robotics*, 7(5):587–600, 2015.
- [25] Xue Bin Peng, Erwin Coumans, Tingnan Zhang, Tsang-Wei Lee, Jie Tan, and Sergey Levine. Learning agile robotic locomotion skills by imitating animals. *arXiv preprint arXiv:2004.00784*, 2020.
- [26] Daniel Holden, Jun Saito, and Taku Komura. A deep learning framework for character motion synthesis and editing. *ACM Transactions on Graphics (TOG)*, 35(4):1–11, 2016.
- [27] Marek Vondrak, Leonid Sigal, Jessica Hodgins, and Odest Jenkins. Video-based 3d motion capture through biped control. *ACM Transactions On Graphics (TOG)*, 31(4):1–12, 2012.
- [28] ZHAO Da, SONG Sifan, SU Jionglong, Zijian JIANG, and Jiaming ZHANG. Learning bionic motions by imitating animals. In *2020 IEEE International Conference on Mechatronics and Automation (ICMA)*, pages 872–879. IEEE, 2020.
- [29] Alexander Mathis, Pranav Mamidanna, Kevin M Cury, Taiga Abe, Venkatesh N Murthy, Mackenzie Weygandt Mathis, and Matthias Bethge. Deeplabcut: markerless pose estimation of user-defined body parts with deep learning. *Nature neuroscience*, 21(9):1281–1289, 2018.
- [30] Estela Bee Dagum and Silvia Bianconcini. *Seasonal adjustment methods and real time trend-cycle estimation*. Springer, 2016.
- [31] Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. *arXiv preprint arXiv:1804.10332*, 2018.
- [32] Konstantinos Bousmalis, Alex Irpan, Paul Wohlhart, Yunfei Bai, Matthew Keelcey, Mrinal Kalakrishnan, Laura Downs, Julian Ibarz, Peter Pastor, Kurt Konolige, et al. Using simulation and domain adaptation to improve efficiency of deep robotic grasping. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 4243–4250. IEEE, 2018.
- [33] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3803–3810. IEEE, 2018.