

Sim-to-Real Transfer for Quadrupedal Locomotion via Terrain Transformer

Hang Lai^{1,2}, Weinan Zhang¹, Xialin He¹, Chen Yu^{2,3}, Zheng Tian⁴, Yong Yu¹, and Jun Wang^{2,5}

Abstract—Deep reinforcement learning has recently emerged as an appealing alternative for legged locomotion over multiple terrains by training a policy in physical simulation and then transferring it to the real world (i.e., sim-to-real transfer). Despite considerable progress, the capacity and scalability of traditional neural networks are still limited, which may hinder their applications in more complex environments. In contrast, the Transformer architecture has shown its superiority in a wide range of large-scale sequence modeling tasks, including natural language processing and decision-making problems. In this paper, we propose Terrain Transformer (TERT), a high-capacity Transformer model for quadrupedal locomotion control on various terrains. Furthermore, to better leverage Transformer in sim-to-real scenarios, we present a novel two-stage training framework consisting of an offline pretraining stage and an online correction stage, which can naturally integrate Transformer with privileged training. Extensive experiments in simulation demonstrate that TERT outperforms state-of-the-art baselines on different terrains in terms of return, energy consumption and control smoothness. In further real-world validation, TERT successfully traverses nine challenging terrains, including sand pit and stair down, which can not be accomplished by strong baselines.

I. INTRODUCTION

Legged locomotion over varied terrains is challenging due to the potential changing dynamics and irregular profiles, surfaces, and obstructions [1]. To tackle this problem, deep reinforcement learning (RL) methods [1–5] have been developed by directly training a controller in simulation with diverse terrains and environmental parameters and then transferring it to the real world using different heuristic adaptation techniques [5–9]. Typically, in the well-known *privileged learning* framework [1], a teacher policy will first encode the privileged information, e.g., heightmaps and physical parameters, which are inaccessible directly in the real world, into a low-dimensional latent vector. Subsequently, a student policy is trained to recover such a latent vector by just using the previous proprioception sequences with conventional sequence models like temporal convolutional network (TCN) [1, 3, 10] or recurrent neural network (RNN) [11, 12].

Though privileged learning and its extension have been widely exploited and achieved remarkable success in quadrupedal locomotion over multiple terrains, it still has some limitations. Firstly, a versatile teacher policy will make full use of the latent vector, denoted as l_t , to represent the environmental property, making it difficult to be estimated



Fig. 1. Application of TERT on A1 robot over multiple challenging terrains.

online. To address this issue, Peng et al. [5] proposed to constrain the mutual information between l_t and the privileged information, which, however, may impede the adaptability of the policy. Secondly, the learned student policy tends to be sensitive to the estimated l_t , and may suffer catastrophic degradation in performance if l_t is out of its training distribution. Therefore, can we avoid the strong dependence on the precision of the estimated vector?

One promising solution is to circumvent the encoded l_t and directly predict the actions of the teacher policy conditioned on the previous proprioception sequences. We argue that this sequence prediction paradigm is exceptionally suitable for the Transformer architecture since it has demonstrated remarkable performance improvement and generalization on a wide range of sequence modeling problems due to its high capacity for modeling and robustness for decoding. In addition, the self-attention mechanism of Transformer could make better use of historical information through credit assignment to capture the characteristics and changes of the environment. Moreover, Transformer has proven to be able to model diverse behaviors [13], which is crucial for transferring to unseen scenarios. Therefore, it is tempting to utilize the modern Transformer models to imitate the teacher policy’s behavior over various terrains using the past trajectories, which can seamlessly take advantage of innovations in sequence modeling literature.

To this end, we propose Terrain Transformer (TERT),

¹Dept. of Computer Sci. and Eng., Shanghai Jiao Tong University, China.

²Digital Brain Lab, Shanghai, China

³School of Info. Sci. and Tech., ShanghaiTech University, China.

⁴School of Creativity and Art, ShanghaiTech University, China.

⁵Centre for Artificial Intelligence, University College London, UK.

a novel Transformer model with privileged learning, and present a two-stage training framework explicitly designed for sim-to-real transfer. More specifically, in the *offline pretraining* stage, the teacher policy interacts with a simulator and collects trajectories for TERT training. This stage resembles the standard Transformer training process in language and offline RL areas [14–16]. After the first stage, Transformer can give relatively accurate predictions conditioned on the teacher’s observation-action sequences. However, the input sequence distribution of the Transformer will shift when deploying. Therefore, in the *online correction* stage, TERT interacts with the simulator while the teacher simultaneously gives actions as the target. Then TERT is trained to fit the teacher’s actions conditioned on its own observation-action sequences.

We compare TERT to the state-of-the-art baselines on different terrains both in simulation and in the real world. From the simulation experiments, TERT surpasses the baselines with higher return and less energy consumption. We then compare TERT and the standard student policy on a real Unitree A1 robot. In the real-world evaluation, TERT successfully traverses nine different terrains, including grass, soil, pebble, slope, sand pit, stair down, etc., as Figure 1 shows. While the student policy accomplished tests on other terrains, it failed on sand pit and stair down, verifying the advantage of Transformer architecture compared to conventional models on challenging terrains. To the best of our knowledge, this is the first work that deals with sim-to-real transfer of quadrupedal locomotion over multiple challenging terrains via Transformer-based sequence modeling, which could become a new paradigm for sim-to-real transfer tasks.

II. RELATED WORK

Sim-to-Real RL for Legged Locomotion. Sim-to-Real reinforcement learning (RL) for legged locomotion has recently gained immense interest due to its potential to eliminate the need for human expertise in controller design [1–5, 17–19]. However, it is non-trivial to transfer the policy trained in simulation to real robots due to the discrepancy between simulation and the real world, known as the *reality gap* [20, 21]. Much effort has been devoted to narrowing this reality gap, either by increasing the simulator’s accuracy [2, 4] or by randomizing the physical parameters such as friction and mass during training (domain randomization) [22–25]. Other works are dedicated to introducing more inductive bias to facilitate the training and transferring of policy. For example, Peng et al. [5] trained a robot by imitating the poses of natural animals. Kumar et al. [3], instead, utilized bioenergetics to help design reward function. Our method is orthogonal to them and can naturally take advantage of the above techniques. Another related line of research is system identification, which tries to estimate the raw environmental parameters [26] or the encoded ones [1, 3] from historical observations for online adaptation.

Transformers. Recent advances in Transformers [14] have led to significant breakthroughs in a number of fields. In natural language processing (NLP) [27, 28], Transformer

has shown great superiority against the traditional LSTM [29] or GRU [11] models. Applications of Transformer have increasingly expanded to other areas, such as computer vision (CV) [30–32] and RL. For example, in the RL area, Decision Transformer [15] utilizes an autoregressive model to generate action sequences conditioned on desired return and past state-action pairs. Trajectory Transformer [16], instead, leverages Transformer to generate entire trajectories by discretizing states, actions, and rewards into tokens. Besides, Wen et al. [33] proposed to model multi-agent RL as a sequence modeling problem and resorted to Transformer to solve it. Furthermore, the Gato model [34] successfully trains and deploys one Transformer model on hundreds of tasks, including robot arm manipulation. Our method draws inspiration from Decision Transformer due to its simplicity and computational efficiency, which are critical in real-world applications.

One closely-related work is Yang et al. [35], which built a Transformer model to fuse information for quadrupedal locomotion. TERT differs from theirs mainly in the following three aspects: i) Different modeling - their Transformer takes the different kinds of information in a *single timestep* as input and models the joint patterns at this moment. In contrast, TERT models quadrupedal locomotion as a sequence prediction problem and takes observations of *multiple timesteps* in the history as input, thus better leveraging the superiority of Transformer in sequence modeling. ii) Different training paradigm - Yang et al. [35] trained the Transformer purely via online RL, while TERT adopts a novel two-stage training framework, which can naturally inherit the merit of privileged learning. iii) Different terrain difficulty - compared with Yang et al. [35], TERT exhibits a more agile and mild behavior and successfully traverses more challenging terrains such as sand pit and stair down, even without visual inputs, which are not included in their work.

III. PRELIMINARIES

A. Reinforcement Learning for Quadrupedal Locomotion

We formulate quadrupedal locomotion as a Partially Observable Markov Decision Process (POMDP) defined by the tuple $(\mathcal{S}, \mathcal{O}, \mathcal{A}, T, r, \gamma)$, where \mathcal{S} , \mathcal{O} , and \mathcal{A} are the state, observation, and action spaces, respectively. $T(s_{t+1} | s_t, a_t)$ is the transition density of state s_t given action a_t , and the reward function is denoted as $r(s_t, a_t)$. $\gamma \in (0, 1)$ is a discount factor. At each timestep t , only the partial observation $o_t \in \mathcal{O}$ can be observed instead of s_t due to the limitation of sensors. The goal of reinforcement learning (RL) is to find the optimal policy $\pi^*: \mathcal{O} \rightarrow \mathcal{A}$ that maximizes the expected return (sum of discounted rewards):

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

B. The Transformer Model

Our model draws upon the Decision Transformer [15]. The difference is that we remove the returns-to-go in the trajectory representation since we find that including the

returns-to-go will slightly degrade the performance of our algorithm in preliminary experiments. To be more specific, we use the GPT [36] architecture consisting of n stacked self-attention layers with causal masking. Give a trajectory: $\tau = (o_1, a_1, o_2, a_2, \dots, o_T, a_T)$ of length T . The Transformer first embeds the inputs into $\{x_i\}_{i=1}^{2T}$ with position embedding. The self-attention layers then map each x_i into the corresponding key, query and value, denoted as k_i , q_i , and v_i , respectively, through linear transformation. The i -th output of the self-attention layer is calculated as:

$$z_i = \sum_{j=1}^i \text{softmax}(\langle \{q_i, k_{j'}\}_{j'=1}^i \rangle_j) \cdot v_j, \quad (2)$$

where $\langle \cdot \rangle$ denotes the dot product operation, and the softmax normalization of the dot product represents the attention weight assigned to the j -th token. Note that only the tokens $j \leq i$ are used when calculating z_i since we should ensure that the Transformer can only access preceding tokens to predict the next ones. The output of the last self-attention layer is then fed into another linear layer to predict the actions conditioned on the previous observation-action sequence:

$$\hat{a}_t \sim \text{Trans}(\cdot \mid o_1, a_1, \dots, o_t), \quad \forall t \in [1, T]. \quad (3)$$

The GPT Transformer is trained to minimize the mean square error between the predicted actions \hat{a}_t and the target actions. For evaluation, the Transformer takes the last T observation-action pairs as input and executes the predicted \hat{a}_T in the environment.

IV. TERRAIN TRANSFORMER

In this section, we introduce Terrain Transformer (TERT), a novel approach that combines high-capacity Transformer models with privileged learning through a two-stage training framework. An overview of the training framework of TERT is illustrated in Figure 2.

A. Teacher Policy Training

Like the standard privileged learning framework [1, 3], our method will first train a teacher policy in simulation with access to the privileged information, which mainly consists of three parts: i) elevation map around the robot base; ii) contact force with the ground; iii) ground-truth physical environmental parameters used for domain randomization, such as friction and mass. Following Kumar et al. [3], let e_t denote the privileged information, which will first be processed into a latent vector l_t by a multi-layer perceptron (MLP) encoder μ . Then, the teacher policy $\bar{\pi}$ gives actions conditioned on l_t and proprioception observation o_t :

$$l_t = \mu(e_t), \quad (4)$$

$$\bar{a}_t = \bar{\pi}(o_t, l_t). \quad (5)$$

The teacher policy and encoder are trained jointly via PPO (Proximal Policy Optimization) algorithm [37].

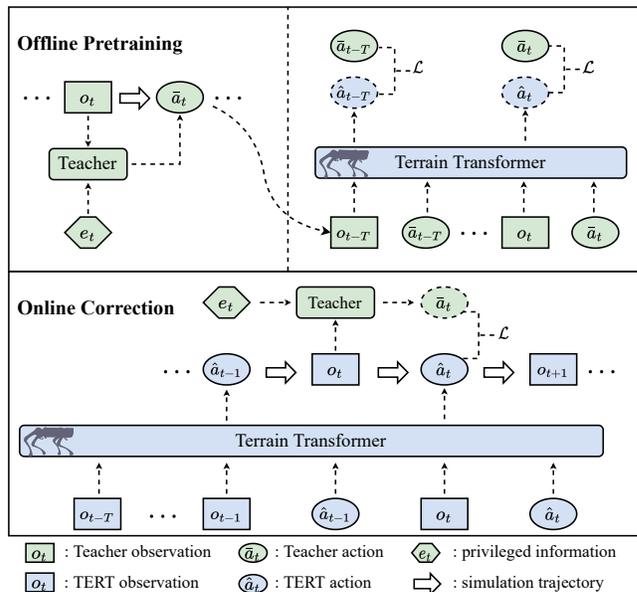


Fig. 2. Illustration of Terrain Transformer training framework. Dashed lines represent actions not executed in the simulator.

B. Terrain Transformer Training

Offline Pretraining. Although teacher policy can achieve near-optimal performance and successfully traverse different terrains in simulation, it can not be directly deployed in the real world since the privileged information is noisy or difficult to obtain [12]. Unlike previous methods that reserve the teacher policy and try to estimate l_t online [1, 3], our method chooses to train a Transformer controller from scratch. Furthermore, we propose a two-stage training framework to better integrate Transformer with privileged learning. Precisely, in the first offline pretraining stage, the teacher policy is executed in simulation over varied terrains with privileged information to collect data (Figure 2 top-left). Then a GPT Transformer is trained on the collected data (observation-action sequences without privileged information) to predict the teacher’s action (Figure 2 top-right):

$$\mathcal{L} = \sum_{t=1}^T (\hat{a}_t - \bar{a}_t)^2, \quad (6)$$

$$(\hat{a}_1, \hat{a}_2, \dots, \hat{a}_T) \sim \text{Trans}(o_1, \bar{a}_1, \dots, o_T, \bar{a}_T),$$

where (o_t, \bar{a}_t) is the observation-action pair in the teacher’s trajectory.

Online Correction. However, a Transformer controller that achieves small prediction losses on teacher’s trajectories does not perform well when tested in the simulator since the input distribution, i.e., the observation-action sequences it encounters, differs from that in the training dataset. Inspired by the DAgger (Data Aggregator) algorithm [38], which queries an expert for action labels of on-policy data, we propose the following online correction stage to address the input distribution shift issue. Specifically, in the online correction stage, the Transformer will then be used to interact with the simulator while the teacher policy gives the target action at the same time. Afterward, the Transformer is trained

to predict the teacher’s action \bar{a}_t conditioned on its own trajectories:

$$\mathcal{L} = \sum_{t=1}^T (\hat{a}_t - \bar{a}_t)^2, \quad (7)$$

$$(\hat{a}_1, \hat{a}_2, \dots, \hat{a}_T) \sim \text{Trans}(o_1, \hat{a}_1, \dots, o_T, \hat{a}_T).$$

Both these two stages play an essential role in the performance of TERT. Without the offline pretraining stage, the initial Transformer controller can hardly give reasonable control commands, and the on-policy data will be restricted and low-quality, thus affecting the efficiency and effectiveness of the algorithm. On the other hand, without the online correction stage, the Transformer model will suffer from the input distribution shift problem, as mentioned before. More experimental comparison of these two variants is provided in Section V-C.

V. EXPERIMENTAL RESULTS

Our experiments aim to answer the following questions:

- How does TERT perform compared with previous state-of-the-art methods in multi-terrain quadrupedal locomotion?
- What are the most valuable components of our overall algorithm?
- Does TERT benefit from the capability and self-attention mechanism of Transformer to model long sequences?
- Could the achievement of TERT in simulation be well transferred to real robots?

To answer these questions, we apply our proposed approach to the Unitree A1 robot [39], which weighs around 12 kg and has 12 motors in total (two for each hip joint and one for each knee joint).

A. Experimental Setup

We implement TERT and baselines based on the open-source codebase in Rudin et al. [40], which leverages the Isaac Gym simulator [41] to support simulation of thousands of robots in parallel and provides multi-terrain simulation, including slopes, stairs, and discrete obstacles. Here we briefly describe some settings that are important for our experiments and defer other details to the original paper.

Observation and Action. Precisely, the proprioception observation $o_t \in \mathbb{R}^{48}$ consists of: base linear and angular velocities, gravity projection, commands, linear and angular velocities of joints, and last action a_{t-1} . The action $a_t \in \mathbb{R}^{12}$ specifies the target positions for each joint, which are then converted to joint torques through a PD controller:

$$\text{torque} = k_p (q_d - q) + k_d (\dot{q}_d - \dot{q}), \quad (8)$$

where q and q_d are the current position and target position for the joint, respectively, and (k_p, k_d) are the hyper-parameters for the PD controller.

Reward Function. We adopt a reward function similar to Rudin et al. [40], which encourages the robot to follow commanded velocities and penalizes velocities along other

axes, large joints torques, accelerations and collisions. We additionally add two reward terms to penalize large action magnitude and encourage smoothness of joint torques [3], which we find helpful for sim-to-real transfer.

Model Architecture. For teacher policy, the privileged information e_t is first encoded by a three-layer MLP with 256 hidden units into a latent vector $l_t \in \mathbb{R}^{12}$. Then l_t is concatenated with o_t before being passed to a three-layer MLP with (512, 256, 128) hidden units. For Transformer, we use 3 self-attention layers with embedding size 256. The dropout rate is set to 0.05 and the trajectory length T is set to 20.

B. Simulation Comparison

We first evaluate our method and baselines in terms of return, control smoothness and energy consumption over different terrains in simulation, where the control smoothness is defined as the average difference between a_t and a_{t-1} . The baselines we compared include:

- **Teacher.** The teacher policy with privileged information serves as an oracle.
- **RMA.** We implement the student policy according to the RMA algorithm [3], which trains a Temporal Convolutional Network (TCN) to estimate l_t using the last 50 steps (on-policy data) and removes the reliance on predefined trajectory generator used in Lee et al. [1].
- **PPO.** A policy directly trained via PPO with o_t as input.
- **StackedPPO.** PPO policy with stacked historical observations and actions as input. Note that PPO can be viewed as a special case of StackedPPO with the stacked length set to 1.
- **GRU.** Replacing the multi-layer perceptron (MLP) in PPO with the GRU [11] model.

To ensure a fair comparison, we use the same teacher policy, i.e., the one listed in Table I, to train both RMA and TERT.

The comparison results are shown in Table I. Our method TERT outperforms other baselines in terms of return on most terrains except rough slope, especially in more challenging scenarios such as stair up and stair down. The gap between TERT and teacher is smaller than that between RMA and teacher, verifying the advantage of Transformer in sequence modeling compared to traditional models such as TCN. Moreover, Stacking past observations and actions or processing them with GRU helps address the partially observable problem but is not comparable to TERT or RMA without the guidance of privileged information. In addition, TERT achieves lower energy consumption and smoother control strategy, though our reward function does not penalize energy consumption explicitly.

C. Empirical Study

In this section, we provide empirical studies to understand the importance of each component and the benefit our method derives from the Transformer architecture.

Ablation Study. We first conduct experiments to evaluate the main components of our algorithm, i.e., the two-stage

TABLE I

COMPARISON RESULTS ON DIFFERENT TERRAINS IN TERMS OF RETURN, SMOOTHNESS, AND ENERGY CONSUMPTION. RESULTS ARE AVERAGED OVER 1000 TRAJECTORIES WITH DIFFERENT DIFFICULTIES. BOLD NUMBERS INDICATE THE BEST SCORES AMONG ALGORITHMS EXCLUDING TEACHER.

	Terrain	Teacher	TERT (ours)	RMA	PPO	StackedPPO	GRU
Return	Smooth Slope	21.33 ± 2.40	21.26 ± 2.51	21.11 ± 2.71	17.30 ± 2.67	19.97 ± 2.69	20.30 ± 2.74
	Rough Slope	20.38 ± 2.64	20.11 ± 3.15	20.21 ± 2.74	16.27 ± 2.89	19.28 ± 2.90	19.54 ± 2.73
	Stair Up	20.49 ± 1.24	19.66 ± 2.26	19.20 ± 2.64	16.43 ± 2.38	18.82 ± 2.13	17.69 ± 3.02
	Stair Down	20.15 ± 3.67	19.26 ± 4.49	18.79 ± 4.78	17.36 ± 4.26	19.08 ± 4.10	19.13 ± 4.04
	Obstacle	22.10 ± 1.98	21.97 ± 2.28	21.71 ± 2.35	17.34 ± 2.72	20.61 ± 1.95	20.55 ± 2.62
	Average	20.89	20.43	20.20	16.94	19.55	19.44
Smooth	Smooth Slope	1.30 ± 0.14	1.26 ± 0.12	1.75 ± 0.15	1.71 ± 0.52	1.55 ± 0.69	1.38 ± 0.25
	Rough Slope	1.35 ± 0.16	1.31 ± 0.14	1.82 ± 0.18	1.77 ± 0.64	1.67 ± 1.03	1.48 ± 0.19
	Stair Up	1.36 ± 0.08	1.38 ± 0.11	1.80 ± 0.16	2.09 ± 0.49	1.70 ± 0.54	1.32 ± 0.28
	Stair Down	1.37 ± 0.22	1.30 ± 0.22	1.85 ± 0.59	1.68 ± 0.48	2.23 ± 1.92	1.57 ± 0.29
	Obstacle	1.25 ± 0.14	1.21 ± 0.10	1.71 ± 0.20	1.72 ± 0.65	1.44 ± 0.51	1.30 ± 0.20
	Average	1.32	1.29	1.78	1.79	1.72	1.41
Energy	Smooth Slope	39.20 ± 5.29	38.93 ± 5.75	38.91 ± 5.69	55.40 ± 10.85	41.67 ± 6.96	42.27 ± 6.29
	Rough Slope	39.57 ± 6.07	38.86 ± 6.78	39.97 ± 6.68	56.41 ± 13.01	43.67 ± 7.51	44.34 ± 5.87
	Stair Up	44.43 ± 4.06	44.57 ± 7.65	41.62 ± 10.54	65.74 ± 14.45	47.69 ± 7.64	40.11 ± 12.01
	Stair Down	37.62 ± 8.03	37.09 ± 10.61	40.37 ± 11.39	55.01 ± 12.98	47.23 ± 14.20	48.21 ± 9.59
	Obstacle	39.43 ± 5.23	38.61 ± 5.59	39.06 ± 5.94	56.33 ± 12.30	40.47 ± 6.59	39.74 ± 5.20
	Average	40.05	39.61	39.98	57.77	44.15	42.93

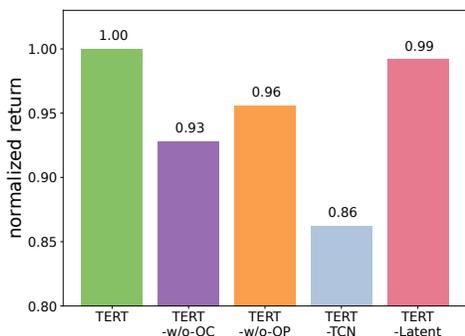


Fig. 3. Average normalized return between $[0, 1]$ of TERT and its ablated variants.

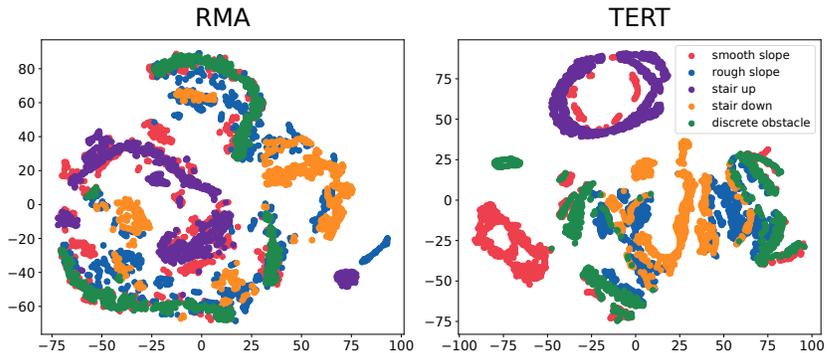


Fig. 4. T-SNE visualization for the last hidden layer of RMA and TERT on different terrains. Points of RMA on five terrains largely overlap while TERT holds more clear boundaries.

training framework, the Transformer model, and the end-to-end imitation scheme. Specifically, we compare the original TERT with four variants: i) **TERT-w/o-OC**. Variant of TERT without online correction (OC) training stage. ii) **TERT-w/o-OP**. Variant of TERT without offline pretraining (OP) training stage. iii) **TERT-TCN**. Replacing Transformer in TERT with TCN model. iv) **TERT-Latent**. Utilizing Transformer to estimate the latent vector like RMA instead of imitating the teacher’s action end-to-end. Results are shown in Figure 3. We find that ablating either the offline pretraining or online correction stage decreases the performance to some extent, validating the effectiveness of our proposed training framework, as discussed in Section IV-B. Furthermore, removing the Transformer causes severe performance degradation, while using Transformer to estimate the latent vector retains most of the advantages of TERT, which further reveals the superiority of Transformer architecture in sequence modeling.

Hidden Layer Visualization. We then unroll TERT and the RMA baseline over different terrains and visualize their last hidden layer in Figure 4. As the t-SNE result shows, the

points of RMA on five terrains largely overlap, while our method holds more clear boundaries. We can also observe a slight overlap for TERT on rough slope and discrete obstacle since both terrains need to tackle uneven ground. This result helps explain why RMA performs excellently on some terrains, such as smooth and rough slopes, but not so well on others. In contrast, TERT can leverage the high-capacity Transformer to better capture the characteristics of the environment and model diverse behaviors, which may help generalize to different terrains.

Attention Weight. The self-attention mechanism endows Transformer models with the ability of context-dependent control [42], which is important for multi-terrain locomotion. We plot the attention weights along a stair-down trajectory of 200 timesteps in Figure 6. As the figure shows, TERT allocates different attention weights over time. Notice that weights assigned to earlier timesteps (e.g., $t - 19$) are comparable to those assigned to later timesteps, indicating that Transformer can make good use of long historical information. Besides, the attention weight shows a cyclical pattern, reflecting the corresponding periodic gait of the

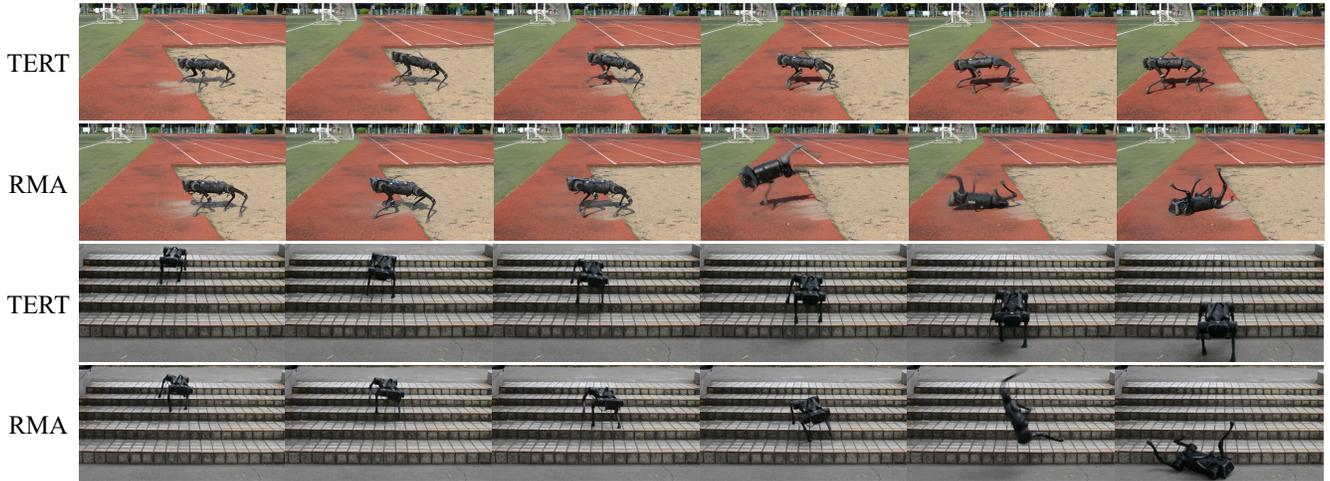


Fig. 5. Snapshots of TERT and RMA traversing across sand pit (top) and stair down (bottom) terrain.

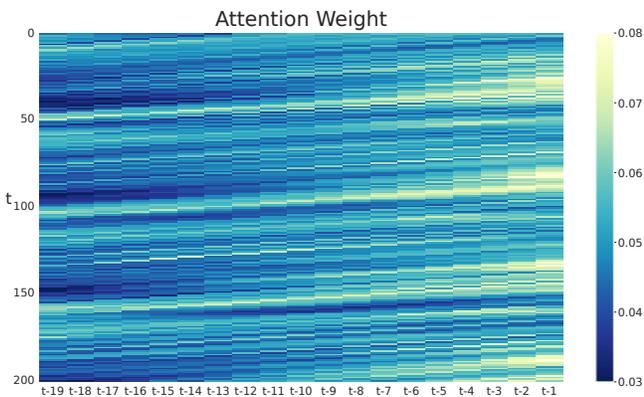


Fig. 6. Heatmap of Transformer attention weight on stair down terrain. Y-axis indicates different timestep t , and x-axis indicates the corresponding history sequences (20 timesteps in our setting). For example, the color of the square at 50-th row and 10-th column represents the attention weight Transformer assigns for (o_{40}, a_{40}) to predict a_{50} .

quadrupedal robot.

D. Real-world Application

We apply TERT and RMA policy to an A1 robot without any fine-tuning in the real world. Both methods are directly run on the A1 hardware without external computing devices. The base linear velocities in o_t are estimated from the accelerometer, and other proprioception information is read from the IMU and other sensors. We set the control frequency to 50Hz and $k_p = 55, k_d = 8$ for both methods.

TABLE II

REAL-WORLD TRAVERSING SUCCESS RATES ON DIFFERENT TERRAINS. RESULTS ARE AVERAGED OVER FIVE TRIALS TO REDUCE DAMAGE TO THE ROBOT HARDWARE.

	Sand Pit	Stair Down	Other Terrains
TERT (ours)	100%	60%	100%
RMA	0%	0%	100%

We select nine different terrains for a comprehensive evaluation, including grass, soil, pebble, smooth floor, damp

wood, slope up, slope down, sand pit and stair down¹. The success rates are listed in Table II. Both methods can pass through the first seven terrains with a 100% success rate, while RMA fails on sand pit and stair down without exception. We attribute this failure to the sensitivity of RMA to the estimated latent vector. In contrast, our method exhibits a more stable control behavior and successfully traverses more challenging terrains, only fails sometimes on stair down. We provide snapshots of trials on sand pit and stair down in Figure 5. Please refer to the supplemental video for detailed comparisons on more terrains.

VI. CONCLUSION

In this paper, we propose Terrain Transformer (TERT), a simple yet effective method to leverage Transformer for quadrupedal locomotion over multiple terrains, including a two-stage training framework to incorporate Transformer with privileged learning. Extensive simulation and real-world experiments demonstrate that TERT outperforms the state-of-the-art baseline on different challenging terrains. Further empirical analyses provide insight that the main superiority of TERT lies in the high capacity for sequence modeling and the self-attention mechanism of Transformer. We believe our work takes an important step towards application of Transformer in quadruped robot locomotion, which could become a new paradigm of sim-to-real transfer for robot control. For future work, it is appealing to integrate elevation information from LiDAR or camera into our framework and investigate its extension to multi-agent control tasks.

ACKNOWLEDGEMENTS

The SJTU Team is supported by Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102) and National Natural Science Foundation of China (62076161). The author Hang Lai is supported by Wu Wen Jun Honorary Doctoral Scholarship, AI Institute, Shanghai Jiao Tong University.

¹We do not include stair up terrain in the comparison since it is too difficult for both methods to accomplish it in the real world without elevation information [3].

REFERENCES

- [1] J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning quadrupedal locomotion over challenging terrain," *Science robotics*, vol. 5, no. 47, eabc5986, 2020.
- [2] 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," *arXiv preprint arXiv:1804.10332*, 2018.
- [3] A. Kumar, Z. Fu, D. Pathak, and J. Malik, "Rma: Rapid motor adaptation for legged robots," *arXiv preprint arXiv:2107.04034*, 2021.
- [4] 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, eaau5872, 2019.
- [5] X. B. Peng, E. Coumans, T. Zhang, T.-W. Lee, J. Tan, and S. Levine, "Learning agile robotic locomotion skills by imitating animals," *arXiv preprint arXiv:2004.00784*, 2020.
- [6] W. Yu, V. C. Kumar, G. Turk, and C. K. Liu, "Sim-to-real transfer for biped locomotion," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2019, pp. 3503–3510.
- [7] W. Yu, J. Tan, Y. Bai, E. Coumans, and S. Ha, "Learning fast adaptation with meta strategy optimization," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 2950–2957, 2020.
- [8] A. Nagabandi, I. Clavera, S. Liu, R. S. Fearing, P. Abbeel, S. Levine, and C. Finn, "Learning to adapt in dynamic, real-world environments through meta-reinforcement learning," *arXiv preprint arXiv:1803.11347*, 2018.
- [9] X. Song, Y. Yang, K. Choromanski, K. Caluwaerts, W. Gao, C. Finn, and J. Tan, "Rapidly adaptable legged robots via evolutionary meta-learning," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, IEEE, 2020, pp. 3769–3776.
- [10] S. Bai, J. Z. Kolter, and V. Koltun, "An empirical evaluation of generic convolutional and recurrent networks for sequence modeling," *arXiv preprint arXiv:1803.01271*, 2018.
- [11] J. Chung, C. Gulcehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *arXiv preprint arXiv:1412.3555*, 2014.
- [12] T. Miki, J. Lee, J. Hwangbo, L. Wellhausen, V. Koltun, and M. Hutter, "Learning robust perceptive locomotion for quadrupedal robots in the wild," *Science Robotics*, vol. 7, no. 62, eabk2822, 2022.
- [13] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *International Conference on Machine Learning*, PMLR, 2021, pp. 8821–8831.
- [14] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [15] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," *Advances in neural information processing systems*, vol. 34, pp. 15 084–15 097, 2021.
- [16] M. Janner, Q. Li, and S. Levine, "Offline reinforcement learning as one big sequence modeling problem," *Advances in neural information processing systems*, vol. 34, pp. 1273–1286, 2021.
- [17] X. Da, Z. Xie, D. Hoeller, B. Boots, A. Anandkumar, Y. Zhu, B. Babich, and A. Garg, "Learning a contact-adaptive controller for robust, efficient legged locomotion," *arXiv preprint arXiv:2009.10019*, 2020.
- [18] Y. Yang, T. Zhang, E. Coumans, J. Tan, and B. Boots, "Fast and efficient locomotion via learned gait transitions," in *Conference on Robot Learning*, PMLR, 2022, pp. 773–783.
- [19] H. Shi, B. Zhou, H. Zeng, F. Wang, Y. Dong, J. Li, K. Wang, H. Tian, and M. Q.-H. Meng, "Reinforcement learning with evolutionary trajectory generator: A general approach for quadrupedal locomotion," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3085–3092, 2022.
- [20] S. Koos, J.-B. Mouret, and S. Doncieux, "Crossing the reality gap in evolutionary robotics by promoting transferable controllers," in *Proceedings of the 12th annual conference on Genetic and evolutionary computation*, 2010, pp. 119–126.
- [21] A. Boeing and T. Bräunl, "Leveraging multiple simulators for crossing the reality gap," in *2012 12th international conference on control automation robotics & vision (ICARCV)*, IEEE, 2012, pp. 1113–1119.
- [22] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel, "Sim-to-real transfer of robotic control with dynamics randomization," in *2018 IEEE international conference on robotics and automation (ICRA)*, IEEE, 2018, pp. 3803–3810.
- [23] Z. Xie, X. Da, M. van de Panne, B. Babich, and A. Garg, "Dynamics randomization revisited: A case study for quadrupedal locomotion," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, IEEE, 2021, pp. 4955–4961.
- [24] W. Zhao, J. P. Queralta, and T. Westerlund, "Sim-to-real transfer in deep reinforcement learning for robotics: A survey," in *2020 IEEE Symposium Series on Computational Intelligence (SSCI)*, IEEE, 2020, pp. 737–744.
- [25] Y. Chebotar, A. Handa, V. Makoviychuk, M. Macklin, J. Issac, N. Ratliff, and D. Fox, "Closing the sim-to-real loop: Adapting simulation randomization with real world experience," in *2019 International Conference on Robotics and Automation (ICRA)*, IEEE, 2019, pp. 8973–8979.
- [26] W. Yu, J. Tan, C. K. Liu, and G. Turk, "Preparing for the unknown: Learning a universal policy with online system identification," *arXiv preprint arXiv:1702.02453*, 2017.
- [27] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, *et al.*, "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [28] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [29] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [30] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [31] K. Lin, L. Wang, and Z. Liu, "End-to-end human pose and mesh reconstruction with transformers," *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.
- [32] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, "Masked autoencoders are scalable vision learners. arxiv 2021," *arXiv preprint arXiv:2111.06377*.
- [33] M. Wen, J. G. Kuba, R. Lin, W. Zhang, Y. Wen, J. Wang, and Y. Yang, "Multi-agent reinforcement learning is a sequence modeling problem," *arXiv preprint arXiv:2205.14953*, 2022.
- [34] S. Reed, K. Zolna, E. Parisotto, S. G. Colmenarejo, A. Novikov, G. Barth-Maron, M. Gimenez, Y. Sulsky, J. Kay, J. T. Springenberg, *et al.*, "A generalist agent," *arXiv preprint arXiv:2205.06175*, 2022.
- [35] R. Yang, M. Zhang, N. Hansen, H. Xu, and X. Wang, "Learning vision-guided quadrupedal locomotion end-to-end with cross-modal transformers," *International Conference on Learning Representation*, 2022.
- [36] A. Radford, K. Narasimhan, T. Salimans, I. Sutskever, *et al.*, "Improving language understanding by generative pre-training," 2018.
- [37] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [38] S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, JMLR Workshop and Conference Proceedings, 2011, pp. 627–635.
- [39] Unitree, *Unitree robotics*, <https://www.unitree.com/>, 2022.
- [40] N. Rudin, D. Hoeller, P. Reist, and M. Hutter, "Learning to walk in minutes using massively parallel deep reinforcement learning," in *Conference on Robot Learning*, PMLR, 2022, pp. 91–100.
- [41] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, *et al.*, "Isaac gym: High performance gpu-based physics simulation for robot learning," *arXiv preprint arXiv:2108.10470*, 2021.
- [42] L. C. Melo, "Transformers are meta-reinforcement learners," in *International Conference on Machine Learning*, PMLR, 2022, pp. 15 340–15 359.

APPENDIX

A. Training Details

TABLE III
DOMAIN RANDOMIZATION RANGE

Parameter	Training Range	Testing Range
Friction	[0.5, 1.25]	[0.1, 2.0]
Added mass (Kg)	[0, 5]	[0, 7]
k_p	[45, 65]	[40, 70]
k_d	[0.7, 0.9]	[0.6, 1.0]

Following Rudin et al. [40], we create 4096 environment instances to collect data in parallel. Each environment is chosen from five types of terrain (smooth slope, rough slope, stair up, stair down, discrete obstacle) with different difficulty levels. In each environment, the robot is initialized with random poses and commanded to walk forward at a speed of 0.4m/s. The robot receives new observations and updates its actions every 0.02 seconds. An episode is terminated if the robot trunk touches the ground or 1000 timesteps (20 seconds in simulation) are collected. We also randomize the physical parameters to further improve the robustness of the policy. The domain randomization range is listed in Table III. We train the teacher policy for 20k iterations; each iteration corresponds to 98304 (4096×24) timesteps. At each stage of Transformer training, we collect a dataset of 20 million timesteps and then train the Transformer for 200k updates with batch size 64.

B. Sequence Length

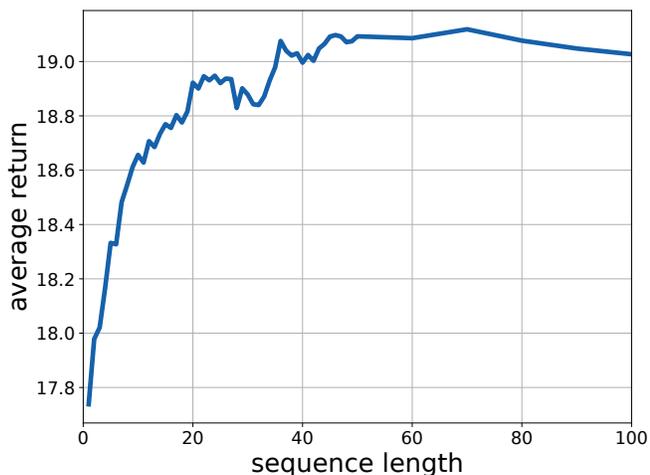


Fig. 7. Average return of TERT after offline pretraining with different sequence length T .

Sequence Length. We further conduct experiments of TERT with different sequence length T . Results are shown in Figure 7. Note that when $T = 1$, TERT degenerates into one-step behavior cloning. We observe that up to a certain level, increasing the sequence length T yields better performance, which reveals the advantage of longer sequence modeling. However, too large T will affect the computation efficiency

of the algorithm, thereby decreasing the control frequency of the robot. Therefore to trade off these two terms, we choose $T = 20$ in our experiments, which achieves acceptable performance according to Figure 7.