

Preference Aligned Diffusion Planner for Quadrupedal Locomotion Control

Xinyi Yuan^{*1}, Zhiwei Shang^{*2}, Zifan Wang², Chenkai Wang⁴, Zhao Shan³,
Meixin Zhu^{†5}, Chenjia Bai^{†3}, Weiwei Wan¹, Kensuke Harada¹, Xuelong Li³

Abstract—Diffusion models demonstrate superior performance in capturing complex distributions from large-scale datasets, providing a promising solution for quadrupedal locomotion control. However, the robustness of the diffusion planner is inherently dependent on the diversity of the pre-collected datasets. To mitigate this issue, we propose a two-stage learning framework to enhance the capability of the diffusion planner under limited dataset (reward-agnostic). Through the offline stage, the diffusion planner learns the joint distribution of state-action sequences from expert datasets without using reward labels. Subsequently, we perform the online interaction in the simulation environment based on the trained offline planner, which significantly diversified the original behavior and thus improves the robustness. Specifically, we propose a novel *weak* preference labeling method without the ground-truth reward or human preferences. The proposed method exhibits superior stability and velocity tracking accuracy in pacing, trotting, and bounding gait under different speeds and can perform a zero-shot transfer to the real Unitree Go1 robots. The project website for this paper is at <https://shangjaven.github.io/preference-aligned-diffusion-legged/>.

I. INTRODUCTION

Learning-based approaches significantly enhance the agility and adaptability of quadrupedal robots to accomplish diverse locomotion tasks [1], [2]. While online learning demonstrates robustness in complex dynamic environments, extensive trial-and-error interactions in simulation are required to learn an effective policy. Thus, learning an online policy can be sample inefficient and requires a meticulously designed reward function. In contrast, offline learning can leverage the advantages of pre-collected offline datasets via model-based controller [3], animal imitation [4], or Reinforcement Learning (RL) policy [5], significantly improving data efficiency and reducing the cost of online interactions [6]. In offline policy learning, diffusion models [7] have shown superior performance in capturing complex action distributions from offline trajectories [8], exhibit promising potential to solve quadruped locomotion tasks with high-dimensional action space and complex action distribution in various terrains. As an example, DiffuseLoco [9] has

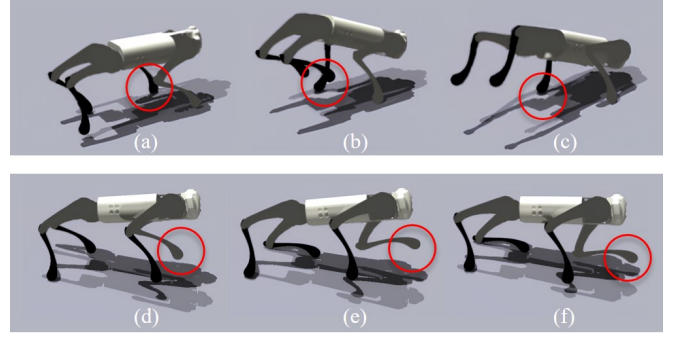


Fig. 1. Failure Cases of the Offline Diffusion Planner under a Single-Source Limited Dataset: (a-c) Bounding Gait Example: The lift delay of the highlighted leg (red circle) deteriorates in subsequent motion sequences, ultimately causing the legged robot to fall. (d-f) Pacing Gait Example: The highlighted leg (red circle) demonstrates a higher lifting compared to the average step height, leading to severe lateral tilting.

recently proposed to train a diffusion planner for quadrupedal locomotion from large-scale offline trajectories.

However, learning robust locomotion planners in a purely offline manner requires collecting a large-scale dataset with wide state coverage [10], especially covering situations that the robot may encounter in the real world. This exposes significant challenges in quadrupedal locomotion scenario, as it's difficult to comprehensively cover gait variations and diverse terrain dynamics. Therefore, it is challenging to learn a robust locomotion policy using the limited offline dataset with inadequate state coverage, as the learned policy may be sensitive to out-of-distribution (OOD) states in the real world. DiffuseLoco [9] addressed the challenge by leveraging a large-scale multi-skill dataset incorporating both quadrupedal and bipedal locomotion strategies from various source policies. The comprehensive state coverage leads to robust offline planner, enabling successful sim-to-real transfer. Rather than extensively increasing the dataset distribution, we aim to address the aforementioned issue by combining offline diffusion modeling with online interactions, where the collection of online trajectories augments the coverage of the state-action space distribution to align with the real world. Notably, unlike previous online refinement algorithm that requires real-time action or reward labels, our proposed framework eliminates the need for external expert labels or reward functions.

In this paper, we propose a two-stage learning framework combining offline and online learning for legged locomotion. In the offline learning stage, the diffusion planner learns the

^{*}Equal Contribution.

[†]Corresponding Authors: Meixin Zhu (e-mail: meixin@ust.hk), Chenjia Bai (e-mail: baichenjia255@gmail.com)

¹ Graduate School of Engineering Science, Osaka University, Japan.

² Robotics and Autonomous Systems Thrust, The Hong Kong University of Science and Technology (Guangzhou), China.

³ Institute of Artificial Intelligence (TeleAI), China Telecom, China

⁴ Department of Statistics and Data Science, Southern University of Science and Technology, China

⁵ School of Transportation, Southeast University, China

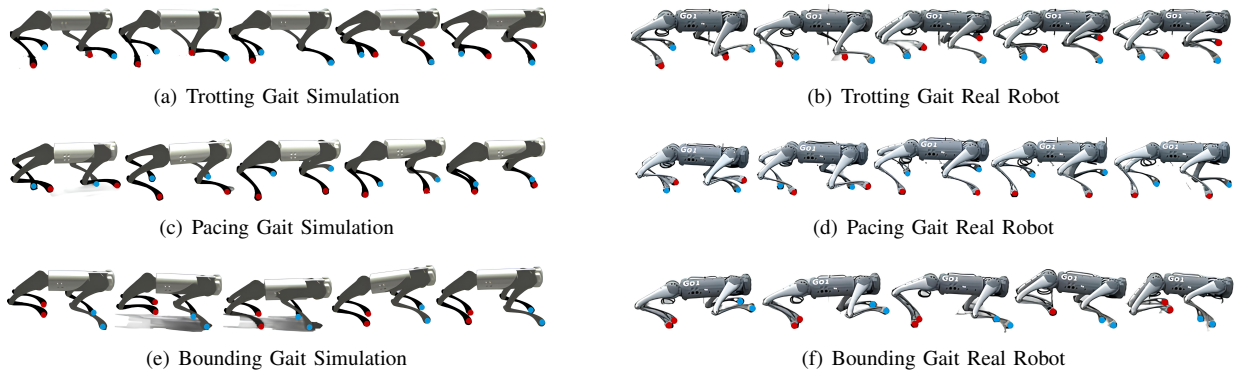


Fig. 2. Video Frames in Simulation and Real World Experiments of the Proposed Architecture: (a-b) Trotting gait simulation and real-world test, (c-d) Pacing gait simulation and real-world test, (e-f) Bounding gait simulation and real-world test

joint distribution of state-action sequence from an limited offline dataset collected by a single pre-trained source policy and does not use reward labels in training. Then, we perform online interactions with the environment based on the diffusion planner, enhancing the diversity of the state distribution and improving robustness via preference alignment. Specifically, we propose a novel online fine-tuning algorithm for the diffusion planner based on preferences, which resembles Direct Preference Optimization (DPO) [11] for diffusion models in text-to-image generation [12]. Importantly, the preference score in our method is measured by distances between the states and the nearest neighbors of expert trajectories, which signifies the optimality of a trajectory and is used as the metric to construct preference pairs. As a result, the preference data can be easily constructed from such a *weak* preference label without the ground-truth reward function or human preferences.

The contributions are summarized as follows: (i) We propose a novel two-stage learning framework that combines online interaction and offline diffusion learning to enhance the diversity and robustness of the diffusion planner under the limited expert dataset; (ii) We give an efficient preference alignment algorithm for offline diffusion planner via DPO and *weak* preference labels; (iii) The resulting diffusion planner exhibits superior performance on stability and velocity tracking accuracy in simulation and can be zero-shot transferred to the real-world Unitree robot.

II. RELATED WORK

A. Learning-Based Approaches for Legged Locomotion

For legged locomotion, learning-based methods automatically capture dynamic behaviors from interacting experiences, largely reducing the need for manual expertise in classical control [13]. Online RL has been widely applied in learning complex locomotion skills in simulation, and adaptive techniques like domain randomization are employed to facilitate the transfer from the simulation to real robot [14], [15], [16], [17], [18]. However, manually designing reward functions and tuning weights can be particularly challenging when dealing with complex tasks. Other methods [19] adopt imitation learning to extract agile locomotion strategies from

real-world animal reference motions, while it requires the motion-capture dataset that is more expensive [20]. Recent studies utilized offline learning on locomotion control with the limited scope confined to gym simulation environments [21], [22]. While DiffuseLoco [9] demonstrated the strong performance of diffusion locomotion planner, it requires large-scale datasets among various source policies (i.e., three different RL-based methods with 14 million data in total).

B. Diffusion Models in Robotics

Diffusion models have demonstrated superior generative capabilities in various robotics tasks such as robotic navigation [23], manipulation [24], [25], and decision-making [26]. For example, Diffuser [27] is a planner that analogizes the planning to the denoising process in diffusion models, demonstrating impressive adaptability in complex long-horizon manipulation tasks. Recent studies [28], [29] represent robot policies also as the diffusion process, where the policy generates joint actions based on multi-modal conditional inputs such as observations or visual information. However, most existing research has been limited to high-level tasks with low-dimensional action spaces, leaving room for exploration in more complex, high-dimensional scenarios such as legged locomotion.

III. TWO-STAGE LEARNING FRAMEWORK

A. System Design

The architecture of the proposed framework is illustrated in Fig. 3. The entire system consists of four stages. Firstly, preference-free offline datasets (i.e., \mathcal{D}) containing various gait patterns are collected in the *Walk-These-Ways* environment. Secondly, behavior cloning is performed by extracting gait-specific diffusion policies from these datasets. Following this, the diffusion policy is fine-tuned using the direct preference optimization method on the constructed preference dataset (i.e., $\mathcal{D}_{\text{pref}}$). Finally, the latest model will be deployed on the Unitree Go1 robot.

B. Offline Dataset Generation

This work utilizes *Walk These Ways* framework [14] as the source locomotion policy and collects offline datasets

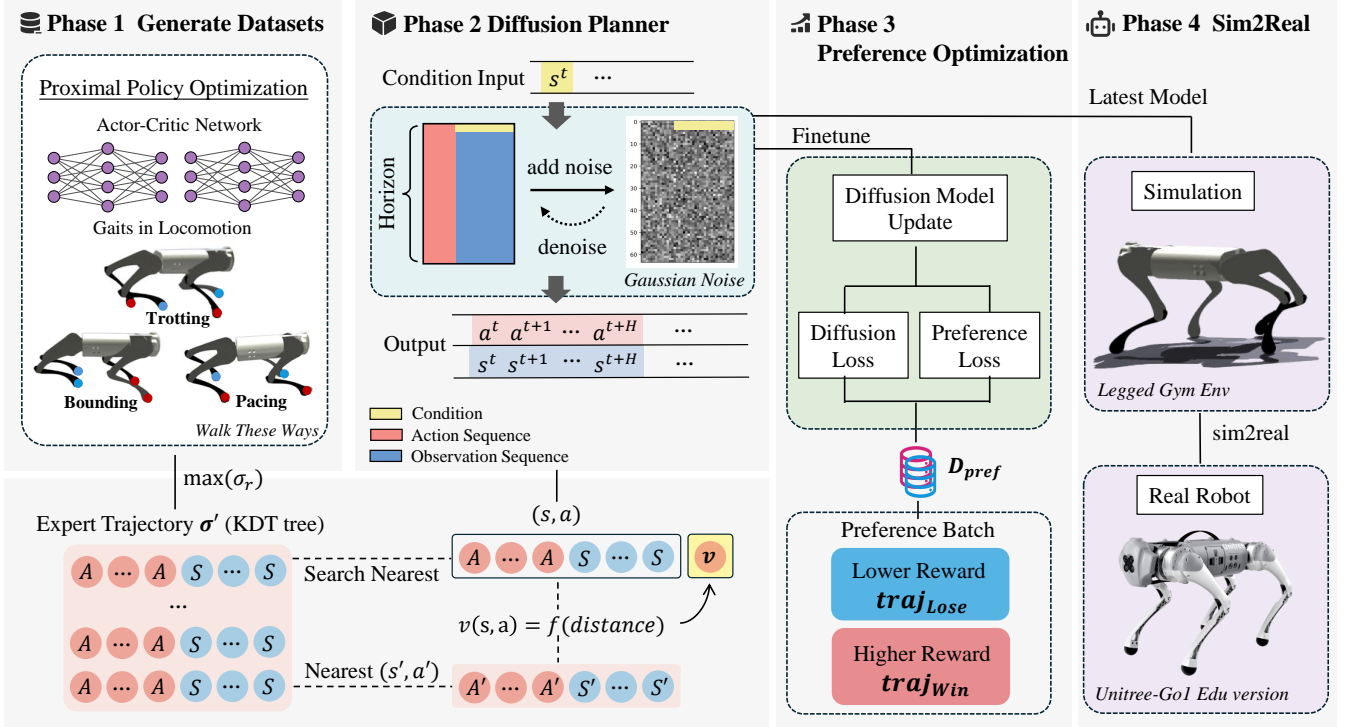


Fig. 3. Proposed Architecture Framework Overview: (1) **Generate Datasets**: the offline datasets among pacing, trotting, and bounding gait are collected through the expert PPO policy in the *walk-these-ways* task. (2) **Behavior Cloning**: given a condition input s^t , the diffusion policy can produce a sequence of states and actions. (3) **Preference Alignment**: Conduct the preference alignment on the offline diffusion planner based on proposed *weak* preference labels. (4) **Sim2Real**: The refined policy is deployed on the Unitree Go1 robot.

for the following gaits: pacing, trotting, and bounding. For each gait, we roll out 2048 episodes with 250 time steps. The data format as $(s_t, \mathbf{a}_t, s_{t+1}, \mathbf{r}_t, \mathbf{d}_t)$, where the \mathbf{d}_t turns True when the episode ends. The behavior policy is denoted as $\pi(\mathbf{a}_t | \mathbf{c}_t, \mathbf{b}_t)$, where \mathbf{c}_t includes commands, and \mathbf{b}_t indicates the behavior parameters. The reward functions of the source locomotion controller include task reward (tracking the command speed), augmented auxiliary (task behavior related), and fixed auxiliary (promoting robot stability).

The Unitree Go1 robot possesses 12 degrees of freedom, with each leg comprising 3 degrees of freedom, corresponding to the hip, thigh, and calf joints. The observation space \mathbf{O}_t of offline diffusion planner consists of the state vector $\mathbf{s}_t = \{\mathbf{v}_t^{cmd}, \mathbf{q}_t, \dot{\mathbf{q}}_t, \mathbf{t}_t\} \in \mathbb{R}^{31}$. Specifically, $\mathbf{v}_t^{cmd} \in \mathbb{R}^3$ incorporates the velocity command along x, y and z axis, $\mathbf{q}_t \in \mathbb{R}^{12}$ represent the joint positions, $\dot{\mathbf{q}}_t \in \mathbb{R}^{12}$ represent the joint velocities, $\mathbf{t}_t \in \mathbb{R}^4$ is the gait and timing related parameter. Each feature of the observations is normalized by a Gaussian distribution before input into the diffusion models.

The action space $\mathbf{A}_t \in \mathbb{R}^{12}$ represents each joint's position targets. We follow the actuator network [2] for actions-to-torques mapping. The legged robot will then execute the resultant torque.

C. Offline Diffusion Planner

We adopt a conditional diffusion model to train an offline diffusion planner for locomotion. The conditional input con-

tains the current observation state \mathbf{s}_0 (ensuring the generated trajectory starts from the current state), and the output is the action sequence $\mathbf{a}_{0:h}$ of length horizon.

The proposed framework uses *Denosing Diffusion Probabilistic Models* (DDPM) [30] with U-net and Transformer backbone. DDPM starts by adding Gaussian noise to the original data stepwise and uses the neural network to learn the inverse denoising process. During the forward process, Gaussian noise will be added gradually to the original data:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad (1)$$

where β_t is the variance scheduler, \mathbf{x}_t and \mathbf{x}_{t-1} are samples from two adjacent diffusion steps.

During the reverse chain, the model starts from the pure Gaussian noise \mathbf{x}_T and gradually extracts the noise to derive $\mathbf{x}_{T-1} \dots \mathbf{x}_0$. DDPM uses $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$ to approximate the conditional denoising distribution:

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t)). \quad (2)$$

For training, we use the simplified objective without the weighting term [30] as

$$\mathcal{L}_t^{\text{simple}}(\theta) = \mathbb{E}_{\mathbf{x}_0, \epsilon, t} [\|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2]. \quad (3)$$

During the training process, the data loader will randomly sample the normalized trajectory segment from the offline dataset with batch size N , and then Gaussian noise will be added iteratively into the trajectories. Subsequently, the diffusion model will denoise the noisy trajectory to reconstruct the

original data. The diffusion loss measures the Mean-Squared Error (MSE) between the true noise and predicted noise:

$$\mathcal{L} = \text{MSE}(\epsilon_k, \epsilon_\theta(\mathbf{a}_k^{t:t+h}, \mathbf{s}_0^t, k)) \quad (4)$$

Notably, ϵ_θ is the noise predictor, the superscript represents the time step in the trajectory (start from time t with the horizon h), while the subscript denotes the diffusion denoising iterations (k). Besides, a fixed mask will be applied on the loss to ignore the deviations in the observation condition (s_0).

During the inference, the Isaac Gym environment will be reset initially, and the current observation serves as the conditional input for trajectory generation. Then, the trained diffusion model will denoise the pure Gaussian noise to derive the action sequence. Specifically, Classifier-Free Guidance (CFG) [31] is employed to blend the conditioned and unconditioned predictions with the weight w :

$$\bar{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{y}) = (1 + w)\epsilon_\theta(\mathbf{x}_t, t, \mathbf{y}) - w\epsilon_\theta(\mathbf{x}_t, t) \quad (5)$$

where $\bar{\epsilon}_\theta$ is the predictor under classifier-free guidance, y represent the conditional information.

After generating the output action sequence of length horizon (h), the robot will execute every action stepwise and record the state transitions. The inference and interaction will continue until the episode length reaches the pre-defined maximum length. To ensure real-time performance, we decrease the diffusion sampling step in inference to 10 steps without sacrificing the locomotion performance. Besides, since the whole generated action sequence will be executed, the required inference frequency within each episode will be reduced compared with only implementing the single action.

IV. PREFERENCE ALIGNED DIFFUSION PLANNER

To enhance the state distribution diversity and policy robustness, we propose the preference alignment method to fine-tune the offline diffusion planner. Firstly, we roll out preference-free datasets \mathcal{D} from the pre-trained offline diffusion planner. When constructing the preference dataset $\mathcal{D}_{\text{pref}}$, we randomly sample two segments from \mathcal{D} without replacement and then assign preference labels based on the following method.

Considering that the ground-truth rewards are sometimes challenging to obtain, we propose a *weak* preference labeling method that allows for constructing preference labels without the reward labels, requiring only a few expert trajectories. The optimal expert trajectory is selected from the given expert trajectories based on their cumulative reward. For each state-action pair (\mathbf{s}, \mathbf{a}) in the preference-free dataset \mathcal{D} , we search for the closest state-action pair $(\mathbf{s}', \mathbf{a}')$ in the *optimal* expert trajectory and calculate their Euclidean distance. The value (v) of specific (\mathbf{s}, \mathbf{a}) is calculated as the equation:

$$v_t = \exp\left(-\frac{\beta \times d_t}{|\mathcal{A}|}\right) \quad (6)$$

where β is a hyper-parameter with a value of 0.5, d_t represents the calculated Euclidean distance, $|\mathcal{A}|$ is the dimension of the action space. Notably, a similar technique was

applied in [32]. However, we have modified this approach by querying only based on the current state and action.

Finally, given two trajectories σ_1 and σ_2 , we compute the corresponding values for each state-action pair within them. The segment with a higher cumulative *value* will be determined as the winning segment (σ_{winning}), as inferred by the following equation:

$$\sigma_{\text{winning}}^1 = \arg \max_{\sigma \in \{\sigma_1, \sigma_2\}} \left(\sum_{t=0}^h v_t^{(1)}, \sum_{t=0}^h v_t^{(2)} \right) \quad (7)$$

To demonstrate the effectiveness of our proposed *weak* preference, we provide the alternative reward-available Preference Label (*strong label*) for experiment comparison. Suppose the reward label can be obtained in the environment, and the segment with higher *cumulative reward* will be preferred, indicated in the following equation:

$$\sigma_{\text{winning}}^2 = \arg \max_{\sigma \in \{\sigma_1, \sigma_2\}} \left(\sum_{t=0}^h r_t^{(1)}, \sum_{t=0}^h r_t^{(2)} \right) \quad (8)$$

According to the Bradley-Terry model, the preference model can be described as (here σ_{winning} is abbreviated as σ^+ , the σ_{losing} is abbreviated as σ^-):

$$P(\sigma^+ > \sigma^-) = \text{Sigmoid}(r(\sigma^+) - r(\sigma^-)) \quad (9)$$

The proposed loss function during the preference alignment incorporates two components: the preference loss (\mathcal{L}_{DPO}) derived from DPO and the regularization term. The preference loss amplifies the difference between the “winning segment” and the “losing segment,” thus improving the diffusion policy’s performance to align with the preference labels and generate more preferred samples. More importantly, the regularization term helps to avoid significant deviations from the original policy. The loss function for the preference alignment stage is as follows:

$$\mathcal{L}(\epsilon_\theta, \mathcal{D}_{\text{pref}}) = \mathcal{L}_{\text{DPO}}(\epsilon_\theta, \mathcal{D}_{\text{pref}}) - \mu \mathbb{E}_{\sigma \in \mathcal{D}} \log \epsilon_\theta(\sigma) \quad (10)$$

The practical approximation and detailed proof can be found in [33].

V. EXPERIMENTAL RESULTS

We conduct the velocity tracking task among three gaits with different speeds to quantify the performance of our proposed framework. The diffusion pipeline is implemented under CleanDiffuser [34], hyperparameters are listed in Table I. For evaluation metrics, we consider the stability and average x -axis speed. We assess an episode as stable when the quadruped robot does not fall during all 250 steps. For every experimental parameter setting, we repeatedly collect 1024 episodes for three random seeds.

We chose three strong baselines for the experiments. The first one is **Conservative Q-learning (CQL)**. CQL mitigates the Q-value overestimation problem by introducing the regularization term in conservative Q-value estimation. CQL is effective in high-dimensional state space and complex environments. The second and third baselines are offline diffusion planners with U-Net backbone and transformer

TABLE I
HYPER-PARAMETERS IN TRAINING, INFERENCE, AND PREFERENCE
ALIGNMENT STAGE

Stage	Hyper-parameter	Value
Offline Diffusion Planner	Batch size	64
	Horizon	64
	Solver	DDPM
	Diffusion steps	20
	Action loss weight	10.0
Inference	w_cg	0.0001
	Sampling steps	10
Preference Alignment	Regularization weight	1.0
	Bias	0
	Temperature	500

backbone, respectively. They are noted as **DDPM-Unet** and **DDPM-Transformer**.

A. Performance of Preference-aligned Diffusion Planner

Table II presents a comparison of the proposed method with baselines across different gaits (pacing, trotting, bounding) and speeds (0.5 m/s and 1.0 m/s) in terms of stability and average x -axis velocity.

CQL consistently failed all the locomotion tasks. We observed in the simulation environment that the quadruped robot either remained stationary on the ground or exhibited subtle irregular jitters. This indicates that CQL struggles with continuous control in complex locomotion tasks.

Stability Performance: The proposed preference-aligned diffusion planner outperformed all baselines in stability across all locomotion tasks. Specifically, for the 1.0 m/s bounding task, the proposed *weak* preference-aligned planner achieves 84.3% stability, much higher than 32.4% of the DDPM-Unet and 46.6% of the DDPM-Transformer. Additionally, in the 0.5m/s trotting task, the stability increased by 19.8% relative to the offline DDPM-Transformer model. These results indicate that our preference-aligned planner exhibits superior stability performance across different speeds among gaits.

Notably, the improvement in stability from the DDPM-Transformer to the preference-aligned planner is attributed to the online alignment stage, which mitigates the inadequate diversity caused by a limited single-source offline dataset and aligns the planner with real-world data distribution. Experiments indicate that the offline diffusion planner in Fig. 3 Phase 2 can generate reasonable action sequences, however, it is extremely sensitive to the OOD states encountered in the simulation environment, examples shown in Fig. 1. This is due to the high-quality expert locomotion trajectories did not enable the diffusion planner to effectively address non-fatal disturbances, the deviated action prediction will gradually be amplified to significant cumulative errors (i.e., falls).

Velocity Tracking Performance: We evaluate the velocity tracking performance by the difference between the average measured x -axis velocity and the target velocity, and a more minor deviation indicates more accurate velocity control. For example, in the 1.0 m/s bounding task, our proposed

planner achieves an average velocity of 0.72 m/s, closely approaching the target speed of 1.0 m/s, significantly outperforming DDPM-Unet (0.42 m/s). In the 0.5 m/s task, the proposed method demonstrates more precise velocity control with minor deviation from the target speed. For example, in the 0.5 m/s trotting task, our method achieved an average velocity of 0.42 m/s, closely matching the target and surpassing DDPM-Unet (0.59 m/s) and DDPM-Transformer (0.36 m/s).

To further present how our proposed method improves the velocity tracking performance, we conduct a case study based on the most challenging 0.5 m/s bounding task with the lowest stability in Fig. 4.

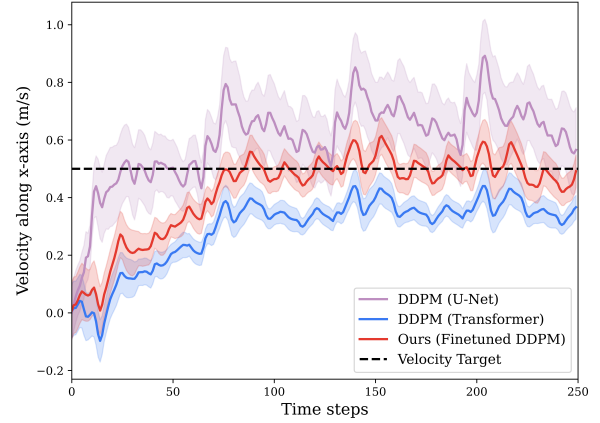


Fig. 4. Velocity tracking result in 0.5 m/s bounding gait between models

Fig. 4 depicts the velocity tracking process within an episode; the real-time velocities are smoothed with a moving average filter. The DDPM-Transformer exhibits a smaller standard deviation while it still deviates from the target velocity. Additionally, DDPM-Unet shows high fluctuation in measured velocity and fails to track the desired velocity. In contrast, the proposed preference-aligned diffusion planner (red line) reaches and maintains close to the target velocity.

In conclusion, the experimental result demonstrates that the proposed method outperforms existing baselines on stability and velocity tracking performance. The quantitative analysis supports the effectiveness and robustness of our proposed two-stage learning framework among different locomotion tasks.

B. Real-world Experiments

To evaluate the robustness and adaptability of the locomotion policy derived by the proposed two-stage learning framework, real-world experiments were conducted on two surfaces with different frictional properties: a rough, high-friction cement surface and a smooth, low-friction PVC floor. Fig. 5 presents key frames during the takeoff and landing phases of the legged robot in the bounding gait, exhibiting foot-terrain interactions across different surfaces. Experiment results indicate that the proposed framework can lead to a robust diffusion-based locomotion policy capable of sim-to-real transfer, even with a limited dataset.

TABLE II
AVERAGE SPEED AND STABILITY COMPARISON AMONG METHODS ON QUADRUPEL LOCOMOTION VELOCITY TRACKING TASKS

	CQL	Offline Planner (DDPM-UNet)		Offline Planner (DDPM-Transformer)		Weak-Preference Aligned Planner		Reward-Preference Aligned Planner	
		Speed (m/s)	Stability (%)	Speed (m/s)	Stability (%)	Speed (m/s)	Stability (%)	Speed (m/s)	Stability (%)
Pacing (0.5 m/s)	fail	0.60	40.9	0.43	72.0	0.46	81.6	0.49	87.3
Trotting (0.5 m/s)	fail	0.59	39.5	0.36	58.7	0.42	78.5	0.50	85.2
Bounding (0.5 m/s)	fail	0.59	40.4	0.21	60.1	0.44	70.8	0.48	72.4
Pacing (1.0 m/s)	fail	0.49	77.4	0.62	81.4	0.64	91.6	0.73	91.8
Trotting (1.0 m/s)	fail	0.48	58.9	0.61	66.6	0.63	84.4	0.67	89.5
Bounding (1.0 m/s)	fail	0.42	32.4	0.78	46.6	0.72	84.3	0.80	89.2

TABLE III
ABLATION STUDIES ON PREFERENCE ALIGNMENT PARAMETERS

		Speed = 1.0 m/s			Speed = 0.5 m/s		
		Pacing	Trotting	Bounding	Pacing	Trotting	Bounding
Preference Number	1024 episodes (*0.5)	89.6	78.8	75.3	69.5	73.0	63.5
	2048 episodes (*1.0)	90.1	77.7	75.9	76.4	79.1	66.4
	3072 episodes (*1.5)	91.8	89.5	89.2	87.3	85.2	72.4
Preference Quality	Weak Preference Label	91.6	84.4	84.3	81.6	78.5	70.8
	Strong Preference Label	91.8	89.5	89.2	87.3	85.2	72.4
Regularization	Without Regularization	—	—	—	—	—	—
	With Regularization	91.8	89.5	89.2	87.3	85.2	72.4

C. Ablation Studies

The ablation studies in Table III investigate the impact of preference dataset size, preference label quality, and regularization methods on stability performance.

Preference Dataset Size: Increasing the preference dataset size from 1024 episodes to 3072 episodes improves the performance across all gaits and speeds. For example, in the 0.5 m/s pacing task, the stability increased by 17.8%. Generally, the sensitivity to the size of the preference dataset varies with the difficulty of different locomotion tasks. In the simplest 1.0 m/s pacing task, 1024 episodes of the preference dataset are sufficient to achieve relatively stable improvement. We observe in the simulation result that a smaller preference dataset may introduce the risk of insufficient coverage on state distribution, leads to deviating from the expected task performance. However, the overall results indicate our proposed preference alignment framework demonstrates stable improvement in the offline diffusion planner under limited preference data.

Preference Label Quality: We compare the performance

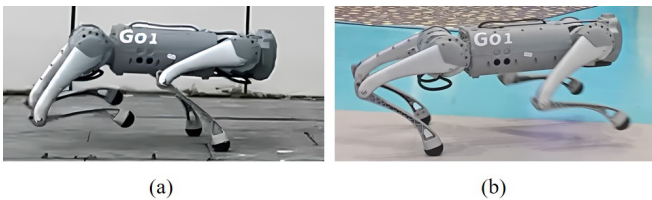


Fig. 5. Bounding Gait Testing on Different Surfaces: (a) Cement Surface which is Rough and High Friction, (b) PVC Floor which is Smooth and Low Friction)

between the reward-based preference label (strong preference) and our proposed reward-unavailable preference label (weak preference). The results in Table III indicate that weak preference labels can achieve satisfactory performance across most locomotion tasks with slight differences compared with strong preference labels. For example, the performance of the weak preference label closely approaches the strong preference label in the 1.0 m/s pacing and 0.5 m/s bounding tasks. Notably, the preference alignment based on weak preference labels demonstrates significant improvement compared with the DDPM-Transformer on all locomotion tasks, indicating the effectiveness of the weak preference labeling method and proposed two-stage learning framework.

Regularization Methods: The ablation results in Table III suggest the critical role of the proposed regularization term in the preference alignment stage. The compared baseline method eliminates the regularization term and solely emphasizes the difference between the winning segment (σ_{winning}) and the losing segment (σ_{losing}) under the current policy and reference policy. The results demonstrate that without the regularization term, the preference alignment will consistently cause *fail* in all locomotion tasks. In comparison, our proposed preference alignment method encourages the generation of winning segments while maintaining the overall likelihood of the winning and losing segments through the regularization term, thus effectively constrains the step size of policy updates to ensure stable and conservative improvements within the online alignment stage.

VI. CONCLUSIONS

This paper presents a two-stage learning framework that integrates offline diffusion learning with online preference

alignment to enhance the diversity and robustness of the diffusion planner. We focus on the challenging scenario with only a single-source limited expert dataset available and validate the effectiveness of the proposed architecture on the locomotion velocity tracking task. We leverage the offline diffusion planner to approximate the complex state-action sequences and further utilize the proposed *weak* preference label to conduct the preference alignment. Experiments indicate that our framework improved the stability and velocity tracking accuracy and can be deployed on Unitree Go1 robots. Future work can incorporate extra modalities, such as vision and language, into the learning framework.

ACKNOWLEDGMENT

This study is supported by the National Natural Science Foundation of China under Grant No.52302379 and Grant No.62306242. Acknowledgment to Computing Lab, Hong Kong University of Science and Technology (Guangzhou) for providing the computational resources utilized in this research.

REFERENCES

- [1] Z. Zhuang, Z. Fu, J. Wang, C. Atkeson, S. Schwertfeger, C. Finn, and H. Zhao, “Robot parkour learning,” *arXiv preprint arXiv:2309.05665*, 2023.
- [2] 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, p. eaau5872, 2019.
- [3] O. Villarreal, V. Barasuol, P. M. Wensing, D. G. Caldwell, and C. Semini, “Mpc-based controller with terrain insight for dynamic legged locomotion,” in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 2436–2442.
- [4] A. Escontrela, X. B. Peng, W. Yu, T. Zhang, A. Iscen, K. Goldberg, and P. Abbeel, “Adversarial motion priors make good substitutes for complex reward functions,” in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2022, pp. 25–32.
- [5] A. Kumar, Z. Fu, D. Pathak, and J. Malik, “Rma: Rapid motor adaptation for legged robots,” *Robotics: Science and Systems XVII*, 2021.
- [6] A. Nie, Y. Flet-Berliac, D. Jordan, W. Steenbergen, and E. Brunskill, “Data-efficient pipeline for offline reinforcement learning with limited data,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 14 810–14 823, 2022.
- [7] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, “Deep unsupervised learning using nonequilibrium thermodynamics,” in *International conference on machine learning*. PMLR, 2015, pp. 2256–2265.
- [8] Z. Zhu, H. Zhao, H. He, Y. Zhong, S. Zhang, Y. Yu, and W. Zhang, “Diffusion models for reinforcement learning: A survey,” *arXiv preprint arXiv:2311.01223*, 2023.
- [9] X. Huang, Y. Chi, R. Wang, Z. Li, X. B. Peng, S. Shao, B. Nikolic, and K. Sreenath, “Diffuselo: Real-time legged locomotion control with diffusion from offline datasets,” *arXiv preprint arXiv:2404.19264*, 2024.
- [10] J. D. Chang, M. Uehara, D. Sreenivas, R. Kidambi, and W. Sun, “Mitigating covariate shift in imitation learning via offline data with partial coverage,” in *Advances in Neural Information Processing Systems*, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan, Eds., 2021. [Online]. Available: <https://openreview.net/forum?id=7PklkyLMRM>
- [11] G. An, J. Lee, X. Zuo, N. Kosaka, K.-M. Kim, and H. O. Song, “Direct preference-based policy optimization without reward modeling,” *Advances in Neural Information Processing Systems*, vol. 36, pp. 70 247–70 266, 2023.
- [12] B. Wallace, M. Dang, R. Rafailov, L. Zhou, A. Lou, S. Purushwalkam, S. Ermon, C. Xiong, S. Joty, and N. Naik, “Diffusion model alignment using direct preference optimization,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8228–8238.
- [13] J. Lee, M. Bjelonic, A. Reske, L. Wellhausen, T. Miki, and M. Hutter, “Learning robust autonomous navigation and locomotion for wheeled-legged robots,” *Science Robotics*, vol. 9, no. 89, p. eadi9641, 2024.
- [14] G. B. Margolis and P. Agrawal, “Walk these ways: Tuning robot control for generalization with multiplicity of behavior,” in *Conference on Robot Learning*. PMLR, 2023, pp. 22–31.
- [15] X. B. Peng, G. Berseth, and M. Van de Panne, “Terrain-adaptive locomotion skills using deep reinforcement learning,” *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, pp. 1–12, 2016.
- [16] L. Schneider, J. Frey, T. Miki, and M. Hutter, “Learning risk-aware quadrupedal locomotion using distributional reinforcement learning,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 11 451–11 458.
- [17] 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, p. eabk2822, 2022.
- [18] 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.
- [19] 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.
- [20] L. Han, Q. Zhu, J. Sheng, C. Zhang, T. Li, Y. Zhang, H. Zhang, Y. Liu, C. Zhou, R. Zhao, *et al.*, “Lifelike agility and play in quadrupedal robots using reinforcement learning and generative pre-trained models,” *Nature Machine Intelligence*, pp. 1–12, 2024.
- [21] S. Levine, A. Kumar, G. Tucker, and J. Fu, “Offline reinforcement learning: Tutorial, review,” and *Perspectives on Open Problems*, vol. 5, 2020.
- [22] F. Torabi, G. Warnell, and P. Stone, “Behavioral cloning from observation,” *arXiv preprint arXiv:1805.01954*, 2018.
- [23] A. Sridhar, D. Shah, C. Glossop, and S. Levine, “Nomad: Goal masked diffusion policies for navigation and exploration,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 63–70.
- [24] U. A. Mishra and Y. Chen, “Reorientdiff: Diffusion model based reorientation for object manipulation,” in *2024 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2024, pp. 10 867–10 873.
- [25] P. M. Scheikl, N. Schreiber, C. Haas, N. Freymuth, G. Neumann, R. Lioutikov, and F. Mathis-Ullrich, “Movement primitive diffusion: Learning gentle robotic manipulation of deformable objects,” *IEEE Robotics and Automation Letters*, 2024.
- [26] A. Ajay, Y. Du, A. Gupta, J. Tenenbaum, T. Jaakkola, and P. Agrawal, “Is conditional generative modeling all you need for decision-making?” *arXiv preprint arXiv:2211.15657*, 2022.
- [27] M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine, “Planning with diffusion for flexible behavior synthesis,” *arXiv preprint arXiv:2205.09991*, 2022.
- [28] T. Pearce, T. Rashid, A. Kanervisto, D. Bignell, M. Sun, R. Georgescu, S. V. Macua, S. Z. Tan, I. Momennejad, K. Hofmann, *et al.*, “Imitating human behaviour with diffusion models,” *arXiv preprint arXiv:2301.10677*, 2023.
- [29] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song, “Diffusion policy: Visuomotor policy learning via action diffusion,” *arXiv preprint arXiv:2303.04137*, 2023.
- [30] J. Ho, A. Jain, and P. Abbeel, “Denosing diffusion probabilistic models,” *Advances in neural information processing systems*, vol. 33, pp. 6840–6851, 2020.
- [31] J. Ho and T. Salimans, “Classifier-free diffusion guidance,” *arXiv preprint arXiv:2207.12598*, 2022.
- [32] J. Lyu, X. Ma, L. Wan, R. Liu, X. Li, and Z. Lu, “Seabo: A simple search-based method for offline imitation learning,” *arXiv preprint arXiv:2402.03807*, 2024.
- [33] Z. Shan, C. Fan, S. Qiu, J. Shi, and C. Bai, “Forward kl regularized preference optimization for aligning diffusion policies,” *arXiv preprint arXiv:2409.05622*, 2024.
- [34] Z. Dong, Y. Yuan, J. Hao, F. Ni, Y. Ma, P. Li, and Y. Zheng, “Clean-diffuser: An easy-to-use modularized library for diffusion models in decision making,” *arXiv preprint arXiv:2406.09509*, 2024.