

TACT-ful: Multi-Channel Terrain Affordance and Compliance Training for Payload-Robust Perceptive Humanoid Locomotion

Anonymous Author(s)

Affiliation

Address

email

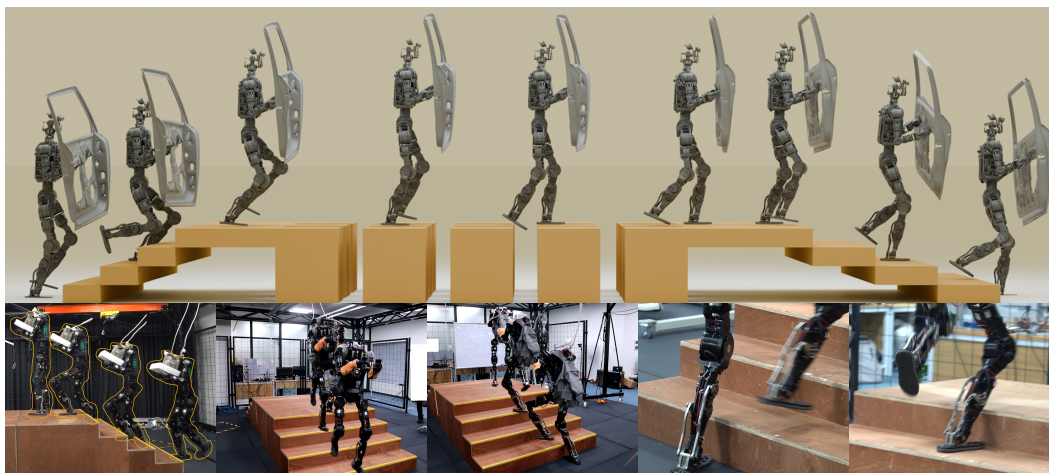


Figure 1: A service humanoid, trained on TACT-ful, traverses structured terrain while carrying heavy payload.

1 **Abstract:** Foothold selection on structured terrain requires explicit reasoning
2 about contact planarity, surface steepness, and kinematic reachability, proper-
3 ties not captured by a single height-based terrain signal. We propose a multi-
4 channel terrain cost combining flatness, steepness, and velocity-aware height fea-
5 sibility, plus a forward climb reward, that simultaneously drives a GPU-parallel
6 divergent component of motion (DCM) foothold planner and shapes a dense
7 per-step affordance reward for an asymmetric actor-critic policy trained with
8 proximal policy optimization (PPO) from depth images. A Bézier swing tra-
9 jectory with adaptive apex bias extends foothold tracking to joint position-and-
10 orientation, using the arc tangent to guide sole orientation through riser cross-
11 ings and tread landings. To support payload tasks, we introduce a lower-body
12 compliance training procedure in which a virtual wrench is injected at a sam-
13 pled load attachment point, generating physically consistent force and moment;
14 wrench-aware compliance targets replace rigid pose penalties, and the policy
15 learns to yield to load-induced perturbations without force sensing. The full sys-
16 tem trains end-to-end with standard PPO, no distillation, and no teacher-student
17 staging, and is deployed on a humanoid directly from simulation with configura-
18 tion changes only. In simulation, the policy reaches 1.0 m/s on stairs with
19 risers up to 0.20 m and improves payload robustness up to ~ 15 kg centered
20 load and for moment-dominated wrist loads without fine-tuning. We also pro-
21 vide a qualitative hardware demonstration on structured terrain. Project website:
22 <https://fai-rl-tech.github.io/tact-locomotion.github.io/>

23 1 Introduction

24 Service humanoids operating in human environments must traverse structured non-flat terrain, in-
25 cluding staircases, stepped platforms, and ramps, while carrying tools or cargo. Foothold placement
26 on such terrain is consequential: a sole that straddles a step edge bears load on a line contact, narrow-
27 ing the friction cone and generating lateral impulses that scale with robot mass. The capture-point [1]
28 divergence rate is independent of mass, but the ground reaction force needed to intercept a diver-
29 gent trajectory scales with total weight, so that early, terrain-quality-aware foothold selection [2]
30 becomes more consequential as platform mass grows. Payload carrying adds a second challenge:
31 a carried load shifts the effective neutral pelvis pose and imposes a time-varying tug on the torso,
32 requiring the locomotion controller to absorb these perturbations without loss of balance.

33 Existing approaches address subsets of this problem. Model-based foothold planners select discrete
34 landing targets subject to stability and reachability constraints [3, 4, 5], but their outputs are not used
35 to shape the reinforcement-learning (RL) policy reward. Perceptive RL methods for humanoids con-
36 sume elevation maps or depth images to adapt gait on rough terrain [6, 7, 8], but typically encode
37 terrain through a single height channel and do not reason explicitly about contact planarity or kine-
38 matic overreach. Compliance under payload is largely treated separately, either through dedicated
39 force controllers, online payload estimators, or manual gait re-tuning, rather than as a behavior the
40 locomotion policy acquires within a single terrain-aware training run.

41 The contributions of this paper are:

- 42 1. **Multi-channel terrain cost with terrain-adaptive swing and tangent-guided foot orienta-**
43 **tion.** Flatness, steepness, velocity-aware height feasibility, and a climb bonus jointly drive a
44 GPU-parallel DCM foothold planner and serve as a dense terrain-affordance learning signal for
45 the RL policy. A Bézier swing trajectory with adaptive apex bias and clearance extends foothold
46 tracking to joint position-and-orientation references, using the arc tangent to guide sole orienta-
47 tion through riser crossings and tread landings. Unlike the closest predecessor, which couples
48 a single-channel DCM planner to a position-only foothold reward under a multi-stage curricu-
49 lum [7], our cost adds explicit flatness, velocity-aware feasibility, and speed-gated climb chan-
50 nels, and the tangent schedule extends the swing reference to sole *orientation*.
- 51 2. **Lower-body compliance training for payload-robust locomotion.** A virtual wrench injected at
52 a sampled load attachment point generates physically consistent force and moment; together with
53 wrench-aware compliance targets for pelvis height and trunk orientation, it teaches the policy to
54 yield to payload-induced perturbations without force sensing, improving robustness up to ~ 15 kg
55 centered load and for moment-dominated wrist loads without payload-specific fine-tuning.
- 56 3. **An integrated end-to-end pipeline with a qualitative hardware demonstration.** The system
57 trains with standard PPO in a single run (no distillation, no teacher-student staging) and com-
58 bines terrain-aware foothold selection with payload compliance, and is deployed zero-shot from
59 simulation on a service humanoid traversing structured terrain while carrying payload.

60 2 Related Work

61 **Foothold planning and terrain-conditioned reinforcement learning.** Model-based foothold
62 planners select discrete landing targets under LIPM capturability constraints: Acosta and Posa [3]
63 formulate MIQP over convex foothold regions, while Xiang et al. [4] jointly optimize step timing
64 and placement under DCM bounds. Kim et al. [5] replace the MIQP with an RL-trained footstep
65 policy mapping elevation maps to 3D targets, decoupled from end-to-end training. Lee et al. [9] use
66 LIP-derived footstep targets as a partial RL reward, and Liu et al. [7] reformulate DCM foothold
67 search as GPU-parallel discrete optimization used as an explicit RL reward with a multi-stage cur-
68 riculum, the most direct predecessor to our approach. Wang et al. [10] show that a purely learned
69 dense foothold reward with a double-critic curriculum can traverse sparse footholds without an ex-
70 plicit planner.

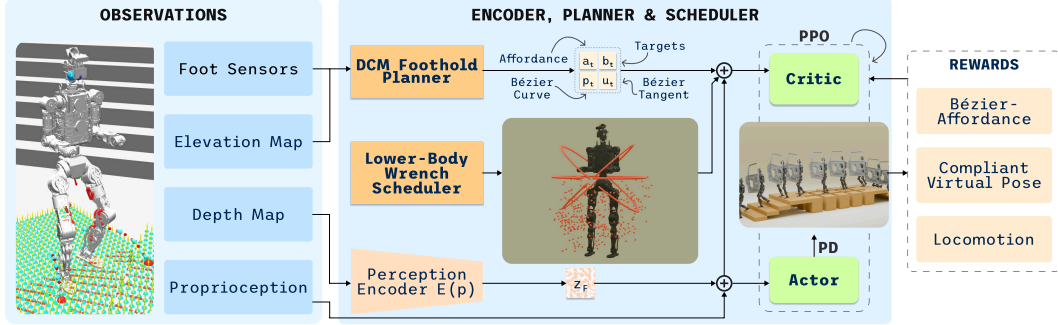


Figure 2: **Overview of the proposed framework.** Two parallel modules feed into the training reward buffer. (i) *DCM foothold planner*: at every control step, a pelvis-mounted elevation map is consumed by the GPU-parallel DCM foothold planner, which selects terrain-optimal landing targets and produces a Bézier swing trajectory reference; these targets define the foothold-tracking and terrain-specific reward terms and are provided as privileged observations to the critic. (ii) *Wrench compliance scheduler*: a virtual wrench is injected at a sampled load attachment point, generating force that shifts the pelvis and moment that tilts the trunk; wrench-aware compliance targets replace rigid pelvis-height and orientation penalties, training the policy to absorb payload-induced perturbations without explicit force sensing.

71 **Perceptive humanoid locomotion and terrain representation.** Agarwal et al. [11], Miki et al.
 72 [12] pioneered the privileged-training-to-depth-distillation paradigm, decoupling policy and percep-
 73 tion learning for zero-shot real-world deployment, which has become standard for perceptive legged
 74 locomotion. Building on this, Song et al. [6] reconstruct a dense egocentric height map from depth
 75 and jointly output joint targets and adaptive gait frequency; Radosavovic et al. [13] train a Trans-
 76 former policy on diverse outdoor terrain without perceptive reward shaping; Long et al. [14] achieve
 77 stair climbing via LiDAR elevation maps; Ben et al. [8] extend perception to full 3D structure via a
 78 voxel-grid representation; Zhang et al. [15] distil terrain-specific experts into a unified transformer
 79 policy with multi-view depth fusion and velocity-scaled features; and Sun et al. [16] train end-to-end
 80 from raw stereo depth images using terrain-specific multi-critic shaping. Vision-based whole-body
 81 policies demonstrate perceptive parkour, stair climbing, gap crossing, and platform jumping, on
 82 humanoids and quadrupeds [17, 18, 19]. Standard terrain representations include probabilistic ele-
 83 vation maps [20], traversability scoring [21], and multi-layer grid maps [22].

84 **Payload-carrying and compliant loco-manipulation.** Large-scale domain randomization yields
 85 implicit payload tolerance [23]; Kumar et al. [24] achieve payload and terrain robustness on
 86 quadrupeds through rapid online system identification; Zhang et al. [25] demonstrate balance re-
 87 covery under large external perturbations via a multi-phase curriculum; Fu et al. [26] make payload
 88 awareness explicit on the Tiangong humanoid via a history-based estimator in a decoupled lower-
 89 body RL architecture, evaluated on flat terrain. Pasricha et al. [27] extend payload limits to $3\times$
 90 nominal via compliant trajectory diffusion. For force-adaptive compliance without force sensors,
 91 Xu et al. [28] train legged robots to absorb collision impulses and sustain payload pulls up to $\frac{2}{3}$ body
 92 weight by tracking a virtual impedance reference; Zhi et al. [29] unify position and force control for
 93 contact-rich loco-manipulation. Terrain property estimation from vision [30] and friction identifica-
 94 tion [31] connect perception to contact mechanics. None of these works couple terrain-quality-aware
 95 foothold selection with compliance training in a single end-to-end RL formulation, leaving the in-
 96 teraction between payload mass, step-height feasibility, and terrain planarity unaddressed.

97 3 Terrain Affordance and Compliance Training

98 Figure 2 Overview the framework: a DCM foothold planner and a wrench compliance scheduler
 99 feed a shared reward for an asymmetric actor-critic policy. All hyperparameters and implementation
 100 details are located in the Appendices.

101 3.1 Capture-Point Stability

102 We use the divergent component of motion (DCM) / capture-point framework [1, 32] with a constant
 103 natural frequency $\omega_0 = \sqrt{g/z_0}$, where g is gravitational acceleration and z_0 is the nominal center
 104 of mass (CoM) height [33]. The DCM is defined as $\xi = x_{\text{CoM}} + \dot{x}_{\text{CoM}}/\omega_0$ and propagates as
 105 $\xi(T) = \xi_0 e^{\omega_0 T}$ over a swing of duration T (step time); applied component-wise in the horizontal
 106 plane this gives the 2D DCM $\xi_T \in \mathbb{R}^2$. The nominal step offset $\mathbf{b}_{\text{nom}} = (b_x, b_y)$ keeps the divergent
 107 mode in steady state [34]; we adopt the linear-CoM-height reduction of Liu et al. [7] and set $\omega = \omega_0$
 108 throughout.

109 3.2 Multi-Channel Terrain Cost

110 The planner selects the cell i that minimizes the total cost

$$111 \mathcal{J}_i = \alpha_{\text{pos}} d_{\text{pos},i} + \alpha_{\text{dcm}} d_{\text{dcm},i} + \alpha_E E_i + \alpha_Q Q_i + \alpha_M M_i - \alpha_{\text{climb}} b_i, \quad (1)$$

111 where $d_{\text{pos},i} = (x_i - v_x T)^2 + \beta (y_i - (v_y T + \text{sign}(f) l_p))^2$ is an asymmetric position residual
 112 from the nominal stride target ($\beta > 1$ penalizes lateral deviation more than sagittal to limit foot
 113 spread; $\mathbf{p}_i = (x_i, y_i)^\top$ is the 2D candidate foothold, (v_x, v_y) the commanded planar velocity, l_p
 114 the nominal lateral inter-foot distance, and $\text{sign}(f) \in \{+1, -1\}$ selects the swinging leg), and
 115 $d_{\text{dcm},i} = \|\xi_T - \mathbf{p}_i - \mathbf{b}_{\text{nom}}\|^2$ is the capture-point stability residual. These are computed from the
 116 elevation map \mathcal{H} in one GPU-batched forward pass.

117 **Flatness cost Q .** A sole straddling a height discontinuity tilts by $\theta_i \approx \arctan(Q_i/L_{\text{foot}})$, reducing
 118 admissible friction to $\mu_{\text{eff},i} \leq \mu_0 \cos \theta_i - \sin \theta_i$, which collapses to zero at $Q_i/L_{\text{foot}} = \mu_0 \approx 0.6$,
 119 precisely the failure regime on standard stairs. We use the elevation range over a footprint-sized
 120 kernel as a GPU-efficient continuous proxy:

$$121 Q_i = \max_{j \in \mathcal{N}_i}(z_j) - \min_{j \in \mathcal{N}_i}(z_j). \quad (2)$$

121 Sampling-based foothold rewards [10] evaluate *contact support coverage*, suited to surface voids but
 122 blind to planarity loss. On structured terrain, a tread-edge landing scores zero under that criterion
 123 yet incurs $Q_i > 0$; the two costs target distinct failure modes and are complementary.

124 **Steepness cost E .** An edge contact on a vertical riser produces a lateral impulse. Sobel-filtered
 125 gradients $g_{x,j}, g_{y,j}$ are max-pooled, deliberately propagating the worst gradient within the footprint:

$$126 E_i = \max_{j \in \mathcal{N}_i} \sqrt{g_{x,j}^2 + g_{y,j}^2} + \varepsilon. \quad (3)$$

126 **Height-feasibility cost M .** A step too tall to reach kinematically or dynamically is worse than a
 127 sub-optimal flat placement. The quadratic penalty

$$128 M_i = \max(|\Delta z_i| - h_{\text{eff}}^*, 0)^2, \quad (4)$$

128 where Δz_i is the stance foot height difference, uses a *velocity-aware effective maximum step height*:

$$129 h_{\text{eff}}^* = h_{\text{min}}^* + (h_{\text{max}}^* - h_{\text{min}}^*) \text{clip}\left(\frac{v_x}{v^*}, 0, 1\right), \quad (5)$$

129 where h_{min}^* and h_{max}^* bound the reachable step height at standstill and rated forward speed v^* ,
 130 respectively. At low speed the robot lacks the leg momentum needed to assist a step-up, so the
 131 effective kinematic reach is reduced; at high speed the dynamic leg drive can clear the full height.
 132 Below a minimum speed threshold a hard override replaces the planned target with the current foot
 133 position, and $h_{\text{eff}}^* \rightarrow h_{\text{min}}^*$ ensures the cost still disfavors high steps in the transition region.

134 **Forward climb bonus b .** To encourage the robot to step *up* onto reachable treads, a speed-gated
 135 capped bonus $b_i = \min(\max(\Delta z_i, 0), h_{\text{eff}}^*) \cdot \mathbf{1}[v_x > v_{\text{min}}]$ rewards upward steps proportionally to
 136 their height, capped at h_{eff}^* .

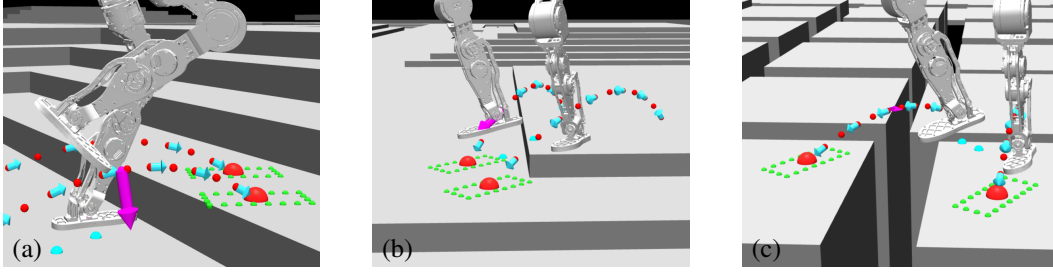


Figure 3: **Bézier swing arcs with adaptive apex bias.** (a) Step-up: apex biased toward the landing target, keeping the peak over the riser face. (b) Step-down: apex biased toward lift-off, extending horizontal travel before descent. (c) Gap: behave analogously, with clearance scaled by $|\Delta z|$.

137 3.3 Bézier Swing Trajectory and Tangent-Guided Foot Orientation

138 After selecting per-foot landing target \mathbf{p}_f^* , the swing foot tracks a quadratic Bézier arc whose tangent
139 simultaneously determines foot orientation:

$$\begin{cases} \mathbf{p}(u) = (1-u)^2 \mathbf{p}_\ell + 2(1-u)u \mathbf{p}_{\text{apex}} + u^2 \mathbf{p}_f^*, \\ \dot{\mathbf{p}}(u) = 2(1-u)(\mathbf{p}_{\text{apex}} - \mathbf{p}_\ell) + 2u(\mathbf{p}_f^* - \mathbf{p}_{\text{apex}}), \end{cases} \quad u \in [0, 1], \quad (6)$$

140 where \mathbf{p}_ℓ is the lift-off position and $\dot{\mathbf{p}}(u)$ points in the direction of instantaneous foot travel. The
141 apex xy is biased toward whichever endpoint is higher:

$$\text{bias} = \text{clip}\left(0.5 + \kappa \frac{\Delta z}{h_{\text{max}}^*}, b_{\text{min}}, b_{\text{max}}\right), \quad \mathbf{p}_{\text{apex},xy} = (1 - \text{bias}) \mathbf{p}_{\ell,xy} + \text{bias} \mathbf{p}_{f,xy}^*, \quad (7)$$

142 For step-up ($\Delta z > 0$), $\text{bias} > 0.5$ places the trajectory peak over the riser face, preventing premature
143 descent onto the vertical surface; for step-down ($\Delta z < 0$), $\text{bias} < 0.5$ extends horizontal travel
144 before descent, reducing the landing impact angle (fig. 3a-b). A phase-conditional schedule derived
145 from $\dot{\mathbf{p}}(u)$ guides sole orientation through riser crossings. Both position proximity and foot heading
146 are jointly optimized during swing via an exponential proximity kernel:

$$r_{\text{foothold}} = \sum_f I_{\text{swing},f} \exp(-\sigma_p \|\mathbf{p}_f - \mathbf{p}_{\text{Béz}}(t)\|^2 - \sigma_d \|\hat{\mathbf{d}}_f - \hat{\mathbf{t}}_f(u)\|^2), \quad (8)$$

147 where $\hat{\mathbf{d}}_f = R_f \hat{\mathbf{e}}_x$ is the foot forward axis and $\sigma_d = 0$ recovers the position-only form when
148 orientation is not used.

149 The apex height adapts to terrain relief:

$$c = \min(c_{\text{min}} + s |\Delta z|, c_{\text{max}}), \quad z_{\text{apex}} = 2(\max(z_\ell, z_f^*) + c) - \frac{1}{2}(z_\ell + z_f^*), \quad (9)$$

150 where s is a clearance scale factor and c_{max} caps the swing peak. The tangent is vertical at

$$u_{\text{peak}} = \frac{z_\ell - z_{\text{apex}}}{z_\ell - 2z_{\text{apex}} + z_f^*}, \quad (10)$$

151 equal to 0.5 for level steps and shifting toward the higher endpoint on transitions. Two
152 windows around u_{peak} define the reference unit vector $\hat{\mathbf{t}}_f(u)$: in the *pre-apex* window
153 $[u_{\text{peak}} - \delta_\ell^-, u_{\text{peak}} + \delta_\ell^+]$, $\dot{\mathbf{p}}$ is rotated 90° in the sagittal plane ($(t_x, t_y, t_z) \mapsto (t_z, t_y, -t_x)$) and nor-
154 malized, yielding a forward-downward direction that guides the sole through the riser and reduces
155 toe-strike risk; in the *post-apex* window $(u_{\text{peak}} + \delta_r^-, u_{\text{peak}} + \delta_r^+]$, $\hat{\mathbf{t}}_f = \dot{\mathbf{p}}/\|\dot{\mathbf{p}}\|$ aligns the foot with
156 the forward travel direction for tread landing; elsewhere $\hat{\mathbf{t}}_f = \mathbf{0}$ (no orientation constraint).

157 3.4 Lower Body Compliance Training

158 A payload exerts a wrench $\mathbf{W} = (\mathbf{F}, \boldsymbol{\tau})$ on the robot, where the moment $\boldsymbol{\tau} = \mathbf{r}_{\text{load}} \times \mathbf{F}$ scales
159 with the CoM-to-load arm \mathbf{r}_{load} independently of mass. $\boldsymbol{\tau}_{\text{ext}}$ tilts the trunk by φ , displacing the
160 CoM by $\Delta x_{\text{CoM}} \approx d_{\text{CoM}} \varphi$; under LIPM this offset grows by $e^{\omega_0 T_{\text{step}}} > 2$ over a typical step, so

161 moment-induced tilt perturbs the capture point more than an equivalent translational offset, and both
 162 components of \mathbf{W} require explicit coverage in the disturbance distribution. We emulate the wrench
 163 by applying a spring-damper at the *virtual load attachment point* $\mathbf{p}_{\text{load}} = \mathbf{p}_{\text{CoM}} + \mathbf{r}_{\text{load}}$:

$$\mathbf{F}_{\text{ext}}(t) = k(\mathbf{p}_a - \mathbf{p}_{\text{load}}(t)) - c\dot{\mathbf{p}}_{\text{load}}(t), \quad (11)$$

164 where \mathbf{p}_a is a fixed world-frame anchor drawn once per episode, with stiffness k and damping c
 165 set so the spring force at maximum displacement R_a reaches a fixed fraction of nominal weight.
 166 Applying the force at \mathbf{p}_{load} rather than the CoM generates $\boldsymbol{\tau}_{\text{ext}} = \mathbf{r}_{\text{load}} \times \mathbf{F}_{\text{ext}}$ automatically. At
 167 each reset, \mathbf{r}_{load} is drawn from a two-component mixture ($\mathbf{r}_{\text{load}} = \mathbf{r}_{\text{body}}$ with prob. ρ_{body} , else \mathbf{r}_{arm})
 168 that jointly covers both carry scenarios: *Body-attached* with $\mathbf{r}_{\text{body}} \sim \text{Ball}(R_{\text{close}})$, direction biased
 169 toward $-\hat{\mathbf{z}}$; and *Arm-extended* with \mathbf{r}_{arm} from a shoulder-frame half-ellipsoid, biased forward. Both
 170 components use cosine-power lobes to concentrate sampling:

$$\begin{cases} p(\hat{\mathbf{d}}) \propto \cos^n \theta, & \theta \in [0, \pi/2] & \text{(body-attached direction),} \\ p(\mathbf{r}_{\text{arm}}) \propto \cos^{nr}(\angle(\mathbf{r}_{\text{arm}}, \hat{\mathbf{x}}_{\text{shoulder}})) & & \text{(arm-extended bias).} \end{cases} \quad (12)$$

171 where $\hat{\mathbf{d}}$ is the unit direction of \mathbf{r}_{body} and θ its polar angle from $-\hat{\mathbf{z}}$. A weight- ϵ isotropic component
 172 blends in lateral and overhead samples. We replace the height and orientation penalties with wrench-
 173 aware compliance targets. For the height,

$$z_{\text{virt}}^*(t) = h_{\text{pelvis}}^{\text{terrain}}(t) + h^* + \alpha_z(z_a(t) - h_{\text{pelvis}}^{\text{terrain}}(t) - h^*), \quad (13)$$

174 where h^* is the commanded base height, $z_a(t) = [\mathbf{p}_{\text{load}}(t)]_z$ is the world-frame z-coordinate of the
 175 virtual attachment point, and $\alpha_z \in [0, 1]$ is the height compliance gain. For the trunk orientation,
 176 the virtual target is derived from the induced moment:

$$\varphi_{\text{virt}}^* = \alpha_\varphi \frac{\tau_y}{k_{\text{rot}}}, \quad \psi_{\text{virt}}^* = \alpha_\psi \frac{\tau_x}{k_{\text{rot}}}, \quad (14)$$

177 where τ_x, τ_y are the roll and pitch components of $\boldsymbol{\tau}_{\text{ext}}$, k_{rot} [N·m/rad] is a rotational stiffness
 178 converting induced moment to a tilt angle, and $\alpha_\varphi, \alpha_\psi \in [0, 1]$ are orientation compliance gains.
 179 The combined compliance reward is

$$r_{\text{comply}} = -(z_{\text{pelvis}} - z_{\text{virt}}^*)^2 - (\varphi_B - \varphi_{\text{virt}}^*)^2 - (\psi_B - \psi_{\text{virt}}^*)^2, \quad (15)$$

180 where φ_B, ψ_B are pelvis pitch and roll from the projected gravity vector. When $\alpha_z = \alpha_\varphi = \alpha_\psi =$
 181 0 the reward collapses to the original rigid height and orientation penalties; at nonzero gains the
 182 policy learns to yield to the full wrench, acquiring wrench-aware compliance in both translation and
 183 rotation as a learned behavior without explicit force or torque sensing.

184 3.5 Policy and Training

185 An asymmetric Actor-Critic Encoder is used. The actor concatenates a convolutional embedding of
 186 a temporally stacked depth image with proprioceptive state and passes the result through a feedfor-
 187 ward multilayer perceptron (MLP). The critic shares the same encoder but also receives privileged
 188 training-only signals. Actions are joint position targets relative to the default pose, converted to
 189 torques via proportional-derivative (PD) control and clipped to mechanical limits. Following the
 190 gait-adaptive extension [6], the policy also outputs a scalar gait frequency f_t that advances an inter-
 191 nal phase clock via an exponential-moving-average (EMA) smoothed update:

$$\hat{f}_t = (1 - \alpha_{\text{EMA}})\hat{f}_{t-1} + \alpha_{\text{EMA}} \text{clip}(f_t, f_{\text{min}}, f_{\text{max}}), \quad \phi_{t+1} = (\phi_t + \Delta t \hat{f}_t) \bmod 1, \quad (16)$$

192 with $(\hat{f}_t, \sin 2\pi\phi_t, \cos 2\pi\phi_t)$ appended to the actor observation. Training uses standard PPO [35]
 193 with adaptive Kullback–Leibler (KL) scheduling, generalized advantage estimation (GAE), and
 194 massively-parallel joint-velocity and terrain-difficulty curricula [36].

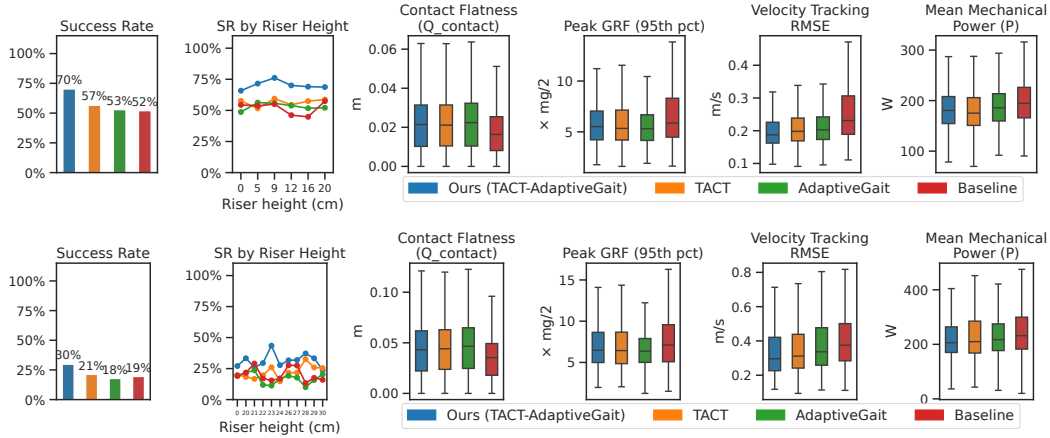


Figure 4: Terrain traversal ablation on **standard terrain** (top) and **hard terrain** (out-of-distribution).

195 4 Results

196 **Experimental setup.** Four variants are compared: *TACT + Adaptive Gait (Ours)*, *TACT-only*,
 197 *Adaptive Gait only*, and *Baseline* (a standard depth-map perceptive policy with no terrain-cost chan-
 198 nels and no privileged elevation-map input to the critic). Each variant is evaluated at iteration 20k
 199 across 4096 environments in MuJoCo [37] on stairs, slopes, and rough terrain; a second, unseen
 200 hard-terrain sweep covers risers 0.20–0.30 m, over which we report the success rate SR_{hard} com-
 201 puted on the strict interior $[0.22, 0.28]$ m.

202 4.1 Terrain Traversal Ablation

203 Ours (TACT + Adaptive Gait) leads on both standard and hard terrain (70 % standard SR and 30 %
 204 SR_{hard}), outperforming TACT-only by 13 and 9 percentage points (pp) and Adaptive Gait only by 17
 205 and 12 pp (Fig. 4); Adaptive Gait only falls below the Baseline on hard terrain (18 % vs. 19 %), con-
 206 firming that frequency re-timing without terrain-quality guidance is actively harmful when risers ap-
 207 proach kinematic limits. The Baseline’s 52 % and 19 % SR coincide with the highest 95th-percentile
 208 ground reaction force (GRF) ($\approx 6.7\times$ and $8.0 \times mg/2$), consistent with edge-biased contacts from
 209 blind foothold placement. On efficiency, Ours achieves lower velocity-tracking root-mean-square
 210 error (RMSE) (0.22 vs. 0.23 m/s) and mechanical power (229 vs. 241 W) than TACT-only on hard
 211 terrain; Q_c is identical between the two (0.023), isolating terrain cost channels as the mechanism
 212 for tread-centered landings. Fig. 5a shows that the SR gap relative to the Baseline is speed-invariant
 213 ($71 / 68 / 69$ % vs. ≈ 20 pp lower at 0.3 / 0.6 / 1.0 m/s), confirming foothold quality, not speed,
 214 drives the difference; Fig. 5b shows that removing all TACT channel weights raises foot-target dis-
 215 tance $2.8\times$ ($0.088 \rightarrow 0.251$ m).

216 4.2 Payload Generalization

217 Fig. 5c evaluates payload generalization across three conditions (pelvis +15 kg, pelvis +20 kg, wrist
 218 +10 kg) on terrain and compares Ours against the Baseline. At moderate centered load (pelvis
 219 +15 kg), Ours achieves 50 % SR versus the Baseline’s 38 % at lower power (247 W vs. 277 W),
 220 consistent with compliance training absorbing the downward wrench and suppressing impulsive
 221 GRF recovery. At high centered load (pelvis +20 kg), both policies degrade substantially to ≈ 37 %
 222 and ≈ 35 % SR respectively, nearly indistinguishable, indicating the load magnitude approaches the
 223 boundary of the training distribution; power rises for both, with Ours remaining 27 W below the
 224 Baseline (293 vs. 320 W), confirming that partial compliance is retained even as balance recovery
 225 deteriorates. Wrist-mounted mass (+10 kg), which generates a large moment rather than a direct
 226 CoM shift, yields the largest margin across all three conditions: Ours achieves 65 % SR versus the

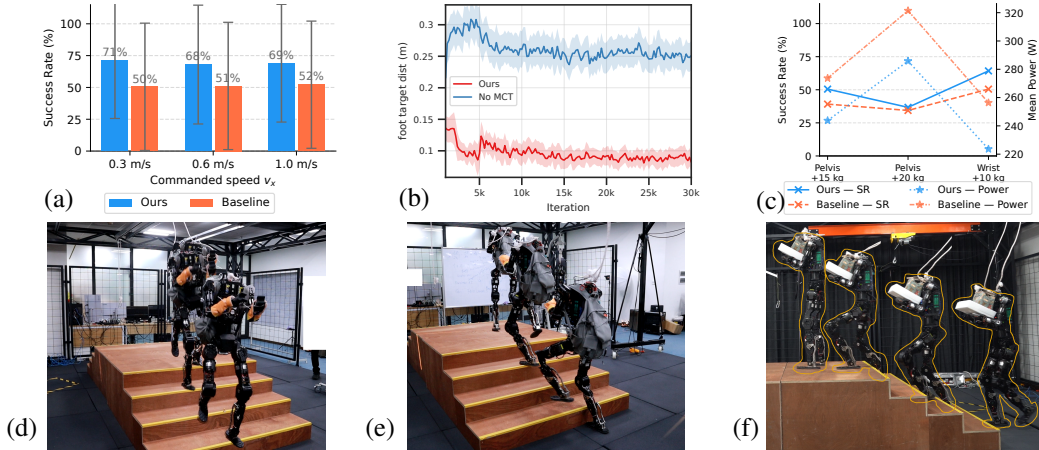


Figure 5: (a) Speed-conditioned SR; (b) foot-target distance; (c) SR (%) and mean power (W). (d–f) Qualitative hardware payload demonstrations, separate from the simulated ablations (d–e: 20 kg backpack; f: 7 kg tray).

227 Baseline’s 50 % at 223 W versus 259 W, consistent with the arm-extended wrench samples in the
 228 disturbance force-field explicitly targeting moment-dominated loads.

229 5 Limitations

230 **Simulation-only evaluation.** All quantitative ablations are conducted in MuJoCo; hardware results
 231 are qualitative demonstrations, with domain randomization as the sole sim-to-real bridge, so absolute
 232 success rates on real structured terrain remain uncharacterized and unmodeled contact phenomena
 233 will shift them. Even in simulation, the best variant reaches only 30 % SR_{hard} on hard terrain
 234 (interior risers 0.22–0.28 m), a ceiling insufficient for reliable deployment on tall staircases.

235 **Compliance and perception scope.** The compliance benefit vanishes near 20 kg of centered load
 236 (Ours $\approx 37\%$ vs. Baseline $\approx 35\%$ SR at pelvis +20 kg): the sampled wrench distribution and leg-
 237 only actuation set an effective payload ceiling, and lateral or distributed loads fall outside the sam-
 238 pled space. The DCM planner depends on accurate elevation-map registration with no fallback under
 239 localization drift, and the forward-facing depth camera gives no lateral coverage. Finally, the ter-
 240 rain cost channels are collectively necessary but individually non-critical, removing any one keeps
 241 foot-target distance within 5 %, and their weights are hand-tuned rather than derived. Extended
 242 analysis of each limitation (adaptive-gait fragility, channel redundancy, upper-body coupling, and
 243 deformable/low-friction terrain) is given in Appendix.

244 6 Conclusion

245 We presented a physics-informed terrain affordance learning system for payload-robust perceptive
 246 humanoid locomotion. A multi-channel terrain cost drives a GPU-parallel DCM foothold plan-
 247 ner and provides a dense terrain-affordance learning signal to the RL policy; a terrain-adaptive
 248 Bézier swing with tangent-guided foot orientation extends foothold tracking to joint position-and-
 249 orientation references. Because the cost is defined from sole geometry rather than learned per
 250 embodiment [12], it is platform-agnostic by construction. Lower-body compliance training via a
 251 cosine-power-weighted virtual spring-damper yields a learned wrench-dependent impedance (con-
 252 ceptually related to virtual impedance tracking [28] but without a separate impedance reference, and
 253 avoiding the online payload estimation of Kumar et al. [24]), improving payload robustness up to
 254 ~ 15 kg centered load and for moment-dominated wrist loads without retraining. The full system
 255 trains with standard PPO in a single run (no teacher-student distillation [11, 15]) and is deployed
 256 zero-shot from simulation on a humanoid; quantitative real-world evaluation and a terrain \times payload
 257 coupling study are the main remaining steps toward reliable deployment.

References

- [1] Jerry Pratt, John Carff, Sergey Drakunov, and Ambarish Goswami. Capture point: A step toward humanoid push recovery. In *2006 6th IEEE-RAS International Conference on Humanoid Robots (Humanoids)*, pages 200–207, 2006. doi:10.1109/ICHR.2006.321385.
- [2] Eric Whitman and Gina Christine Fay. Terrain aware step planning system. U.S. Patent Application Publication US20200117198A1, assigned to Boston Dynamics, Inc., April 2020. URL <https://patents.google.com/patent/US20200117198A1/en>. Published Apr. 16, 2020; granted as US11287826B2.
- [3] Brian Acosta and Michael Posa. Perceptive mixed-integer footstep control for underactuated bipedal walking on rough terrain. *IEEE Transactions on Robotics*, 41:4518–4537, 2025. doi:10.1109/TRO.2025.3587998.
- [4] Zhaoyang Xiang, Upama Pant, and Ayonga Hereid. Perceptive variable-timing footstep planning for humanoid locomotion on disconnected footholds, 2026. URL <https://arxiv.org/abs/2603.07400>.
- [5] Minku Kim, Brian Acosta, Pratik Chaudhari, and Michael Posa. Learning a vision-based footstep planner for hierarchical walking control. *2025 IEEE-RAS 24th International Conference on Humanoid Robots (Humanoids)*, pages 1–8, 2025. URL <https://arxiv.org/abs/2508.06779>.
- [6] Haolin Song, Hongbo Zhu, Tao Yu, Yan Liu, Mingqi Yuan, Wengang Zhou, Hua Chen, and Houqiang Li. Gait-adaptive perceptive humanoid locomotion with real-time under-base terrain reconstruction. *IEEE Robotics and Automation Letters*, 11(4):4969–4976, 2026. doi:10.1109/LRA.2026.3664167.
- [7] Yan Liu, Tao Yu, Haolin Song, Hongbo Zhu, Nianzong Hu, Yuzhi Hao, Xiuyong Yao, Xizhe Zang, Hua Chen, and Jie Zhao. FastStair: Learning to run up stairs with humanoid robots, 2026. URL <https://arxiv.org/abs/2601.10365>.
- [8] Qingwei Ben, Botian Xu, Kailin Li, Feiyu Jia, Wentao Zhang, Jingping Wang, Jingbo Wang, Dahua Lin, and Jiangmiao Pang. Gallant: Voxel grid-based humanoid locomotion and local-navigation across 3D constrained terrains, 2025. URL <https://arxiv.org/abs/2511.14625>.
- [9] Ho Jae Lee, Seungwoo Hong, and Sangbae Kim. Integrating model-based footstep planning with model-free reinforcement learning for dynamic legged locomotion. In *2024 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 11248–11255, 2024. doi:10.1109/IROS58592.2024.10801468.
- [10] Huayi Wang, Zirui Wang, Junli Ren, Qingwei Ben, Tao Huang, Weinan Zhang, and Jiangmiao Pang. BeamDojo: Learning agile humanoid locomotion on sparse footholds. In *Proceedings of Robotics: Science and Systems*, Los Angeles, CA, USA, June 2025. doi:10.15607/RSS.2025.XXI.068. URL <https://www.roboticsproceedings.org/rss21/p068.html>.
- [11] Ananye Agarwal, Ashish Kumar, Jitendra Malik, and Deepak Pathak. Legged locomotion in challenging terrains using egocentric vision, 2022. URL <https://arxiv.org/abs/2211.07638>.
- [12] Takahiro Miki, Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learning robust perceptive locomotion for quadrupedal robots in the wild. *Science Robotics*, 7(62):eabk2822, 2022. doi:10.1126/scirobotics.abk2822. URL <https://doi.org/10.1126/scirobotics.abk2822>.
- [13] Ilija Radosavovic, Sarthak Kamat, Trevor Darrell, and Jitendra Malik. Learning humanoid locomotion over challenging terrain, 2024. URL <https://arxiv.org/abs/2410.03654>.

- 304 [14] Junfeng Long, Junli Ren, Moji Shi, Zirui Wang, Tao Huang, Ping Luo, and Jiangmiao Pang.
305 Learning humanoid locomotion with perceptive internal model, 2024. URL <https://arxiv.org/abs/2411.14386>.
306
- 307 [15] Yuanhang Zhang, Younggyo Seo, Juyue Chen, Yifu Yuan, Koushil Sreenath, Pieter Abbeel,
308 Carmelo Sferrazza, Karen Liu, Rocky Duan, and Guanya Shi. RPL: Learning robust humanoid
309 perceptive locomotion on challenging terrains, 2026. URL <https://arxiv.org/abs/2602.03002>.
310
- 311 [16] Wandong Sun, Yongbo Su, Leoric Huang, Alex Zhang, Dwyane Wei, Mu San, Daniel Tian,
312 Ellie Cao, Baoshi Cao, Yang Liu, Finn Yan, Ethan Xie, and Zongwu Xie. Now You See That:
313 Learning end-to-end humanoid locomotion from raw pixels, 2026. URL <https://arxiv.org/abs/2602.06382>.
314
- 315 [17] Ziwen Zhuang, Shenzhe Yao, and Hang Zhao. Humanoid parkour learning, 2024. URL
316 <https://arxiv.org/abs/2406.10759>.
- 317 [18] Zhen Wu, Xiaoyu Huang, Lujie Yang, Yuanhang Zhang, Xi Chen, Pieter Abbeel, Rocky Duan,
318 Angjoo Kanazawa, Carmelo Sferrazza, Guanya Shi, and C. Karen Liu. Perceptive Humanoid
319 Parkour: Chaining dynamic human skills via motion matching, 2026. URL <https://arxiv.org/abs/2602.15827>.
320
- 321 [19] David Hoeller, Nikita Rudin, Dhionis Sako, and Marco Hutter. Anymal parkour: Learning agile
322 navigation for quadrupedal robots, 2023. URL <https://arxiv.org/abs/2306.14874>.
- 323 [20] Péter Fankhauser, Michael Bloesch, and Marco Hutter. Probabilistic terrain mapping for mobile
324 robots with uncertain localization. *IEEE Robotics and Automation Letters*, 3(4):3019–
325 3026, 2018. doi:10.1109/LRA.2018.2849506. URL <https://doi.org/10.1109/LRA.2018.2849506>.
326
- 327 [21] David D. Fan, Kyohei Otsu, Yuki Kubo, Anushri Dixit, Joel Burdick, and Ali-akbar Agha-
328 mohammadi. STEP: Stochastic traversability evaluation and planning for risk-aware off-road
329 navigation. In *Proceedings of Robotics: Science and Systems*, Virtual, July 2021. doi:
330 10.15607/RSS.2021.XVII.021. URL [https://www.roboticsproceedings.org/rss17/](https://www.roboticsproceedings.org/rss17/p021.html)
331 [p021.html](https://www.roboticsproceedings.org/rss17/p021.html).
- 332 [22] Péter Fankhauser and Marco Hutter. A universal grid map library: Implementation and use
333 case for rough terrain navigation. In Anis Koubaa, editor, *Robot Operating System (ROS):*
334 *The Complete Reference (Volume 1)*, volume 625 of *Studies in Computational Intelligence*,
335 chapter 5, pages 99–120. Springer, Cham, 2016. doi:10.1007/978-3-319-26054-9_5. URL
336 https://doi.org/10.1007/978-3-319-26054-9_5.
- 337 [23] Ilija Radosavovic, Tete Xiao, Bike Zhang, Trevor Darrell, Jitendra Malik, and Koushil
338 Sreenath. Real-world humanoid locomotion with reinforcement learning, 2023. URL <https://arxiv.org/abs/2303.03381>.
339
- 340 [24] Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. RMA: Rapid motor adaptation
341 for legged robots. In *Proceedings of Robotics: Science and Systems*, Virtual, July 2021. doi:
342 10.15607/RSS.2021.XVII.011. URL [https://www.roboticsproceedings.org/rss17/](https://www.roboticsproceedings.org/rss17/p011.html)
343 [p011.html](https://www.roboticsproceedings.org/rss17/p011.html).
- 344 [25] Tong Zhang, Boyuan Zheng, Ruiqian Nai, Yingdong Hu, Yen-Jen Wang, Geng Chen, Fanqi
345 Lin, Jiongye Li, Chuye Hong, Koushil Sreenath, and Yang Gao. HuB: Learning extreme
346 humanoid balance, 2025. URL <https://arxiv.org/abs/2505.07294>.
- 347 [26] Lequn Fu, Yijun Zhong, Xiao Li, Yibin Liu, Zhiyuan Xu, Jian Tang, and Shiqi Li. Load-
348 aware locomotion control for humanoid robots in industrial transportation tasks, 2026. URL
349 <https://arxiv.org/abs/2603.14308>.

- 350 [27] Anuj Pasricha, Joewie Koh, Jay Vakil, and Alessandro Roncone. Dynamics-compliant trajec-
351 tory diffusion for super-nominal payload manipulation, 2025. URL [https://arxiv.org/
352 abs/2508.21375](https://arxiv.org/abs/2508.21375).
- 353 [28] Botian Xu, Haoyang Weng, Qingzhou Lu, Yang Gao, and Huazhe Xu. Facet: Force-adaptive
354 control via impedance reference tracking for legged robots, 2025. URL [https://arxiv.
355 org/abs/2505.06883](https://arxiv.org/abs/2505.06883).
- 356 [29] Peiyuan Zhi, Peiyang Li, Jianqin Yin, Baoxiong Jia, and Siyuan Huang. Learning a uni-
357 fied policy for position and force control in legged loco-manipulation, 2025. URL [https:
358 //arxiv.org/abs/2505.20829](https://arxiv.org/abs/2505.20829).
- 359 [30] Jiaqi Chen, Jonas Frey, Ruyi Zhou, Takahiro Miki, Georg Martius, and Marco Hutter. Ident-
360 ifying terrain physical parameters from vision - towards physical-parameter-aware loco-
361 motion and navigation. *IEEE Robotics and Automation Letters*, 9(11):9279–9286, 2024. doi:
362 [10.1109/LRA.2024.3455788](https://doi.org/10.1109/LRA.2024.3455788). URL <https://doi.org/10.1109/LRA.2024.3455788>.
- 363 [31] Hajun Kim, Dongyun Kang, Min Gyu Kim, Gijeong Kim, and Hae Won Park. On-
364 line friction coefficient identification for legged robots on slippery terrain using smoothed
365 contact gradients. *IEEE Robotics and Automation Letters*, 10(4):3150–3157, 2025. doi:
366 [10.1109/LRA.2025.3541428](https://doi.org/10.1109/LRA.2025.3541428). URL <https://doi.org/10.1109/LRA.2025.3541428>.
- 367 [32] Johannes Engelsberger, Christian Ott, and Alin Albu-Schäffer. Three-dimensional bipedal walk-
368 ing control based on divergent component of motion. *IEEE Transactions on Robotics*, 31(2):
369 355–368, 2015. doi:[10.1109/TRO.2015.2405592](https://doi.org/10.1109/TRO.2015.2405592).
- 370 [33] Majid Khadiv, Alexander Herzog, S. Ali. A. Moosavian, and Ludovic Righetti. Walking con-
371 trol based on step timing adaptation. *IEEE Transactions on Robotics*, 36(3):629–643, 2020.
372 doi:[10.1109/TRO.2020.2982584](https://doi.org/10.1109/TRO.2020.2982584).
- 373 [34] Twan Koolen, Tomas De Boer, John Rebula, Ambarish Goswami, and Jerry Pratt. Cap-
374 turrability-based analysis and control of legged locomotion, part 1: Theory and application
375 to three simple gait models. *The International Journal of Robotics Research*, 31:1094–1113,
376 07 2012. doi:[10.1177/0278364912452673](https://doi.org/10.1177/0278364912452673).
- 377 [35] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
378 policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017. URL [https://
379 arxiv.org/abs/1707.06347](https://arxiv.org/abs/1707.06347).
- 380 [36] Nikita Rudin, David Hoeller, Philipp Reist, and Marco Hutter. Learning to walk in minutes
381 using massively parallel deep reinforcement learning. In Aleksandra Faust, David Hsu, and
382 Gerhard Neumann, editors, *Proceedings of the 5th Conference on Robot Learning*, volume
383 164 of *Proceedings of Machine Learning Research*, pages 91–100. PMLR, 08–11 Nov 2022.
384 URL <https://proceedings.mlr.press/v164/rudin22a.html>.
- 385 [37] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based
386 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
387 pages 5026–5033, 2012. doi:[10.1109/IROS.2012.6386109](https://doi.org/10.1109/IROS.2012.6386109).

388 Appendix

389 A Implementation Details

390 **DCM derivation (§3.1).** Liu et al. [7] show that for a linear CoM height profile $z(t) = k_z t + z_0$
 391 during swing, the Variable-Height Inverted Pendulum dynamics reduce to $\dot{\xi} \approx \omega(t) \xi$ when $a =$
 392 $1 + k_z / (2\sqrt{gz_0}) \approx 1$; we adopt this result and set $\omega = \omega_0$. The steady-state nominal step offset that
 393 keeps the divergent mode bounded [34] is $b_x = v_x T / (e^{\omega_0 T} - 1)$ and $b_y = \text{sign}(f) l_p / (1 + e^{\omega_0 T})$,
 394 where l_p is the nominal lateral inter-foot distance and $\text{sign}(f) \in \{+1, -1\}$ selects the swinging leg.

Table 1: DCM foothold planner configuration.

Parameter	Value	Notes
z_0	0.90 m	Nominal CoM height
$T_{\text{swing}} (= T \text{ in } \S 3.1)$	0.45 s	Half gait period
l_p	0.20 m	Lateral inter-foot distance
κ	0.4	Bézier apex bias gain
b_{\min}, b_{\max}	0.25, 0.75	Bézier bias clamp bounds
c_{\min}	0.05 m	Base swing clearance
s	0.5	Clearance scale (m per m step height)
c_{\max}	0.20 m	Maximum swing clearance
$\delta_{\ell}^-, \delta_{\ell}^+$	0.30, 0.05	Lift phase: swing fractions before/after apex where tangent orientation is blended
δ_r^-, δ_r^+	0.05, 0.25	Reach phase: swing fractions before/after apex where landing orientation is blended
Elevation map	$37 \times 25, 0.05$ m	1.80×1.20 m at pelvis
Foot footprint kernel	0.25×0.10 m	For Q_i, E_i pooling
Search window	± 0.30 m (x and y)	Around nominal stride; fallback to nominal at mean valid height if empty
h_{\min}^*	0.05 m	Effective step height at standstill
h_{\max}^*	0.28 m	Effective step height at rated speed
v^*	0.5 m/s	Rated forward speed for h_{eff}^* scaling (Eq. (5))
v_{\min}	0.05 m/s	Standing override: planned target replaced by current foot position
Lateral penalty	$\beta = 2.5$	Lateral vs. sagittal asymmetry
$\alpha_{\text{pos}}, \alpha_{\text{dcm}}$	1.0, 0.5	Position and capture-point stability
$\alpha_E, \alpha_Q, \alpha_M, \alpha_{\text{climb}}$	0.6, 4.0, 6.0, 1.5	Terrain channels (Eq. (1))
<i>Gait-adaptive extension</i>		
f_{\min}, f_{\max}	1.0, 1.5 Hz	Gait frequency bounds
α_{EMA}	0.2	Frequency smoothing weight

395 B Policy Learning Details

396 **Observation space (actor).** Base angular velocity, projected gravity vector, velocity command,
 397 joint positions and velocities relative to default, last action, binary foot contacts, gait phase ($\sin \phi$,
 398 $\cos \phi$), and a stacked depth image (four temporally delayed frames at 30 Hz, blurred, range-clipped,
 399 and normalized). Proprioceptive signals are maintained as a 5-frame history at the 50 Hz control
 400 frequency; the depth image stack is updated at the camera pipeline rate (30 Hz).

401 **Observation space (critic, training only).** Actor observations plus: true CoM velocity, noiseless
 402 elevation map (ray misses filled with a sentinel value), ground-truth foot contact forces and heights,
 403 foot air times, and planner landing targets \mathbf{p}_f^* for each foot.

404 **Training hyperparameters.** PPO with clip parameter $\epsilon = 0.2$, learning rate 10^{-3} with adaptive
 405 KL-based scheduling, GAE $\lambda = 0.95$, discount $\gamma = 0.99$, 5 learning epochs per rollout, 4 mini-
 406 batches, 24 steps per environment per rollout.

Table 2: Lower-body compliance training parameters (§3.4).

Parameter	Value	Notes
<i>Spring-damper</i>		
R_a	0.40 m	Anchor ball radius
k	200 N/m	Virtual spring stiffness
c	20 N-s/m	Virtual damping
<i>Wrench sampling</i>		
ρ_{body}	0.6	Body-attached component probability
R_{close}	0.10 m	Body-attached ball radius
n	2	Cosine-power exponent (bias toward $-\hat{z}$)
n_r	3	Cosine-power exponent (arm-extended forward bias)
$(a_{\text{fwd}}, a_{\text{lat}}, a_{\text{vert}})$	(0.50, 0.40, 0.30) m	Arm-extended half-ellipsoid semi-axes
ϵ	0.1	Isotropic mixture weight
<i>Compliance targets (Eq. (13), (14))</i>		
k_{leg}	300 N/m	Effective leg stiffness; gives $\alpha_z = k/(k + k_{\text{leg}}) = 0.40$
α_z	$k/(k + k_{\text{leg}}) = 0.40$	Height compliance gain
k_{rot}	50 N-m/rad	Rotational stiffness (moment-to-tilt conversion)
$\alpha_\varphi, \alpha_\psi$	0.5	Orientation compliance gains (pitch, roll)

Table 3: Reward terms. Terrain-specific terms are highlighted; all other terms are standard locomotion rewards. I_f denotes foot contact, F_f contact force, L_c centroidal angular momentum.

Term	w	Definition
<i>Velocity tracking</i>		
$r_{\text{vel},xy}$	+3.5	$\exp(-\ e_{xy}\ ^2/0.25)$
$r_{\text{vel},z}$	+3.0	$\exp(-e_z^2/0.5)$
<i>Gait quality</i>		
r_{gait}	+2.0	Phase match: stance/swing vs. gait clock
r_{airtime}	+2.0	Single-stance duration reward
r_{orient}	-7.0	$\ \mathbf{g}_{xy}\ ^2$ (base tilt)
r_{pelvis}	-3.0	$\ \mathbf{g}_{xy}\ ^2$ (pelvis)
r_{action}	-0.8	$\ a_t - a_{t-1}\ ^2$
$r_{\text{collision}}$	-2.0	Self-collision force > 10 N
r_{limits}	-1.0	Joint-limit violation
r_{slip}	-0.4	Tangential contact velocity
r_L	-0.001	$\ L_c\ ^2$ angular momentum
<i>Terrain-specific (contribution)</i>		
r_{foothold}	+2.1	Eq. (8); $\sigma_p = 10$, $\sigma_d \in \{0, 5.0\}$
$r_{\text{clearance}}$	-0.5	$\sum_f h_f^{\text{foot}} - h_f^{\text{terrain}} - d^* \cdot \ v_{f,xy}\ $, $d^* = 0.06$ m
r_{stumble}	-0.5	$\mathbf{1}[\ F_f^{xy}\ > 4 F_f^z \wedge \ v_f^{xy}\ > 0.15$ m/s]
r_{comply}	-1.5	Eq. (15); $h^* = 0.90$ m, $\alpha_z = k/(k + k_{\text{leg}})$

Table 4: Domain randomization ranges applied at episode startup.

Parameter	Range	Mode
Waist-link mass offset	$[-2, +2]$ kg	add
Other body masses	$\times [0.95, 1.05]$	scale
Foot friction	$[0.3, 1.6]$	absolute
PD gains (K_p, K_d)	$\times [0.7, 1.1]$	scale
Joint stiffness / damping	$\times [0.7, 1.3]$	scale
Joint armature	$\times [0.2, 5.0]$	scale
CoM offset (pelvis)	± 0.05 m	add
Encoder bias	± 0.015 rad per joint	add
Depth camera pose	pos $x, y \pm 0.01$ m, $z (-0.03, +0.01)$ m; rot roll/yaw $\pm 2^\circ$, pitch $\pm 10^\circ$	add

407 **Curriculum.** A two-stage velocity curriculum starts at $[0, 0.5]$ m/s and expands to $[0, 1.0]$ m/s
408 after 120,000 steps. The terrain curriculum uses stairs with 0.05–0.20 m risers and 0.30–0.60 m
409 treads, open-width stairs, and flat sections.

Table 5: Actor network specification.

Component	Specification
Depth encoder	Conv2d: channels [32, 32], kernel 3×3 , max-pool FC: [128, 64], output embedding $\in \mathbb{R}^{64}$
Actor MLP	Input: embed (64) \oplus proprioception; hidden (1024, 512, 256, 128), ELU activation
Action	Joint position increments, clipped to limits
History	5-frame obs history at 50 Hz

410 **Gait-adaptive frequency action.** The policy outputs a scalar gait frequency $f_t \in [f_{\min}, f_{\max}]$
411 as an additional action dimension appended after joint targets. Each physics step the gait phase
412 advances as $\phi_t = (\phi_{t-1} + \Delta t_{\text{phys}} \cdot f_t) \bmod 1.0$. The raw policy output is clipped to $[f_{\min}, f_{\max}]$
413 and smoothed by an EMA before use: $f_t = (1 - \alpha_{\text{EMA}}) f_{t-1} + \alpha_{\text{EMA}} \text{clip}(f_t^{\text{raw}}, f_{\min}, f_{\max})$. The
414 phase and its sinusoidal encoding ($\sin \phi_t, \cos \phi_t$) are included in the actor observation; f_{\min}, f_{\max} ,
415 and α_{EMA} are listed in the DCM planner table (§A).

416 **Auxiliary reward terms.** *Clearance:* $\sum_f |h_f^{\text{foot}} - h_f^{\text{terrain}} - d^*| \cdot \|v_{f,xy}\|$ penalizes swing proximity
417 to the terrain, velocity-weighted so edge collisions at high swing speed incur proportionally larger
418 cost. *Stumble:* fires when $\|F_f^{xy}\| > 4|F_f^z|$ and $\|v_f^{xy}\| > 0.15$ m/s; the criterion $\mu_{\text{eff}} > 4$ reliably
419 distinguishes riser edge-strikes from high-force landings. *Compliant base height:* $(z_{\text{pelvis}} - z_{\text{virt}}^*)^2$
420 tracks a virtual compliant target (eq. (13)); at zero compliance gain it reduces to the standard terrain-
421 relative penalty, preserving the LIPM eigenfrequency $\omega_0 = \sqrt{g/z_0}$ through elevation changes.

422 C Experimental Details

423 **Training setup.** All policies are trained on a single NVIDIA H100 SXM5 80 GB GPU with 8192
424 parallel MuJoCo environments. Each iteration collects 24 simulation steps per environment. The
425 PPO rollout and network update run on the same GPU; physics simulation is CPU-parallelized
426 across the host node’s available cores. All training uses the mjlab framework with MuJoCo 3.8.0
427 as the physics backend. Total wall-clock training time per variant is approximately 10–12 h.

428 **Training terrain configurations.** Figure 6 shows the eight terrain types used during curriculum
429 training. Stair variants (a–d) are the primary training signal for the multi-channel terrain cost and
430 foothold reward; slope and rough terrain (e–h) provide gradient and contact-stability diversity. The
431 policy is never trained with riser heights above 0.20 m; the hard-terrain evaluation extends risers to
432 0.30 m and is fully out-of-distribution.

433 **Evaluation setup.** All ablation variants are evaluated at iteration 20 k with 4096 environments us-
434 ing different random seeds in MuJoCo on stairs, slopes, and rough terrain. Standard evaluation ter-
435 rain spans the same riser range as training: stair flights with risers 0.05 m to 0.20 m and treads 0.25 m
436 to 0.55 m, including 3-step and 5-step flights, open-width stairs (no side walls), slopes up to 23° , and
437 random rough height fields 0–0.15 m. Hard-terrain evaluation uses riser heights 0.20–0.30 m with
438 configurations unseen during training. An episode succeeds if the robot traverses $D_{\text{target}} = 10$ m
439 without triggering fall termination (base tilt $> 70^\circ$, self-collision, or joint-limit violation). SR_{hard}
440 is computed over the strict interior riser range [0.22, 0.28] m, excluding the boundary cases at the
441 edges of the hard sweep.

442 **Ablation variants.** All variants share the base training setup; only cost weights, reward terms, and
443 observation inputs differ. Each is evaluated at iteration 20 k with 4096 environments over different
444 random seeds.

445 *TACT + Adaptive Gait (Ours):* full system with all four terrain-cost channels ($\alpha_E, \alpha_Q, \alpha_M, \alpha_{\text{climb}} >$
446 $0; \beta > 1$ controls lateral asymmetry in d_{pos}), foothold tracking reward, and gait-adaptive frequency
447 head [6].

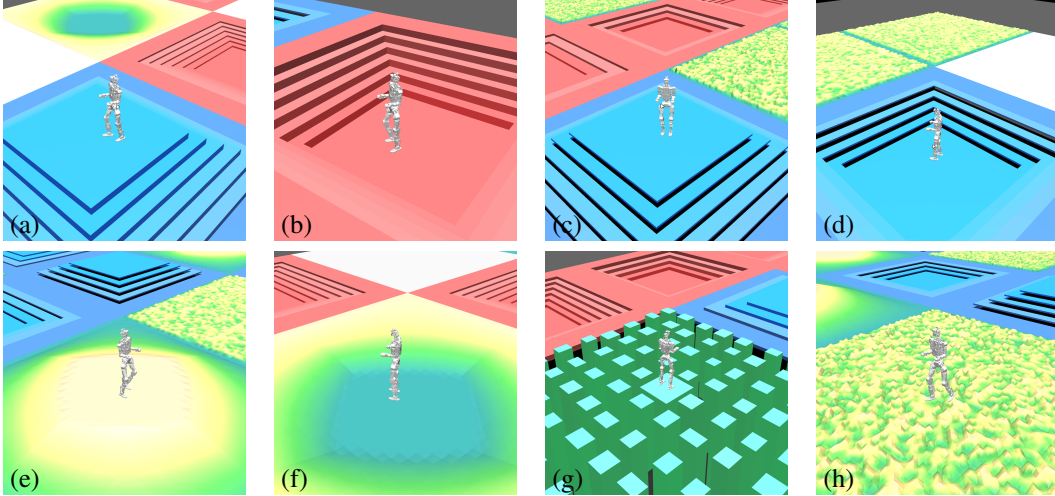


Figure 6: **Training terrain configurations.** (a) Pyramid stairs, ascending. (b) Pyramid stairs, descending. (c) Open-width stairs, ascending (no side walls). (d) Open-width stairs, descending. (e) Pyramid slope, ascending. (f) Pyramid slope, descending. (g) Stepping stones. (h) Gravel (random rough height field). Stair risers span 0.05 m to 0.20 m with treads 0.25 m to 0.55 m; slopes up to 23°; rough field height 0–0.15 m. Hard-terrain evaluation extends risers to 0.30 m (out-of-distribution).

448 *TACT-only*: all terrain-cost channels and foothold tracking reward active; gait-adaptive frequency
 449 head removed.

450 *Adaptive Gait only*: gait-adaptive frequency head added; all terrain-cost channel weights and
 451 foothold tracking reward set to zero ($\alpha_Q = \alpha_E = \alpha_M = \alpha_{\text{climb}} = 0$, $w_{\text{foothold}} = 0$); no privileged height-
 452 map input to the critic.

453 *Baseline*: standard depth-map perceptive policy; neither TACT channels nor gait adaptation; no
 454 privileged critic input.

455 **Payload generalization (terrain).** Three attachment conditions are evaluated: pelvis-mounted
 456 +15 kg and +20 kg (CoM offset; tests balance under downward wrench) and wrist-mounted +10 kg
 457 (large moment arm; tests trunk-tilt compliance). Mass is added at test time with no policy modifi-
 458 cation. A no-mass-DR baseline is trained identically to the full policy but without waist-link mass
 459 offset randomization. Each condition uses 100 episodes at $v_x = 0.5 \text{ m s}^{-1}$ on a 5-step staircase with
 460 0.20 m risers.

461 **Planned evaluations (future work).** The following targeted protocols are specified for complete-
 462 ness; their quantitative results are deferred to future work and not claimed here: non-nominal pelvis-
 463 height tracking ($h^* \in \{h_0 - 6, h_0 - 3, h_0 + 3\}$ cm, recording h^* -RMSE), and a maximum sustained
 464 downward pull sweep (F_{max} at SR > 90 %).

465 **Metric definitions.**

- 466 • SR: success rate (%), fraction of episodes that traverse D_{target} .
- 467 • SR_{hard}: SR restricted to riser heights $\Delta z \in [0.22, 0.28]$ m; isolates performance on the most
 468 demanding subset of the hard sweep (strict interior, excluding boundary cases at 0.20 and 0.30 m).
- 469 • Q_c : mean flatness cost Q_i (Eq. (2)) at the elevation-map cell under each foot at touchdown,
 470 averaged over all touchdown events.
- 471 • F^{95} : $\text{pct}_{95}(\|F_f\|/(mg/2))$, 95th-percentile normalized ground reaction force at touchdown
 472 events.
- 473 • E_v : velocity-tracking RMSE, $\sqrt{N^{-1} \sum_{t=1}^N \|v_{xy,t} - v_{\text{cmd},t}\|^2}$ (m/s).

- 474 • P : mean absolute mechanical power, $N^{-1} \sum_{t=1}^N \|\boldsymbol{\tau}_t \odot \dot{\mathbf{q}}_t\|_1$ (W); lower values indicate more
475 efficient gait [16].
- 476 • **Foot-target distance**: mean planar Euclidean distance between the executed touchdown centroid
477 and the nearest DCM-planned foothold at swing initiation, averaged over all touchdown events
478 in the episode (the “foot-target distance” reported in §4.1).

479 D Hardware and Deployment

480 **Sensor architecture.** The actor policy receives a single perceptual stream at runtime: a noisy
481 depth image from a forward-facing RGB-D (color-and-depth) camera. A raycaster-based elevation
482 map mounted at the pelvis is available during training only and is used exclusively by the privileged
483 critic; it is not available at deployment. Additional per-foot proximity rays and a pelvis height
484 sensor supply scalar terrain signals used exclusively in reward computation and privileged critic
485 observations. The simulation models each sensor with geometry-matched raycasters; the depth
486 camera adopts the intrinsics of a RealSense D435 and injects pose domain randomization to bridge
487 the sim-to-real gap.

488 **Depth camera processing pipeline.** Raw depth frames are captured at 30 Hz, cropped to remove
489 border artifacts, resized to 30×30 px, Gaussian-blurred to suppress raycaster aliasing, and normal-
490 ized to $[0, 1]$. Four temporally spaced frames are stacked into a history buffer and fed to the depth
491 encoder at each 50 Hz control step. Camera intrinsics are matched to a RealSense D435; extrinsic
492 pose domain randomization (table 4) covers mounting uncertainty. The identical pipeline runs at
493 deployment with no modification.

494 **Deployment.** The trained actor (depth encoder + policy MLP) is exported to ONNX and executed
495 at 50 Hz on the robot’s onboard computer via a joint-position PD interface. No real-world fine-
496 tuning is performed; the sim-to-real gap is closed entirely by domain randomization during training.

497 E Additional Results

498 Flat-Terrain Payload Generalization

499 Fig. 7 evaluates payload generalization across three conditions (pelvis +15 kg, pelvis +25 kg, wrist +15 kg) on flat terrain and compares Ours against the Baseline, isolating compliance behavior from terrain traversal difficulty. At moderate load (pelvis +15 kg, wrist +15 kg), Ours maintains 76–79% SR against the Baseline’s 67–76% while consuming 9–20% less power, consistent with compliance training suppressing impulsive GRF recovery. At high centered load (pelvis +25 kg), Ours drops to 47.5% (below the Baseline’s 60.0%), indicating that a large CoM shift exceeds the compliance training distribution and causes over-compensation; power remains lower (257 vs. 283 W), confirming the policy still reduces impulsive contacts but can no longer maintain balance. Wrist-mounted mass, which generates a large moment rather than a direct CoM offset, is handled comparably in SR but with a clear power advantage for Ours (186 vs. 231 W), consistent with the arm-extended wrench samples in the disturbance force-field.

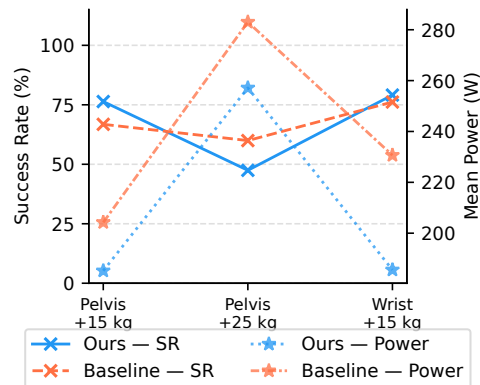


Figure 7: SR (%) and mean power (W) vs. payload mass on flat terrain.

500 **Cross-Embodiment Generalization**

501 Fig. 8 shows the ablation evaluated on two platforms (Platform-A (H1-2 class) and Unitree G1)
 502 using identical reward weights and hyperparameters, with only the robot asset and joint configura-
 503 tion changed. The relative ordering of variants is preserved across embodiments, supporting the
 504 platform-agnostic claim in §6.

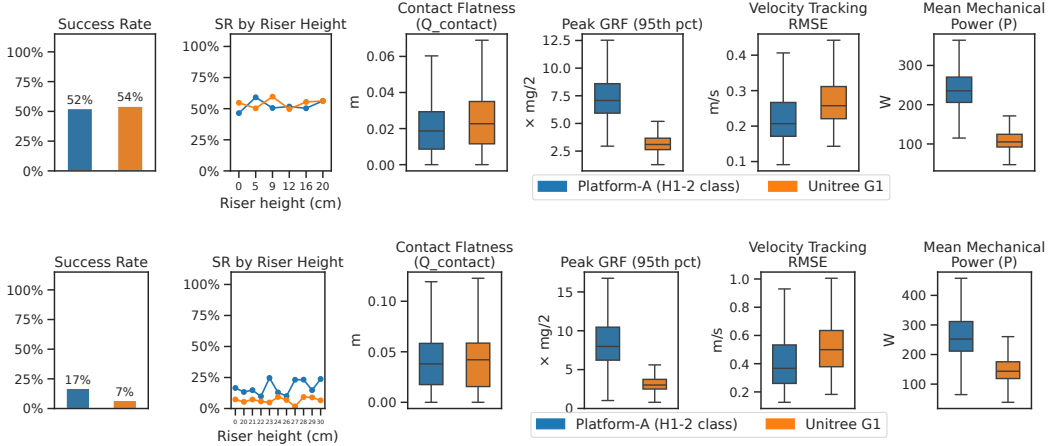


Figure 8: Cross-embodiment ablation on Platform-A (H1-2 class) and Unitree G1: **standard terrain** (top) and **hard terrain, out-of-distribution** (bottom). Identical reward weights and hyperparameters; only robot asset and joint configuration differ.

505 **F Extended Limitations**

506 **Adaptive-gait fragility under kinematic stress.** The *Adaptive Gait* only variant falls below the
 507 Baseline on hard terrain (18 % vs. 19 % SR): EMA-smoothed frequency re-timing without terrain-
 508 quality guidance modifies stance-to-swing departure timing with no corresponding adjustment to
 509 the landing target, producing foot-target mismatches that the foothold reward would otherwise sup-
 510 press. Adaptive gait frequency is therefore an auxiliary capability contingent on accurate foothold
 511 guidance, not an independent locomotion primitive, and its benefit is confined to regimes where
 512 kinematic margins are non-trivial. A natural remedy is to co-condition the frequency head on the
 513 terrain-cost signal so that re-timing decisions are aware of contact quality, or to gate frequency
 514 updates on a kinematic feasibility check before committing to the new landing target.

515 **Terrain-cost channel redundancy.** Removing any single terrain cost channel leaves foot-target
 516 distance within 5 % of the full system (0.085–0.093 m vs. 0.088 m), while removing all three raises
 517 it 2.8× to 0.251 m. The channels are thus collectively necessary but individually non-critical, and
 518 the manually tuned weighting ($\alpha_Q=4.0$, $\alpha_E=0.6$, $\alpha_M=6.0$) does not follow from principled mini-
 519 mality. A compressed two-channel formulation and a weight-sensitivity sweep across platforms are
 520 left to future work. More broadly, replacing hand-tuned weights with a differentiable meta-learning
 521 or Bayesian optimization pass over the cost coefficients could yield principled, platform-adaptive
 522 weightings without requiring per-robot re-tuning.

523 **Compliance distribution and upper-body coupling.** Near 20 kg of centered CoM shift, two fac-
 524 tors compound. First, the spring-damper sampling concentrates on body-attached downward-biased
 525 and arm-extended forward-biased scenarios; lateral carry, two-point distributed loads, and dynami-
 526 cally swinging payloads produce wrenches outside this sampled space. Second, the policy controls
 527 leg joints only: at high payload masses, upper-body inertia shifts the effective CoM in ways the
 528 lower-body policy cannot observe, as the pelvis inertial measurement unit (IMU) and joint encoders
 529 provide no direct signal of trunk angular momentum. Together these set the effective payload ceiling.
 530 Both failure modes have tractable remedies: broadening the wrench distribution to include lateral

531 and multi-point loads during training, and extending the controlled joints to include the torso or
532 arm would directly address the sampling gap and the observability gap respectively. An online pay-
533 load estimator feeding a residual compliance target could further extend the effective range without
534 requiring a larger wrench training distribution.

535 **Tangent-guided orientation failure on slopes.** The Bézier swing reference orients the sole using
536 the arc tangent of the trajectory at the landing point, which on discrete stair steps closely approxi-
537 mates the riser-face normal and drives the foot to land with the sole parallel to the tread. On smooth
538 slopes this mechanism fails: the swing trajectory approaches the inclined surface from above along
539 a path whose touchdown tangent is approximately horizontal, independent of terrain slope angle.
540 The resulting orientation reference does not rotate the sole to match the surface normal, so the foot
541 arrives level on a tilted surface. On ascending grades the surface rises toward the direction of travel,
542 placing the toe geometrically closer to the ground than the heel at the moment of touchdown; the
543 foot therefore contacts toe-first, concentrating GRF on the toe edge. On descending grades the
544 inverse holds, producing heel-first contact with the toe elevated. Both deviate from the designed
545 flat-sole landing, narrow the friction cone, and raise F^{95} ; on grades above roughly 15° the resulting
546 ankle-roll moment is frequently sufficient to trigger fall termination.

547 A complementary failure arises in the foothold planner. The flatness cost Q_i measures height vari-
548 ance within a local neighbourhood around each candidate cell: a uniformly inclined surface pro-
549 duces near-zero inter-cell variance, so $Q_i \approx 0$ regardless of slope angle, and the planner assigns
550 high contact quality to steep cells. The steepness cost E_i is intended to compensate, but with the
551 current weighting ($\alpha_E = 0.6$ vs. $\alpha_M = 6.0$) it is the weakest channel and is insufficient to steer
552 footholds away from high-gradient cells when flatness and reachability are simultaneously satisfied.
553 The combination of incorrect sole orientation and inadequate slope-avoidance pressure jointly ac-
554 counts for the observed degradation on slope terrain relative to stairs of comparable height gain per
555 step. A corrected planner would supplement Q_i with a surface-normal tilt term and extend the swing
556 orientation reference from the trajectory tangent to the estimated surface normal at the target cell,
557 analogous to normal-aligned foot placement used in model-based planners.

558 Additional Bézier-specific failures emerge on slopes. The apex height is set relative to the take-
559 off elevation; on ascending grades the landing point is higher than take-off by $\Delta z = l_{\text{stride}} \sin \theta$,
560 reducing effective clearance over the landing surface by the same amount. At steep grades and
561 long strides the foot can contact the slope during mid-swing before the apex is reached, a form of
562 early collision absent on flat terrain and stairs. On descending grades the inverse produces excess
563 apex clearance, which extends swing duration and can desynchronize foot arrival with the DCM-
564 timed stance transition. Addressing both the orientation and clearance failures requires replacing
565 the trajectory-tangent reference with a surface-normal estimate at the target cell (available from the
566 elevation map) and setting the apex height relative to the landing elevation rather than the take-off
567 elevation, a straightforward modification to the existing Bézier parameterization.

568 **Elevation-map dependency and terrain scope.** The DCM planner consumes a pelvis-mounted
569 elevation map requiring accurate extrinsic odometry; localization drift degrades foothold selection
570 monotonically with no fallback planner, and the forward-facing depth camera provides no lateral
571 or rearward coverage. Highly deformable surfaces (sand, gravel, foam) violate the rigid-terrain
572 assumption underlying both the Bézier clearance computation and the flatness cost Q_i , causing
573 systematic under-clearance, and wet or icy contacts fall outside the domain-randomization friction
574 range ($\mu \in [0.3, 1.6]$). Terrain with moving obstacles is not modeled. Robustness to localization drift
575 could be improved by incorporating a fallback stance-hold controller that activates when odometry
576 confidence is low, and by fusing the depth camera into a local elevation map that does not depend
577 on global pose. Extending the friction domain randomization range and adding deformable-terrain
578 simulation (e.g., via compliant contact models in MuJoCo) are the most direct paths to covering
579 low-friction and soft-surface deployments. Handling dynamic obstacles would require integrating a
580 short-horizon obstacle prediction model into the foothold search window.