

COMPLIANTWBC: Whole-Body Compliance for Heavy Humanoids via Force Latent Estimation and Residual Impedance Targets

Anonymous Author(s)

Affiliation

Address

email

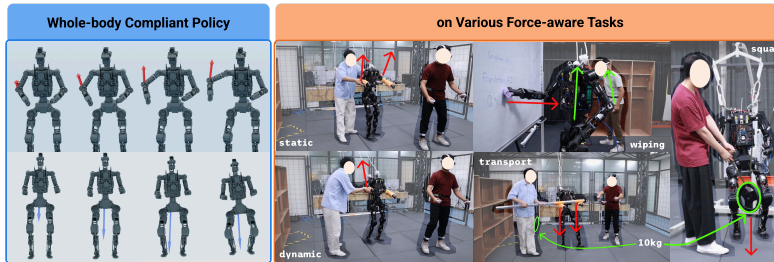


Figure 1: COMPLIANTWBC pairs a force-aware base policy with a bounded residual on the impedance equilibrium, trained under a Phong-weighted force-origin sampler that extends compliance to the whole body. We demonstrate it on a range of force-aware tasks on a real heavy humanoid.

1 **Abstract:** Whole-body compliant control is essential for deploying heavy hu-
2 manoids under high payload in human-centric environments. Most prior force-
3 aware learning-based pipelines focus on end-effector resistance, per-link upper-
4 body springs, or end-effector stiffness modulation, leaving *arbitrary-site pertur-*
5 *bations on heavy platforms with lower-body engagement* largely unaddressed.
6 We close this gap with COMPLIANTWBC comprising: (1) A multi-site whole-
7 body impedance *reference controller* generalizes classical Cartesian impedance
8 to any controlled link via null-space projection and a multi-contact balancing
9 quadratic programming (QP) and supplies per-link virtual targets that supervise
10 an RL policy through a compliance-fidelity reward. (2) A *Force Latent Encoder*,
11 co-trained with the base policy behind a gradient barrier and shaped by wrench
12 reconstruction, supervised-contrastive clustering, a KL bottleneck, and temporal
13 smoothness, yields a wrench-grounded, class-clustered latent; a bounded *Resid-*
14 *ual Policy* then edits the impedance equilibrium to absorb wrench-estimation er-
15 ror, preserving *local* classical-impedance passivity under fixed stiffness. (3) A
16 *Phong-weighted force-origin sampler* with an axis-decoupled pelvis anchor in-
17 duces lower-body-inclusive compliance via two interpretable parameters instead
18 of a hand-crafted schedule. We evaluate COMPLIANTWBC in simulation against
19 compliant and stiff baselines and demonstrate it on a real heavy humanoid across
20 static/dynamic force reaction, board wiping, squat under payload, and cooperative
21 payload transport. Project website: <https://compliantwbc.github.io/>

22 **Keywords:** Whole-Body Compliant Control, Latent Force Modeling, Residual
23 Policy Learning

24 1 Introduction

25 Heavy humanoid loco-manipulation requires whole-body compliance. For instance, pelvis bracing
26 against a wall, hip translation against a pulled cart, or double-support redistribution while lifting
27 cannot be expressed by upper-body springs alone. Classical operational-space methods [1, 2, 3, 4, 5,
Under submission. Do not distribute.

28 6, 7] cover this regime analytically, but require accurate dynamics, hand-tuned task hierarchies, and
29 contact-mode reasoning that are brittle on heavy platforms. Learning-based pipelines relax these
30 requirements; however, the closest predecessors restrict compliance to the end-effector or to the
31 upper-body chain [8, 9, 10, 11, 12], leaving arbitrary-site multi-contact compliance with lower-body
32 impedance participation largely unaddressed.

33 In the learning-based context, the policy observes only proprioception and a motion command (the
34 external wrench \mathbf{f}_{ext} is unobserved on hardware) and must output torques that *yield* compliantly to
35 \mathbf{f}_{ext} at arbitrary contact sites while sustaining balance and command tracking. Transferring classical
36 compliance into a learning pipeline with reinforcement learning (RL) then faces two challenges.
37 First, the classical impedance formulation considers a single end-effector, so a principled *whole-*
38 *body* heuristic must compose centroidal-momentum balance with per-link impedance at arbitrary
39 sites and distribute the resulting wrench to a multi-contact support set. Second, training must actu-
40 ally exercise the whole body (not just the end-effectors) under force perturbations, and the learned
41 policy must *encode* compliance rather than override it. Consequently, residual schemes that only
42 edit motor commands or tracked motions [13, 14] lose the analytical controller’s physical bounds.
43 On the other hand, some whole-body trackers and teleoperation pipelines [15, 16, 17, 18, 19] expose
44 no impedance interface during training, so compliance must be injected separately while running the
45 full analytical pipeline at deployment inherits its runtime cost and contact-mode fragility on heavy
46 platforms.

47 We address both challenges with COMPLIANTWBC (Figure 1). For the first, we derive an analytical
48 multi-site reference controller (§3) composing centroidal-momentum balance with per-link Carte-
49 sian impedance at arbitrary sites. Rather than cloning its torque, we use its per-link virtual targets
50 to supervise a base RL policy through a compliance-fidelity reward [12, 20], so the deployed policy
51 avoids the QP’s runtime cost and contact-mode fragility. For the second, a Phong-weighted force-
52 origin sampler (§4.2) drives perturbations across the whole body—including pelvis and lower-body
53 sites that upper-body-only curricula never excite—so the policy *encodes* compliance rather than
54 overriding it. A variational *Force Latent Encoder*, jointly trained with the base policy behind a
55 gradient barrier from PPO and shaped by wrench reconstruction, contrastive clustering [21], a KL
56 bottleneck, and temporal smoothness, supplies a class-clustered, wrench-grounded estimate of the
57 unobserved contact state [22]. Finally, a bounded residual on the frozen base edits the impedance
58 *equilibrium* rather than motor commands or gains [23, 24], so its effect is upper-bounded by the
59 fixed stiffness and the closed loop is *locally* passive by the classical Cartesian-impedance argu-
60 ment (§F) [1, 25]. We study these behaviors on a heavy humanoid platform (70 kg, versus the 35 kg
61 Unitree G1 [26] used in most prior learning-based compliant humanoid work; specifications in §B)
62 to tackle a wider range of contact-rich tasks.

63 Our contributions can be summarized as follows: (1) An **analytical multi-site whole-body compli-**
64 **ance reference controller** (§3, §A) composing centroidal-momentum balance with per-link Carte-
65 sian impedance at arbitrary sites, whose per-link virtual targets serve as an RL compliance-fidelity
66 reward rather than a behavior-cloning target. (2) A **Force Latent Encoder with residual-on-**
67 **equilibrium architecture** (§4.1) whose latent is grounded in privileged wrench supervision and
68 clustered by perturbation class, and whose bounded residual edits the impedance equilibrium, pre-
69 serving analytic compliance under sim-to-real adaptation. (3) A **Phong-weighted force-origin sam-**
70 **pler with an axis-decoupled pelvis anchor** (§4.2) that induces multi-contact, lower-body-inclusive
71 compliance, replacing a hand-crafted per-site frequency/magnitude schedule, validated on a real
72 heavy humanoid (§5).

73 2 Related Work

74 2.1 Compliance Control for Humanoids

75 Classical compliance strategies, including task-space impedance and admittance control [1, 2, 3,
76 4, 5, 6, 7], regulate interaction forces via virtual mass–spring–damper dynamics, with operational-
77 space extensions that handle multiple frames through hierarchical task-priority projection and a

Table 1: Humanoid whole-body-control landscape. *Lower-body compliance*: yields with the legs/pelvis (\checkmark), only in a quadruped setting (quadruped), or not at all ($-$). *Multi-site*: where compliance is realized—end-effector only (EE), upper body (upper only), center-of-mass reference (CoM), or arbitrary body links (whole-body). *Force latent*: an explicit, policy-conditioning force/contact embedding (\checkmark), an implicit one (implicit), or none ($-$). *Residual architecture*: where a learned correction is applied.

Method	Lower-body Compliance	Multi-site	Force latent	Residual architecture	Humanoid platform
FALCON [8]	$-$	EE	$-$	$-$	G1/Booster
GentleHumanoid [11]	$-$	upper only	implicit	$-$	G1
FACET [10]	quadruped	CoM	$-$	impedance ref.	Go2/G1
SoftMimic [12]	\checkmark	whole-body	$-$	$-$	G1
CHIP [9]	$-$	EE	$-$	hindsight goal	G1
ResMimic [14]	$-$	$-$	$-$	on motion	G1
ASAP [13]	$-$	$-$	$-$	on action	G1
CompliantWBC (ours)	\checkmark	\checkmark	\checkmark	on equilibrium	in-house, 70 kg

78 balancing QP that distributes the centroidal wrench to support contacts. These pipelines provide rig-
79 orous stability guarantees but require accurate dynamics, hand-crafted task hierarchies, and explicit
80 contact-mode reasoning that is brittle on heavy humanoid platforms; time-varying stiffness addition-
81 ally breaks passivity unless an explicit storage mechanism is enforced [27, 28]. Recent whole-body
82 tracking controllers [29, 30, 31] may achieve agile loco-manipulation via motion imitation but do
83 not expose an impedance interface, so compliance must be injected separately.

84 To relax these requirements, recent learning-based methods train policies to exhibit compliant beh-
85 avior through reward shaping, perturbation curricula, per-link virtual springs, or end-effector stiff-
86 ness modulation [8, 11, 10, 9]. Most restrict compliance to the end-effector or the upper-body
87 chain, leaving lower-body bracing, hip translation, and double-support redistribution—the regimes
88 that dominate heavy loco-manipulation—unaddressed. SoftMimic [12] is a notable exception that
89 achieves genuine whole-body compliance, but tracks a single motion rather than an arbitrary con-
90 figuration. A parallel line targets heavy-payload loco-manipulation with multi-policy or trajectory-
91 optimized references [32, 33] but similarly does not expose a per-link impedance interface. Our
92 work composes centroidal-momentum balance and per-link Cartesian impedance at arbitrary sites
93 inside an RL loop, admitting end-effector-only and upper-body-only formulations as special cases
94 (Table 1).

95 2.2 Residual Policies and Latent Force Estimation

96 Residual reinforcement learning [23, 34, 35] composes a learned correction on top of a hand-crafted
97 base controller, retaining analytical stability while absorbing unmodeled dynamics. The idea of
98 learning residuals on controller parameters with RL was pioneered by Buchli et al. [36], who
99 optimized both reference trajectories and gain schedules on hardware. Our work operates in the
100 same spirit but residualizes on the impedance *equilibrium* rather than on gains. Recent humanoid
101 pipelines instantiate this residual idea differently, focusing on motor commands, on tracked motions,
102 or on goals via hindsight relabeling [13, 14, 9]. However, none of these approaches edit a per-link
103 impedance interface with an explicit force bound. Consequentially, the residual’s effect on induced
104 contact force is implicit rather than physically bounded. In contrast, we compose per-link impedance
105 at arbitrary sites and place a *bounded* residual on the per-link *equilibrium* over a *frozen* base.

106 A complementary line learns compact latent embeddings of unobserved context from proprioceptive
107 history, with supervised or contrastive objectives shown to prevent collapse to motion-phase memo-
108 rization [37, 38, 39, 40, 21]. Recent work specializes such latents to estimate external wrenches from
109 proprioception without dedicated force sensors [41, 42, 43]. Relative to this RMA-style lineage, we
110 infer a wrench-grounded, class-clustered latent from history observations behind a gradient barrier,
111 which drives a *bounded residual on the per-link impedance equilibrium*. Because the link stiffness
112 \mathbf{K}_ℓ is held fixed, the residual’s effect is upper-bounded by an explicit virtual-target shift times the
113 stiffness, so the closed loop is *locally* passive by the classical Cartesian-impedance argument (§F).

114 3 Whole-Body Impedance Heuristic

115 A humanoid has configuration $\mathbf{q} \in \text{SE}(3) \times \mathbb{R}^{n_j}$, generalized velocity $\boldsymbol{\nu}$, and floating-base dynamics
 116 $M\dot{\boldsymbol{\nu}} + C\boldsymbol{\nu} + \mathbf{g} = \mathbf{S}^\top \boldsymbol{\tau} + \sum_c \mathbf{J}_c^\top \boldsymbol{\lambda}_c + \sum_i \mathbf{J}_{\ell_i}^\top \mathbf{f}_{\text{ext},i}$ with support-contact reaction wrenches $\boldsymbol{\lambda}_c$
 117 and external wrenches $\mathbf{f}_{\text{ext},i}$ at links ℓ_i . Classical Cartesian impedance at a single end-effector [1]
 118 prescribes $M_d \ddot{\mathbf{x}} + D_d \dot{\mathbf{x}} + K_d(\mathbf{x} - \mathbf{x}_d) = \mathbf{f}_{\text{ext}}$, with $\mathbf{x} \in \mathbb{R}^3$. We extend this to a hierarchy of
 119 frames coupled by balance.

120 **Reference controller (summary).** For controlled links \mathcal{L} (hands, elbows, knees, pelvis, torso) and
 121 support contacts \mathcal{C}_s , the analytical reference controller composes (i) a centroidal-momentum balance
 122 task whose wrench is distributed to the support contacts by a balancing QP over friction cones and
 123 centers of pressure [5], (ii) per-link Cartesian impedance at arbitrary sites projected into the balance
 124 null space [3], and (iii) a joint-space posture task. It produces a reference torque $\boldsymbol{\tau}_c$ and per-link
 125 virtual targets $\mathbf{x}_{\ell,d}^c$. We give the full controller, including the energy-tank passivity treatment for
 126 time-varying stiffness, in §A, and use it strictly as a *reference*: of its outputs, *only* the per-link virtual
 127 targets $\mathbf{x}_{\ell,d}^c$ enter the policy reward (§4). The QP-distributed torque $\boldsymbol{\tau}_c$ is never cloned or rewarded,
 128 so the deployed policy inherits neither the QP’s runtime cost nor its contact-mode fragility on heavy
 129 platforms. Admitting knees and the pelvis as controlled and support sites is the mechanism by which
 130 the heuristic extends two-footed, end-effector-only settings to multi-contact, lower-body-inclusive
 131 postures.

132 **Force-aware virtual targets.** The single signal carried into training is the per-link virtual target.
 133 For a perturbed link $\ell \in \mathcal{L}$ under an external wrench \mathbf{f}_{ext} applied at body site p ,

$$\mathbf{x}_{\ell,d} = \mathbf{x}_{\ell}^{\text{ref}}(t) + \mathbf{K}_{\ell}^{-1} \mathbf{f}_{\ell}(t), \quad \mathbf{f}_{\ell} = \text{Ad}_{p \rightarrow \ell}^{-\top} \mathbf{f}_{\text{ext}}, \quad (1)$$

134 where \mathbf{f}_{ℓ} is the external wrench transported to link ℓ ’s frame (equal to \mathbf{f}_{ext} when applied at ℓ) and
 135 \mathbf{K}_{ℓ}^{-1} is the anisotropic Cartesian compliance, so the correction has units of displacement. This gen-
 136 eralizes GentleHumanoid’s scalar spring to arbitrary sites and anisotropic stiffness: restricting \mathcal{L} to
 137 upper-body links with $\mathbf{K}_{\ell} = K_p \mathbf{I}_3$, $K_p \sim \mathcal{U}(5, 250)$ N/m and $\mathbf{D}_{\ell} = 2\sqrt{m_v K_p} \mathbf{I}_3$ ($m_v = 0.1$ kg) re-
 138 covers GentleHumanoid’s upper-body Cartesian-spring response [11] (up to its motion-data-driven
 139 guiding contacts), and restricting the perturbed set to the end-effectors reduces to FALCON’s end-
 140 effector perturbation model [8] as an analytical special case. Time-varying \mathbf{K}_{ℓ} would break pas-
 141 sivity [28]; we hold \mathbf{K}_{ℓ} fixed (§4.1), which makes the closed loop locally passive by the classical
 142 Cartesian-impedance argument (§F), and defer the energy-tank treatment of the general time-varying
 143 case to §A.

144 4 Method: Force-Aware Base Policy and Residual on Impedance Targets

145 4.1 Pipeline Architecture

146 The analytical controller of §3 assumes rigid contacts, perfect torque control, and accurate dynamics,
 147 and assumes the external wrench \mathbf{f}_{ext} is observed; none holds on a heavy humanoid. We therefore
 148 train a two-stage pipeline (Figure 2): (i) a *force-aware base policy* that jointly trains with a Force La-
 149 tent Encoder to induce compliance behavior, and (ii) a *residual policy* on top of the frozen base that
 150 compensates for the discrepancy between the base’s wrench estimate and the true external wrench.

151 **Base policy.** We first train a force-aware whole-body policy $\pi_{\text{base}}(\mathbf{o}, \mathbf{c}, \mathbf{z}_F) \rightarrow \boldsymbol{\tau}$ that takes as input a
 152 proprioceptive observation \mathbf{o} , a motion command \mathbf{c} , and the latent wrench embedding \mathbf{z}_F produced
 153 by the Force Latent Encoder, and outputs the joint torque. The policy is trained end-to-end with
 154 PPO [44] against a compliance-fidelity reward

$$r_t = r_{\text{track}} - w_c \|\mathbf{x}_{\ell} - \mathbf{x}_{\ell,d}^c\|^2 - w_e \|\boldsymbol{\tau}\|^2, \quad (2)$$

155 whose term $w_c \|\cdot\|^2$ drives the controlled links toward the analytical controller’s virtual-target de-
 156 formation $\mathbf{x}_{\ell,d}^c$ (Eq. (1)) under the ground-truth perturbation. This incentivizes the policy to learn
 157 multi-site compliance without inheriting the balancing QP’s runtime cost at deployment. The full
 158 reward formulations (Table 3) and the two-stage training loop (Algorithm 1) are given in §B.

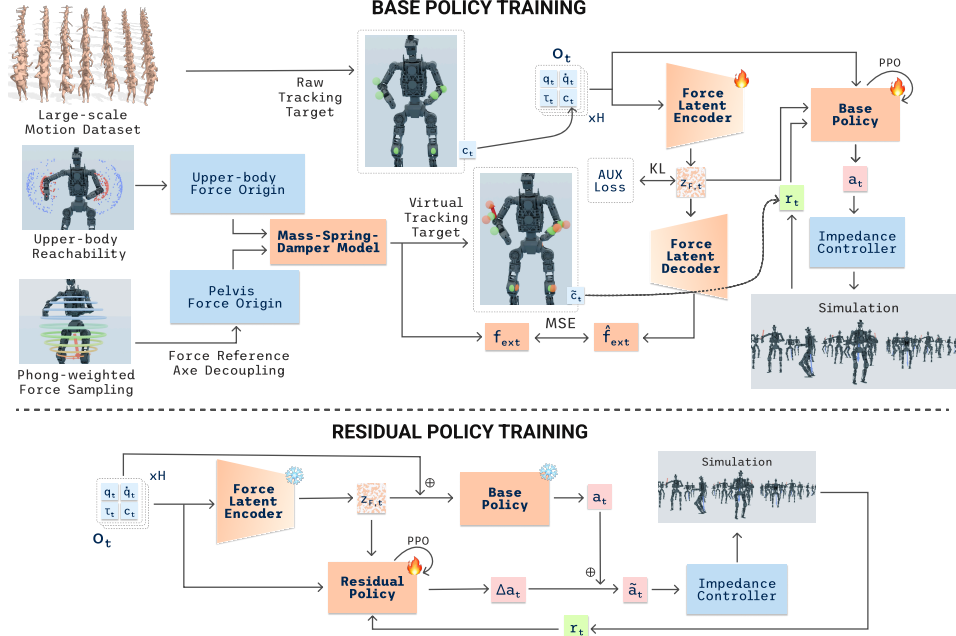


Figure 2: Two-stage training of COMPLIANTWBC: (i) a base policy jointly trained with the Force Latent Encoder to induce whole-body compliance; (ii) a bounded residual on the impedance equilibrium that compensates for physics discrepancies.

159 **Force Latent Encoder** follows a variational encoder architecture that maps an observation history
 160 window (proprioceptions, velocities, last torques, and motion command) to a force latent $z_F \in \mathbb{R}^{d_F}$.
 161 The encoder is paired with two auxiliary heads — a *wrench decoder* ψ that regresses the per-site
 162 external wrench on all controlled links, and a *projection head* g on the unit sphere dedicated to the
 163 contrastive loss. The encoder is shaped *only* by the auxiliary objective in Eq. (3) while the posterior
 164 mean is detached before it is fed to π_{base} and π_{res} . This forms a *gradient barrier* preventing PPO
 165 gradients from flowing through the encoder, so that the latent cannot be co-opted away from its
 166 wrench-grounded objective toward reward-hacking features.

167 The auxiliary objective \mathcal{L}_ϕ combines four terms: a per-site wrench-reconstruction loss (the main
 168 term) grounding z_F in the true contact wrench through ψ ; a supervised contrastive loss [21] over
 169 perturbation-class labels (which limb, which direction); a KL/VIB bottleneck [45] suppressing di-
 170 mensions that carry neither wrench- nor class-relevant information; and a temporal-smoothness term
 171 on consecutive posterior means:

$$\mathcal{L}_\phi = \lambda_w \mathcal{L}_{\text{wrench}} + \lambda_s \mathcal{L}_{\text{supcon}} + \lambda_k \mathcal{L}_{\text{kl}} + \lambda_m \mathcal{L}_{\text{smooth}}. \quad (3)$$

172 Together they make z_F a wrench-grounded, class-clustered, low-bandwidth embedding of the priv-
 173 ileged contact state (full forms in §B).

174 **Residual policy** is a small MLP added on top of the frozen base policy and frozen Force Latent
 175 Encoder. The residual reads $(o_t, c_t, z_F, \hat{f}_{\text{ext}})$, with the frozen wrench estimate $\hat{f}_{\text{ext}} = \psi(z_F)$, and
 176 outputs a bounded delta on the impedance equilibrium:

$$\tilde{x}_{\ell,d} = x_{\ell,d} + \epsilon_x \tanh(\Delta x_{\ell,d}), \quad (4)$$

177 with positional bound $\epsilon_x = 5$ cm. When the frozen wrench estimate $\hat{f}_{\text{ext}} = \psi(z_F)$ is biased or noisy
 178 with respect to the true external wrench, the residual edits the impedance *equilibrium* to absorb the
 179 discrepancy without changing the base policy and without modulating gains. Editing the equilibrium
 180 rather than the torque preserves the base controller’s analytic compliance: a 5 cm shift in $x_{\ell,d}$
 181 induces at most $K_\ell \cdot 5$ cm of force, an explicit physical bound the residual cannot violate. We train
 182 the residual in a separate stage rather than jointly with the base: under a single PPO loss the residual
 183 would absorb compliance the base could itself learn, turning its bound into a constraint on the joint
 184 system rather than an interpretable discrepancy compensator. More details in §B.

185 **4.2 Compliant virtual targets via external-force sampling**

186 **Phong-weighted sampling.** For perturbed upper-body links (wrists, elbows) we sample force origins from a kinematic reachability cloud and drag the link toward them [11]. For the pelvis—which endures much larger forces over a limited range of motion—we instead draw a force origin from a continuous *Phong-weighted* distribution inside a ball of radius $R_a = 0.5$ m centered on the pelvis, with direction \hat{n} following a Phong lobe [46] concentrated on the downward axis $-\hat{e}_z$ with exponent $n=2$:

$$p(\hat{n}) = (1 - \epsilon) \frac{n + 1}{2\pi} \max(-\hat{n} \cdot \hat{e}_z, 0)^n + \epsilon \frac{1}{4\pi}, \quad (5)$$

192 with isotropic mixing weight $\epsilon = 0.1$, sampled in closed form by inverse CDF (§C) to give the offset $\delta = r(\sin \theta \cos \phi, \sin \theta \sin \phi, -\cos \theta)$. The downward bias emphasizes gravity on a carried load while lateral and upward samples cover reaching and lifting. Unlike a fixed cloud the distribution is gap-free, needs no forward-kinematics precomputation, and exposes its bias through two interpretable parameters (n, ϵ) (Figure 3).

197 **Axis-decoupled pelvis anchor and compliant height target.** Anchoring the offset δ rigidly to the moving pelvis yields a self-chasing constant load, whereas a world-fixed anchor lets the force grow without bound as the robot walks away. We instead decouple the anchor by axis: the horizontal component co-moves with the floating base (leaving locomotion a bounded, body-relative load), while the vertical component—the axis we want the policy to yield along—is pinned to the *commanded* base height h^* , giving the pelvis force-origin anchor $\mathbf{a}_{\text{pelvis}}$:

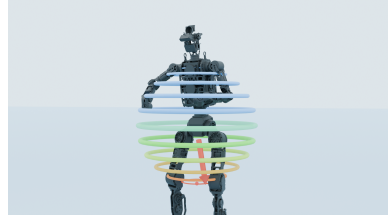


Figure 3: Phong-weighted sampling generates downward-biased force origins, decoupled by axis.

$$\mathbf{a}_{\text{pelvis}}(t) = \underbrace{\mathbf{x}_{\text{pelvis}}^{xy}(t) + \delta^{xy}}_{\text{base-relative (horizontal)}} + \underbrace{(h^* + \delta^z) \hat{e}_z}_{\text{gravity-anchored (vertical)}}, \quad (6)$$

208 so the anchor height $z_a = h^* + \delta^z$ is a per-chunk gravity-aligned datum that does not sink with the robot. We then evaluate the base-height term of r_{track} against a compliance-aware target $z_{\text{virt}}^*(t) = h^* + \alpha_z(z_a(t) - h^*)$, where the gain $\alpha_z = k/(k + k_{\text{leg}})$ is set by a two-spring static balance between the pelvis-anchor stiffness k and the effective vertical leg stiffness k_{leg} (§C): a stiff anchor pulls the comfortable pose toward the load, a stiff stance resists it. Penalizing $(z_{\text{pelvis}} - z_{\text{virt}}^*)^2$ thus interpolates from rigid tracking ($\alpha_z = 0$) to full anchor-following ($\alpha_z = 1$), so compliant lower-body sinking emerges without an explicit compliance offset.

215 **5 Experiments**

216 **Evaluation goal.** We ask three questions: (i) does *whole-body* compliance from multi-site (pelvis, torso, elbow, wrist) force sampling improve over aggressive tracking and end-effector-only compliance, especially for forces away from the hands? (ii) does the residual on virtual targets improve force-aware tracking without sacrificing yielding? and (iii) how does COMPLIANTWBC behave in the real world on contact-rich tasks? We address (i)–(ii) in §5.2 and (iii) in §5.3.

221 **5.1 Settings and Metrics**

222 We compare against two external baselines—a stiff whole-body tracker **TWIST2** [15] (nominal command \mathbf{x}^{ref} , no force-aware targets) and an embodiment-specific re-implementation of upper-body-only **GentleHumanoid** [11]—and two ablations: **Ours w/o pelvis force sampling** (upper-body perturbations only, approximating GentleHumanoid-style compliance) and **Ours w/o residual** (Stage-1 base policy with the Force Latent Encoder, no Stage-2 residual). **Ours (full)** adds the bounded residual on virtual targets (Eq. (4)) with per-link stiffness fixed by the curriculum. All variants share an identical training recipe, environment budget, and paired random seeds.

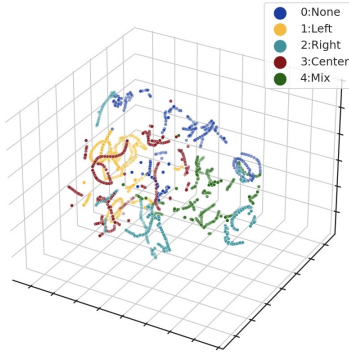


Figure 4: Force latent embedding projections via t-SNE of different force profiles across contact sites.

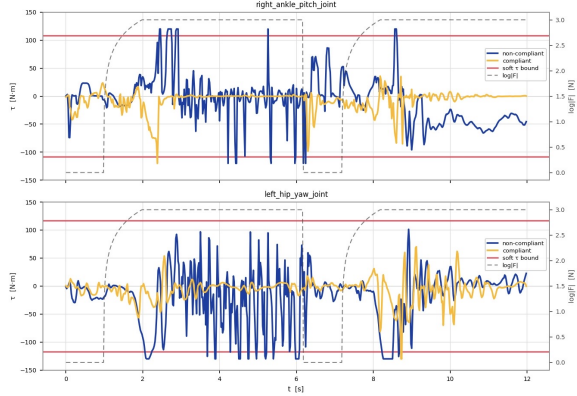


Figure 5: Compliant and Non-compliant controller reaction under pelvis perturbations.

Table 2: Performance of COMPLIANTWBC versus compliant and non-compliant baselines over 100 paired rollouts. Arrows indicate the preferred direction; bracketed columns are reported in units of 10^{-2} . $E_{\text{cmd}}^{\text{free}}/E_{\text{cmd}}^{\text{force}}$: wrist command-tracking RMSE without/with external force; ρ_{τ} : fraction of near-saturated joints (lower is more compliant); R_{LB} : fraction of torque deviation borne by the lower body; S : success (no fall) rate.

Controller	Tracking quality		WB compliance and robustness		
	$E_{\text{cmd}}^{\text{free}} \downarrow$ [$\times 10^{-2}$]	$E_{\text{cmd}}^{\text{force}} \downarrow$ [$\times 10^{-2}$]	$\rho_{\tau} \downarrow$ [$\times 10^{-2}$]	$R_{\text{LB}} \uparrow$	$S \uparrow$
TWIST2	2.18 \pm 0.31	3.14 \pm 0.42	2.74 \pm 0.38	0.16 \pm 0.04	0.59
GENTLEHUMANOID	5.22 \pm 0.47	7.86 \pm 0.61	2.12 \pm 0.29	0.14 \pm 0.03	0.91
OURS w/o pelvis force sampling	4.12 \pm 0.38	5.89 \pm 0.51	1.87 \pm 0.24	0.16 \pm 0.03	0.93
OURS w/o residual policy	4.29 \pm 0.35	5.43 \pm 0.47	1.56 \pm 0.21	0.20 \pm 0.03	0.94
OURS (full)	4.65 \pm 0.41	6.57 \pm 0.53	0.93 \pm 0.15	0.31 \pm 0.04	0.98

229 **Metrics.** Apart from comparing **Task success** S and **Tracking performance** $E_{\text{cmd}}^{\text{free}}$, we also report
 230 metrics that characterize the compliant response to external forces, including **Torque reaction** ρ_{τ} ,
 231 **Lower-body participation** R_{LB} , and **Force-aware tracking** $E_{\text{cmd}}^{\text{force}}$ upon external forces. **Torque**
 232 **reaction** ρ_{τ} reflects the fraction of near-saturated joints. A stiff controller fights the force (high ρ_{τ});
 233 a compliant one yields, keeping it low. **Lower-body participation** R_{LB} is the fraction of the total
 234 perturbation-induced torque deviation (relative to a matched no-perturbation rollout) borne by the
 235 lower-body joints; we treat it as a behavioral *diagnostic* of whether the policy engages the legs, not
 236 as a standalone success criterion. The metric formulations are given in the appendix.

237 5.2 Simulation Experiments

238 **Protocol.** Each controller runs 100 paired rollouts (a force trapezoid matching training) over a fixed
 239 site pool (wrist, elbow, torso, pelvis, hip, knee), split evenly between static (standing) and dynamic
 240 (stepping, locomotion) motions, plus a perturbation-free pass per (*controller, motion*) pair for the
 241 free-space and lower-body-participation metrics.

242 **Results.** Table 2 lays out the tracking–compliance trade-off. TWIST2 tracks best in free space but
 243 succeeds on only 59% of trials, fighting every perturbation; GENTLEHUMANOID reaches $S = 0.91$
 244 by relaxing the upper-body limbs but shows the lowest whole-body engagement ($R_{\text{LB}} = 0.14$).
 245 OURS (FULL) sits on the favorable frontier, with the highest lower-body participation and the low-
 246 est joint saturation ρ_{τ} at the highest success rate, at only a modest free-tracking gap to TWIST2.
 247 The ablations isolate each component: removing pelvis sampling collapses R_{LB} toward the GEN-
 248 TLEHUMANOID regime, while removing the residual recovers slight tracking precision at the cost
 249 of compliance and robustness, demonstrated by worst joint saturation, lower-body participation,
 250 success rate. We provide further qualitative insight into the learned behavior of the Force Latent En-
 251 coder and the compliant response of the policy. Figure 4 shows the Force Latent Encoder organizes
 252 its embeddings into distinct per-site clusters (left/right limb, pelvis, mixed, and no-perturbation)—
 253 a force-origin-aware representation rather than a collapsed force/no-force mode, a prerequisite for

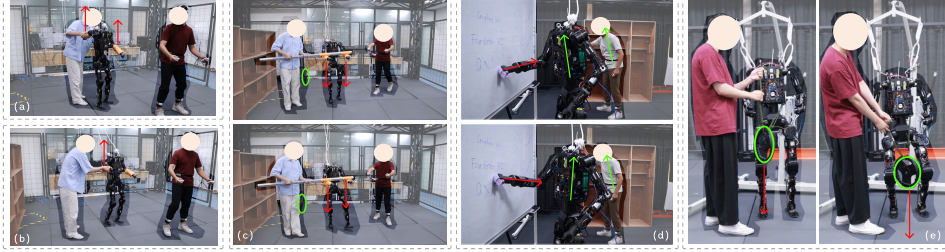


Figure 6: We deploy COMPLIANTWBC in the real world on five force-aware tasks: (a) Static Force Reaction; (b) Dynamic Force Reaction; (c) Cooperative Payload Transport; (d) Board Wiping; and (e) Squat Under Payload.

254 site-dependent reactions. Figure 5 contrasts compliant and non-compliant joint torques under a
 255 pelvis perturbation: the compliant policy keeps lower-body torques much smaller and stays sta-
 256 ble through the contact, consistent with its lower ρ_τ in Table 2 and concretizing the resist-and-fail
 257 behavior behind TWIST2’s low success rate.

258 5.3 Real-World Experiments

259 We deploy our trained policy on a real humanoid via teleoperation [47, 48, 49] on five tasks: *Static*
 260 *Force Reaction*, *Dynamic Force Reaction*, *Cooperative Payload Transport*, *Board Wiping*, and *Squat*
 261 *Under Payload* (Figure 6). The first two mirror the simulation perturbation protocol; the remaining
 262 three apply sustained, structured contacts that no single-site end-effector formulation can absorb,
 263 directly stressing the multi-site, whole-body assumptions of our pipeline (§4). We report these as
 264 *qualitative* demonstrations; quantitative on-robot evaluation is left to future work.

265 Under large pushes while tracking a dynamic walking motion (Figure 6b,c) the robot stays
 266 balanced—the regime where the stiff TWIST2 collapses in simulation ($S = 0.59$). In *Cooperative*
 267 *Payload Transport* the operator hands over a 100 N bar; the robot yields at the wrists, redistributes
 268 the load through torso and hips, and re-establishes a stable gait—the yield-then-reconcile behav-
 269 ior of the impedance equilibrium (Eq. (1)). *Squat Under Payload* stresses the same mechanism
 270 vertically, spreading sustained pelvis load across the chain rather than saturating the hip and knee.
 271 *Board Wiping* is the clearest whole-body case: as the operator presses a board against the robot’s
 272 hand, the humanoid leans its trunk back and shifts its CoM over the support polygon—which an
 273 upper-body-only baseline cannot do, since no perturbation reaches its lower body in training. These
 274 demonstrations are qualitatively consistent with the simulation ablation.

275 6 Conclusion

276 We presented COMPLIANTWBC, a two-stage RL pipeline that defines a multi-site whole-body
 277 impedance reference controller and uses it as a compliance-fidelity reward signal. Stage 1 jointly
 278 trains the Force Latent Encoder and a force-aware base policy against this reward while stage 2
 279 trains a bounded residual on the impedance equilibrium that compensates for the discrepancy be-
 280 tween the frozen wrench estimator and the true external wrench, with per-link stiffness held fixed so
 281 that passivity follows from the classical Cartesian-impedance argument. We validate the result on
 282 heavy-payload loco-manipulation tasks that require lower-body and multi-contact compliance and
 283 observe compliance behavior with high lower-body participation both in simulation and in the real
 284 world.

285 **Limitations & Future Work.** Our passivity result is *local* (fixed stiffness, near the equilibrium, as-
 286 suming the base realizes the analytical torque; §F), and the balancing QP (§A) can become infeasible
 287 under large simultaneous perturbations. The Force Latent Encoder is trained on simulated ground-
 288 truth wrenches, so hardware relies on proprioceptive generalization; a force–torque or current-based
 289 estimator [43, 41] would help. Finally, experiments use a single heavy platform and the real-world
 290 study is qualitative, so isolating the heaviness advantage and adding quantitative on-robot metrics
 291 remain key next steps, alongside bimanual sustained-contact manipulation and end-to-end stiffness
 292 learning under a differentiable energy-tank gate.

References

- 293
- 294 [1] Neville Hogan. Impedance control: An approach to manipulation: Parts I–III. *Journal of*
295 *Dynamic Systems, Measurement, and Control*, 107:1–24, 1985.
- 296 [2] Christian Ott, Alin Albu-Schäffer, Andreas Kugi, and Gerd Hirzinger. On the passivity-based
297 impedance control of flexible joint robots. *IEEE Transactions on Robotics*, 24(2):416–429,
298 2008.
- 299 [3] Oussama Khatib. A unified approach for motion and force control of robot manipulators: The
300 operational space formulation. *IEEE Journal of Robotics and Automation*, 3(1):43–53, 1987.
- 301 [4] Luis Sentis and Oussama Khatib. Synthesis of whole-body behaviors through hierarchical
302 control of behavioral primitives. *International Journal of Humanoid Robotics*, 2(4):505–518,
303 2005.
- 304 [5] Bernd Henze, Máximo A. Roa, and Christian Ott. Passivity-based whole-body balancing for
305 torque-controlled humanoid robots in multi-contact scenarios. *The International Journal of*
306 *Robotics Research*, 35(12):1522–1543, 2016.
- 307 [6] Alexander Dietrich. *Whole-Body Impedance Control of Wheeled Humanoid Robots*, vol-
308 ume 116 of *Springer Tracts in Advanced Robotics*. Springer, 2016. doi:10.1007/
309 978-3-319-40557-5.
- 310 [7] Alin Albu-Schäffer, Christian Ott, and Gerd Hirzinger. A unified passivity-based control
311 framework for position, torque and impedance control of flexible joint robots. *The Interna-*
312 *tional Journal of Robotics Research*, 26(1):23–39, 2007. doi:10.1177/0278364907073776.
- 313 [8] Yuanhang Zhang, Yifu Yuan, Prajwal Gurunath, Ishita Gupta, Shayegan Omidshafiei, Ali-
314 akbar Agha-mohammadi, Marcell Vazquez-Chanlatte, Liam Pedersen, Tairan He, and Guanya
315 Shi. FALCON: Learning force-adaptive humanoid loco-manipulation. *arXiv preprint*
316 *arXiv:2505.06776*, 2025.
- 317 [9] Sirui Chen, Zi-Ang Cao, Zhengyi Luo, Fernando Castañeda, Chenran Li, Tingwu Wang,
318 Ye Yuan, Linxi Fan, C. Karen Liu, and Yuke Zhu. CHIP: Adaptive compliance for humanoid
319 control through hindsight perturbation. *arXiv preprint arXiv:2512.14689*, 2025.
- 320 [10] Botian Xu, Haoyang Weng, Qingzhou Lu, Yang Gao, and Huazhe Xu. FACET: Force-adaptive
321 control via impedance reference tracking for legged robots. *arXiv preprint arXiv:2505.06883*,
322 2025.
- 323 [11] Qingzhou Lu, Yao Feng, Baiyu Shi, Michael Pisen, Zhenan Bao, and C. Karen Liu. Gen-
324 tleHumanoid: Learning upper-body compliance for contact-rich human and object interaction.
325 *arXiv preprint arXiv:2511.04679*, 2025.
- 326 [12] Gabriel B. Margolis, Michelle Wang, Nolan Fey, and Pulkit Agrawal. SoftMimic: Learning
327 compliant whole-body control from examples. *arXiv preprint arXiv:2510.17792*, 2025.
- 328 [13] Tairan He, Jiawei Gao, Wenli Xiao, Yuanhang Zhang, Zi Wang, Jiashun Wang, Zhengyi Luo,
329 Guanqi He, Nikhil Sobanbab, Chaoyi Pan, et al. Asap: Aligning simulation and real-world
330 physics for learning agile humanoid whole-body skills. *arXiv preprint arXiv:2502.01143*,
331 2025.
- 332 [14] Siheng Zhao, Yanjie Ze, Yue Wang, C. Karen Liu, Pieter Abbeel, Guanya Shi, and Rocky
333 Duan. ResMimic: From general motion tracking to humanoid whole-body loco-manipulation
334 via residual learning. *arXiv preprint arXiv:2510.05070*, 2025.
- 335 [15] Yanjie Ze, Siheng Zhao, Weizhuo Wang, Angjoo Kanazawa, Rocky Duan, Pieter Abbeel,
336 Guanya Shi, Jiajun Wu, and C. Karen Liu. Twist2: Scalable, portable, and holistic humanoid
337 data collection system. *arXiv preprint arXiv:2511.02832*, 2025.

- 338 [16] Yixuan Li, Yutang Lin, Jieming Cui, Tengyu Liu, Wei Liang, Yixin Zhu, and Siyuan Huang.
339 Clone: Closed-loop whole-body humanoid teleoperation for long-horizon tasks, 2025.
- 340 [17] Chenhao Lu, Xuxin Cheng, Jialong Li, Shiqi Yang, Mazeyu Ji, Chengjing Yuan, Ge Yang, Sha
341 Yi, and Xiaolong Wang. Mobile-television: Predictive motion priors for humanoid whole-body
342 control, 2025. URL <https://arxiv.org/abs/2412.07773>.
- 343 [18] Zipeng Fu, Qingqing Zhao, Qi Wu, Gordon Wetzstein, and Chelsea Finn. Humanplus: Hu-
344 manoid shadowing and imitation from humans, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2406.10454)
345 [2406.10454](https://arxiv.org/abs/2406.10454).
- 346 [19] Yanjie Ze, Zixuan Chen, João Pedro Araújo, Zi ang Cao, Xue Bin Peng, Jiajun Wu, and
347 C. Karen Liu. Twist: Teleoperated whole-body imitation system, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2505.02833)
348 [2505.02833](https://arxiv.org/abs/2505.02833).
- 349 [20] Zewen He, Chenyuan Chen, Dilshod Azizov, and Yoshihiko Nakamura. Cotap: Compliant
350 task pipeline and reinforcement learning of its controller with compliance modulation, 2025.
351 URL <https://arxiv.org/abs/2509.25443>.
- 352 [21] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola,
353 Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *Advances*
354 *in Neural Information Processing Systems (NeurIPS)*, 2020.
- 355 [22] Arnaud Demont, Mehdi Benallegue, Abdelaziz Benallegue, Pierre Gergondet, Antonin Dal-
356 lard, Rafael Cisneros, Masaki Murooka, and Fumio Kanehiro. The kinetics observer: A tightly
357 coupled estimator for legged robots, 2024. URL <https://arxiv.org/abs/2406.13267>.
- 358 [23] Tobias Johannink, Shikhar Bahl, Ashvin Nair, Jianlan Luo, Avinash Kumar, Matthias Loskyll,
359 Juan Aparicio Ojea, Eugen Solowjow, and Sergey Levine. Residual reinforcement learning for
360 robot control. In *IEEE International Conference on Robotics and Automation (ICRA)*, pages
361 6023–6029, 2019. doi:10.1109/ICRA.2019.8794127.
- 362 [24] Mark J. Balas and Susan A. Frost. Adaptive control of linear modal systems using residual
363 mode filters and a simple disturbance estimator. In *Proceedings of the 2011 American Control*
364 *Conference*, pages 2338–2343, 2011. doi:10.1109/ACC.2011.5991208.
- 365 [25] Yu Zhang, Long Cheng, Xiuzhe Xia, and Haoyu Zhang. Learning variable impedance skills
366 from demonstrations with passivity guarantee, 2024. URL [https://arxiv.org/abs/2306.](https://arxiv.org/abs/2306.11308)
367 [11308](https://arxiv.org/abs/2306.11308).
- 368 [26] Unitree Robotics. Unitree G1 Humanoid Robot. <https://www.unitree.com/cn/g1/>,
369 2024. Accessed: 2026-05-29.
- 370 [27] Federica Ferraguti, Cristian Secchi, and Cesare Fantuzzi. A tank-based approach to impedance
371 control with variable stiffness. In *IEEE International Conference on Robotics and Automation*
372 *(ICRA)*, pages 4948–4953, 2013. doi:10.1109/ICRA.2013.6631284.
- 373 [28] Klas Kronander and Aude Billard. Stability considerations for variable impedance control.
374 *IEEE Transactions on Robotics*, 32(5):1298–1305, 2016. doi:10.1109/TRO.2016.2593492.
- 375 [29] Tairan He, Wenli Xiao, Toru Lin, Zhengyi Luo, Zhenjia Xu, Zhenyu Jiang, Jan Kautz, Changliu
376 Liu, Guanya Shi, Xiaolong Wang, Linxi Fan, and Yuke Zhu. Hover: Versatile neural whole-
377 body controller for humanoid robots, 2025. URL <https://arxiv.org/abs/2410.21229>.
- 378 [30] Tairan He, Zhengyi Luo, Xialin He, Wenli Xiao, Chong Zhang, Weinan Zhang, Kris Ki-
379 tani, Changliu Liu, and Guanya Shi. Omnih2o: Universal and dexterous human-to-humanoid
380 whole-body teleoperation and learning. *arXiv preprint arXiv:2406.08858*, 2024.

- 381 [31] Zhengyi Luo, Jinkun Cao, Kris Kitani, Weipeng Xu, et al. Perpetual humanoid control for real-
382 time simulated avatars. In *IEEE/CVF International Conference on Computer Vision (ICCV)*,
383 pages 10895–10904, 2023.
- 384 [32] Kaiyan Xiao, Zihan Xu, Zhe Cheng, Chengju Liu, and Qijun Chen. Kinematics-aware multi-
385 policy reinforcement learning for force-capable humanoid loco-manipulation, 2025. URL
386 <https://arxiv.org/abs/2511.21169>.
- 387 [33] Hao Zhang, Yves Tseng, Ding Zhao, and H. Eric Tseng. Interaction-aware whole-body control
388 for compliant object transport, 2026. URL <https://arxiv.org/abs/2603.03751>.
- 389 [34] Tom Silver, Kelsey Allen, Josh Tenenbaum, and Leslie Kaelbling. Residual policy learning.
390 *arXiv preprint arXiv:1812.06298*, 2018.
- 391 [35] Xiang Zhang, Changhao Wang, Lingfeng Sun, Zheng Wu, Xinghao Zhu, and Masayoshi
392 Tomizuka. Efficient sim-to-real transfer of contact-rich manipulation skills with online ad-
393 mittance residual learning. In *7th Annual Conference on Robot Learning*, 2023. URL
394 <https://openreview.net/forum?id=gFXVysXh48K>.
- 395 [36] Jonas Buchli, Freek Stulp, Evangelos Theodorou, and Stefan Schaal. Learning variable
396 impedance control. *I. J. Robotic Res.*, 30:820–833, 06 2011. doi:10.1177/0278364911402527.
- 397 [37] Ashish Kumar, Zipeng Fu, Deepak Pathak, and Jitendra Malik. RMA: Rapid motor adaptation
398 for legged robots. In *Robotics: Science and Systems (RSS)*, 2021.
- 399 [38] Ashish Kumar, Zhongyu Li, Jun Zeng, Deepak Pathak, Koushil Sreenath, and Jitendra Malik.
400 Adapting rapid motor adaptation for bipedal robots. In *2022 IEEE/RSJ International Confer-
401 ence on Intelligent Robots and Systems (IROS)*, pages 1161–1168. IEEE, 2022.
- 402 [39] Joonho Lee, Jemin Hwangbo, Lorenz Wellhausen, Vladlen Koltun, and Marco Hutter. Learn-
403 ing quadrupedal locomotion over challenging terrain. *Science Robotics*, 5(47):eabc5986, 2020.
- 404 [40] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive
405 predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- 406 [41] Daegy Lim, Myeong-Ju Kim, Junhyeok Cha, Donghyeon Kim, and Jaeheung Park. Pro-
407 prioceptive external torque learning for floating base robot and its applications to humanoid
408 locomotion, 2023. URL <https://arxiv.org/abs/2309.04138>.
- 409 [42] Peiyuan Zhi, Peiyang Li, Jianqin Yin, Baoxiong Jia, and Siyuan Huang. Learning unified force
410 and position control for legged loco-manipulation, 2025. URL [https://arxiv.org/abs/
411 2505.20829](https://arxiv.org/abs/2505.20829).
- 412 [43] Haochen Shi, Songbo Hu, Yifan Hou, Weizhuo Wang, Karen Liu, and Shuran Song. Minimalist
413 compliance control. *arXiv preprint arXiv:2603.00913*, 2026.
- 414 [44] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal
415 policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 416 [45] Alexander A Alemi, Ian Fischer, Joshua V Dillon, and Kevin Murphy. Deep variational infor-
417 mation bottleneck. In *International Conference on Learning Representations (ICLR)*, 2017.
- 418 [46] Bui Tuong Phong. Illumination for computer generated pictures. *Commun. ACM*, 18(6):
419 311–317, June 1975. ISSN 0001-0782. doi:10.1145/360825.360839. URL [https://doi.
420 org/10.1145/360825.360839](https://doi.org/10.1145/360825.360839).
- 421 [47] Zhigen Zhao, Liuchuan Yu, Ke Jing, and Ning Yang. Xrobotoolkit: A cross-platform frame-
422 work for robot teleoperation, 2025. URL <https://arxiv.org/abs/2508.00097>.

- 423 [48] Joao Pedro Araujo, Yanjie Ze, Pei Xu, Jiajun Wu, and C. Karen Liu. Retargeting matters:
424 General motion retargeting for humanoid motion tracking, 2025. URL [https://arxiv.org/
425 abs/2510.02252](https://arxiv.org/abs/2510.02252).
- 426 [49] PICO Immersive Pte. Ltd. PICO Motion Tracker. [https://www.picoxr.com/global/
427 products/pico-motion-tracker](https://www.picoxr.com/global/products/pico-motion-tracker), 2023. Accessed: 2026-05-29.
- 428 [50] Emanuel Todorov, Tom Erez, and Yuval Tassa. MuJoCo: A physics engine for model-based
429 control. In *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*,
430 pages 5026–5033, 2012. doi:10.1109/IROS.2012.6386109.

431 A Analytical whole-body reference controller

432 This appendix gives the full analytical compliance controller summarized in §3. It is used *only*
 433 as a reference that generates the per-link virtual targets $\mathbf{x}_{\ell,d}^c$ of Eq. (1) used in the policy reward;
 434 the reference torque $\boldsymbol{\tau}_c$ and the energy tank below are derivation/reference quantities and are *not*
 435 executed at deployment.

436 A humanoid has configuration $\mathbf{q} \in \text{SE}(3) \times \mathbb{R}^{n_j}$, generalized velocity $\boldsymbol{\nu}$, and floating-base dynamics
 437 $M\dot{\boldsymbol{\nu}} + C\boldsymbol{\nu} + \mathbf{g} = \mathbf{S}^\top \boldsymbol{\tau} + \sum_c \mathbf{J}_c^\top \boldsymbol{\lambda}_c + \sum_i \mathbf{J}_{\ell_i}^\top \mathbf{f}_{\text{ext},i}$ with support-contact reaction wrenches $\boldsymbol{\lambda}_c$
 438 and external wrenches $\mathbf{f}_{\text{ext},i}$ at links ℓ_i . Classical Cartesian impedance at a single end-effector [1]
 439 prescribes $M_d \ddot{\mathbf{x}} + D_d \dot{\mathbf{x}} + K_d(\mathbf{x} - \mathbf{x}_d) = \mathbf{f}_{\text{ext}}$; we lift this to a hierarchy of frames coupled by
 440 balance.

441 **Reference compliant torque.** For controlled links \mathcal{L} (hands, elbows, knees, pelvis, torso), support
 442 contacts \mathcal{C}_s , and centroidal momentum \mathbf{h} , the reference torque is

$$\boldsymbol{\tau}_c = \underbrace{\mathbf{J}_h^\top \mathbf{w}_h^{\text{imp}}}_{\text{centroidal balance}} + \underbrace{\sum_{\ell \in \mathcal{L}} N_{\text{bal}}^\top \mathbf{J}_\ell^\top [\mathbf{K}_\ell(\mathbf{x}_{\ell,d} - \mathbf{x}_\ell) + D_\ell(\dot{\mathbf{x}}_{\ell,d} - \dot{\mathbf{x}}_\ell)]}_{\text{per-link Cartesian impedance}} + \underbrace{N_{\text{task}}^\top [\mathbf{K}_q(\mathbf{q} - \mathbf{q}_{\text{nom}}) + D_q \dot{\mathbf{q}}]}_{\text{joint-space posture}} \quad (7)$$

443 where the centroidal-impedance wrench is $\mathbf{w}_h^{\text{imp}} = \mathbf{K}_h(\mathbf{h}_d - \mathbf{h}) + D_h(\dot{\mathbf{h}}_d - \dot{\mathbf{h}}) + m\mathbf{g}$, $N_{\text{bal}} =$
 444 $\mathbf{I} - \mathbf{J}_h^\top (\mathbf{J}_h M^{-1} \mathbf{J}_h^\top)^{-1} \mathbf{J}_h M^{-1}$ is the dynamically consistent null-space projector of the balance
 445 task [3], and N_{task} projects orthogonal to both balance and per-link tasks. Each per-link Cartesian
 446 wrench is first mapped to joint torque by \mathbf{J}_ℓ^\top and then projected by N_{bal}^\top , so every term acts in the
 447 n_j -dimensional generalized-force space. The desired centroidal wrench is realized by the support
 448 contacts via a balancing QP in the spirit of [5]:

$$\min_{\{\boldsymbol{\lambda}_c\}} \left\| \mathbf{w}_h^{\text{imp}} - \sum_c \mathbf{G}_c \boldsymbol{\lambda}_c - \sum_i \mathbf{G}_{\text{ext},i} \mathbf{f}_{\text{ext},i} \right\|^2 + \gamma_{\text{reg}} \sum_c \|\boldsymbol{\lambda}_c\|^2 \quad \text{s.t.} \quad \boldsymbol{\lambda}_c \in \mathcal{F}_{\mu_c}, \text{CoP}_c \in \mathcal{S}_c \quad (8)$$

449 with grasp maps \mathbf{G}_c , linearized friction cones \mathcal{F}_{μ_c} , support polygons \mathcal{S}_c , and regularizer weight γ_{reg} .
 450 Admitting knees into \mathcal{C}_s adds support sites and generalizes the two-footed setting to multi-contact,
 451 lower-body-involved postures.

452 **Energy tank (general time-varying-stiffness case).** Time-varying $\mathbf{K}_\ell(t)$ with null-space projec-
 453 tion injects power and breaks passivity [27, 28]. For that general case one may add a scalar tank
 454 energy $T(t) \in [T_{\text{min}}, T_{\text{max}}]$,

$$\dot{T} = \sum_\ell \dot{\mathbf{x}}_\ell^\top D_\ell \dot{\mathbf{x}}_\ell - P_{\text{inject}}(\dot{\mathbf{K}}_\ell, \mathbf{N}_{\text{bal}}, \mathbf{N}_{\text{task}}), \quad (9)$$

455 filled by dissipated power and drained by non-passive control, freezing gains as $T \rightarrow T_{\text{min}}$. In our
 456 deployed pipeline \mathbf{K}_ℓ is held fixed and the residual edits only the equilibrium, so this tank is inert
 457 at runtime and is retained only as scaffolding for the time-varying-stiffness extension (cf. §F).

458 B Architecture details

459 **Base policy.** 3-layer MLP, hidden width 512, trained with PPO [44]. Observation \mathbf{o} comprises
 460 joint configurations, velocities, projected gravity, the last action, and the motion command \mathbf{c} . The
 461 latent \mathbf{z}_F enters as an additional input.

462 **Force Latent Encoder.** The encoder takes a history window $\mathbf{h}_t =$
 463 $(\mathbf{q}_{t-H:t}, \boldsymbol{\nu}_{t-H:t}, \boldsymbol{\tau}_{t-H:t-1}, \mathbf{c}_{t-H:t})$ of $H = 16$ steps (160 ms at 100 Hz), flattens it, and
 464 passes it through a 3-layer MLP (widths $512 \rightarrow 512 \rightarrow 256$, LayerNorm on the last hidden
 465 layer) emitting $(\boldsymbol{\mu}_t, \log \boldsymbol{\sigma}_t^2)$ over a latent of dimension $d_F = 32$. A reparameterized sample
 466 $\mathbf{z}_t = \boldsymbol{\mu}_t + \boldsymbol{\sigma}_t \odot \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, drives the wrench head $\psi : \mathbf{z}_F \rightarrow \hat{\mathbf{f}}_{\text{ext}} \in \mathbb{R}^{3|\mathcal{L}|}$ on all $|\mathcal{L}|$ controlled
 467 links and the projection head $g : \mathbf{z}_F \rightarrow \mathbb{S}^{d_g-1}$ with $d_g = 64$.

Auxiliary losses.

$$\mathcal{L}_{\text{wrench}} = \|\psi(\mathbf{z}_t) - \mathbf{f}_{\text{ext}}\|_2^2, \quad (10)$$

$$\mathcal{L}_{\text{supcon}} = \mathcal{L}^{\text{sup}}(g(\mathbf{z}_t), y_t), \quad (11)$$

$$\mathcal{L}_{\text{kl}} = \text{D}_{\text{KL}}(\mathcal{N}(\boldsymbol{\mu}_t, \boldsymbol{\sigma}_t^2) \parallel \mathcal{N}(\mathbf{0}, \mathbf{I})), \quad (12)$$

$$\mathcal{L}_{\text{smooth}} = \|\boldsymbol{\mu}_t - \boldsymbol{\mu}_{t-1}\|_2^2. \quad (13)$$

Weights: $(\lambda_w, \lambda_s, \lambda_k, \lambda_m) = (1.0, 0.3, 10^{-3}, 0.1)$. \mathcal{L}^{sup} is the supervised contrastive loss of [21] with perturbation-class labels y_t indicating which limb is perturbed and in which direction.

Residual policy. 3-layer MLP, hidden width 256, trained with PPO against the same reward of Eq. (2). Output is bounded as in Eq. (4) with $\epsilon_x = 5$ cm; the per-link stiffness \mathbf{K}_ℓ is held fixed by the curriculum throughout Stage 2.

C Force sampling details

Phong-weighted force-origin sampling. The inverse-CDF expressions for the Phong lobe of Eq. (5) are

$$\cos \theta = u^{1/(n+1)}, \quad \phi \sim \mathcal{U}[0, 2\pi), \quad r = R_a v^{1/3}, \quad (14)$$

with $u, v \sim \mathcal{U}[0, 1]$, θ measured from $-\hat{e}_z$. The cube-root radial warp yields uniform density in volume; with probability ϵ the direction is drawn isotropically instead.

The pelvis stiffness (k, c) is set so the spring force at maximum displacement R_a reaches a fixed fraction of the robot’s nominal weight: the stiffness term yields a coherent load that grows as the policy drifts from the anchor and admits a well-defined rest position against the leg stiffness, while the damping term reproduces the velocity-dependent tug of an elastically attached inertia. Per chunk we draw hold time $T_{\text{hold}} \sim \mathcal{U}(0.2, 2.0)$ s, ramp-down $\mathcal{U}(0.25, 1.0)$ s, and upper-body stiffness $K_p \sim \mathcal{U}(5, 250)$ N/m, with upper-body type sampled from {none, single-arm, both arms}; knees and feet are never perturbed directly. The shared downward bias of Eq. (5) makes the dominant carried-load regime coherent without an explicit pelvis-to-upper-body direction-inheritance heuristic.

Force law. Upper-body anchors apply a unilateral spring that auto-ramps as the body approaches and releases [11], while the pelvis is driven by a mass–spring–damper coupling to the sampled anchor to model a heavy object held through a compliant grasp. The two laws share the same anchor-based form and are gated on link type:

$$\mathbf{f}_\ell(t) = \begin{cases} k(\mathbf{a}_{\text{pelvis}}(t) - \mathbf{x}_{\text{pelvis}}(t)) - c\dot{\mathbf{x}}_{\text{pelvis}}(t), & \ell = \text{pelvis}, \\ K_p(\mathbf{a}_\ell - \mathbf{x}_\ell(t)) \mathbf{1}[\mathbf{a}_\ell \text{ on active side}], & \ell \in \mathcal{L}_{\text{upper}}. \end{cases} \quad (15)$$

D Algorithm

Reward. The trained objective is the compliance-fidelity reward of Eq. (2), where $x_{\ell,d}^c$ is the analytical controller’s impedance-implied trajectory under the realized (ground-truth) perturbation, available in simulation. The term $w_c \|\cdot\|^2$ rewards following this compliant deformation rather than the original motion command—generalizing GentleHumanoid’s [11] reference-dynamics reward to the whole body. The encoder’s wrench-reconstruction MSE does *not* appear in r_t : it shapes ϕ only through \mathcal{L}_ϕ (Eq. (3)), behind the gradient barrier, so the policy reward never confuses “latent fits the wrench” with “robot behaves compliantly.”

Hyperparameters. $N_1 = 5 \times 10^8$ environment-steps for Stage 1; $N_2 = 5 \times 10^8$ environment-steps for Stage 2; 8192 parallel environments in MjWarp (MuJoCo-based [50]); PPO with clip 0.2, discount $\gamma = 0.99$, GAE $\lambda = 0.95$; control rate 100 Hz; reward weights $w_c = 0.5$, $w_e = 10^{-3}$; force-latent auxiliary weights $(\lambda_w, \lambda_s, \lambda_k, \lambda_m)$ as in Eq. (3).

Algorithm 1 COMPLIANTWBC training pipeline

Require: Motion dataset \mathcal{M} , force-sampling distribution \mathcal{C} (§4.2), simulator \mathcal{S}

Stage 1: Force-aware base policy (PPO) + Force Latent Encoder (aux).

- 1: Initialize π_{base} , encoder ϕ , wrench head ψ , projection head g randomly.
- 2: **for** $i = 1, \dots, N_1$ **do**
- 3: Sample $(m, \mathbf{f}_{\text{ext}}, y, \mathcal{C}_s) \sim \mathcal{M} \times \mathcal{C}$ (y : perturbation class).
- 4: **for** $t = 1, \dots, T$ **do**
- 5: $(\boldsymbol{\mu}_t, \log \sigma_t^2) \leftarrow \phi(\mathbf{q}_{t-H:t}, \boldsymbol{\nu}_{t-H:t}, \boldsymbol{\tau}_{t-H:t-1}, \mathbf{c}_{t-H:t}); \mathbf{z}_F \leftarrow \boldsymbol{\mu}_t.\text{detach}()$.
- 6: $\boldsymbol{\tau}_t \leftarrow \pi_{\text{base}}(\mathbf{o}_t, \mathbf{c}_t, \mathbf{z}_F)$; step \mathcal{S} .
- 7: Roll the analytical controller on the same state to obtain $\mathbf{x}_{\ell,d}^c$; compute reward (2).
- 8: **end for**
- 9: PPO update on π_{base} [44].
- 10: Aux update on $(\phi, \psi, g): \nabla \mathcal{L}_\phi$ from Eq. (3) using $(\mathbf{f}_{\text{ext}}, y)$.
- 11: **end for**

Stage 2: Residual on impedance targets (PPO).

- 12: Freeze $\pi_{\text{base}}, \phi, \psi$. Initialize the residual policy π_{res} randomly.
- 13: **for** $i = 1, \dots, N_2$ **do**
- 14: Sample $(m, \mathbf{f}_{\text{ext}}, \mathcal{C}_s) \sim \mathcal{M} \times \mathcal{C}$.
- 15: **for** $t = 1, \dots, T$ **do**
- 16: $\mathbf{z}_F \leftarrow \phi(\cdot); \hat{\mathbf{f}}_{\text{ext}} \leftarrow \psi(\mathbf{z}_F)$.
- 17: Form analytical target $\mathbf{x}_{\ell,d}$ via Eq. (1) using $\hat{\mathbf{f}}_{\text{ext}}$.
- 18: $\Delta \mathbf{x}_{\ell,d} \leftarrow \pi_{\text{res}}(\mathbf{o}_t, \mathbf{c}_t, \mathbf{z}_F, \hat{\mathbf{f}}_{\text{ext}})$.
- 19: Compose $\tilde{\mathbf{x}}_{\ell,d}$ via (4); \mathbf{K}_ℓ held fixed by curriculum.
- 20: $\boldsymbol{\tau}_t \leftarrow \pi_{\text{base}}(\mathbf{o}_t, \mathbf{c}_t, \mathbf{z}_F)$ executed against $(\tilde{\mathbf{x}}_{\ell,d}, \mathbf{K}_\ell)$; step \mathcal{S} ; compute reward (2).
- 21: **end for**
- 22: PPO update on π_{res} [44].
- 23: **end for**
- 24: **return** $\pi_{\text{base}}, \phi, \psi, \pi_{\text{res}}$.

502 E Metrics Formulation

503 **Task success.** A trial is successful if the robot remains valid throughout the force perturbation roll-
504 out without triggering the evaluation termination conditions. In our evaluation, these terminations
505 are restricted to fall/instability events: root linear velocity exceeding 5 m/s, root height dropping be-
506 low 0.5 m, or root height rising above 1.0 m. Let $f_i \in \{0, 1\}$ denote whether any such termination
507 is triggered during trial i . The success indicator is

$$s_i = \mathbf{1}[f_i = 0], \quad (16)$$

508 and the reported success rate is

$$S = \frac{1}{N} \sum_{i=1}^N s_i. \quad (17)$$

509 **Command tracking error.** For the commanded body frames \mathcal{K}_{cmd} , currently the left and right
510 wrist frames, we measure Cartesian tracking error relative to the reference motion. Let $\mathcal{T}_{\text{free}}$ denote
511 timesteps for which the external-force envelope is inactive, and let $\mathcal{T}_{\text{force}}$ denote timesteps for which
512 it is active, i.e., the normalized applied-force envelope $a_t > 0.05$ (the trapezoid amplitude of §4.2).
513 For a time set \mathcal{T} , the command tracking RMSE is

$$E_{\text{cmd}}(\mathcal{T}) = \sqrt{\frac{1}{|\mathcal{T}| |\mathcal{K}_{\text{cmd}}|} \sum_{t \in \mathcal{T}} \sum_{k \in \mathcal{K}_{\text{cmd}}} \|\mathbf{p}_{k,t} - \mathbf{p}_{k,t}^{\text{ref}}\|_2^2}. \quad (18)$$

514 We report the free-space error $E_{\text{cmd}}^{\text{free}} = E_{\text{cmd}}(\mathcal{T}_{\text{free}})$ from the no-force pass and the force-window er-
515 ror $E_{\text{cmd}}^{\text{force}} = E_{\text{cmd}}(\mathcal{T}_{\text{force}})$ over the active perturbation window, which quantifies tracking relaxation
516 under external contact.

Table 3: Full reward bundle for COMPLIANTWBC (both stages), grouped by purpose with per-term weights; the compliance-fidelity term of Eq. (2) corresponds to the *Compliance* group.

Term	Form	Weight
<i>Whole-body tracking</i>		
tracking_joint_dof	$\exp(-0.15 \sum_j (q_j - q_{j,d})^2)$	+2.0
tracking_joint_vel	$\exp(-0.01 \sum_j (\dot{q}_j - \dot{q}_{j,d})^2)$	+0.2
tracking_root_rotation	$\exp(-5 \theta_q^2)$	+1.0
tracking_root_linear_vel	$\exp(-\ \mathbf{v}_b - \mathbf{v}_{b,d}\ ^2)$	+1.0
tracking_root_angular_vel	$\exp(-\ \boldsymbol{\omega}_b - \boldsymbol{\omega}_{b,d}\ ^2)$	+1.0
tracking_keybody_pos	$\exp(-10 \sum_k \ \Delta_b^k - \Delta_{b,d}^k\ ^2)$	+2.0
tracking_keybody_pos_global	$\exp(-10 \sum_k \ \mathbf{x}_k - \mathbf{x}_{k,d}\ ^2)$	+0.1
<i>Compliance — admittance force & virtual target</i>		
force_reward	$\exp(-0.0625 \langle \ \mathbf{F}_\ell^{\text{app}} - \mathbf{F}_\ell^{\text{exp}}\ \rangle_\ell) \mathcal{K}[\text{not exc.}]$	+2.0
force_target_tracking	$\exp(-25 \langle \ \mathbf{x}_\ell - \bar{\mathbf{x}}_\ell\ \rangle_\ell^2)$	+2.0
force_target_vel_tracking	$\exp(-\langle \ \dot{\mathbf{x}}_\ell - \dot{\bar{\mathbf{x}}}_\ell\ \rangle_\ell^2)$	+1.0
force_ext_penalty	$\langle \mathcal{K}[\ \mathbf{F}_\ell^{\text{app}}\ > F_{\text{max},\ell} \wedge \ \mathbf{F}_\ell^{\text{app}}\ > \ \mathbf{F}_\ell^{\text{exp}}\ + \frac{\delta}{2}] \rangle_\ell$	-6.0
keypoint_tracking_imp	$\exp(-10 \sum_{k \in \mathcal{K}} \ \mathbf{x}_k - \bar{\mathbf{x}}_k\ ^2), \bar{\mathbf{x}}_k = \mathbf{x}_\ell^{\text{virt}}$ on \mathcal{L}	+2.0
<i>Lower-body yielding under load</i>		
lower_keypoint_tracking	$\exp(-10 \sum_{k \in \mathcal{K}_{10}} \ \mathbf{x}_k - \bar{\mathbf{x}}_{k,10}\ ^2)$	+2.0
tracking_compliant_root_translation_z	$\exp(-5 (z_{\text{root}} - z_{\text{virt}})^2)$	+1.0
squat_reward	$\exp(-25 (\Delta z - c_z \sum_\ell \ \mathbf{F}_\ell\ ^2), c_z = 3 \times 10^{-4} \text{ m/N})$	+1.0
<i>Safety — joint, torque, contact, collision</i>		
alive	$\mathcal{K}[\text{alive}]$	+0.5
dof_pos_limits	$\sum_j [\text{ReLU}(q_j - q_j^{\text{max}}) + \text{ReLU}(q_j^{\text{min}} - q_j)]$	-5.0
joint_limit	same form as dof_pos_limits (additional)	-10.0
dof_torque_limits	$\sum_j \text{ReLU}(\tau_j / \tau_{j,\text{max}} - 0.95)$	-1.0
self_collisions	$\mathcal{K}[\ \mathbf{F}^{\text{self}}\ > 10 \text{ N}]$ summed over contact pts	-10.0
feet_stumble	$\mathcal{K}[\exists f : \ \mathbf{F}_f^{xy}\ > 4 F_f^z]$	-1.25
feet_contact_forces	$\sum_f \text{ReLU}(F_f^z - 500 \text{ N})$	-5×10^{-4}
feet_slip	$\sum_f \sqrt{\ \dot{\mathbf{x}}_f^{xy}\ } \mathcal{K}[F_f^z > 5]$	-0.1
<i>Regularization — smoothness & locomotion</i>		
dof_vel	$\sum_j \dot{q}_j^2$	-10^{-4}
dof_acc	$\sum_j \ddot{q}_j^2$	-5×10^{-8}
ankle_dof_vel	$\sum_{j \in \mathcal{J}_{\text{ank}}} \dot{q}_j^2$	-2×10^{-4}
ankle_dof_acc	$\sum_{j \in \mathcal{J}_{\text{ank}}} \ddot{q}_j^2$	-10^{-7}
action_rate_l2	$\ \mathbf{a}_t - \mathbf{a}_{t-1}\ ^2$	-0.1
ang_vel_xy	$\omega_{b,x}^2 + \omega_{b,y}^2$	-0.01
feet_air_time	$\sum_f \min(0, t_f^{\text{air}} - 0.5\text{s}) \mathcal{K}[\text{first contact}] \mathcal{K}[\ \mathbf{v}_{b,d}^{xy}\ > 0.05]$	+5.0

517 **Torque reaction.** We report the fraction of near-saturated joints, ρ_τ :

$$\rho_\tau = \frac{1}{T n_j} \sum_{t=1}^T \sum_{j=1}^{n_j} \mathbf{1}[|\tau_{j,t}| > 0.9 \tau_j^{\text{max}}], \quad (19)$$

518 where τ_j^{max} is the torque limit of joint j . A stiff controller fights the force (high ρ_τ); a compliant
519 one yields, keeping it low.

520 **Lower-body participation.** To verify that the controller actually uses the lower body under multi-
521 site contacts, we compare the torque deviation from a matched no-perturbation rollout:

$$R_{\text{LB}} = \frac{\sum_t \sum_{j \in \mathcal{J}_{\text{lower}}} |\tau_{j,t} - \tau_{j,t}^0|}{\sum_t \sum_{j \in \mathcal{J}_{\text{all}}} |\tau_{j,t} - \tau_{j,t}^0| + \epsilon}. \quad (20)$$

522 Here τ^0 is the torque sequence produced by the same controller on the same command trajectory
523 without external perturbation. A high R_{LB} during hip, pelvis, torso, or knee perturbations indicates
524 that the policy redistributes effort through the legs instead of treating force response as an upper-
525 body-only problem.

526 F Local passivity by construction

527 **Proposition 1** (Bounded interaction force and local passivity). *Let the closed-loop system satisfy:*
528 *(i) the base policy realizes τ_c of Eq. (7) at the equilibrium $(\mathbf{q}^*, \mathbf{0})$; (ii) the residual output is bounded*

529 as in Eq. (4); (iii) the per-link stiffness \mathbf{K}_ℓ is constant in time. Then (a) the residual shifts each
 530 link equilibrium by at most ϵ_x , so the induced restoring force is bounded by $\|\mathbf{K}_\ell\| \epsilon_x$; and (b) if
 531 additionally (iv) the residual-edited equilibrium $\tilde{\mathbf{x}}_{\ell,d}$ varies slowly enough that its power injection is
 532 dominated by the link damping \mathbf{D}_ℓ , the storage function $S = \frac{1}{2} \boldsymbol{\nu}^\top \mathbf{M} \boldsymbol{\nu} + \frac{1}{2} \sum_\ell (\mathbf{x}_\ell - \tilde{\mathbf{x}}_{\ell,d})^\top \mathbf{K}_\ell (\mathbf{x}_\ell -$
 533 $\tilde{\mathbf{x}}_{\ell,d})$ satisfies $\dot{S} \leq \sum_i \mathbf{f}_{\text{ext},i}^\top \dot{\mathbf{x}}_{\text{ext},i}$, where $\mathbf{x}_{\text{ext},i}$ is the displacement of contact point i ; i.e. the closed
 534 loop is passive with respect to the external-wrench port.

535 *Sketch.* Bound (a) is immediate from Eq. (4): the residual moves the equilibrium by at most ϵ_x , so
 536 with fixed \mathbf{K}_ℓ the induced restoring force changes by at most $\|\mathbf{K}_\ell\| \epsilon_x$. For (b), differentiating S and
 537 substituting the closed-loop dynamics gives $\dot{S} = \sum_i \mathbf{f}_{\text{ext},i}^\top \dot{\mathbf{x}}_{\text{ext},i} - (\text{damping dissipation}) - \sum_\ell (\mathbf{x}_\ell -$
 538 $\tilde{\mathbf{x}}_{\ell,d})^\top \mathbf{K}_\ell \dot{\tilde{\mathbf{x}}}_{\ell,d}$. With \mathbf{K}_ℓ constant the $\dot{\tilde{\mathbf{x}}}_{\ell,d}$ term that would otherwise inject non-passive power van-
 539 ishes identically; the remaining moving-equilibrium term is dominated by the damping dissipation
 540 under assumption (iv), yielding the stated inequality. For the general time-varying-stiffness case the
 541 energy tank of §A would absorb the $\dot{\tilde{\mathbf{x}}}_{\ell,d}$ term; in our pipeline stiffness is fixed and the residual does
 542 not modulate gains, so the tank is reference scaffolding rather than a runtime mechanism.

543 **Scope and what this does and does not buy.** The result is *local* (around the equilibrium) and
 544 *by construction*: it follows from the residual bound and fixed stiffness, not from the learned com-
 545 ponents. Global stability under adversarial external forces is not claimed and remains open. Two
 546 consequences are worth stating explicitly. (i) Action-residual baselines (ASAP, ResMimic) edit the
 547 torque directly and do not admit this storage function; their effect on induced contact force is not
 548 physically bounded by construction. (ii) The proposition is a safety property for human–robot con-
 549 tact, not a performance property: it does not guarantee tracking, success rate, or sim-to-real transfer.
 550 Those are settled empirically (§5).