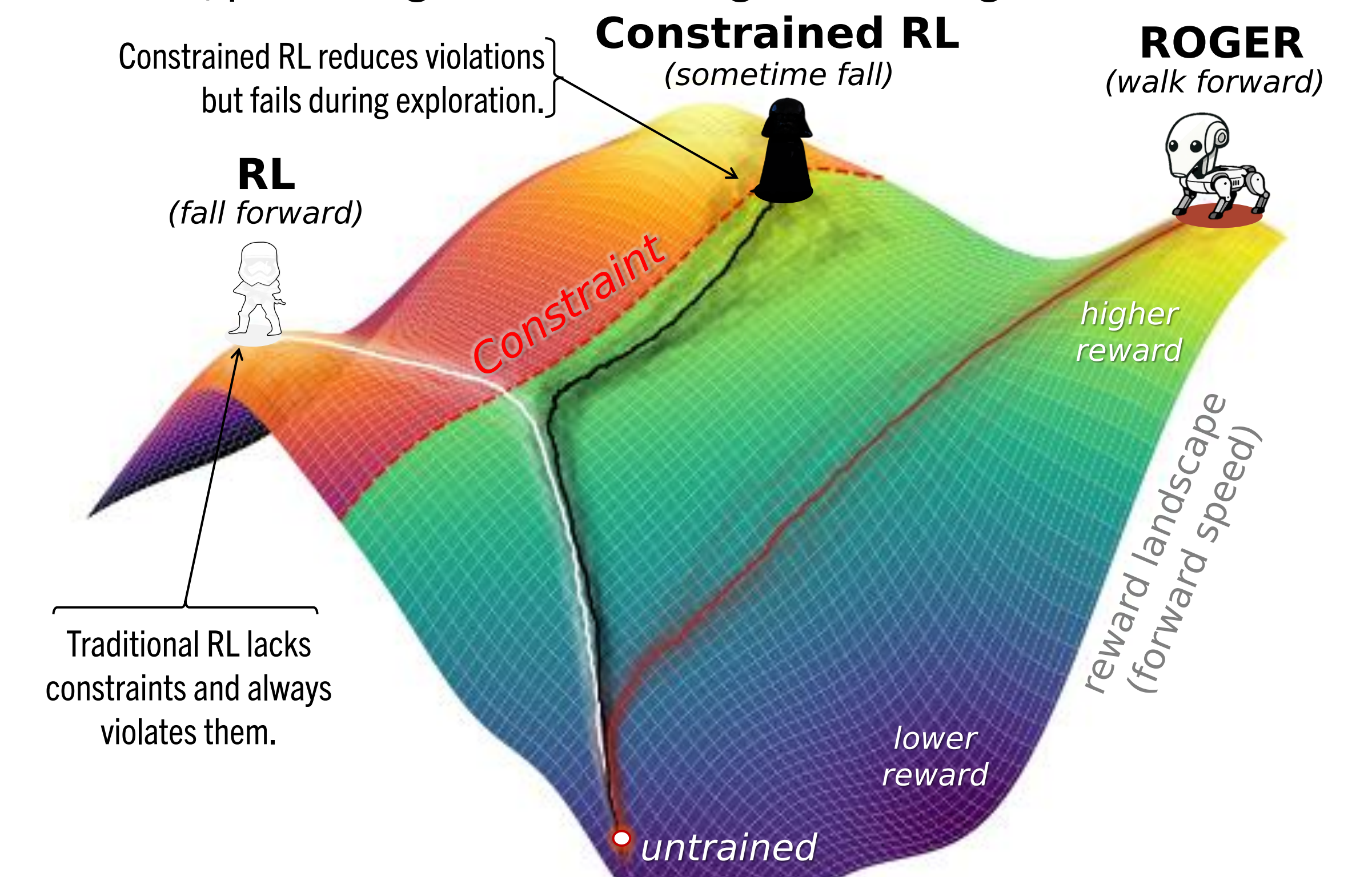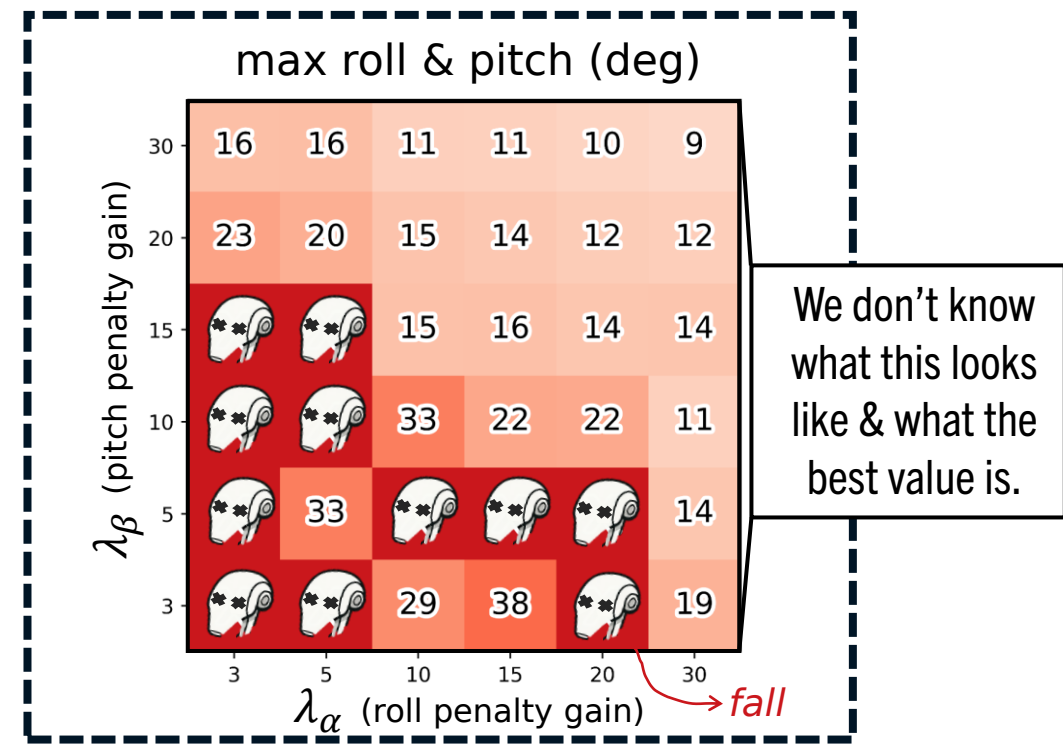# Gain Tuning Is Not What You Need:
## Reward Gain Adaptation for Constrained Locomotion Learning

*Arthicha Srisuchinnawong and Poramate Manoonpong (VISTEC, Thailand)*

BRAIN VISTEC
ROBOTICS SCIENCE AND SYSTEMS

"HELLO THERE"
"ROGER ROGER"
GENERAL REWADIOUS, AFTER HAVING AN ARMY OF ROGER.
"I'VE FOUND A TRICK; SIMPLE, BUT POWERFUL! YOUR REWARD WILL MAKE A FINE ADDITION TO MY EXP."

## Introduction

- In general, Reinforcement Learning (RL) is **sample inefficient** and **cannot satisfy constraints** during learning, making real-world RL impossible.

- **Finding a good reward function** is also **exhaustive**, **sensitive**, and **unpredictable**, with no clear way to determine which gain values guarantee constraint satisfaction.

max roll & pitch (deg)

We don't know what this looks like & what the best value is.

- This work introduces <u>ROGER that uses embodied interaction and intuitive constraint thresholds</u> to recompute new weighting gains online, preventing violations throughout learning.

**Constrained RL** *(sometime fall)*
**ROGER** *(walk forward)*
Constrained RL reduces violations but fails during exploration.
**RL** *(fall forward)*
higher reward
reward landscape (forward speed)
Constraint
lower reward
Traditional RL lacks constraints and always violates them.
untrained

"YOU WERE THE CHOSEN ONE! SAID TO BRING LEARNING TO THE REAL WORLD, NOT GET STUCK IN ENDLESS TUNING! TO ENFORCE CONSTRAINTS, NOT BREAK THEM!"
"I DON'T LIKE SIM. IT OVERFITS, AND DOESN'T TRANSFER."
RL-KIN, AFTER POLICY DEPLOYMENT.
MASTER CANNOT-BE, AFTER EVERY FAILED RUN.
NOTHING. JUST HIGH GROUND.

## Constrained RL

- **Fixed-weighting** constrained RL enforces constraints by manually tuning penalty weights or shaping reward functions.
- **Adaptive-weighting** constrained RL adjusts penalty weights dynamically using primal-dual or switching mechanisms.

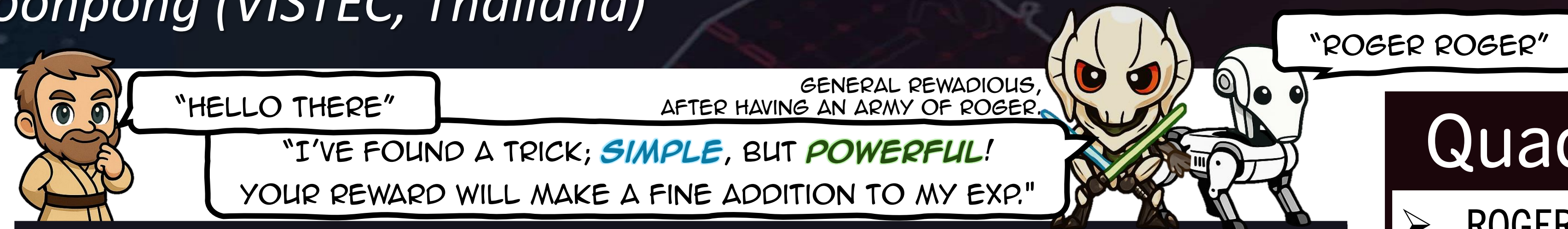| Techniques | Reward Gain ($\lambda_t$) | Penalty Gain ($\lambda_{it}$) |
|---|---|---|
| Fixed-weighting & Control Barrier Function (CBF) | 1.0 (fixed) | tuned and fixed |
| Primal-Dual Optimization (PDO) | 1.0 (fixed) | $[\lambda_{it} + \eta_\lambda(\tilde{R}_{it} - (\tau_i - \delta_i))]_+$ |
| Online Learning of Auxiliary Loss (OL-AUX) | 1.0 (fixed) | $\lambda_{it} + \eta_\lambda \nabla \mathbb{E}[R_{vt}]\nabla \mathbb{E}[\tilde{R}_{it}]$ |
| Constrained Rectified Policy Optimization (CRPO) | 0.0 if (exist $\tilde{R}_{it} > \tau_i - \delta_i$) else 1.0 | 1.0 if ($\tilde{R}_{it} > \tau_i - \delta_i$) else 0.0 |

**References**
[Constrained RL] A review of safe reinforcement learning, PAMI, 2024.
[CRPO] Crpo: Safe reinforcement learning with convergence guarantee, ICML, 2021.
[OL-AUX] Adaptive auxiliary task weighting for reinforcement learning. NeurIPS, 2019.
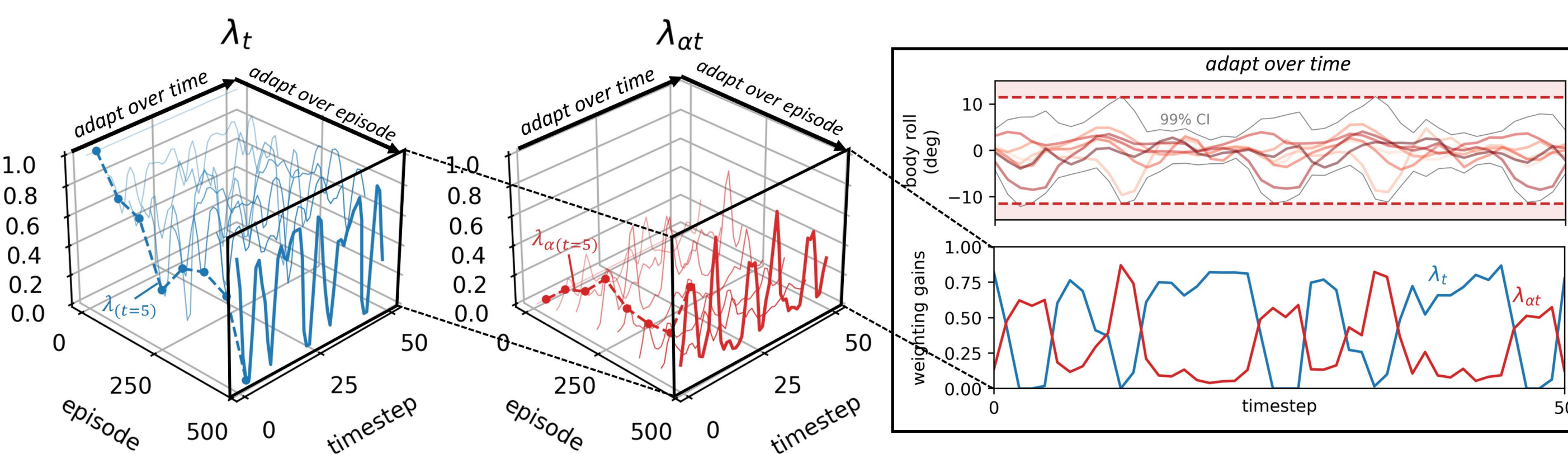[SME-AGOL] Interpretable neural control with adaptable online learning for sample efficient, TNNLS, 2025.
[GOLLUM] Growable neural control with online learning for continual locomotion learning, IJRR, 2025.
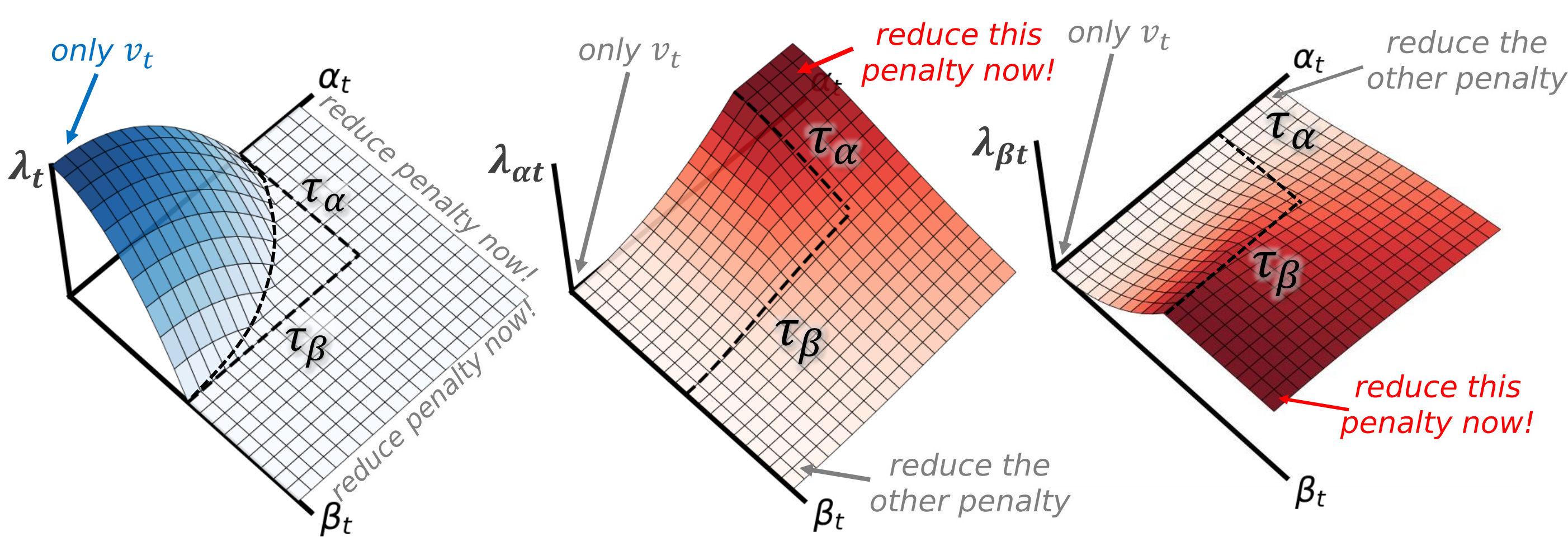
## ROGER *(Reward-Oriented Gains via Embodied Regulation)*

- **Key Idea**: ROGER prioritizes constraint satisfaction when penalties approach thresholds and automatically shifts focus back to task rewards once safe behavior is restored.

$\lambda_t$    $\lambda_{\alpha t}$    adapt over time

adapt over time    adapt over episode

- ROGER uses **fixed manifolds** below to map current penalty values to weighting gains on-the-fly. Unlike learning-based methods (e.g., PDO, OL-AUX), ROGER has **no update delay**, **no tuning**, and **minimal computation** time (~0.46 ms).

only $v_t$    $\alpha_t$
reduce this penalty now!
reduce the other penalty
$\lambda_t$    $\tau_\alpha$    $\lambda_{\alpha t}$    $\tau_\alpha$    $\lambda_{\beta t}$    $\tau_\alpha$
$\tau_\beta$    $\tau_\beta$    $\tau_\beta$    $\beta_t$
reduce this penalty now!

- Only 3 steps to construct these manifolds: (1) estimate constraint proximities using thresholds, (2) compute penalty ratios, and (3) map to weighting gains.

intuitive constraint thresholds (max roll, max pitch, max torque)

$$\delta_{\alpha t} = \frac{|\alpha_t|}{\tau_\alpha}$$
$$\delta_{\beta t} = \frac{|\beta_t|}{\tau_\beta}$$

$$\Delta_t = min\{\delta_{\alpha t}^2 + \delta_{\beta t}^2, 1.0\}$$
$$r_{\alpha t} = (\delta_{\alpha t}^2)/(\delta_{\alpha t}^2 + \delta_{\beta t}^2)$$
$$r_{\beta t} = (\delta_{\beta t}^2)/(\delta_{\alpha t}^2 + \delta_{\beta t}^2)$$

$$\begin{bmatrix} \lambda_t \\ \lambda_{\alpha t} \\ \lambda_{\beta t} \end{bmatrix} = \begin{bmatrix} 1 - \Delta_t \\ r_{\alpha t}\Delta_t \\ r_{\beta t}\Delta_t \end{bmatrix}$$

constraint proximities    penalty ratios    adaptive weighting gains

- ROGER is **Lyapunov-stable** and guarantees improvement of the primary objective/reward after learning.
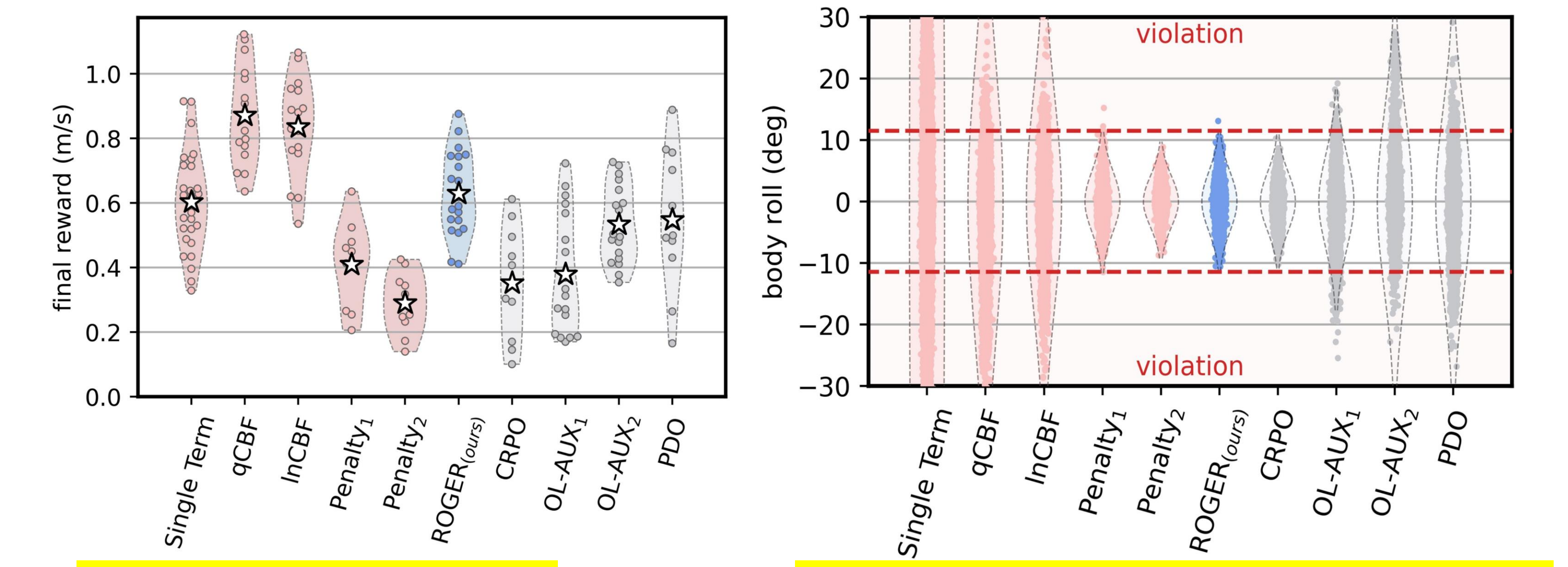
**Algorithm 1** Reinforcement Learning with ROGER
1: Perform exploration and collect trajectory $\tau$.
2: Compute estimated penalties ($|\tilde{\alpha}_t|$ and $|\tilde{\beta}_t|$) using: $\begin{bmatrix} |\tilde{\alpha}_t| \\ |\tilde{\beta}_t| \end{bmatrix} = \sum \gamma^i \begin{bmatrix} |\alpha_{t+i}| \\ |\beta_{t+i}| \end{bmatrix} / \sum \gamma^i.$
3: Compute weighting gains ($\lambda$s) using ROGER: $\begin{bmatrix} \lambda_t \\ \lambda_{\alpha t} \\ \lambda_{\beta t} \end{bmatrix} = \begin{bmatrix} 1 - \Delta(|\tilde{\alpha}_t|, |\tilde{\beta}_t|) \\ r_{\alpha t}\Delta(|\tilde{\alpha}_t|, |\tilde{\beta}_t|) \\ r_{\beta t}\Delta(|\tilde{\alpha}_t|, |\tilde{\beta}_t|) \end{bmatrix}.$
4: Update policy using combined reward or advantage:
$R_t = \lambda_t v_t - \lambda_{\alpha t}|\alpha_t| - \lambda_{\beta t}|\beta_t|$  or  $A_t = \lambda_t A_{vt} - \lambda_{\alpha t}A_{\alpha t} - \lambda_{\beta t}A_{\beta t}.$
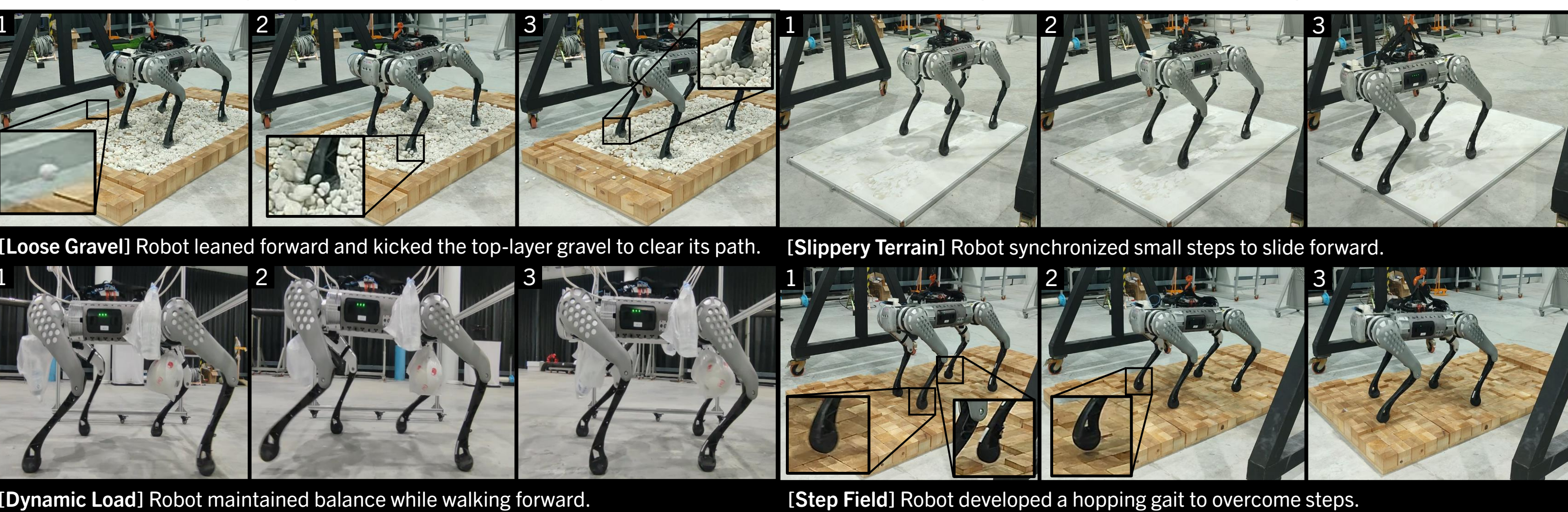
ROGER-
POORLY TUNED RL    WELL-TUNED CONSTRAINED RL

## Quadruped Robot

- ROGER achieves constraint satisfaction with **<5 violations over 50,000 training timesteps**, even when using small training samples (only ≈50 timesteps/update).
- ROGER yields **50% higher final reward** compared to other state-of-the-art techniques.

final reward (m/s)    body roll (deg)    violation    violation

Single Term, qCBF, lnCBF, Penalty₁, Penalty₂, ROGER(ours), CRPO, OL-AUX₁, OL-AUX₂, PDO

- **Physical locomotion learning** was achieved **within 30 mins from scratch without any falls** on challenging conditions (loose gravel, step field, slippery terrain, and with dynamic load).

[Loose Gravel] Robot leaned forward and kicked the top-layer gravel to clear its path. [Slippery Terrain] Robot synchronized small steps to slide forward.
[Dynamic Load] Robot maintained balance while walking forward. [Step Field] Robot developed a hopping gait to overcome steps.

## MuJoCo Benchmark

- ROGER can be applied to **other types of robots and control**, e.g., FCNN + PPO.
- Under **restricted constraints**, ROGER prioritizes reducing penalties first and then optimizes the main objective once the constraints are satisfied.

ROGER: less oscillation, upright body
default: more oscillation, oscillating body
CRPO: pitch upward, less oscillation, upright body

| Robot/Environment | | Final Values | | | #Violation (throughout) | | |
|---|---|---|---|---|---|---|---|
| | | Default | CRPO | ROGER | Default | CRPO | ROGER |
| Ant | Distance (m; ∝ reward) | 106.02 | 57.50 | **117.45** | n/a | n/a | n/a |
| | Torque (Nm), $\tau = 1.0$ | **0.25** | 0.93 | 0.39 | $10^{-10}$ | 0.82 | $10^{-6}$ |
| | Height (m), $\tau = 0.2$ | 0.00 | 0.00 | 0.00 | $10^{-5}$ | 0.0 | $10^{-5}$ |
| | Roll (°), $\tau = 45$ | 12.37 | 16.04 | **10.00** | $10^{-4}$ | 0.02 | $10^{-5}$ |
| | Pitch (°), $\tau = 45$ | 8.27 | 10.88 | **5.35** | $10^{-8}$ | $10^{-8}$ | $10^{-7}$ |
| Cheetah | Distance (m; ∝ reward) | 46.10 | 82.48 | **92.25** | n/a | n/a | n/a |
| | Torque (Nm), $\tau = 1.0$ | 0.76 | 0.98 | **0.55** | $10^{-12}$ | 0.88 | $10^{-5}$ |
| | Pitch (°), $\tau = 10$ | 6.13 | **3.66** | 5.53 | 0.08 | $10^{-5}$ | 4.7 |
| Hopper | Distance (m; ∝ reward) | 5.23 | 6.42 | **6.57** | n/a | n/a | n/a |
| | Torque (Nm), $\tau = 1.0$ | 0.78 | 0.80 | **0.33** | 0.01 | 0.01 | $10^{-7}$ |
| | Pitch (°), $\tau = 10$ | 5.10 | 4.38 | **2.06** | 0.14 | 0.01 | $10^{-11}$ |
| Walker | Distance (m; ∝ reward) | 5.99 | 1.91 | **11.65** | n/a | n/a | n/a |
| | Torque (Nm), $\tau = 1.0$ | 0.91 | 0.94 | **0.40** | $10^{-7}$ | 0.01 | $10^{-8}$ |
| | Pitch (°), $\tau = 45$ | 21.00 | 7.62 | 8.31 | 0.07 | 8.07 | $10^{-8}$ |

blue text: best performance (highest reward, lowest penalty); red text: %violation > 0.01

Ant    Cheetah    Hopper    Walker

**In conclusion, use ROGER when:**
- You're doing **high-stakes** real-world **fine-tuning** or **continual learning**.
- You're **tired of reward tuning**.

Video    GitHub    https://github.com/Arthicha/ROGER_ROGER_public

"TRY OR TRY NOT, THERE IS NO DUE"
DA'YO, AFTER YOU READ THIS POSTER.