

PASSIVE WALKER RL: FROM FSM TO JAX/BRAX

EXPERT CONTROLLERS & SCALABLE RL

YUNUS EMRE DANABAŞ

MECHATRONIC ENGINEERING
SABANCI UNIVERSITY

JULY 17, 2025

Sabancı
Universitesi

- Introduction & Motivation
- Methodology
- Results & Discussion
- Impact & Ethical Issues
- Project Management & Timeline
- Conclusion & Future Work
- Appendix

- Legged locomotion remains a challenging control problem:
 - ▶ High-dimensional dynamics, contact events, underactuation
 - ▶ Sample-inefficient exploration in standard RL
- Expert demonstrations can bootstrap learning and ensure safety
- JAX/Brax offers GPU-accelerated, vectorised simulation for fast RL
- Goal: combine rule-based experts with modern RL for efficient, robust passive walking

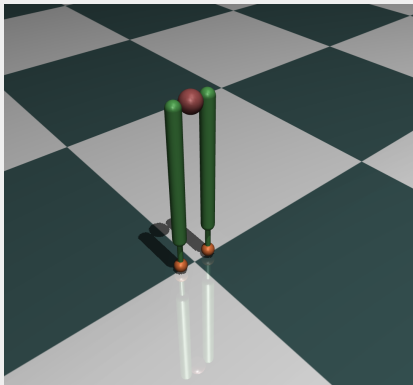
PROBLEM STATEMENT

- Design a control pipeline for a planar passive walker:
 - ▶ Leverage a finite-state expert for stable baseline gait
 - ▶ Learn a differentiable policy via behaviour cloning (BC)
 - ▶ Fine-tune with Proximal Policy Optimization (PPO)
- Achieve both sample efficiency and high final performance
- Scale up training using a vectorised JAX/Brax implementation
- Systematically explore hyperparameters to identify robust configurations

KEY CONTRIBUTIONS

1. **Expert-to-RL Pipeline:** FSM → Behaviour Cloning → BC-seeded PPO for “walk-from-day-one” learning
2. **Imitation-Regularised PPO:** Clipped surrogate loss augmented with decaying BC term for stability
3. **Vectorised JAX/Brax Pipeline:** High-throughput, GPU-accelerated simulation enabling large-scale experiments
4. **Large-Scale Hyperparameter Sweep:** 120-job grid over reward scaling, learning rate, network capacity
5. **Open-Source Release:**
github.com/yunusdanabas/passive_walker_rl

PHYSICS MODEL & SIMULATION SETUP



- **MuJoCo model:** 5 bodies, 7 DoFs (planar slide x, z , yaw hinge, hip hinge, two prismatic knees)
- **Virtual ramp:** gravity tilted 11.5° downhill
- **Simulation:** physics timestep = 1 ms; control at 200–1 000 Hz

JOINTS & ACTUATORS

Joint	Type
slide_x	prismatic
slide_z	prismatic
yaw	hinge
hip	hinge
left_knee	prismatic
right_knee	prismatic

Ranges: slide x, z, yaw, hip
unbounded; knees ± 0.30 m.

Actuator	k_p
hip_act	5
left_knee_act	1000
right_knee_act	1000

Control range: hip ± 0.5 rad; knees
 ± 0.3 m.

EXPERT FINITE-STATE CONTROLLER

FSM LOGIC (DETAILED CONDITIONS)

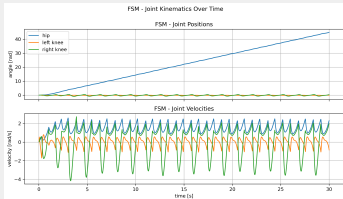
■ Hip Controller:

- ▶ *Phase*: one leg swings, the other supports.
- ▶ *Switch when*:
 - Swing foot contacts ground ($z_{\text{foot}} < 0.05 \text{ m}$),
 - And trunk returns upright ($\text{pitch} < 0$).
- ▶ *Action*: desired hip angle toggles between $\pm 0.3 \text{ rad}$.

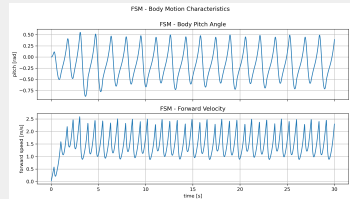
■ Knee Controller:

- ▶ *States*: **Stance** (extended) vs. **Retraction** (bent).
- ▶ *To Retract*:
 - Opposite foot has just landed ($z_{\text{other foot}} < 0.05 \text{ m}$),
 - And local thigh upright ($\text{pitch} < 0$).
- ▶ *To Extend*:
 - Swing leg has moved forward past neutral ($\text{hip pitch} > \text{threshold}$).
- ▶ *Positions*: retracted = $+0.30 \text{ m}$, stance = 0 m .

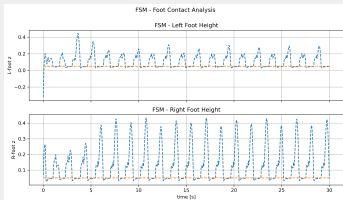
FSM KINEMATICS SUMMARY



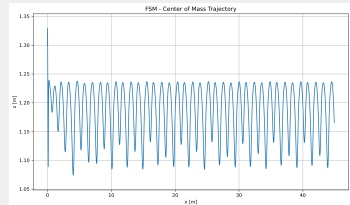
(a) Joint Angles & Velocities



(b) Torso Pitch & Speed



(c) Foot-Ground Contact



(d) Center of Mass Path

Figure: Expert FSM reference trajectories: (top) joint kinematics and (bottom) foot contact timing and CoM trajectory.

■ Variants explored:

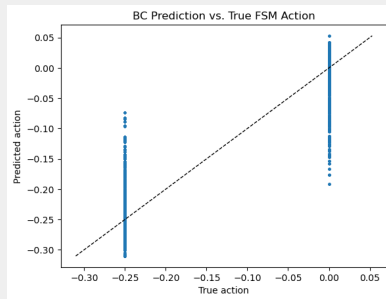
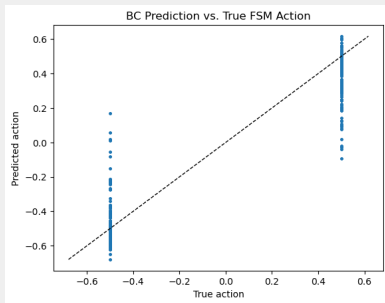
- ▶ Hip-only BC
- ▶ Knees-only BC
- ▶ Full BC (hip + both knees) with four loss functions: MSE, Huber, L1, Combined

■ **Data:** $\sim 10^5$ expert steps at $\sim 10^3$ Hz from FSM \rightarrow observations (11 dims) expert actions (3 dims)

■ **Model:** 2-layer ReLU MLP (256 hidden units)

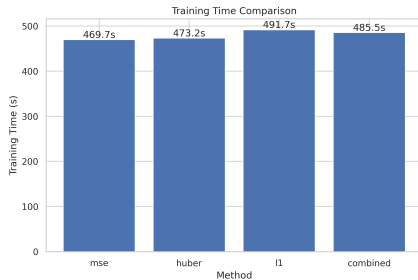
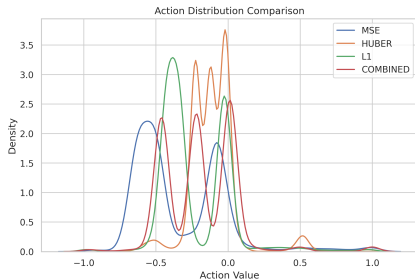
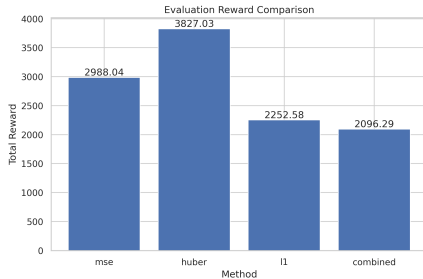
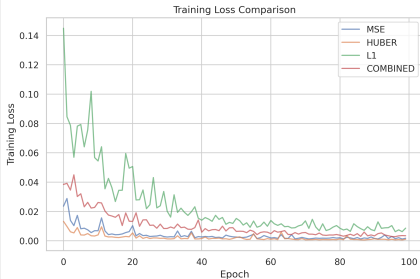
■ **Training:** $\mathcal{O}(10^5)$ samples, $\sim 10^2$ epochs, batch ~ 64 , $\alpha \sim 10^{-4}$

BC PREDICTION VS. TRUE FSM ACTIONS



- Scatter of BC predictions versus FSM labels shows tight clustering around $y = x$ for both joints.
- Indicates low bias and accurate reproduction of expert commands across the gait cycle.
- Used 100,000 samples and three loss variants (MSE, Huber, L1, combined).

BC LOSS VARIANTS: VISUAL COMPARISON



BC LOSS VARIANTS: KEY INSIGHTS

■ Training Loss Convergence:

- ▶ All four loss functions (MSE, Huber, L1, Combined) converge over 100 epochs.
- ▶ Huber and L1 exhibit slightly smoother decay and robustness to outliers.

■ Evaluation Reward:

- ▶ Huber-trained policy achieves the highest total reward when deployed.
- ▶ Combined loss performs comparably but with greater variance.

■ Action Distribution:

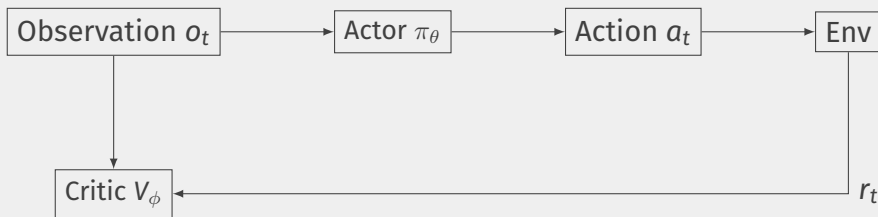
- ▶ MSE yields tightly concentrated actions around the expert mean.
- ▶ Combined loss produces a broader distribution, indicating exploratory behavior.

■ Training Efficiency:

- ▶ MSE and Huber are fastest per epoch.
- ▶ L1 and Combined incur minor extra cost due to additional absolute/huber computations.

PPO FINE-TUNING STRATEGIES

PPO FINE-TUNING OVERVIEW



- **Actor** (π_θ): small MLP, one hidden layer
- **Critic** (V_ϕ): two-layer MLP

- **Rollouts**: on-policy data + GAE
- **Updates**: PPO clip + $\beta(t)$ -regularisation

PPO UPDATE ALGORITHM

Require: θ, ϕ, β_0

```
1: for iteration = 1...N do
2:    $\mathcal{B} \leftarrow \text{COLLECTTRAJECTORIES}(\pi_\theta)$ 
3:    $\{\hat{A}, \hat{R}\} \leftarrow \text{GAE}(\mathcal{B}.\text{rewards}, \mathcal{B}.\text{dones}, V_\phi(\mathcal{B}.\text{obs}))$ 
4:   for epoch = 1...K do
5:     for all mini-batch  $b \subset \mathcal{B}$  do
6:        $L_{\text{clip}} \leftarrow \text{PPOCLIP}(b, \pi_\theta^{\text{old}}, \hat{A})$ 
7:        $L_{\text{bc}} \leftarrow \|\pi_\theta(b.\text{obs}) - b.\text{bc\_targets}\|^2$ 
8:        $L \leftarrow L_{\text{clip}} + \beta(t) L_{\text{bc}}$ 
9:        $\theta \leftarrow \theta - \alpha \nabla_\theta L$ 
10:    end for
11:  end for
12:   $\phi \leftarrow \phi - \alpha_v \nabla_\phi \|V_\phi(\mathcal{B}.\text{obs}) - \hat{R}\|^2$ 
13:   $\beta(t+1) \leftarrow \text{Decay}(\beta(t))$ 
14: end for
```

PPO + BC LOSS FUNCTIONS

Clipped Surrogate Objective

$$L_{\text{PPO}} = -\mathbb{E}_t \left[\min(r_t \hat{A}_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t) \right], \quad r_t = \frac{\pi_{\theta}(a_t | o_t)}{\pi_{\theta_{\text{old}}}(a_t | o_t)}$$

Total Loss with BC Regularisation

$$L_{\text{total}} = L_{\text{PPO}} + \beta(t) \mathbb{E}_t [\|\pi_{\theta}(o_t) - a_t^{\text{BC}}\|^2], \quad \beta(o) = \beta_o, \beta(N) \rightarrow o.$$

IMITATION WEIGHT ANNEALING



- $\beta(0) = \beta_0 > 0$: strong imitation early
- $\beta(t) \downarrow 0$: pure PPO later
- Balances stability vs. exploration

BC-SEEDED VS SCRATCH PPO

BC-seeded PPO

- Policy initialised from BC weights
- Imitation term $\beta(t) > 0$ early
- Faster convergence to walking gait
- Smaller batch sizes, fewer env steps

Scratch PPO

- Random policy initialization
- No imitation regularisation ($\beta(t) \equiv 0$)
- Requires more exploration
- Longer training to reach stability

TRAINING FLOW RECAP

1. **Rollouts:** collect on-policy trajectories under π_θ
2. **Advantage Estimation:** compute \hat{A}, \hat{R} via GAE
3. **Policy Update:** clipped PPO + $\beta(t)$ BC loss
4. **Critic Update:** regress V_ϕ to returns \hat{R}
5. **Annealing:** linearly decay imitation weight $\beta(t)$

Repeat for N iterations

VECTORISED JAX/BRAX PIPELINE

PIPELINE OVERVIEW

- **MuJoCo → Brax:** Parse your passive walker XML via `brax.io.mjcf.load_model` into a JAX-native System.
- **Environment Wrapper:** `BraxPassiveWalker` inherits `PipelineEnv`, implements `reset` / `step` in pure JAX.
- **PD Control & Reward:** $\tau = K_p(q_{\text{targ}} - q) - K_d \dot{q}$, `reward = Δx` , terminates on low torso height or large pitch.
- **Vectorisation & JIT:** Compile once with `jax.jit`, run `vmap` over N envs in parallel.
- **Train with PPO:** Call `brax.training.agents.ppo.train` on batched data for maximum throughput.

BRAXPASSIVEWALKER: reset

```
1 def reset(self, rng):
2     # 1) sample initial noise on joints
3     rng, sub = jax.random.split(rng)
4     noise = (2*jax.random.uniform(sub,(3,))-1)*self.reset_noise
5     q0 = self.sys.init_q.at[self.act_idx].add(noise*self.action_scale)
6     qd0 = jnp.zeros_like(self.sys.init_qd)
7
8     # 2) initialize pipeline state
9     ps = self.pipeline_init(q0, qd0)
10
11     # 3) return Brax State
12     return State(
13         pipeline_state=ps,
14         obs=self._get_obs(ps),
15         reward=0.0,
16         done=0.0,
17         metrics={"prev_x": ps.x.pos[0,0]}
18     )
```

- Randomises hip and knees within $\pm \text{reset_noise}$
- Builds initial Brax pipeline state \rightarrow position + velocity
- Packs into State with obs, reward, done, metrics

BRAXPASSIVEWALKER: step

```
1 def step(self, state, action):
2     # 1) PD controller
3     act_scaled = jnp.clip(action, -1, 1) * self.action_scale
4     q, qd      = state.pipeline_state.q[self.act_idx], state.pipeline_state.qd[self.
        act_idx]
5     tau        = self.kp * (act_scaled - q) - self.kd * qd
6
7     # 2) forward simulation
8     ps_next    = self.pipeline_step(state.pipeline_state, tau)
9
10    # 3) reward & done
11    dx          = ps_next.x.pos[0, 0] - state.metrics["prev_x"]
12    height_ok   = ps_next.x.pos[0, 2] > 0.5
13    pitch_ok    = abs(quat_to_euler(ps_next.x.rot[0])[1]) < 0.8
14    done_flag   = jnp.logical_not(height_ok & pitch_ok)
15
16    # 4) pack next state
17    return state.replace(
18        pipeline_state=ps_next,
19        obs=self._get_obs(ps_next),
20        reward=dx,
21        done=done_flag.astype(jnp.float32),
22        metrics={"prev_x": ps_next.x.pos[0, 0]}
23    )
```

- **PD torques:** $\tau = K_p(a_{\text{target}} - q) - K_d \dot{q}$
- **Reward:** forward progress Δx
- **Termination:** torso height < 0.5 m or pitch || 0.8 rad
- **State replace:** updates obs, reward, done, and prev_x

HIGH-THROUGHPUT TRAINING

- **Massive sweep:** 120 PPO jobs \times 24M steps each (few hours on RTX 4060 Ti).
- **Parallel simulators:** 128 envs in lock-step via vmap $\rightarrow \times 100$ speed-up vs. loop.
- **Compile once:** JAX JIT fixes the compute graph, reuses across all envs iterations.
- **Toolchain:** Brax for physics, Flax/Equinox for networks, Optax for optimizers.
- **Reproducible logging:** msgpack payloads, deterministic seeds, helper scripts for aggregate plots.

EXPERIMENTAL SWEEP DESIGN

EXPERIMENTAL SWEEP DESIGN

■ Goals:

- ▶ Find the single best hyper-parameter configuration
- ▶ Analyse trends over reward scale, learning rate and network size

■ **Grid:** 3 seeds \times 2 reward scales \times 4 LRs \times 5 architectures = 120 jobs

■ **Metrics:** final episode reward (mean \pm std over seeds), wall-clock time, stable completions

Dimension	Values (count)
Seeds	{0,1,2} (3)
Reward scale	{0.5,1.0} (2)
Learning rate	{1e-3,5e-4,1e-4,1e-5} (4)
Architectures	{tiny,small,medium,deep,deepXL} (5)
Total runs	120

Table: Sweep grid dimensions (values in braces).

SWEEP RESULTS: REWARD SCALE & LEARNING RATE TRENDS

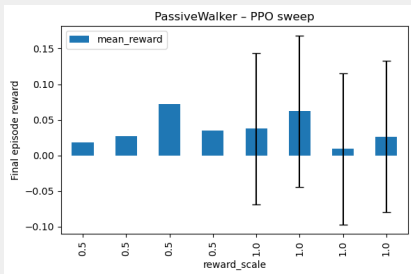


Figure: Mean \pm std reward vs. reward scale

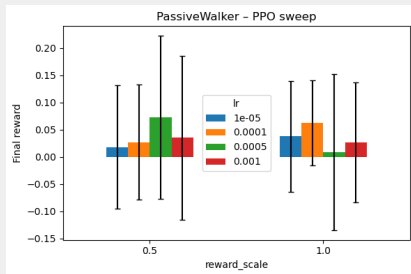


Figure: Mean \pm std reward vs. learning rate

Best Hyper-Parameter Configuration

S2 | R = 0.5 | LR = 1e-3 | Arch = medium

- Highest mean reward across seeds
- Shows that moderate reward scaling and a higher learning rate perform best
- Typical wall-clock: ~2–10 min per job on RTX 4060 Ti
- All 120 jobs completed successfully (no simulator crashes)

Guideline for future experiments: start with
medium+LR=1e-3+scale=0.5.

RESULTS & DISCUSSION

BEHAVIOUR CLONING PERFORMANCE

- **BC variants:** hip-only & knee-only (context), focus on full BC (hip + knees)
- **Losses compared:** MSE, Huber, L1, Combined

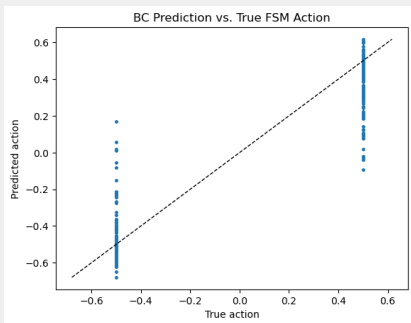


Figure: Hip: BC vs. FSM labels

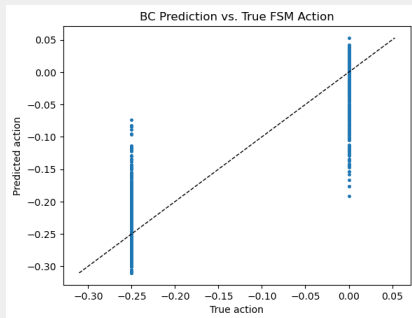


Figure: Knee: BC vs. FSM labels

BC LOSS VARIANT COMPARISON

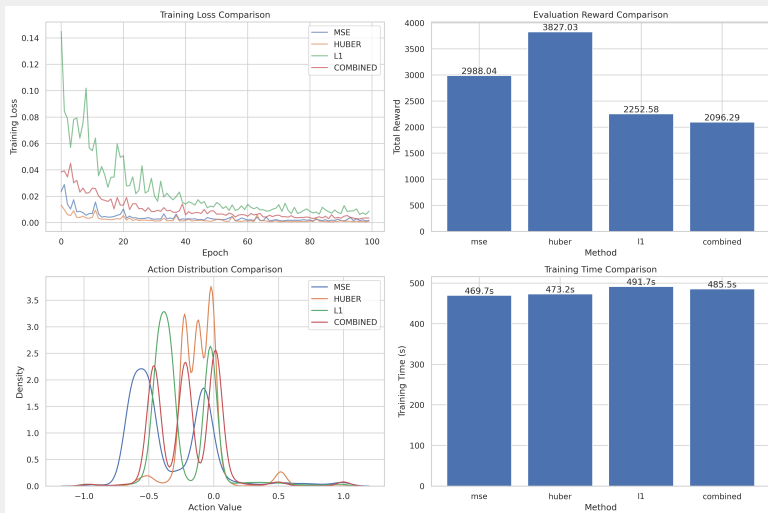


Figure: Training loss, evaluation reward, action-density and runtime for each BC objective

PPO FINE-TUNING PERFORMANCE

- **BC-seeded PPO** converges in $\sim 50\%$ of the iterations vs. scratch PPO
- **Scratch PPO** achieves similar final gait but needs $\times 2 - \times 3$ more samples
- Decaying imitation weight $\beta(t)$ stabilises early training

HYPERPARAMETER SWEEP: AGGREGATED TRENDS

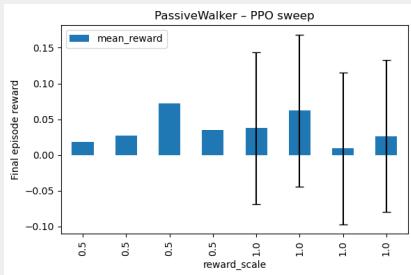


Figure: Mean \pm std reward vs. reward scale

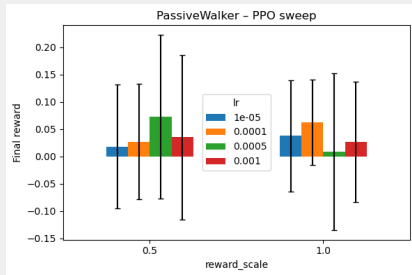


Figure: Mean \pm std reward vs. learning rate

DISCUSSION & KEY TAKEAWAYS

- **Behaviour Cloning:** Huber loss offers best trade-off of bias vs. robustness
- **Imitation-seeded PPO:** drastically cuts sample complexity and speeds convergence
- **Sweep guidelines:** medium networks, $LR=1e-3$, $scale=0.5$ for strong performance
- **High-throughput pipeline:** vectorised JAX/Brax makes such large sweeps practical

IMPACT & ETHICAL ISSUES

- **Accelerated Research:** “Walk-from-day-one” imitation jump-starts RL in legged locomotion.
- **Open-Source Toolkit:** JAX/Brax pipeline + scripts for large-scale sweeps released on GitHub.
- **Broader Applications:** Assistive robotics, search-and-rescue, exploratory platforms.
- **Efficiency Gains:** Vectorised simulation achieves $\times 100$ speed-up vs. naive loops.

- **Real-World Safety:** Simulator-to-robot gap demands rigorous hardware safety checks.
- **Bias in Demonstrations:** FSM expert encodes narrow gait patterns—limits generalisation.
- **Compute Footprint:** 120-job sweep deep training consume GPU hours—track budgets.
- **Dual-Use Risks:** Legged controllers may be repurposed for surveillance or military.
- **Transparency & Reproducibility:** Full code, configs, logs published; encourage peer verification.

PROJECT MANAGEMENT & TIMELINE

- **Solo Effort:** All design, implementation, experiments by the author
- **Learning Curve:** JAX/MuJoCo/Brax resources sparse—many tools mastered from scratch
- **Challenges:** XML→Brax compatibility, PD tuning, large-scale sweep orchestration
- **Resources:** Initial work on CPU laptop; final deep runs on RTX 4060 Ti workstation
- **Deliverables:** Well-organized GitHub repo with code, data, scripts, and documentation

TIMELINE

Phase		Activities
FSM & Data Collection		Expert FSM design, MuJoCo demos, dataset logging
Behaviour Cloning		MLP architecture, loss variants, BC training
PPO	Fine-Tuning	BC-seeded & scratch PPO, ablation studies
Brax Port & Hyperparameter Sweep		MJCF→Brax conversion, 120-job PPO grid
Write-up & Release		Report chapters, presentation, open-source push

CONCLUSION & FUTURE WORK

CONCLUSION

- **Pipeline Success:** FSM→BC→PPO→Brax achieved robust passive walking “from day one.”
- **Sample Efficiency:** Imitation regularisation halved RL training steps vs. scratch.
- **Scalability:** JAX/Brax vectorisation made 24 M-step sweeps feasible in hours.
- **Open Science:** All code, models, results publicly available for community reuse.

FUTURE WORK

- **Complex Terrain:** Extend to uneven ground, stairs, and variable slopes.
- **Hardware Validation:** Transfer policies to real robot—study sim-to-real gaps and safety.
- **Domain Randomisation:** Improve robustness to mass, friction, and sensor noise variations.
- **Advanced Architectures:** Graph-based or attention-powered controllers for multi-limb coordination.
- **Energy Efficiency:** Incorporate power/regret into reward for practical deployment.

ADDITIONAL RESOURCES

- **Code & Models:** https://github.com/yunusdanabas/passive_walker_rl.git
- **Contact:** yunusdanabas@sabanciuniv.edu