

Step 1: IK Computes Position Error

$$\Delta p_{ee} = p_{ee}^{goal} - p_e^{current}$$

(EE goal position - Current EE position)



Step 2: Error Transformed to Reward

$$r_{arm} = \exp\left(-\frac{\|\Delta p_{ee}\|}{2\sigma}\right) \times 0.8$$

- Small error → High reward (≈ 0.8)
- Large error → Low reward (≈ 0)



Step 3: Total Reward Drives Policy

$$r_{total} = r_{loco} + r_{arm}$$

$$\nabla_{\theta} J = \mathbb{E} [\nabla_{\theta} \log \pi(\mathbf{a}|\mathbf{o}) \cdot r_{total}]$$

(Locomotion + Manipulation reward)



Step 4: Locomotion "Compromises"

Policy discovers:

- Pure walking → Large error → Low r_{total}

Policy learns:

- Adjust legs → Support arm → High r_{total}



Emergent Whole-Body Control

Legs serve global goal (EE tracking)
instead of independent locomotion