

Policy composition in reinforcement learning via multi-objective policy optimization

Shruti Mishra^{1,2} Ankit Anand² Jordan Hoffmann² Nicolas Heess²
Martin Riedmiller² Abbas Abdolmaleki² Doina Precup^{2,3,4}

¹Sony AI ²Google DeepMind ³McGill University ⁴Mila

A multi-objective approach to leverage experts for locomotion

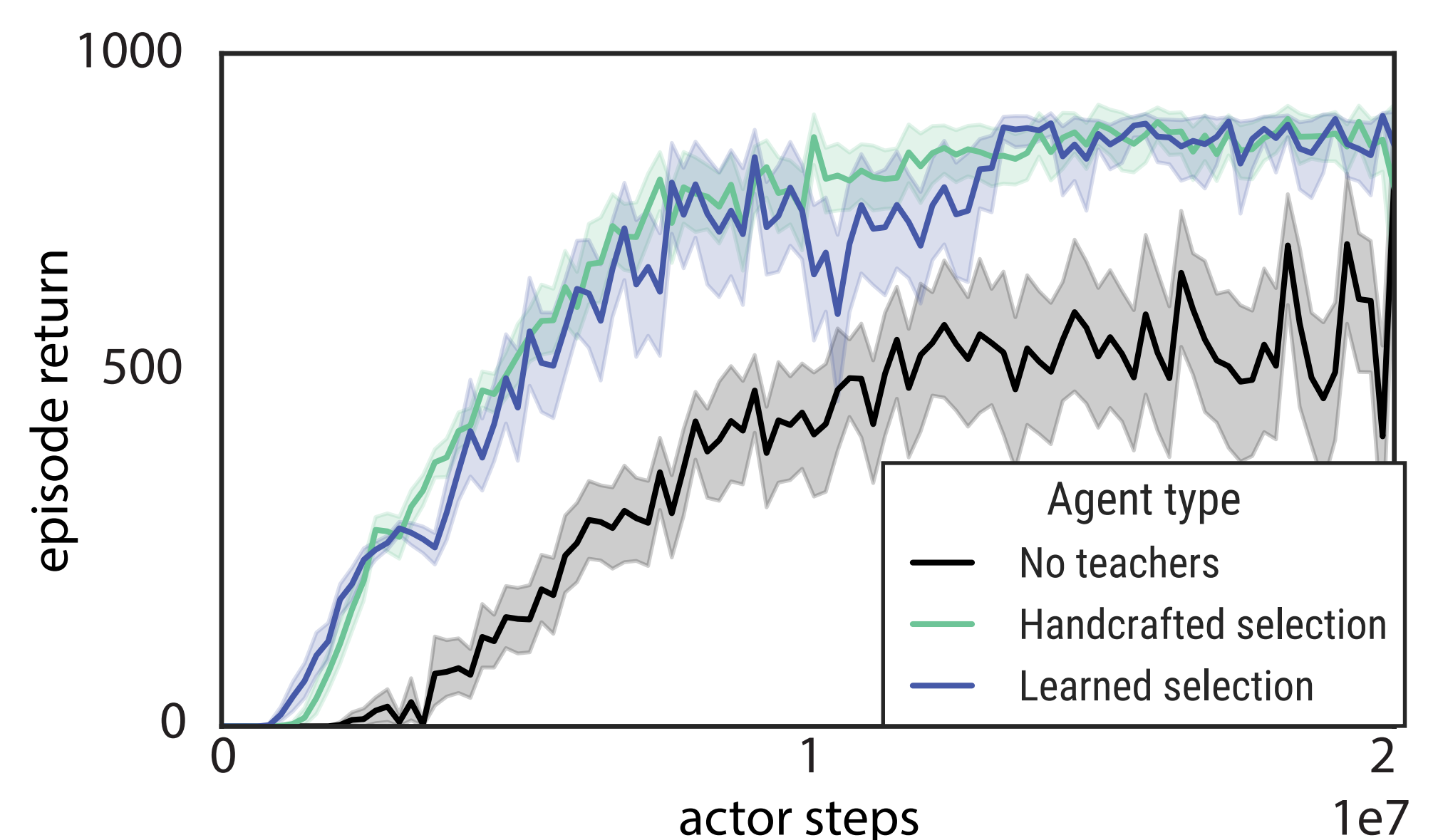
Motivation

- Pre-existing expertise can be useful for agents in locomotion domains.
- Improvements to learning efficiency can enable more general-purpose agents in environments.

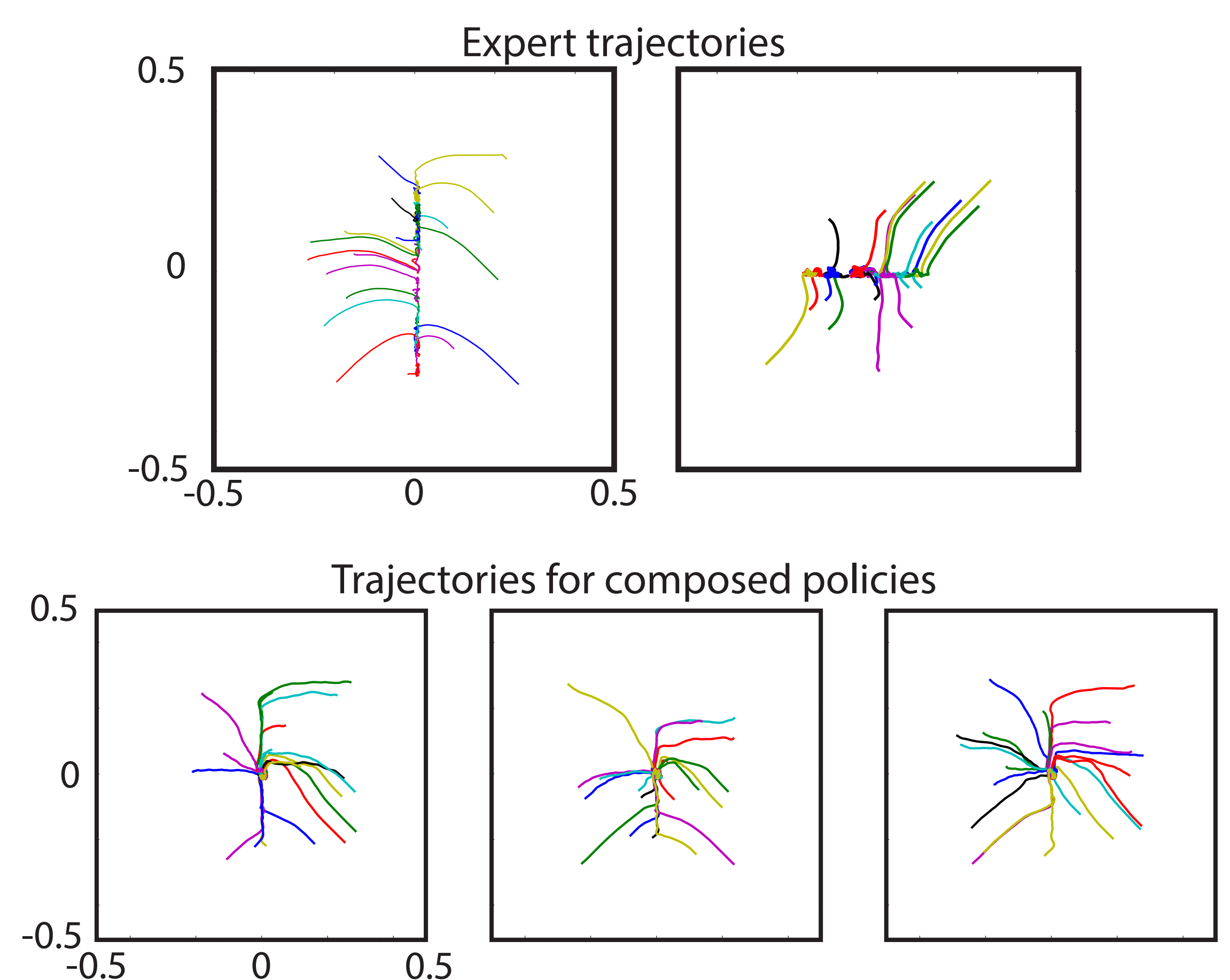


Outcomes

- Agents successfully choose how to leverage pre-existing policies.
- This can overcome limitations of reward shaping.

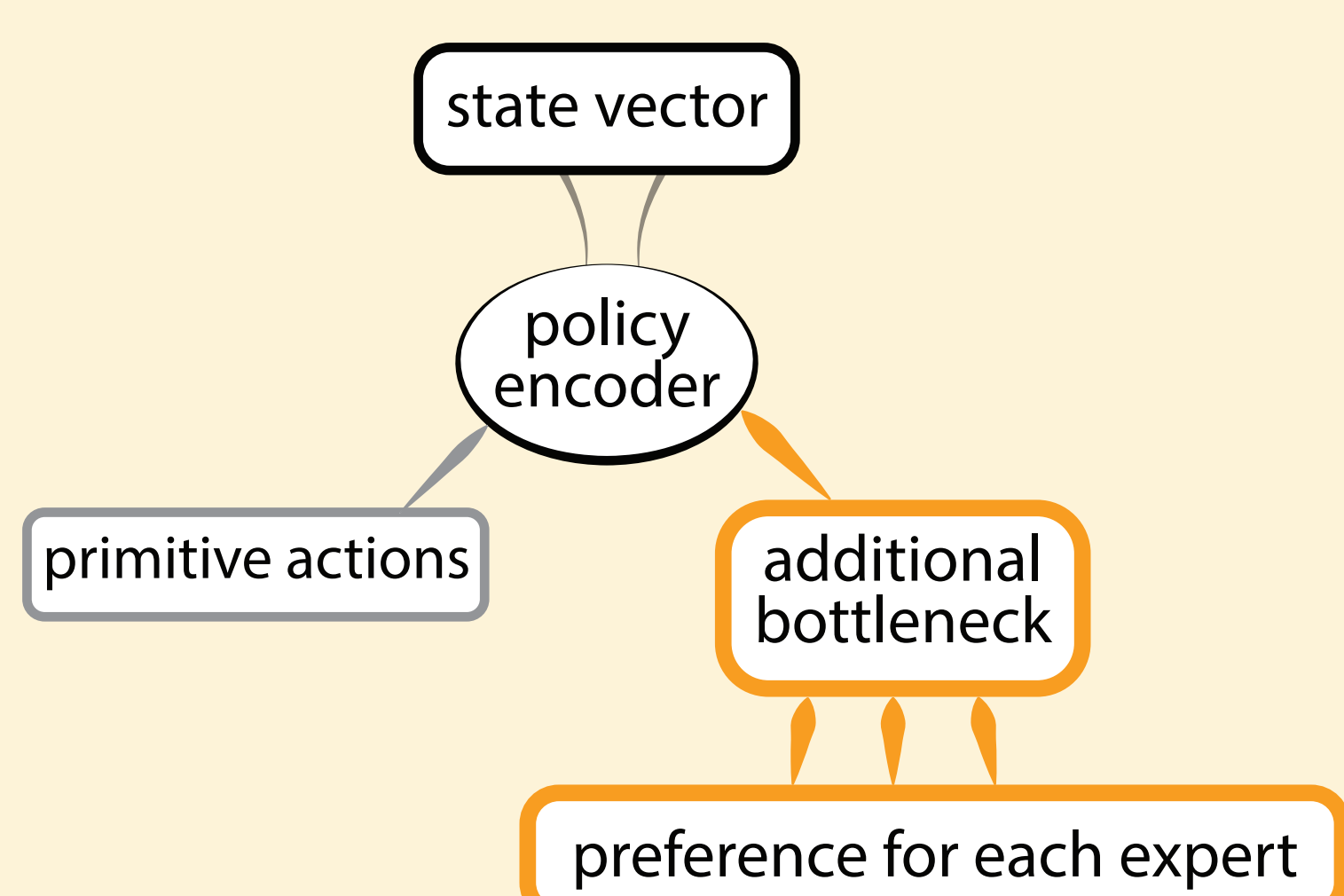


- Policies are shaped by experts.



Method

1. Reinforcement learning agents in locomotion domains have access to pre-existing, sub-optimal policies.
2. Adherence to teacher policies is framed as additional objectives in a multi-objective RL setup.
3. Agents control the adherence to pre-existing policies using the task objective.



contact: shruti.mishra10@alumni.imperial.ac.uk