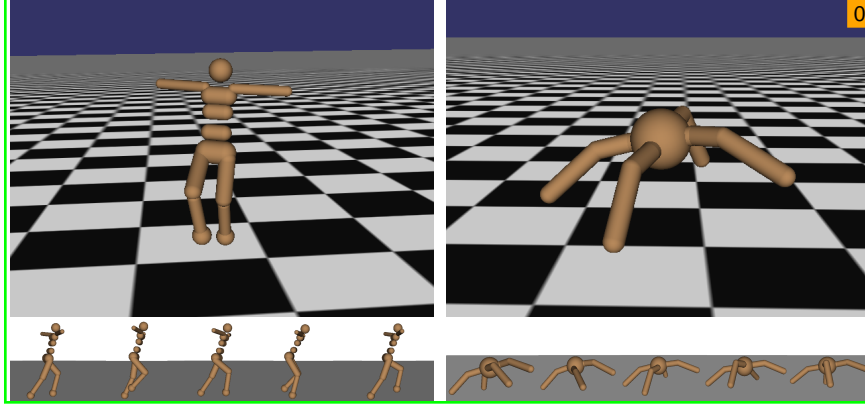


figure (0.98)



figure\_caption (0.95)

Figure 1: Top figures: robot models used for 3D locomotion. Bottom figures: a sequence of frames from the learned gaits. Videos are available at <https://sites.google.com/site/gaepapersupp>.

title (0.92)

## 6.2.2 TASK DETAIL 2

plain\_text (0.95)

For the cart-pole balancing task, we collected 20 trajectories per batch, with a maximum length of 1000 timesteps, using the physical parameters from Barto et al. (1983).

plain\_text (0.96)

The simulated robot tasks were simulated using the MuJoCo physics engine (Todorov et al., 2012). The humanoid model has 33 state dimensions and 10 actuated degrees of freedom, while the quadruped model has 29 state dimensions and 8 actuated degrees of freedom. The initial state for these tasks consisted of a uniform distribution centered on a reference configuration. We used 50000 timesteps per batch for bipedal locomotion, and 200000 timesteps per batch for quadrupedal locomotion and bipedal standing. Each episode was terminated after 2000 timesteps if the robot had not reached a terminal state beforehand. The timestep was 0.01 seconds.

plain\_text (0.97)

The reward functions are provided in the table below.

Task	Reward
3D biped locomotion	$v_{\text{fwd}} - 10^{-5} \ u\ ^2 - 10^{-5} \ f_{\text{impact}}\ ^2 + 0.2$
Quadruped locomotion	$v_{\text{fwd}} - 10^{-6} \ u\ ^2 - 10^{-3} \ f_{\text{impact}}\ ^2 + 0.05$
Biped getting up	$-(h_{\text{head}} - 1.5)^2 - 10^{-5} \ u\ ^2$

plain\_text (0.95)

Here,  $v_{\text{fwd}}$  := forward velocity,  $u$  := vector of joint torques,  $f_{\text{impact}}$  := impact forces,  $h_{\text{head}}$  := height of the head.

plain\_text (0.96)

In the locomotion tasks, the episode is terminated if the center of mass of the actor falls below a predefined height: .8 m for the biped, and .2 m for the quadruped. The constant offset in the reward function encourages longer episodes; otherwise the quadratic reward terms might lead to a policy that ends the episodes as quickly as possible.

title (0.90)

## 6.3 EXPERIMENTAL RESULTS 9

plain\_text (0.98)

All results are presented in terms of the cost, which is defined as negative reward and is minimized. Videos of the learned policies are available at <https://sites.google.com/site/gaepapersupp>. In plots, “No VF” means that we used a time-dependent baseline that did not depend on the state, rather than an estimate of the state value function. The time-dependent baseline was computed by averaging the return at each timestep over the trajectories in the batch.

title (0.90)

### 6.3.1 CART-POLE 11

plain\_text (0.97)

The results are averaged across 21 experiments with different random seeds. Results are shown in Figure 2, and indicate that the best results are obtained at intermediate values of the parameters:  $\gamma \in [0.96, 0.99]$  and  $\lambda \in [0.92, 0.99]$ .