



For each  $t$ :

network loss,  
network layers

reward  $r_t$

Next layer  
from pool

- (1) info. characterizing  
environment's network
- (2) info. characterizing  
task

state  $\mathbf{s}_t$

REML

agent

$\pi_{\text{PPO}}$

$\mathbf{a}_t$

action

layer pool

