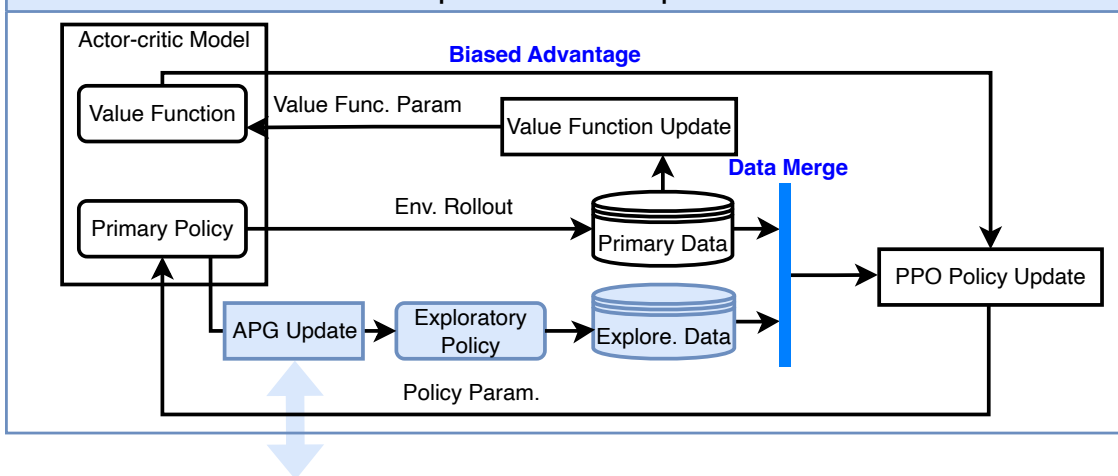
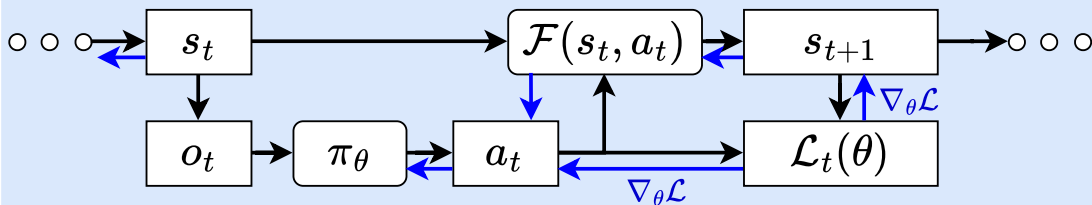


## PPO Update with Biased Exploration



## Policy Update with Analytical Policy Gradient



$s_t$ State	$\mathcal{L}_t(\theta)$ Loss/Return	$\pi_\theta$ Observation	$\rightarrow$ Forward Process
$o_t$ Observation	$s_{t+1}$ Next State	$\mathcal{F}$ Differentiable Dynamics	$\leftarrow$ Backpropagation