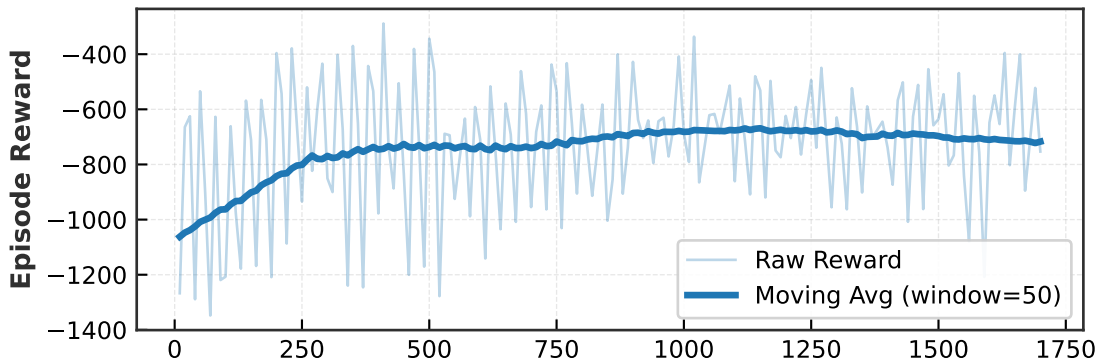


(a) Reward Convergence



(b) Training Stability

