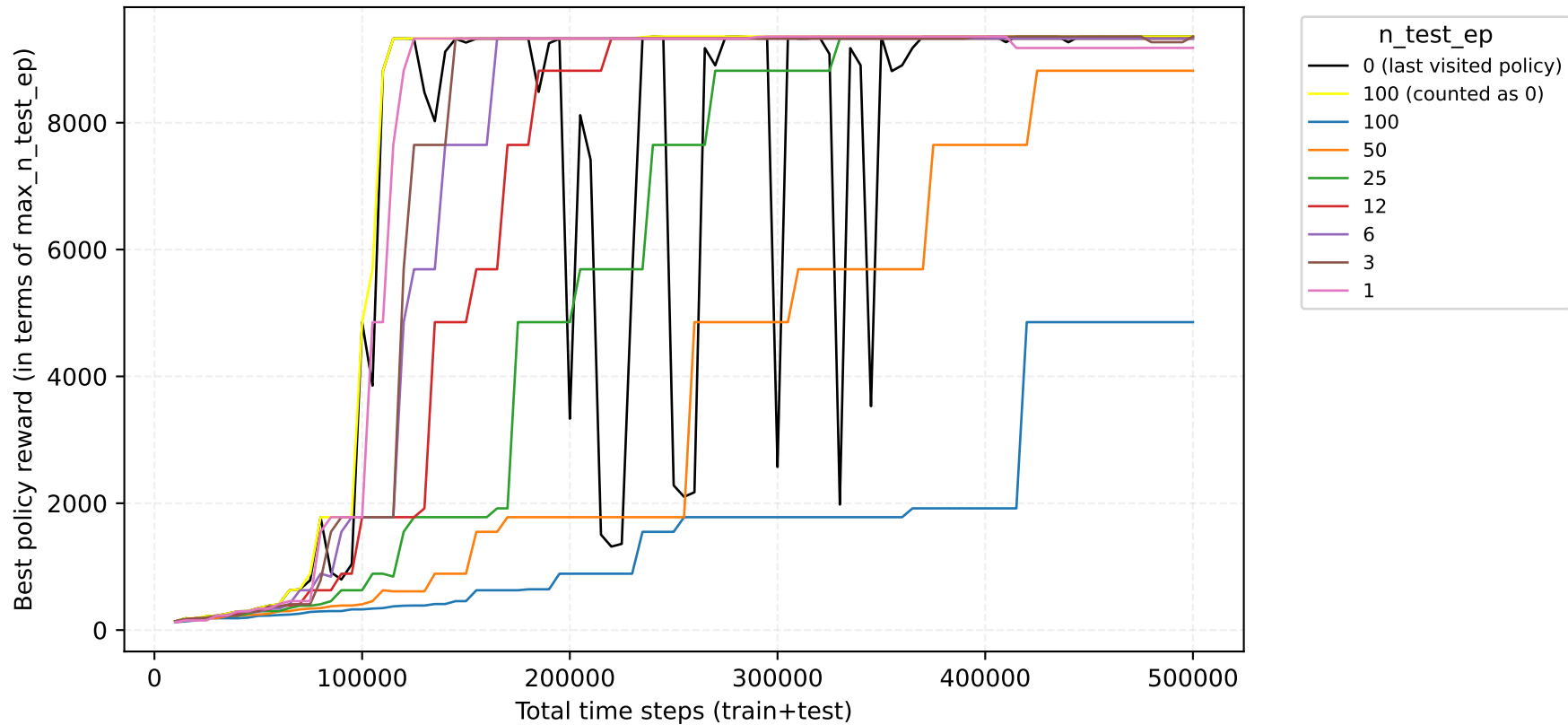


Learning-curve: best policy found
(iteratively validating policies with $n_{\text{test_ep}}$ test episodes)



(2048 train timesteps = training 1 policy = 1 trajectory)