

Abstract

We study an online continual lifelong learning scenario where the underlying dynamics of the world can change, and the agent has no ability to reset, meaning that mistakes compound catastrophically into the future.

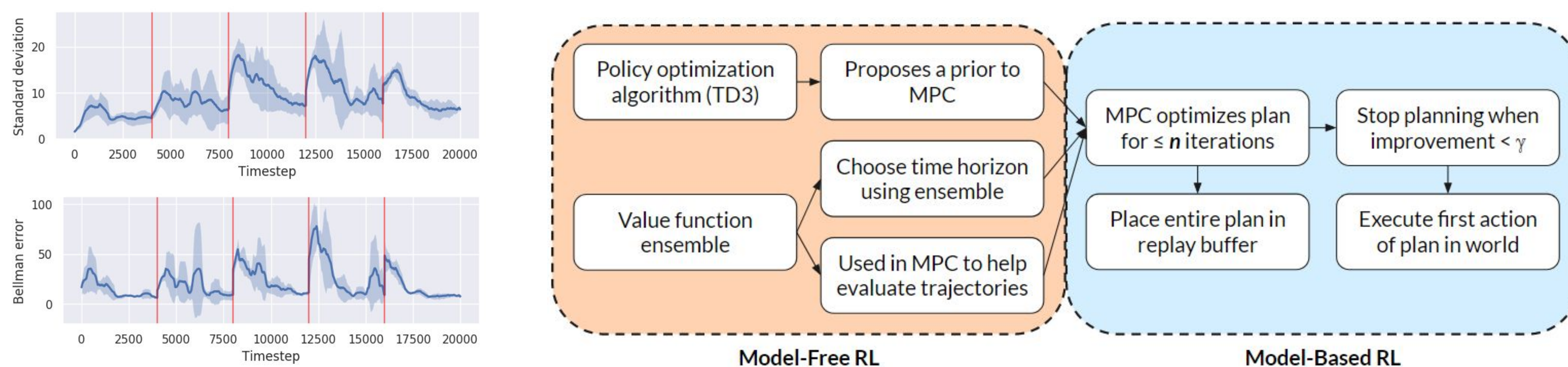
Model-based planning is effective immediately given the model, but is prohibitively expensive and can be biased by the planning horizon. Model-free learning is slow, leading to crucial mistakes that can land the agent into areas where learning is difficult, and can struggle to adapt to world changes.

What do we need to learn effectively in this setting with limited computation?

Adaptive Online Planning

AOP augments MPC with the model-free learning of value and policy networks. Model-based planning allows for effective control and rapid adaptation early in life, while model-free learning allows for reduced computation, stable behavior, and higher asymptotic performance later in life.

We quantify uncertainty as the disagreement of an ensemble of value functions. Early in life, standard deviation is high. At world changes (red), the value function weakly predicts the value of the current state (Bellman error). When these are below thresholds, AOP uses a lower planning horizon.

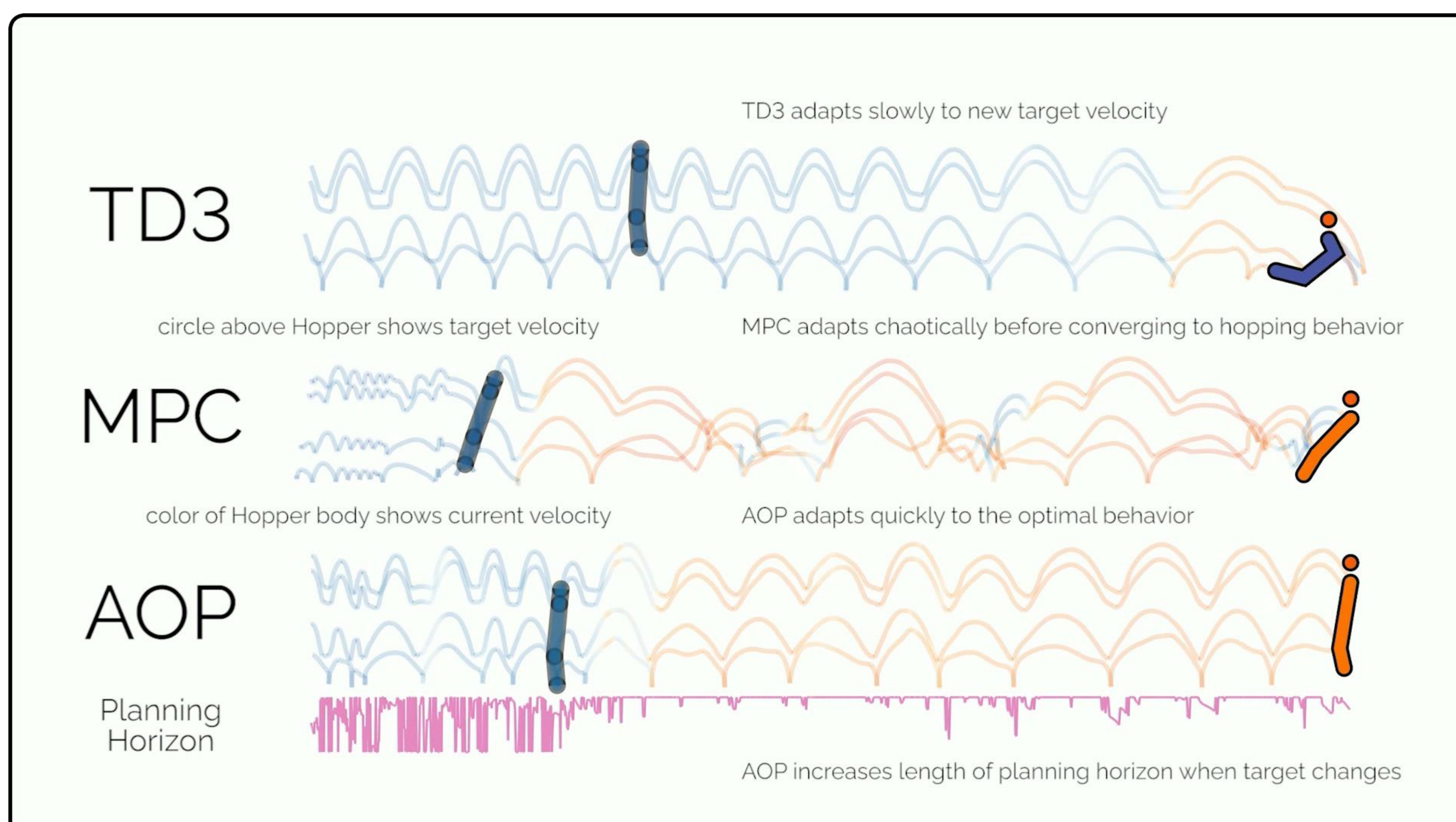


AOP executes more planning iterations when planning iterations improve the planned trajectory, using the value function to guide planning and exploration. The policy is preferred when it is strong in comparison to the planned trajectory.

Fast Adaptation

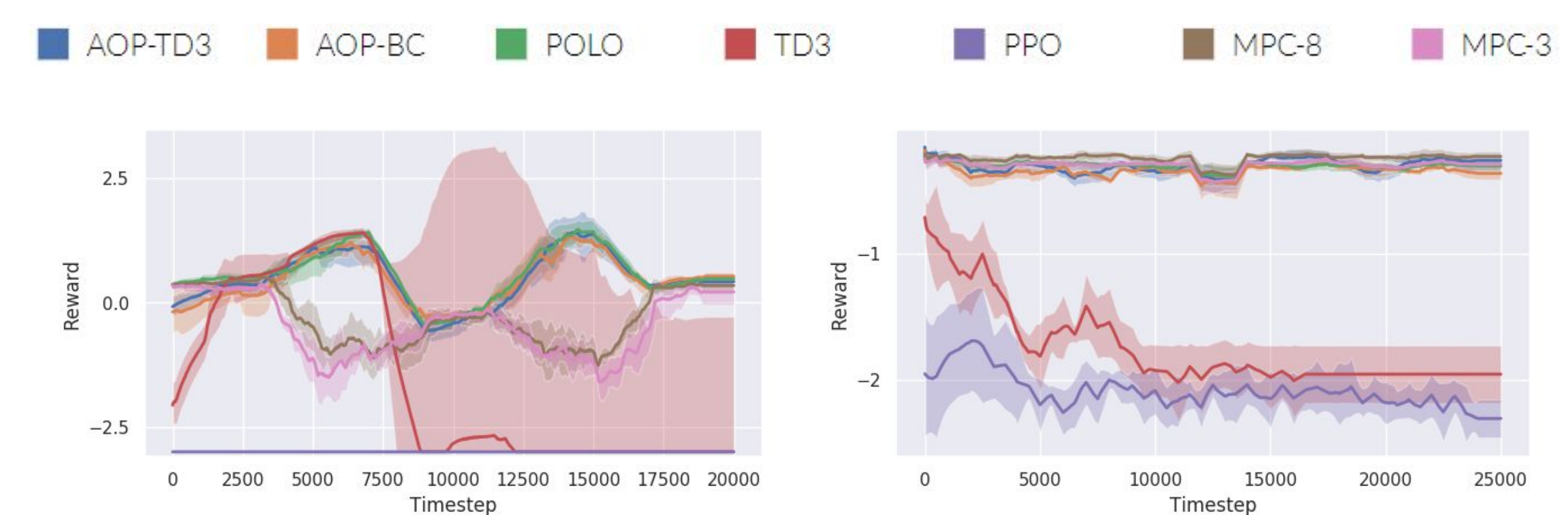
The hopper tries to achieve a target forward velocity. The blue outline of the hopper shows when the agent is tasked with switching from a slow (blue) to a fast (orange) target velocity.

TD3 adapts slowly and falls over, while MPC acts chaotically, leading to suboptimal behavior. AOP quickly adapts to the new world setting, increasing its level of planning accordingly to do so.



Changing Worlds

AOP achieves the performance of more expensive model-based planners while avoiding the performance degradation over time of model-free methods.



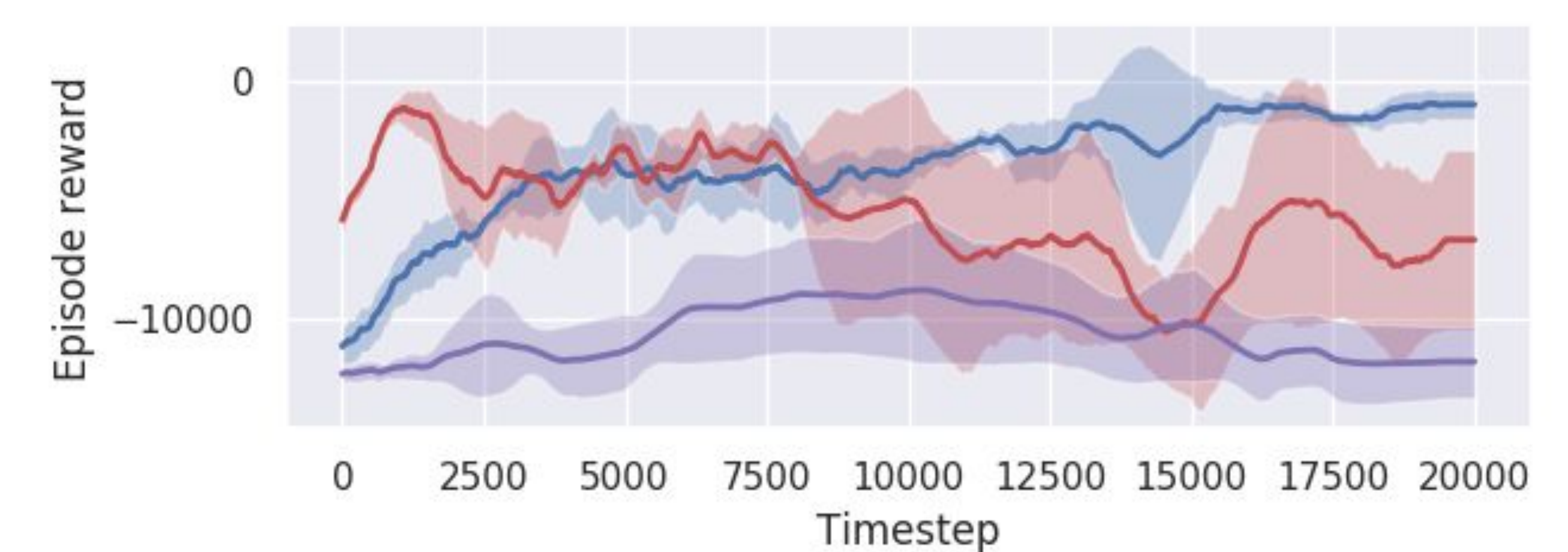
Novel States

Furthermore, AOP demonstrates learning; model-based planning only stays static in performance. Left: AOP learns to generalize to unseen worlds. Right: AOP performs better in mazes it sees many times (red lines show new worlds).



Policy Degradation

Unlike the model-free algorithms, the policy learned by AOP does not degrade over time. Below shows the policy learned by the algorithm ran in a standard episode from the starting state at the beginning of the agent's lifetime.



Planning

AOP uses significantly less planning than comparable settings for other algorithms over a lifetime.

AOP-TD3	AOP-BC	POLO	MPC-8	MPC-3
11.39%	11.40%	37.50%	100%	37.50%

More Information

Website: bit.ly/aop_neurips

arXiv: arxiv.org/abs/1912.01188

Code: github.com/kzl/aop