

Surprise!

Agent orange and agent blue are trained with...

1. **The same off-policy algorithm (DDPG).**
2. **The same dataset.**

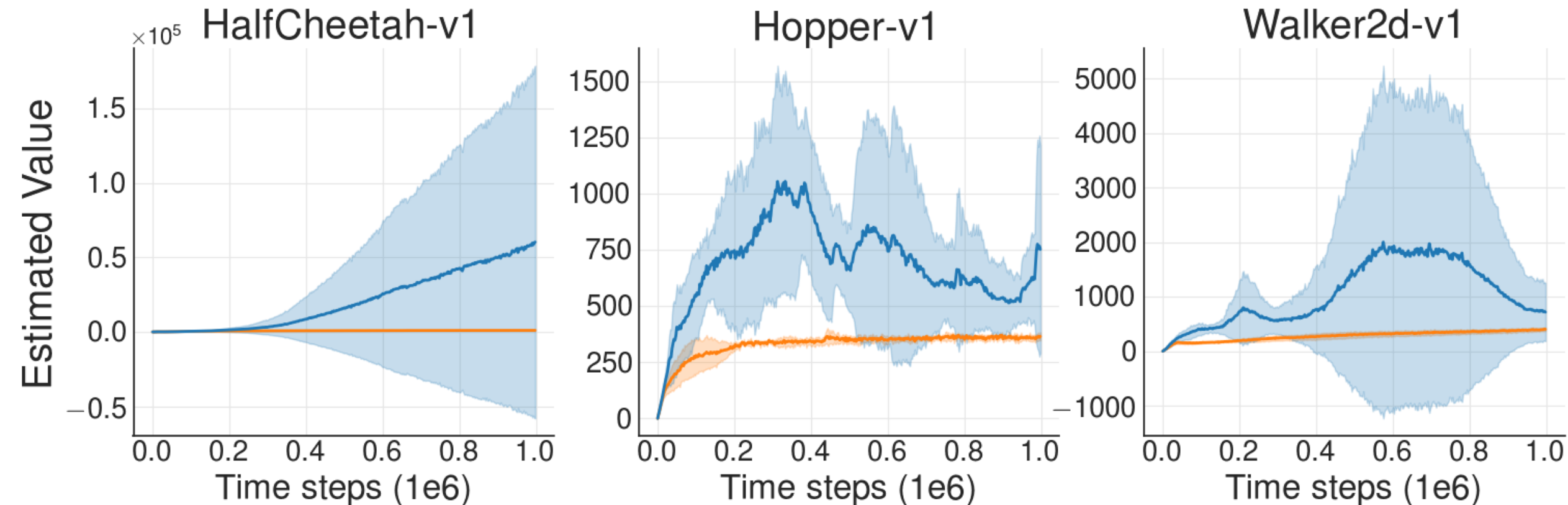
The Difference?

1. **Agent orange:** Interacted with the environment.
 - Standard RL loop.
 - Collect data, store data in buffer, train, repeat.
2. **Agent blue:** Never interacted with the environment.
 - Trained with data collected by agent orange concurrently.

1. Trained with the same off-policy algorithm.
2. Trained with the same dataset.
3. One interacts with the environment. One doesn't.

Off-policy deep RL fails when **truly off-policy**.

Value Predictions



Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

The diagram illustrates the Bellman optimality equation $Q(s, a) \leftarrow r + \gamma Q(s', a')$. It features two labels at the bottom: "GIVEN" in red and "GENERATED" in blue. Red arrows point from "GIVEN" to the state-action pair (s, a) and the reward r . A red arrow points from the previous state-action pair (s', a') to the discounted value term $\gamma Q(s', a')$. A blue arrow points from "GENERATED" to the next state-action pair (s', a') .

Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

1. $(s, a, r, s') \sim \text{Dataset}$
2. $a' \sim \pi(s')$

Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$
 $\rightarrow Q(s, a) = \mathbf{bad}$

Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$$
$$\rightarrow Q(s, a) = \mathbf{bad}$$

Extrapolation Error

$$Q(s, a) \leftarrow r + \gamma Q(s', a')$$

$$(s', a') \notin \text{Dataset} \rightarrow Q(s', a') = \mathbf{bad}$$
$$\rightarrow Q(s, a) = \mathbf{bad}$$

Extrapolation Error

Attempting to evaluate π without (sufficient) access to the (s, a) pairs π visits.

Batch-Constrained Reinforcement Learning

Only choose π such that we have access to the (s, a) pairs π visits.

Batch-Constrained Reinforcement Learning

1. $a \sim \pi(s)$ such that $(s, a) \in Dataset$.
2. $a \sim \pi(s)$ such that $(s', \pi(s')) \in Dataset$.
3. $a \sim \pi(s)$ such that $Q(s, a)$ is maxed.

Batch-Constrained Deep Q-Learning (BCQ)

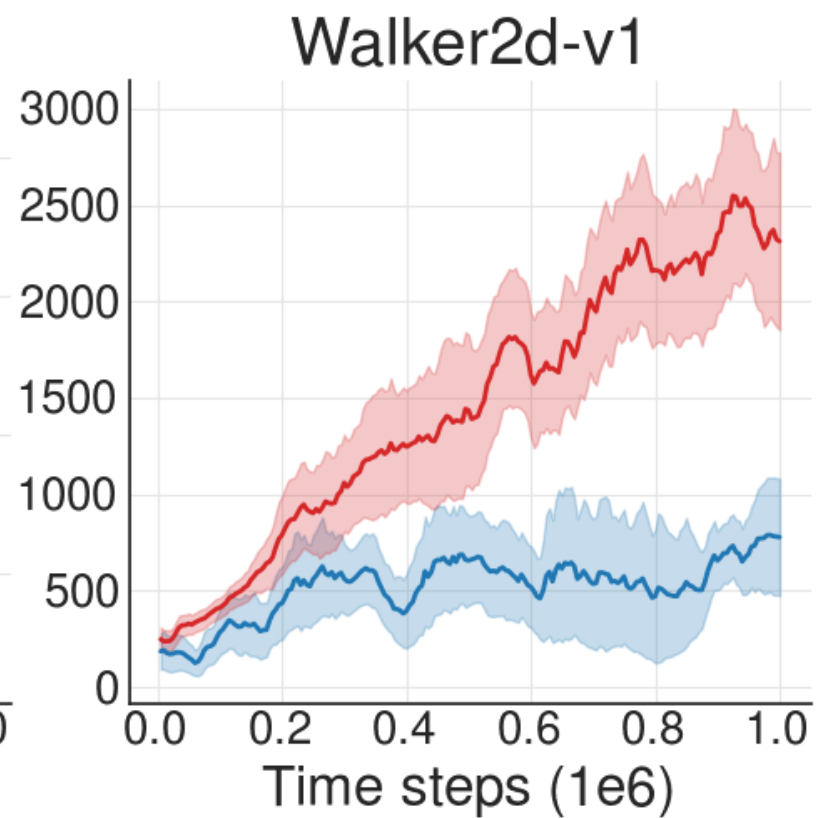
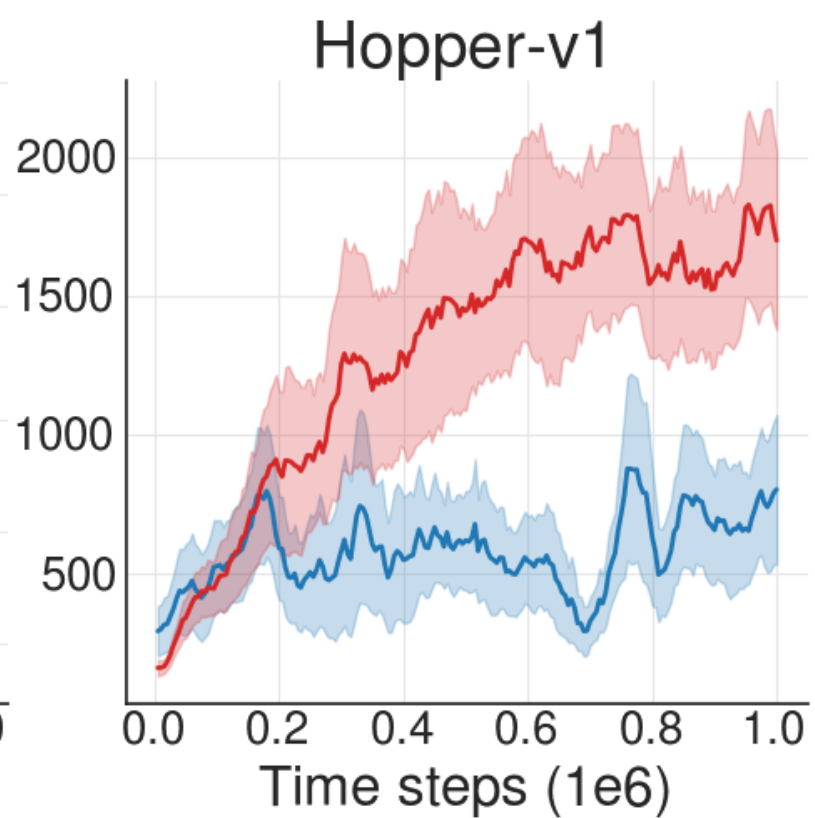
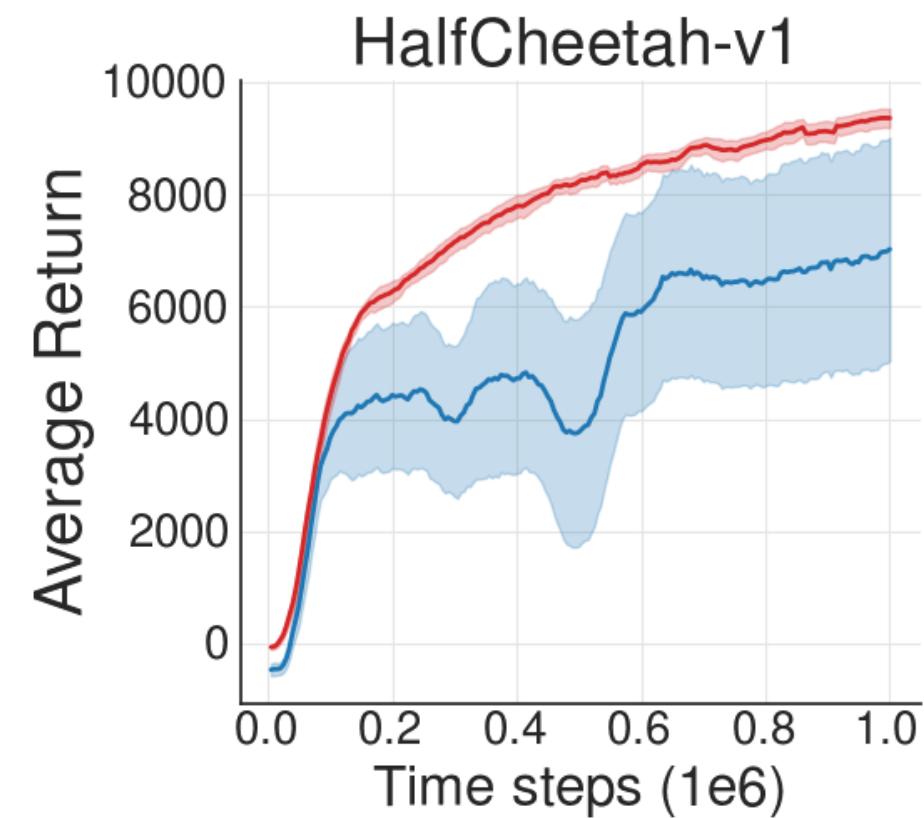
First imitate dataset via generative model:

$$G(a|s) \approx P_{Dataset}(a|s).$$

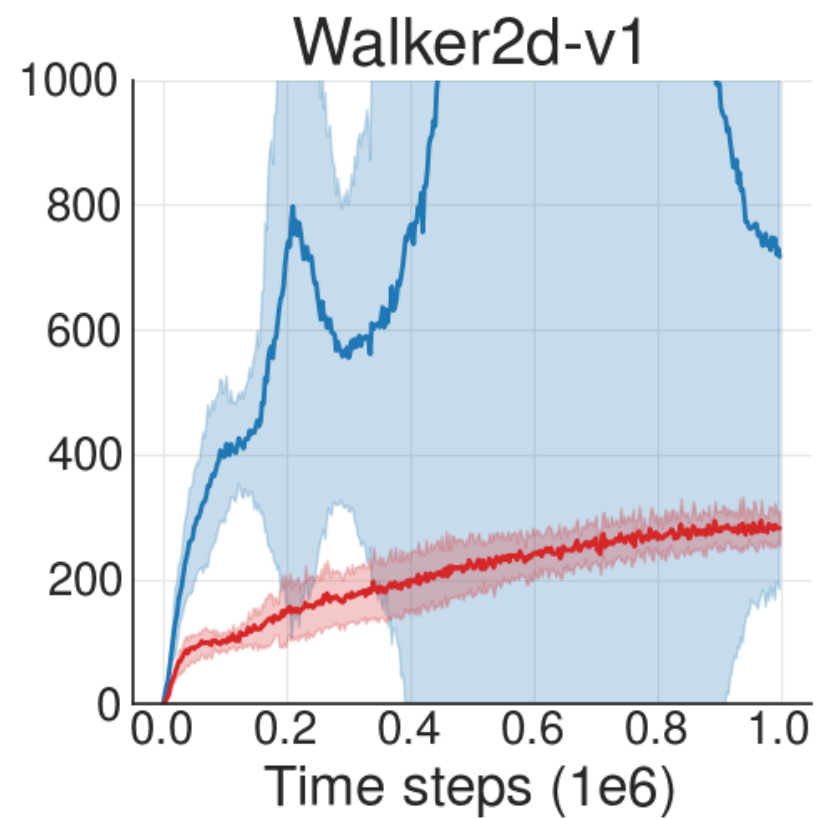
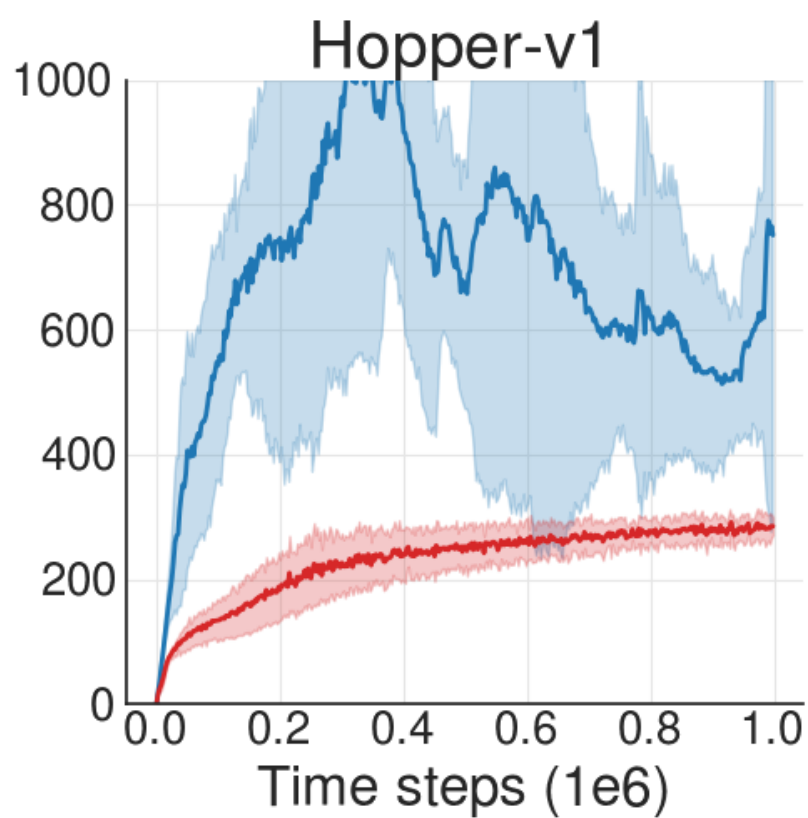
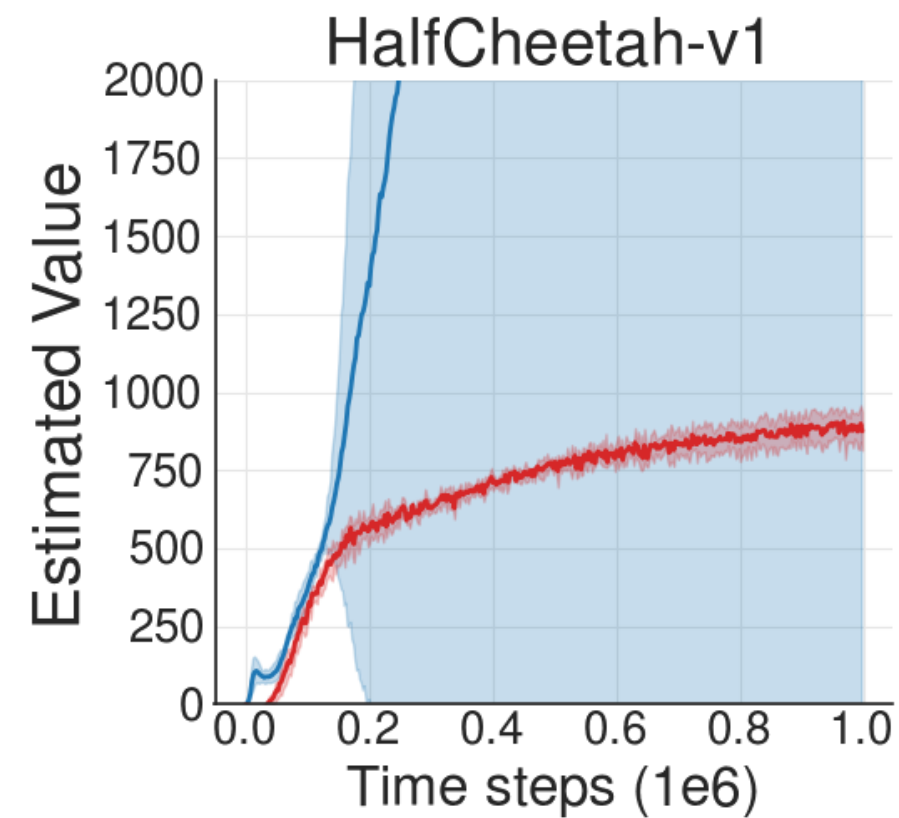
$$\pi(s) = \operatorname{argmax}_{a_i} Q(s, a_i), \text{ where } a_i \sim G$$

(I.e. select the best action that is likely under the dataset)

(+ some additional deep RL **magic**)



■ BCQ ■ DDPG



■ BCQ ■ DDPG