

# GENERAL POLICY EVALUATION AND IMPROVEMENT BY LEARNING TO IDENTIFY FEW BUT CRUCIAL STATES

FRANCESCO FACCIO, ADITYA RAMESH, VINCENT HERRMANN, JEAN HARB AND

JÜRGEN SCHMIDHUBER

{francesco, aditya, vincent.herrmann, juergen}@idsia.ch, jean.merheb-harb@mail.mcgill.ca

## PROBLEM AND MOTIVATION

- **Reinforcement Learning (RL):** Find optimal policy  $\pi^*$
- **Policy optimization:** Given a class of policies, find the policy parameters maximizing  $J(\pi_\theta)$  (Sutton et al., 1999):  

$$J(\pi_\theta) = \int_S \mu_0(s) V^{\pi_\theta}(s) ds = \int_S \mu_0(s) \int_{\mathcal{A}} \pi_\theta(a|s) Q^{\pi_\theta}(s, a) da ds$$
- **Problem:** Value functions are defined for a *single policy*. During optimization, the information on previous policies is potentially lost
- **Solution idea:** Learn a single value function able to evaluate many policies

## PARAMETER-BASED VALUE FUNCTIONS

- **Parameter-Based Value Functions** (Faccio et al., 2021) generalize over multiple policies by incorporating the policy parameters as an additional input
- **PSVF:** Parameter based state-value function  

$$V(s, \theta) := \mathbb{E}[R_t | s_t = s, \theta]$$
- **PAVF:** Parameter based action-value function  

$$Q(s, a, \theta) := \mathbb{E}[R_t | s_t = s, a_t = a, \theta]$$
- **PSSVF:** Parameter based start-state-value function  

$$V(\theta) := \mathbb{E}[R_0 | \theta]$$
- The PSSVF models  $J(\theta)$  directly as a differentiable function  $V(\theta)$ , which is the expectation of  $V(s, \theta)$  over the initial states

$$V(\theta) := \mathbb{E}_{s \sim \mu_0(s)}[V(s, \theta)] = \int_S \mu_0(s) V(s, \theta) ds = J(\pi_\theta).$$

If we can learn  $V_w(\theta)$ , we can improve the policy by simply taking gradient ascent steps

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta V_w(\pi_\theta)$$

- **Problem:** How can we give the policy parameters as input to the value function when  $\pi_\theta$  is a neural network?

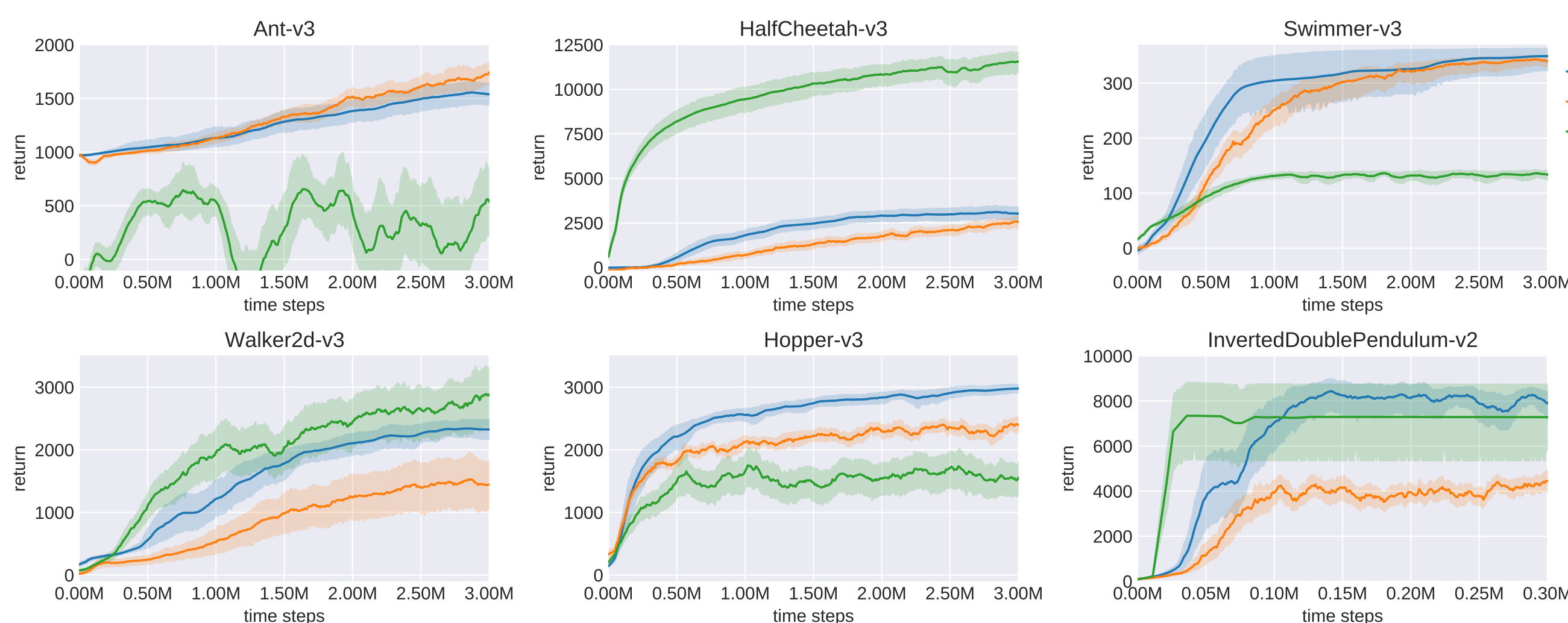
## POLICY FINGERPRINTING

- **Policy fingerprinting** (Harb et al., 2020) creates a lower-dimensional policy representation
- It learns a set of  $K$  **probing states**  $\{\tilde{s}_k\}_{k=1}^K$  and an evaluation function  $V_\phi$
- To evaluate a policy  $\pi_\theta$ , it computes the ‘probing actions’  $\tilde{a}_k$  that the policy produces in the probing states. Then the concatenated vector of these actions is given as input to  $V_\phi$  and mapped to the return
- Setting  $w = \{\phi, \tilde{s}_1, \dots, \tilde{s}_K\}$ , we optimize  $V_w$  minimizing MSE:

$$\begin{aligned} \min_w \mathcal{L}_V &:= \min_w \mathbb{E}_{(\pi_\theta, r) \in B} [(V_w(\theta) - r)^2] \\ &= \min_{\phi, \tilde{s}_1, \dots, \tilde{s}_K} \mathbb{E}_{(\pi_\theta, r) \in B} [(V_\phi([\pi_\theta(\tilde{s}_1), \dots, \pi_\theta(\tilde{s}_K)]) - r)^2] \end{aligned}$$

## MAIN RESULTS

- Comparison with DDPG (Lillicrap et al., 2015) and ARS (Mania et al., 2018) using deep deterministic policies (2 hidden layers, 256 neurons per layer)



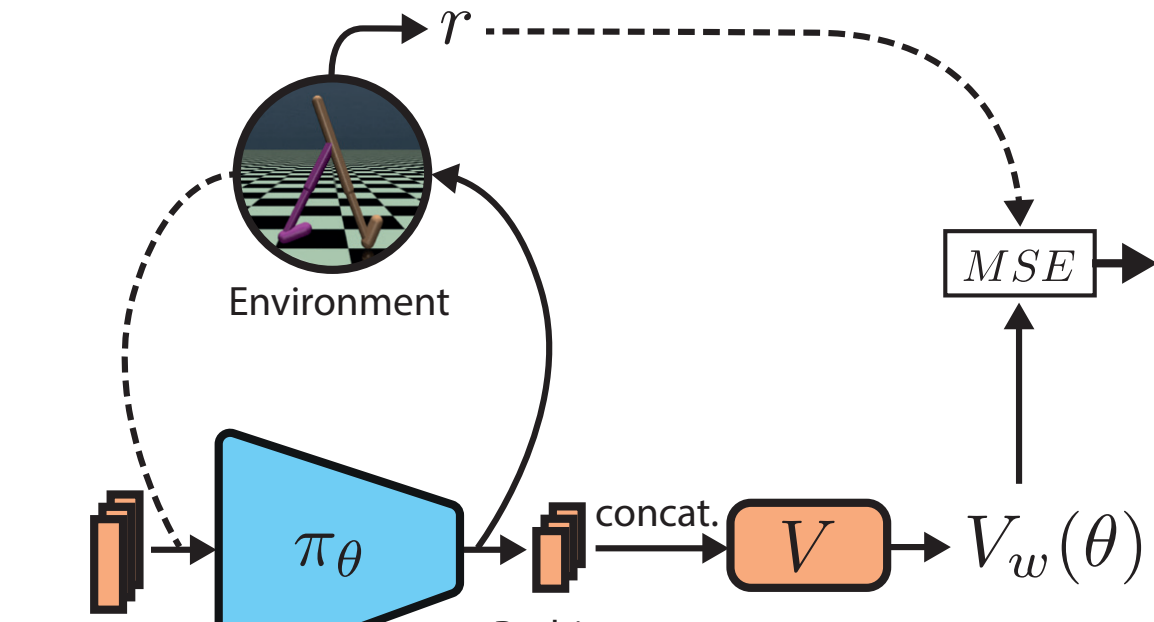
## REFERENCES

- F. Faccio, L. Kirsch, and J. Schmidhuber. Parameter-based value functions. In *9th International Conference on Learning Representations, ICLR*, 2021.
- J. Harb, T. Schaul, D. Precup, and P.-L. Bacon. Policy evaluation networks. *arXiv preprint arXiv:2002.11833*, 2020.
- T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- H. Mania, A. Guy, and B. Recht. Simple random search of static linear policies is competitive for reinforcement learning. In *Advances in Neural Information Processing Systems, NeurIPS*, 2018.
- R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Proceedings of the 12th International Conference on Neural Information Processing Systems*, 1999.

## ACTOR CRITIC ALGORITHM

### Off-policy actor-critic with PSSVF

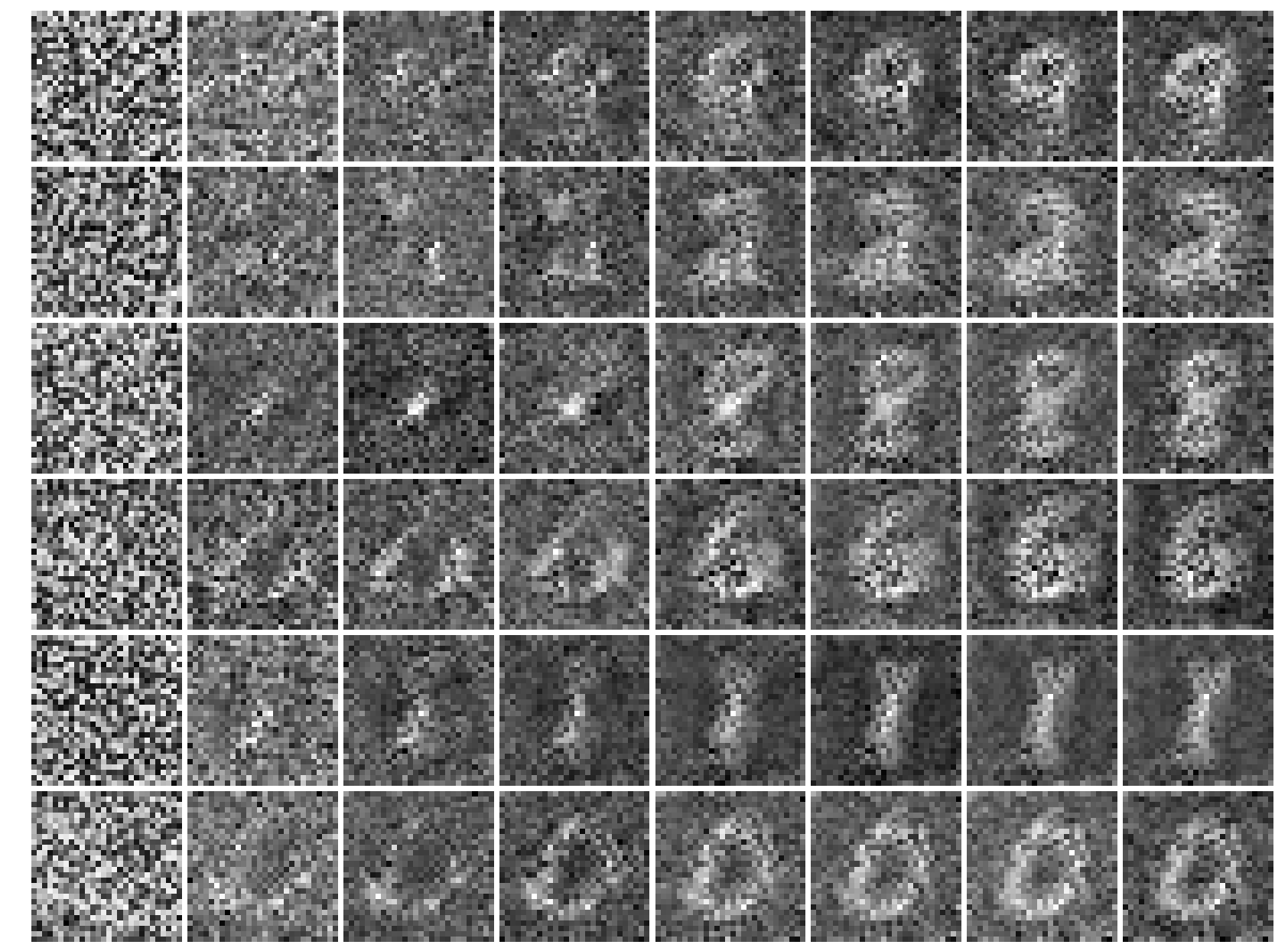
Given the behavioral  $\pi_b$ , find  $\pi_\theta$  maximizing  $J(\theta)$ :



1. Collect data with  $\pi_b$  (expensive in RL)
2. Use data to train  $V(\theta)$
3. Improve  $\pi_\theta$  following  $\nabla_\theta J(\pi_\theta)$  (offline optimization)
4. Set new behavioral  $\pi_\theta \leftarrow \pi_b$
5. Repeat until convergence

## DEMONSTRATION ON MNIST

- We apply our algorithm to MNIST classification. The return is the negative supervised loss. We obtain 87% accuracy. Learned probing states are digits

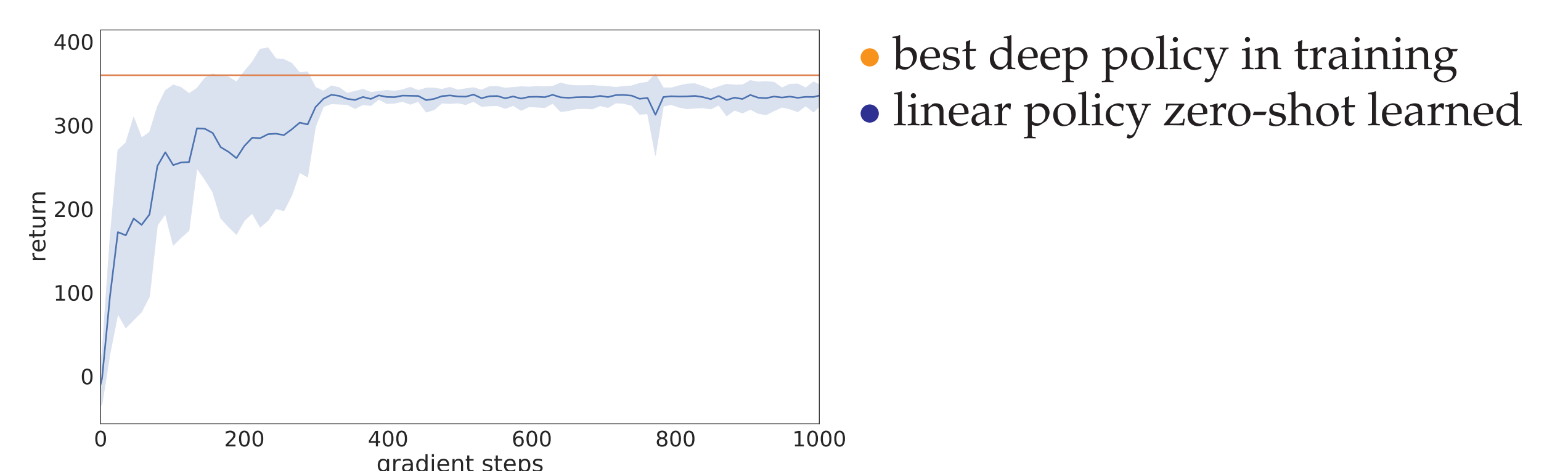


### Offline policy improvement:

- With an offline dataset  $\{\pi_{\theta_i}, l_i\}_{i=1}^N$  of randomly initialized CNNs and their losses (maximum accuracy 12%), we train  $V_w$  to evaluate such policies
- We randomly initialize a new CNN and take many steps of gradient ascent through the fixed value function, obtaining a final CNN whose accuracy is around 65% on the test set

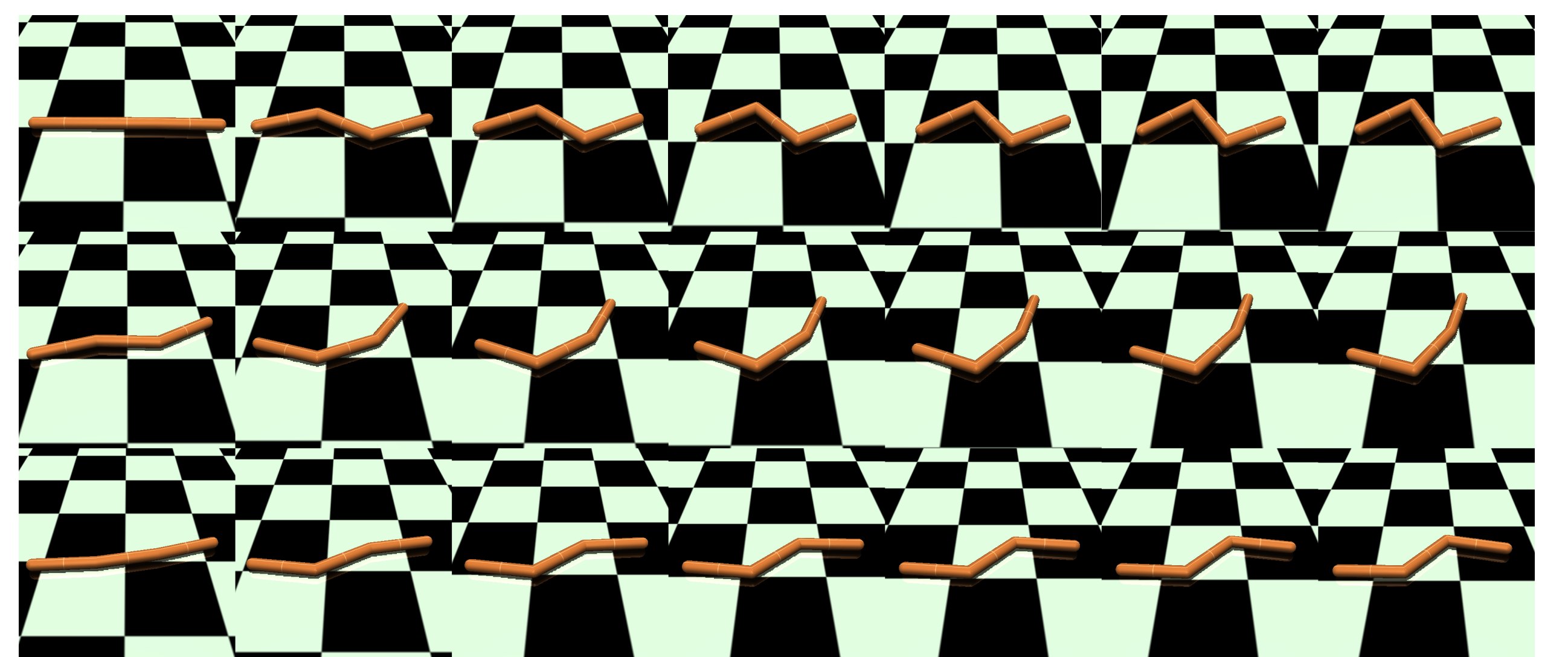
## MORE RESULTS

- A PSSVF trained using deep deterministic policies zero-shot learns a linear policy with similar performance in Swimmer



- The method extracts crucial abstract knowledge about the environment in form of very few learned abstract states sufficient to fully specify the behavior of many policies

- A randomly initialized policy can learn optimal behaviors in Swimmer by knowing how to act only in the following 3 crucial learned states (similar results for Hopper with 5 states):



- Videos of learned probing states for the RL experiments:

