

# Sample-Efficient Automated Deep Reinforcement Learning

Jörg K.H. Franke<sup>1</sup>, Gregor Köhler<sup>2</sup>, André Biedenkapp<sup>1</sup>, Frank Hutter<sup>1,3</sup>



<sup>1</sup>University of Freiburg, Germany  
<sup>2</sup>German Cancer Research Center (DKFZ), Heidelberg, Germany  
<sup>3</sup>Bosch Center for Artificial Intelligence, Renningen, Germany

# Problem Setting

- We want to train an off-policy Reinforcement Learning Agent on a new task.
- We don't know the optimal hyperparameters.
- We assume, environment interactions are expensive or limited.

# Our Approach:

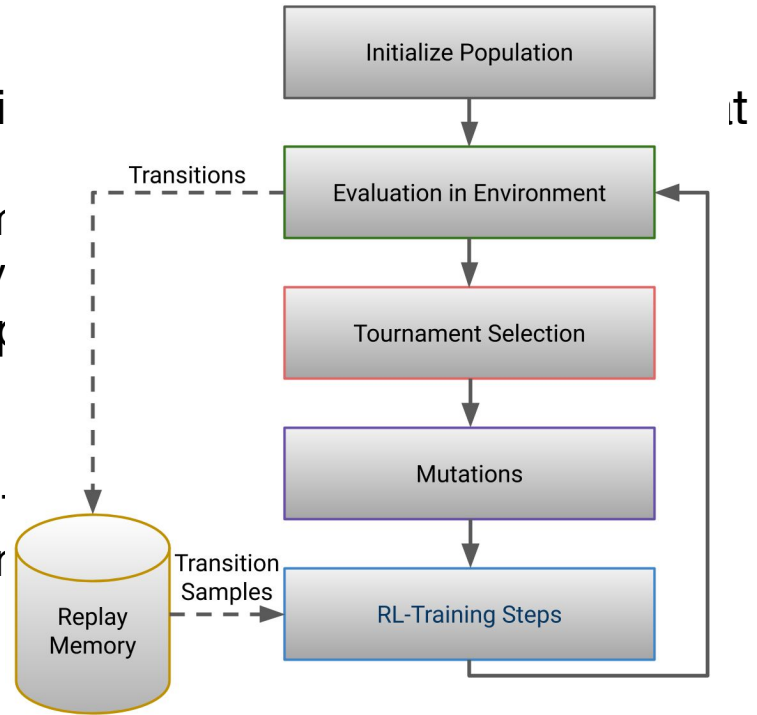
Sample-Efficient Automated Deep Reinforcement Learning (SEARL)

## Key idea:

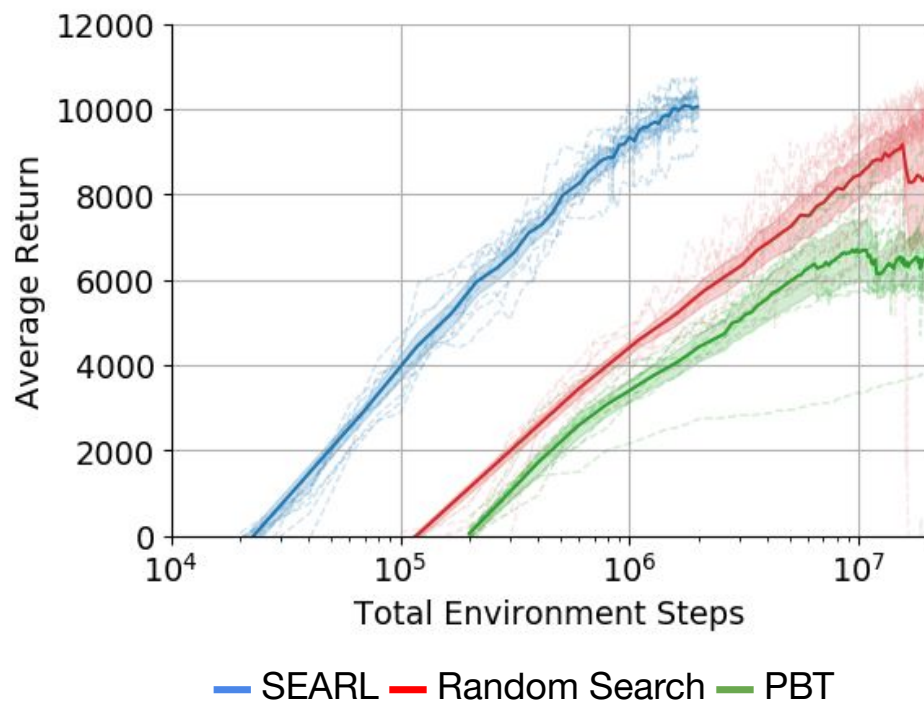
Use an evolutionary routine with shared replay memory to do RL training and population-based hyperparameter and neural architecture optimization.

# SEARL in a nutshell

- Neuroevolution for RL-training and online at the same time
  - Neuroevolution allows to find the optimal HP
  - Agent evaluations fill the shared replay memory
  - Agents are trained using the shared replay memory
- Automatically finds optimal HP during training
- Dynamic adaptation instead of fixed optimal HP
- Saves environment interactions

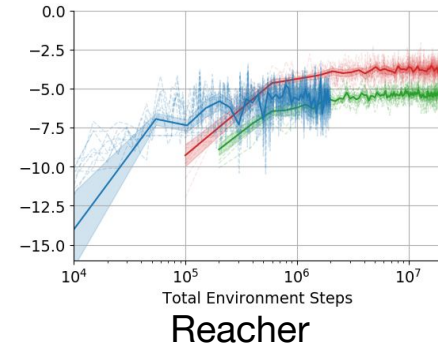
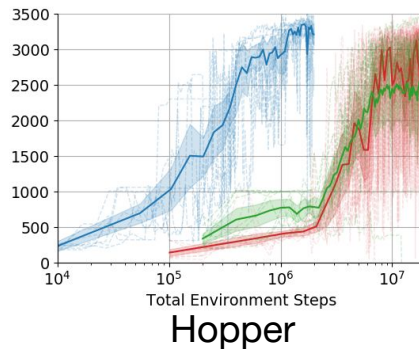
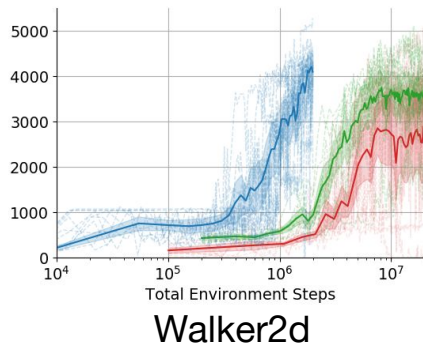
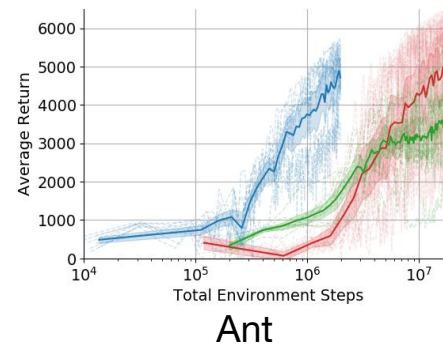
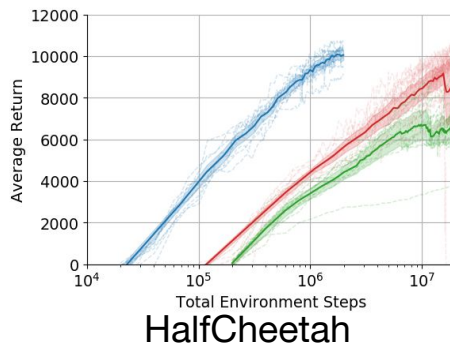


# Results - SEARL for TD3 in MuJoCo “HalfCheetah”



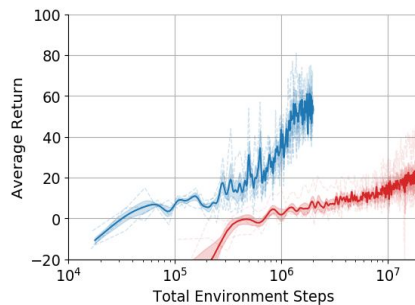
# Results - SEARL for TD3 with MuJoCo Suite

- SEARL
- Random Search
- PBT

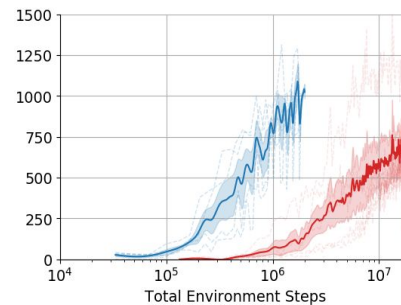


# Results - SEARL for DQN with Atari environments

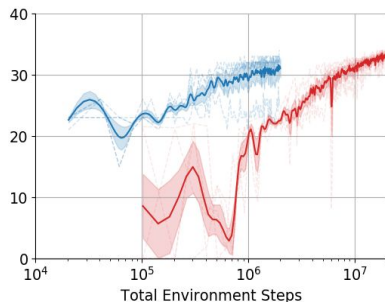
— SEARL  
— Random Search



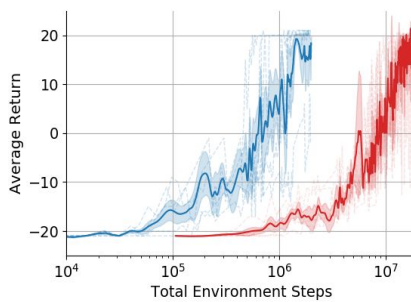
Boxing



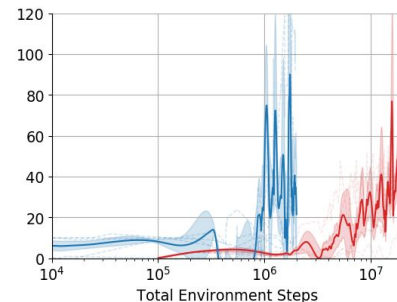
Enduro



Freeway



Pong



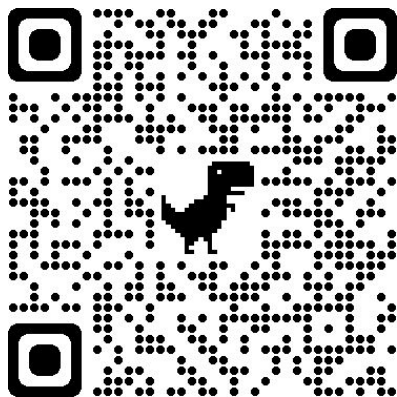
RoadRunner

# Conclusion

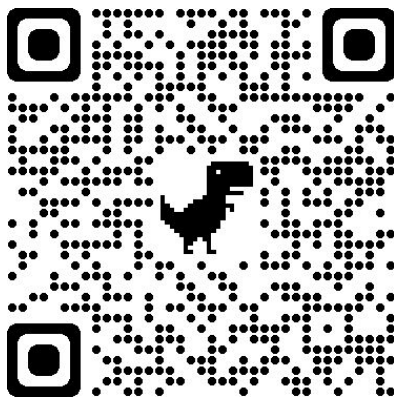
- SEARL **jointly trains** an RL off-policy agent and **optimizes** the **hyperparameters** including the neural network **architecture**.
- SEARL is based on **neuroevolution** and a **shared replay memory**.
- The population in SEARL benefits from a **diverse set of experience** in the shared replay memory.
- SEARL requires **up to ten times fewer environment interactions** as random search or Population Based Training (PBT) and nearly the same amount of environment steps as a normal RL agent training.



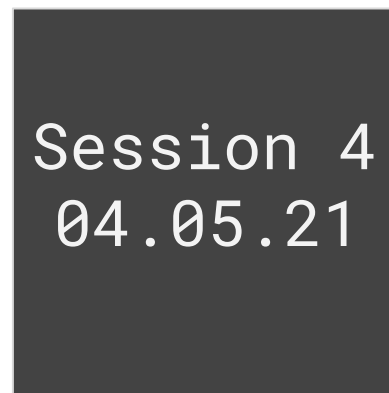
# Thanks for your attention!



Paper



Code



Poster