

Model-based Evolutionary Policy Search for Skill Learning in Continuous Domains

Jan Hendrik Metzen, University Bremen, Robotics Group, Robert-Hooke-Str. 5, 28359 Bremen, Germany

1 Abstract

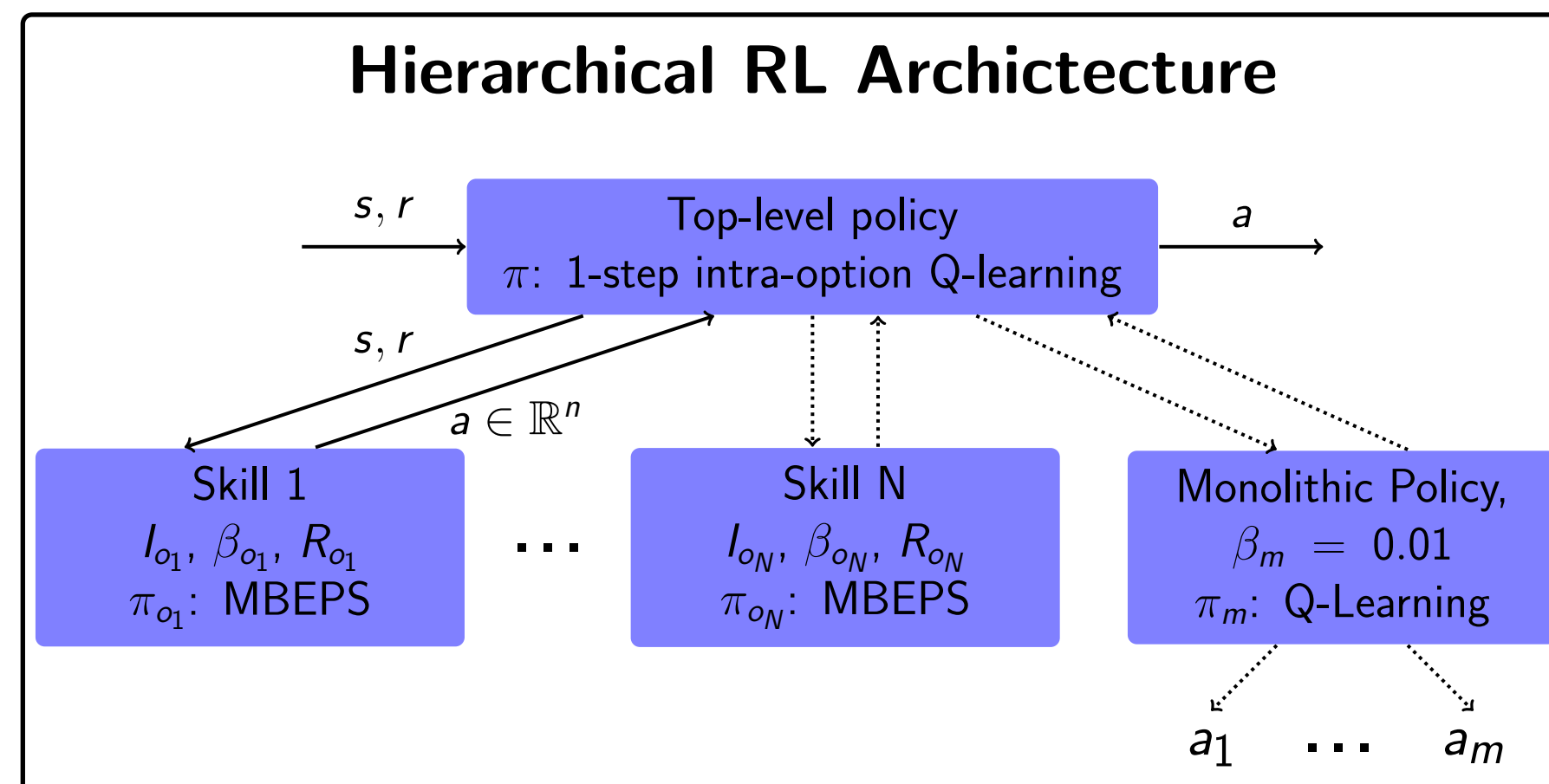
We investigate the utility of evolutionary policy search (EPS) for skill learning in a Hierarchical Reinforcement Learning architecture. While EPS is well suited for domains with **continuous state and action spaces**, it is susceptible to **non-stationarities** caused by concurrent learning of several skills and higher layers of the architecture. We hypothesize that (if **model-learning** is feasible) EPS should be used instead for **planning** based on trajectory sampling in a model which “smooths out” non-stationarities.

2 Evolutionary Policy Search (EPS)

- Parametrized policy representation $\pi(\theta)$
- Population-based approach where each individual encodes parameter vector θ
- **Expected return of $\pi(\theta)$:**
 $J(\theta) = E[R(s_0)|s_0 \sim S_0, P_{ss'}^a, R_{ss'}^a, \pi(\theta)]$
- Fitness function $f(\theta)$ is sample-based approximation of $J(\theta)$
- Evolution Strategy used to find θ which maximizes $f(\theta)$

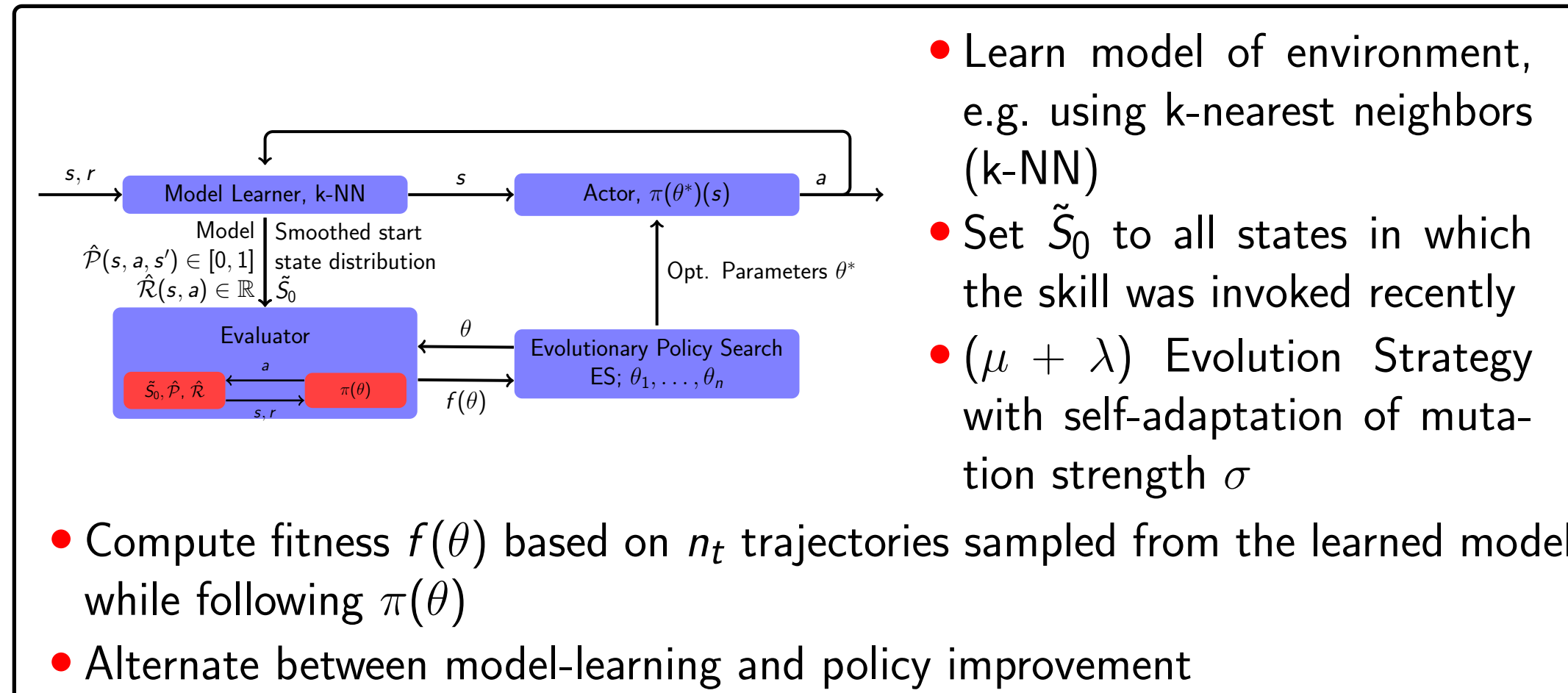
3 Using EPS for Skill Learning

- **Motivation:** EPS is well-suited for problems with continuous state and action spaces [1, 2, 4]
- **Challenge:** Structured problems may require sophisticated policies with high dimensional parameter vector
- **Approach:** Hierarchical RL with autonomous skill discovery allows to split complex problems into simpler subproblems. For these subproblems, simple (e.g. linear) policies with a small number of parameters may suffice

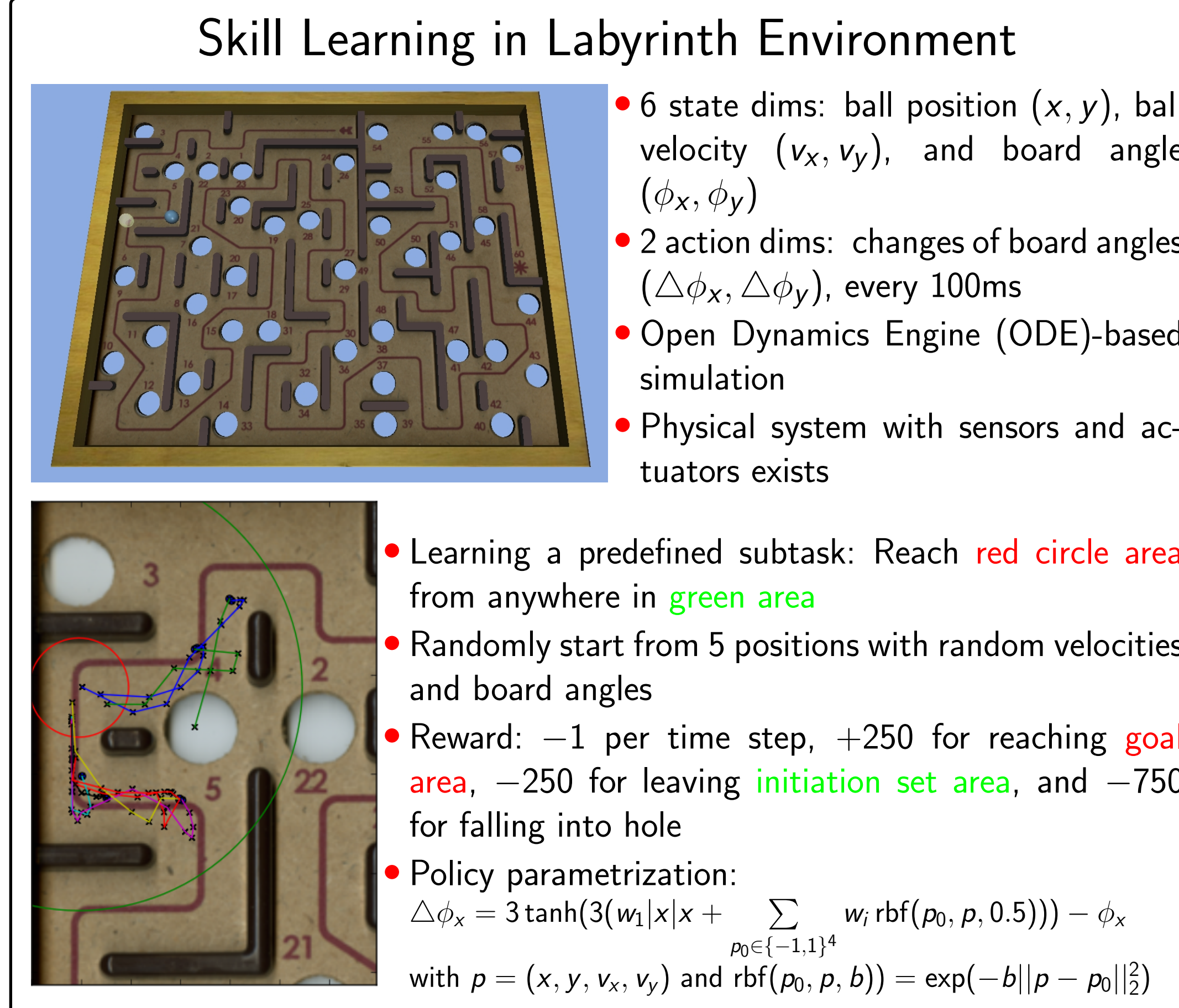
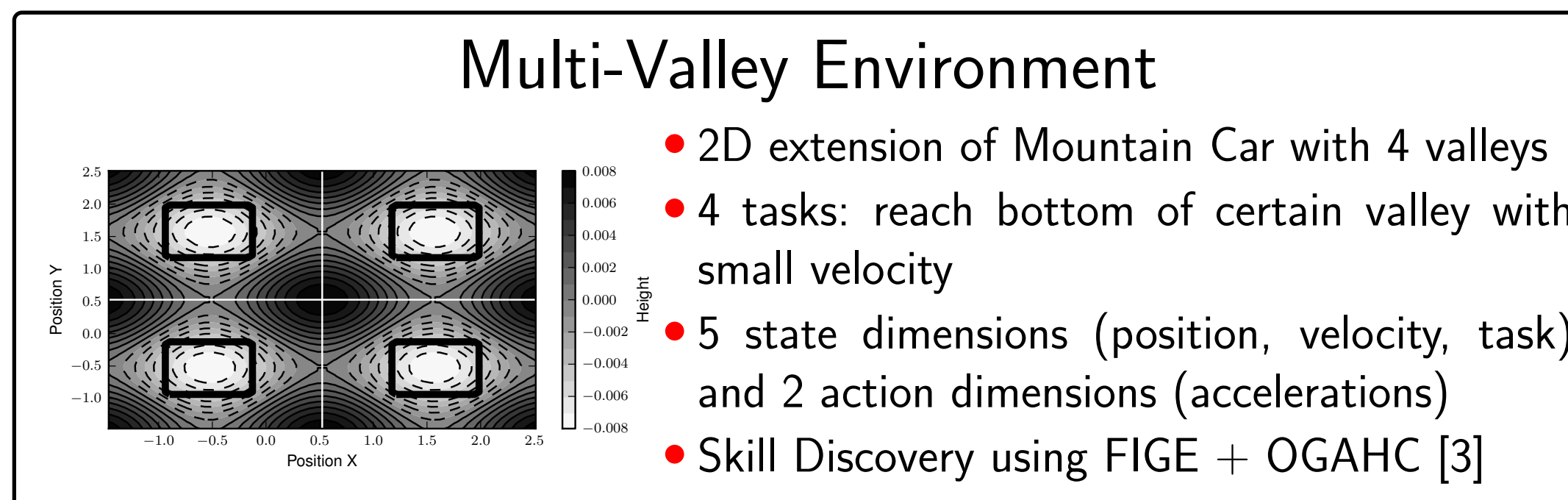


- **Challenge:** Concurrent learning of several skills and higher layers of the architecture makes option's start state distribution S_0 non-stationary
- **Approach:** Derive policy by planning (trajectory sampling) in learned model of environment in option's initiation set. Use EPS for planning in model; the start state distribution can be kept stationary during one planning iteration (i.e. one generation)

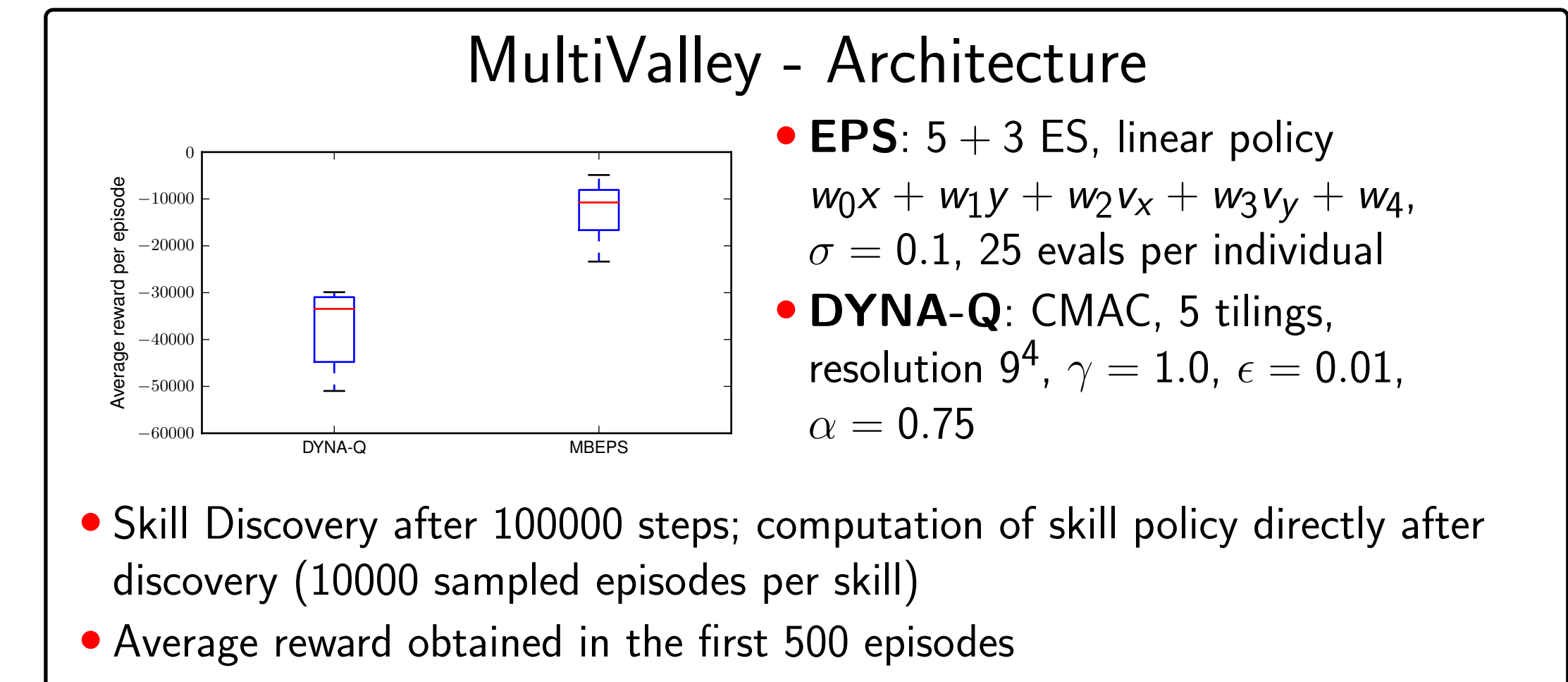
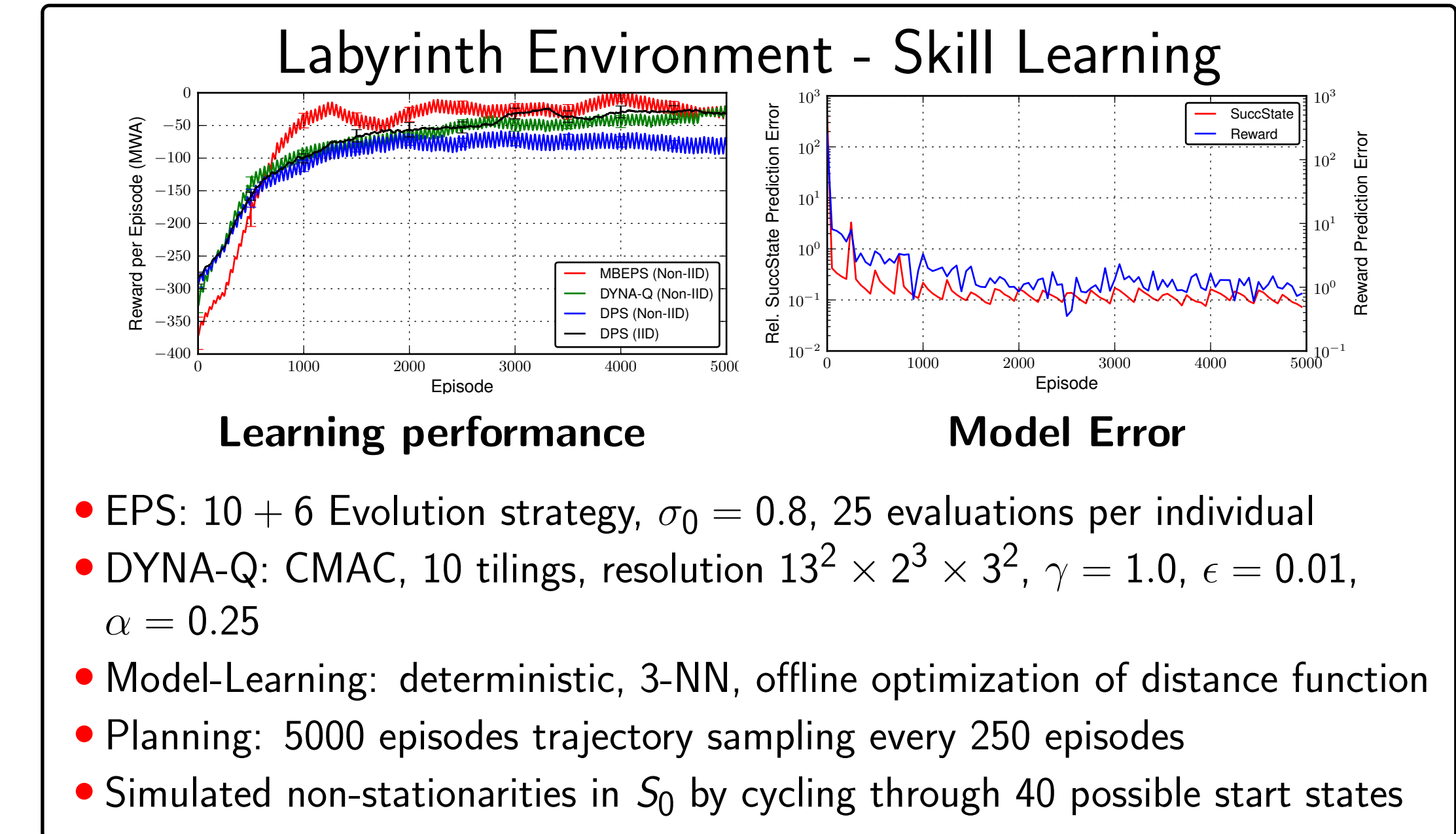
4 Model-based EPS



5 Scenarios



6 Results



7 Outlook

- Comparison with other policy search approaches
- Using other model learning approaches; systematic exploration (RMax-like)
- Integrate Skill Discovery, Skill Learning, and Compositional Learning and evaluate in entire Labyrinth Environment

8 References

- [1] Verena Heidrich-Meisner and Christian Igel. Variable metric reinforcement learning methods applied to the noisy mountain car problem. In *8th European Workshop on Reinforcement Learning (EWRL 2008)*, pages 136–150. Springer-Verlag, 2008.
- [2] Shivaram Kalyan Krishnan and Peter Stone. An empirical analysis of value function-based and policy search reinforcement learning. In *The Eighth International Conference on Autonomous Agents and Multiagent Systems (AAMAS)*, May 2009.
- [3] Jan Hendrik Metzen. Online skill discovery using graph-based clustering. In *10th European Workshop on Reinforcement Learning (EWRL 2012)*, 2012.
- [4] Julian Togelius, Tom Schaul, Daan Wierstra, Christian Igel, Faustino Gomez, and Jürgen Schmidhuber. Ontogenetic and phylogenetic reinforcement learning. *Künstliche Intelligenz*, (3/09):30–33, September 2009.