

Branch Dueling Deep Q-Networks for Robotics Applications

Master's Thesis

Baris Yazici

Msc. Mahmoud Akl

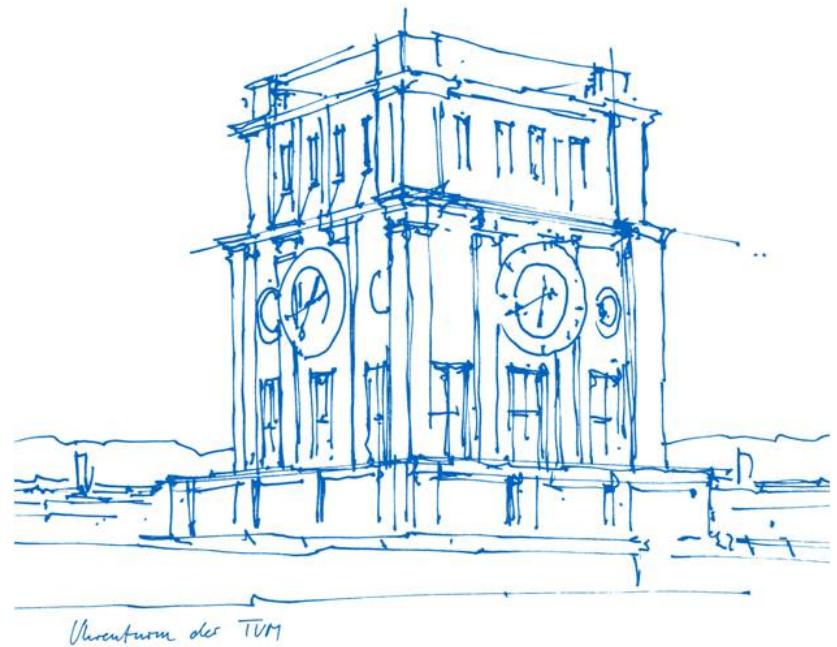
Technical University of Munich

Computer Science Faculty

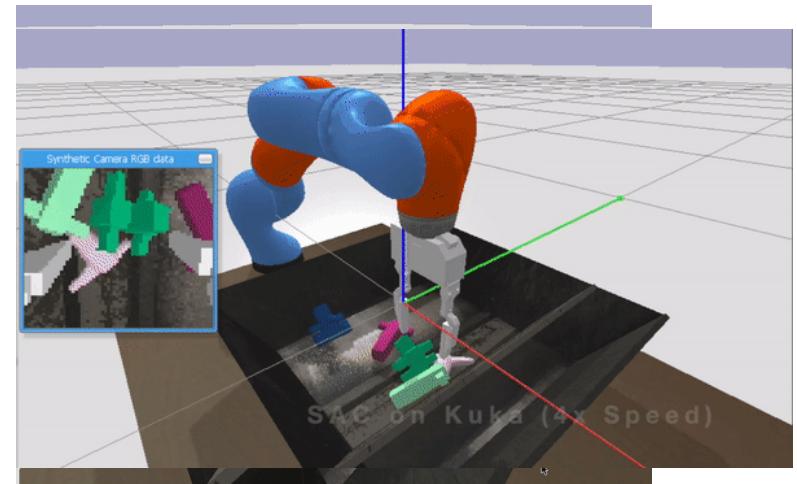
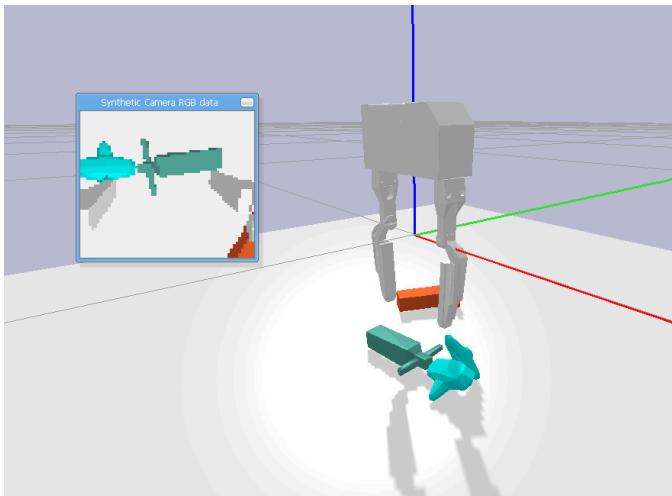
Chair for Robotics, Artificial Intelligence

and Realtime Systems

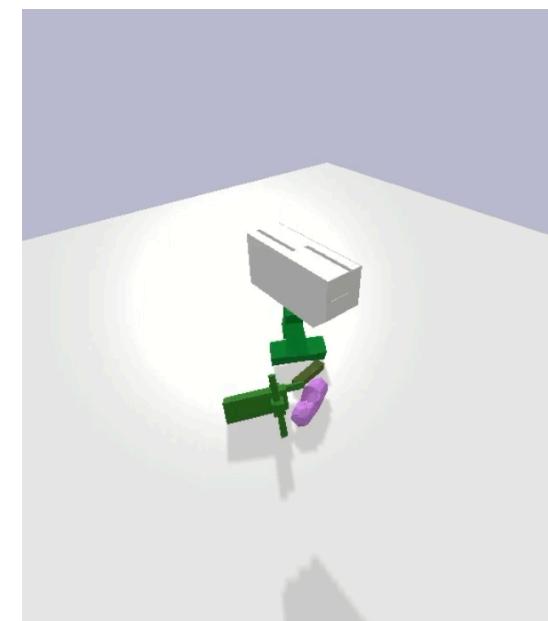
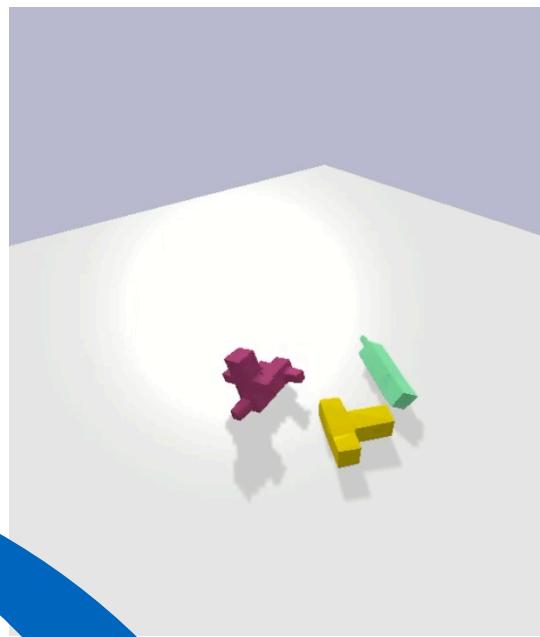
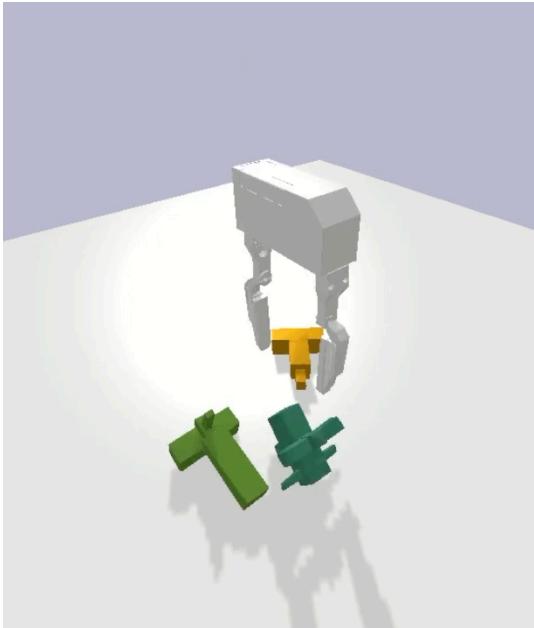
Munich, 05.10.2020



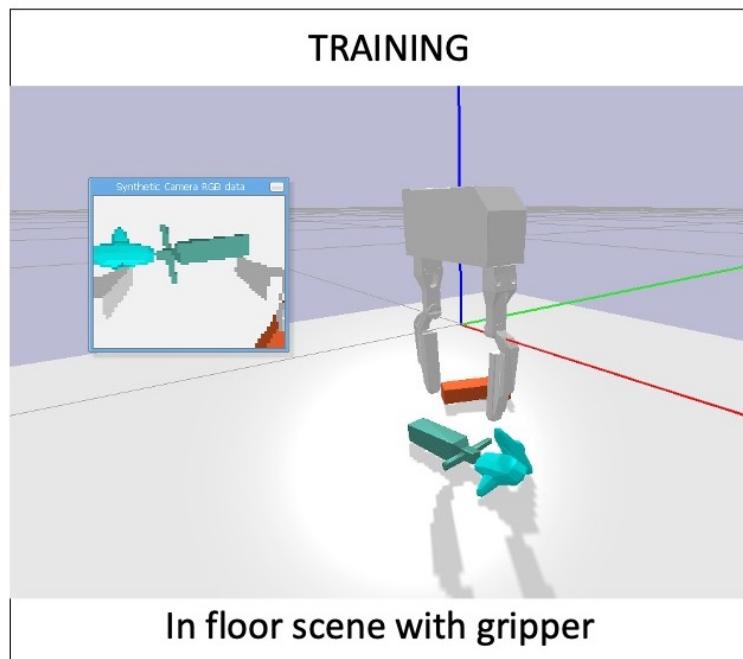
Adaptive Grasping with Domain Transfer



Train one grasping model for different **scenes** and **environments**

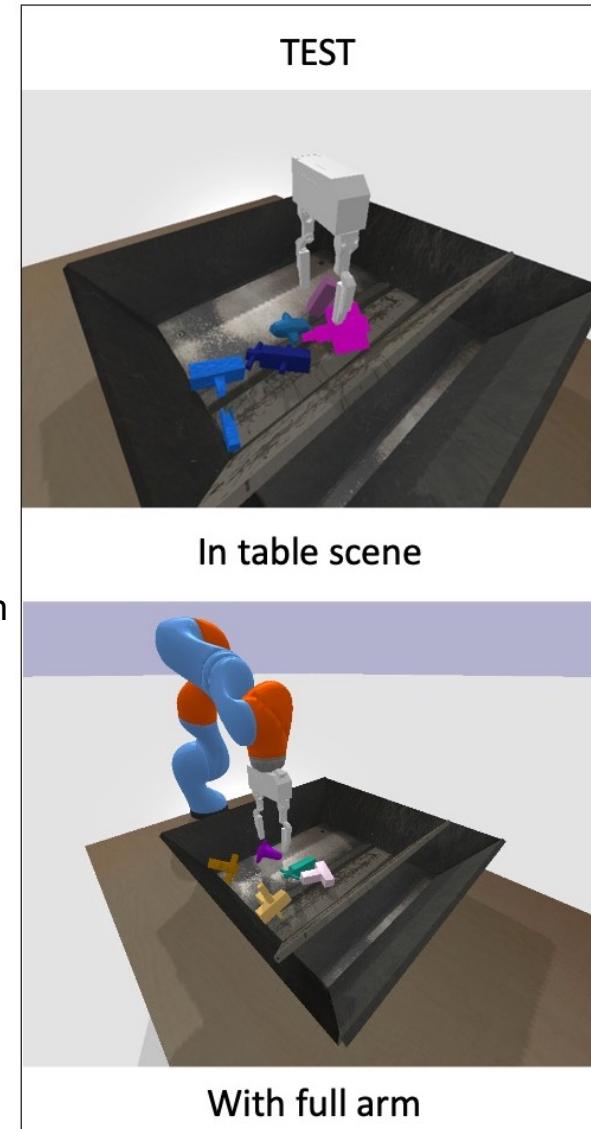


Train-Test Pipeline

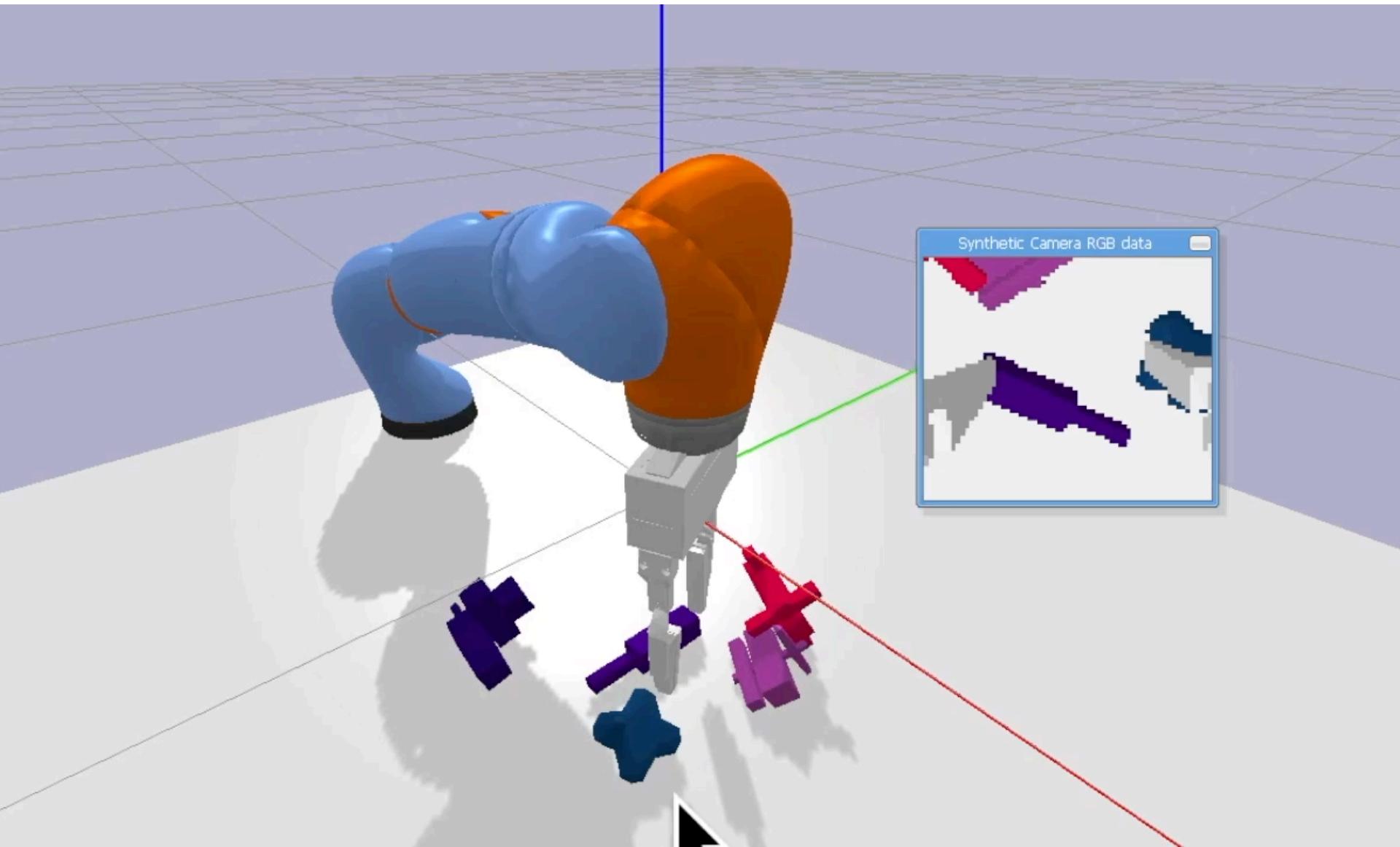


New scene
New Domain

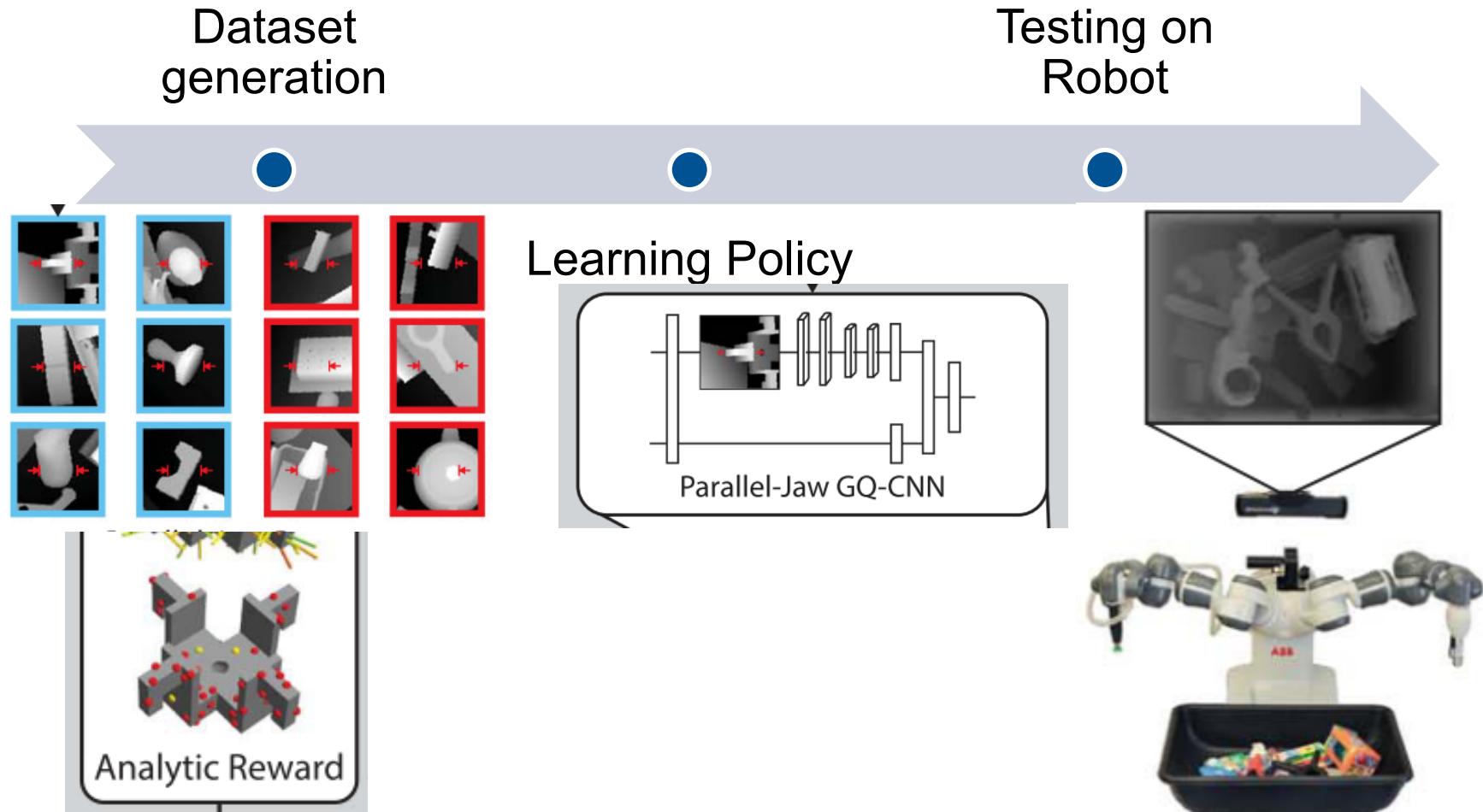
Two black arrows originate from the bottom right of the TRAINING section. The top arrow points towards the top part of the TEST section, labeled "New scene". The bottom arrow points towards the bottom part of the TEST section, labeled "New Domain".



Why Adaptive Robust Grasping?

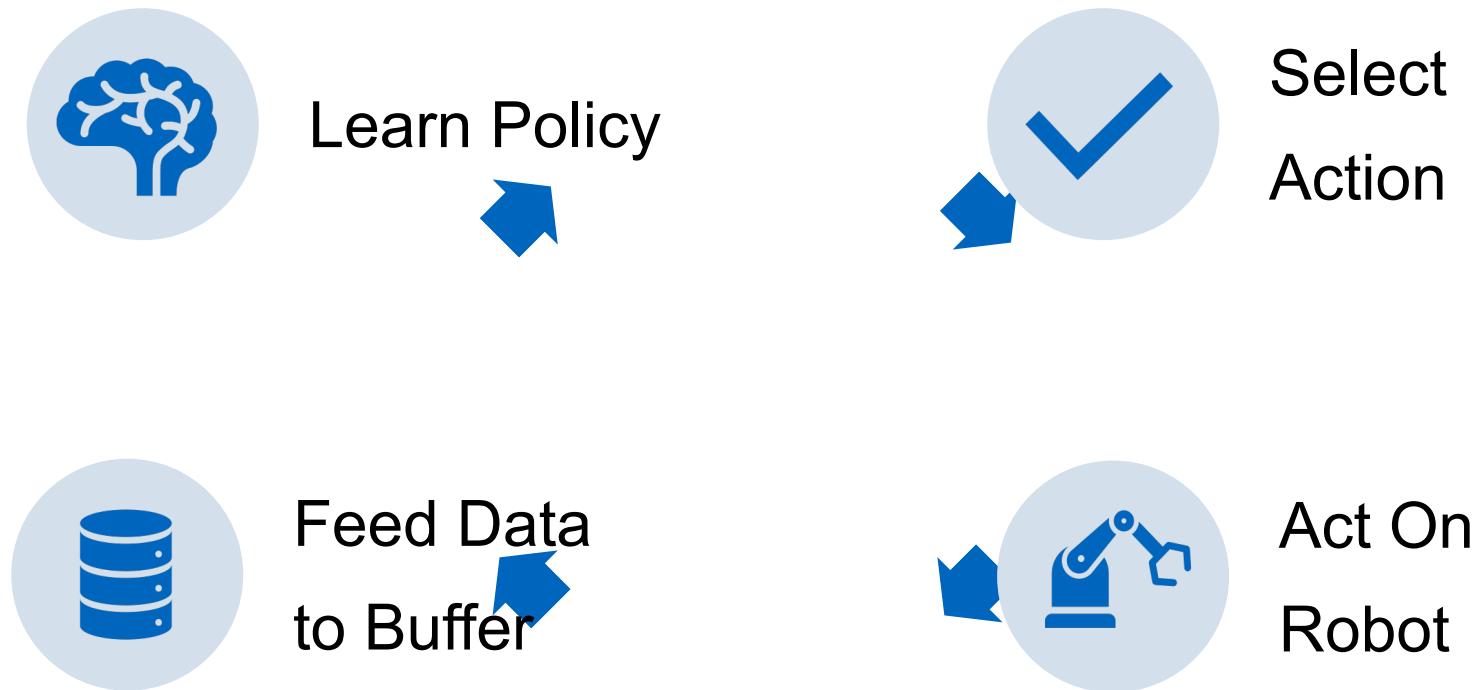


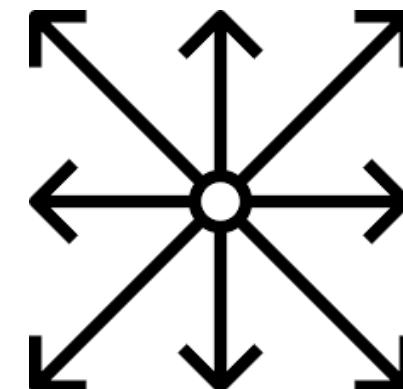
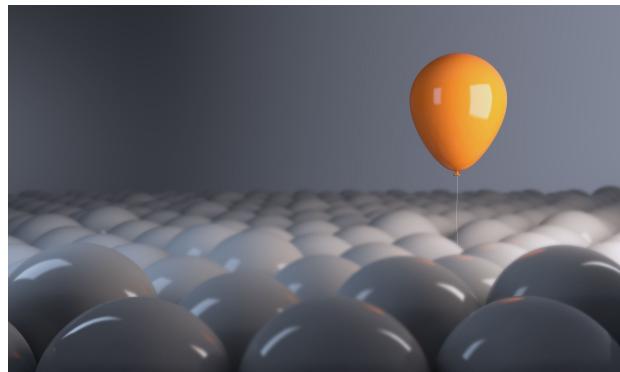
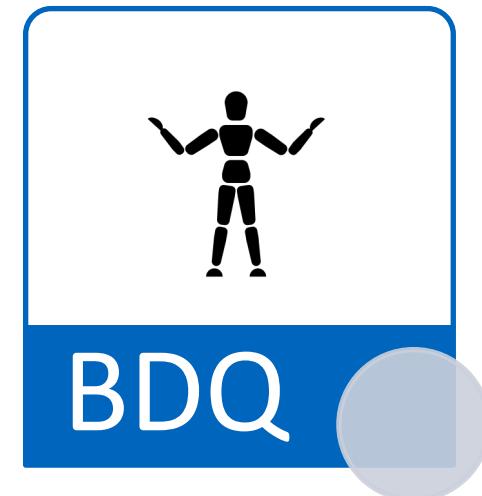
Current State of Art

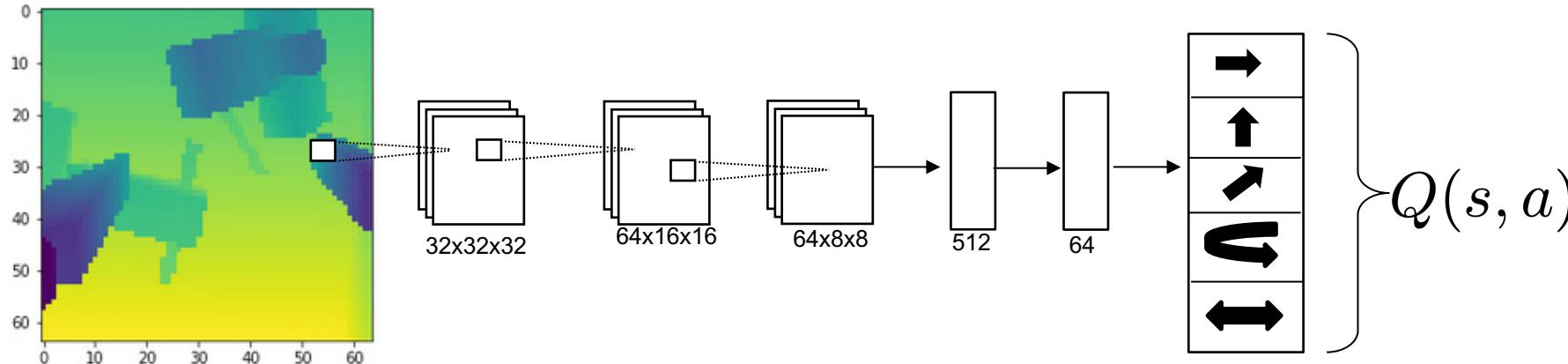
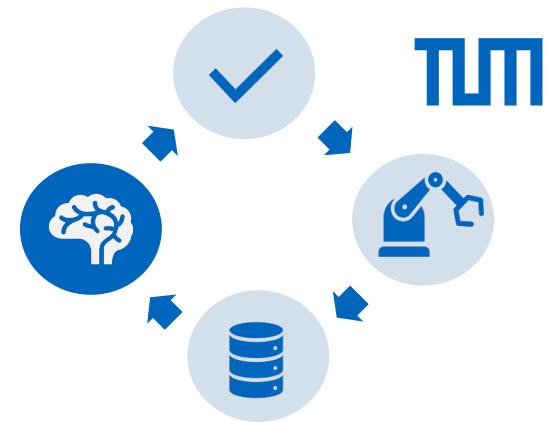


Mahler et al., 2019

Our Approach







$$L(\theta) = E_{(s, a, r, s') \sim U(D)} \left[Q_t(s', a; \theta^-) - Q(s, a; \theta) \right]^2$$

$$Q^*(s, a) = \mathbb{E}_{s' \sim \mathcal{E}} \left[r + \gamma \max_{a'} Q^*(s', a') \middle| s, a \right]$$



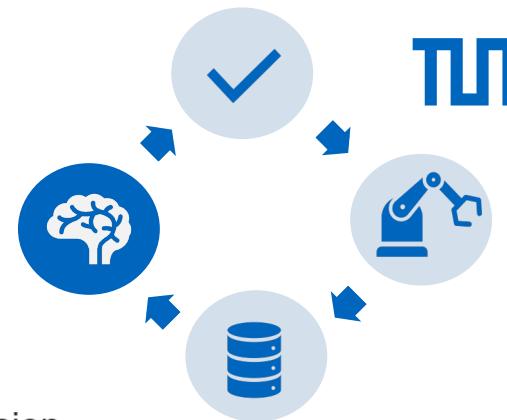
SAC



DQN

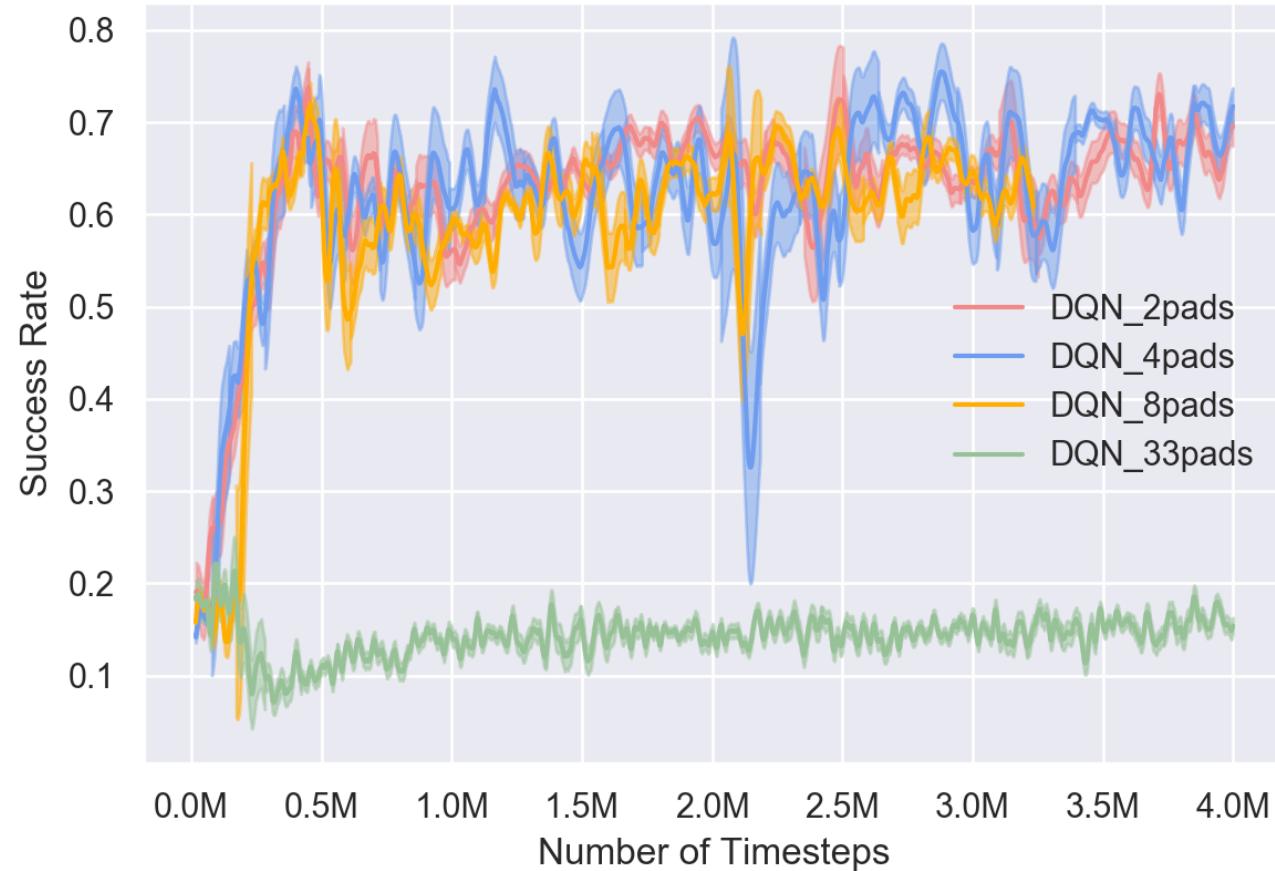


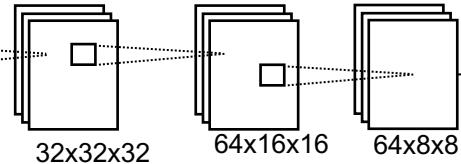
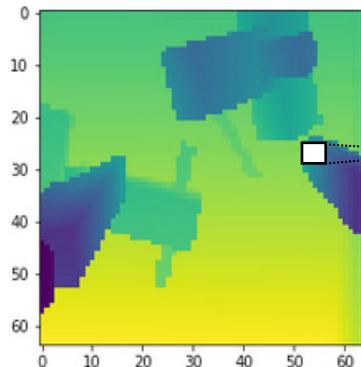
BDQ



TUM

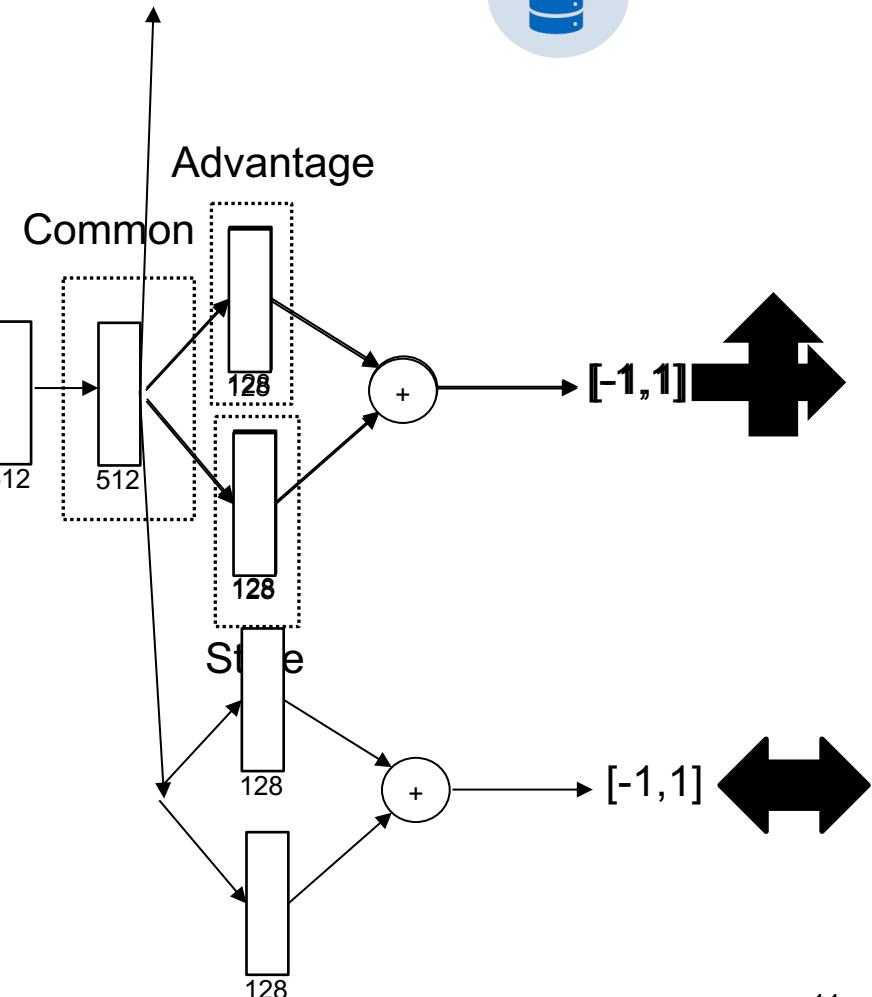
DQN simplified env. increasing action dimension

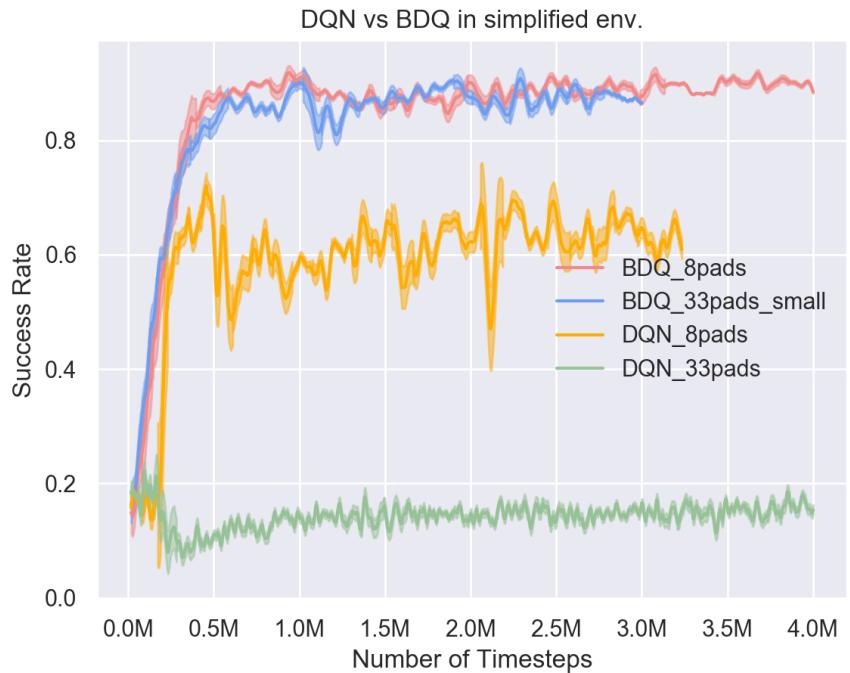
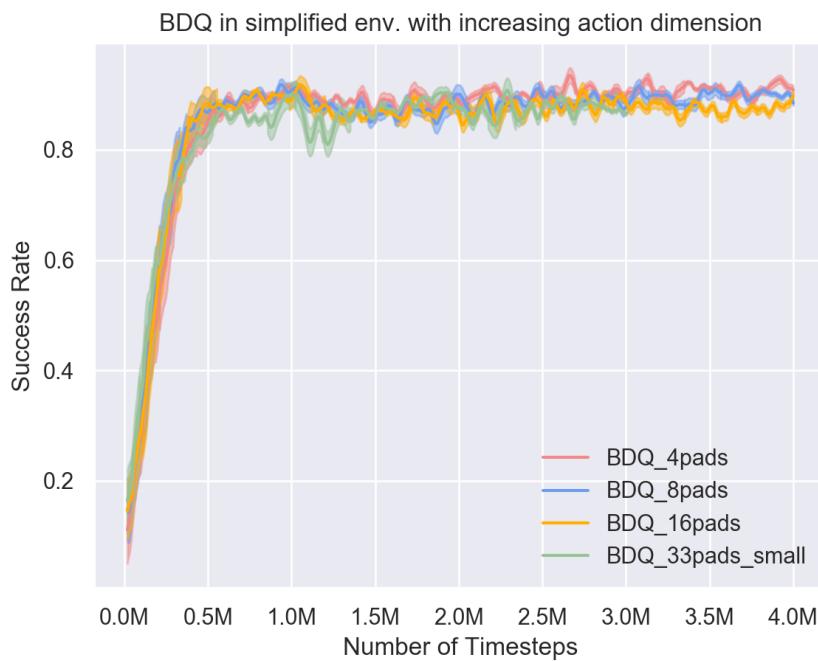
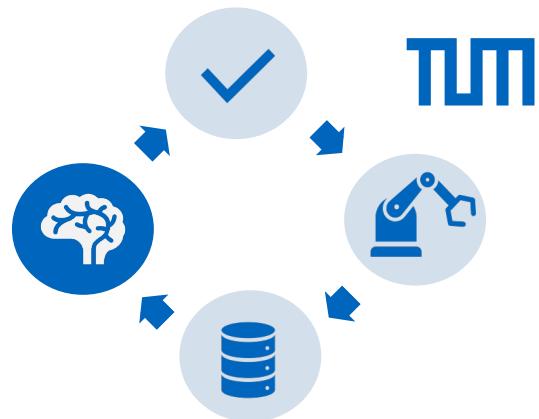




$$L = \mathbb{E}_{(s, a, r, s') \sim D} \left[\frac{1}{N} \sum_d (y_d - Q_d(s, a_d))^2 \right]$$

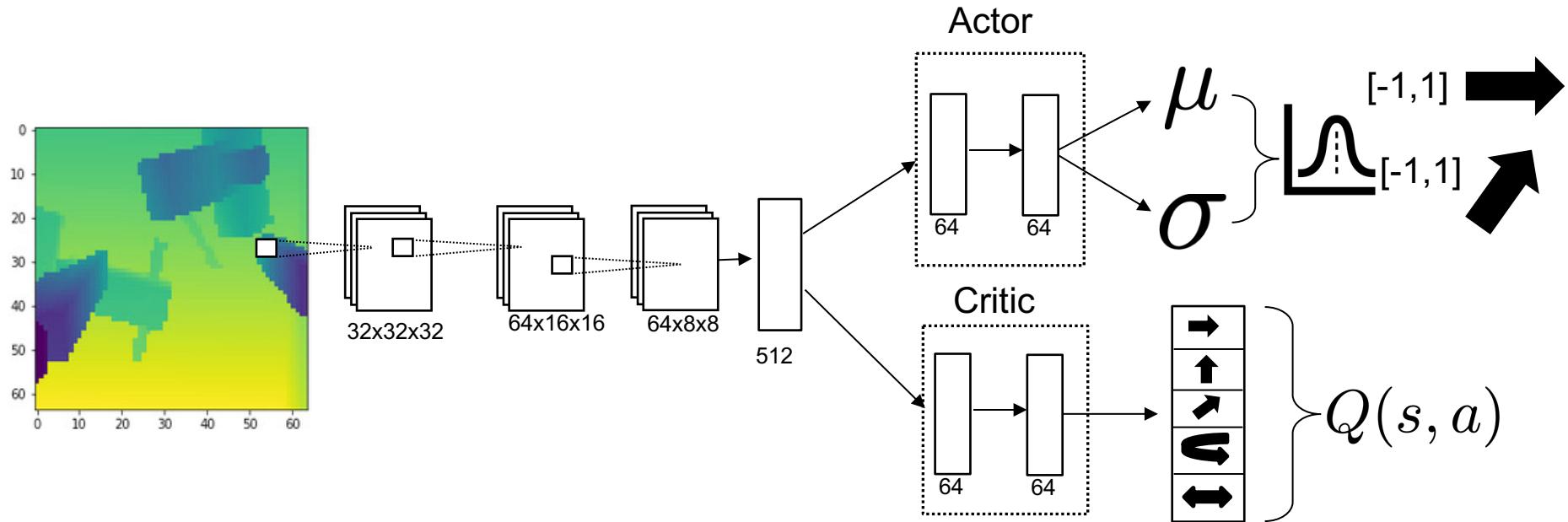
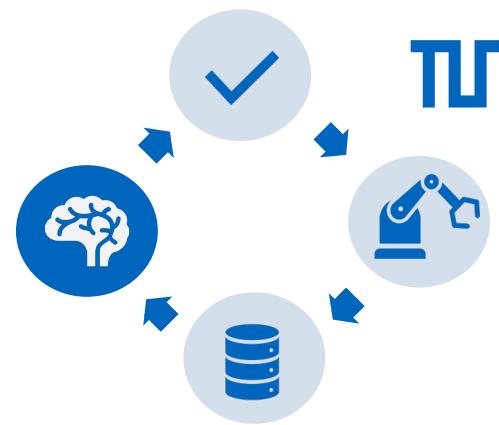
$$y = r + \gamma \frac{1}{N} \sum_d Q_d^-(s', \arg \max_{a'_d \in A_d} Q_d(s', a'_d))$$







TUM



$$J(\pi) = \sum^T E_{(s_t, a_t)}[r(s_t, a_t) + \alpha H(\pi(\cdot | s_t)]$$

$$V^\pi(s) = {}^t E_{a \sim \pi}[Q^\pi(s, a) - \alpha \log \pi(a | s)]$$



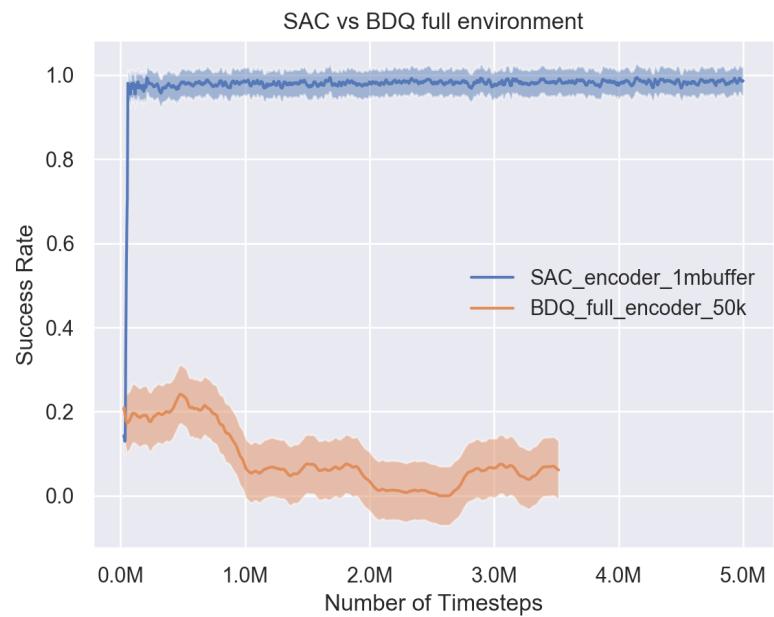
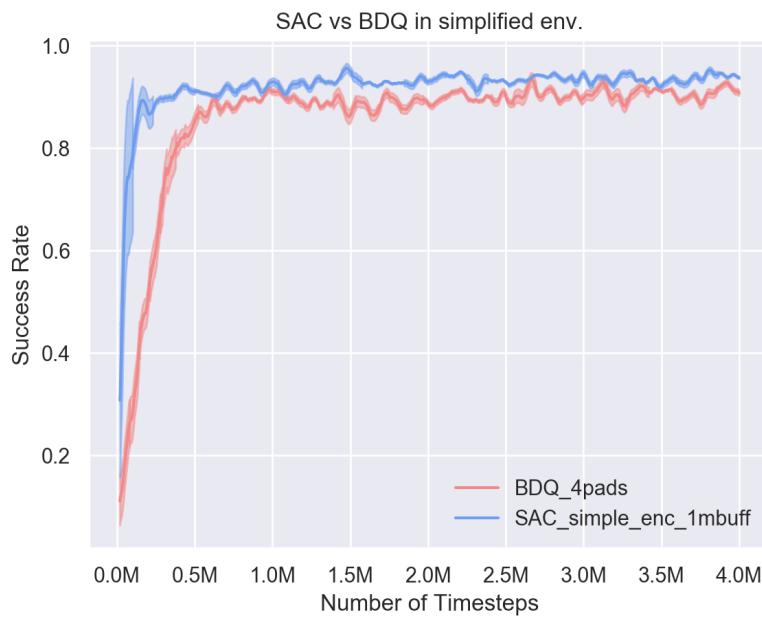
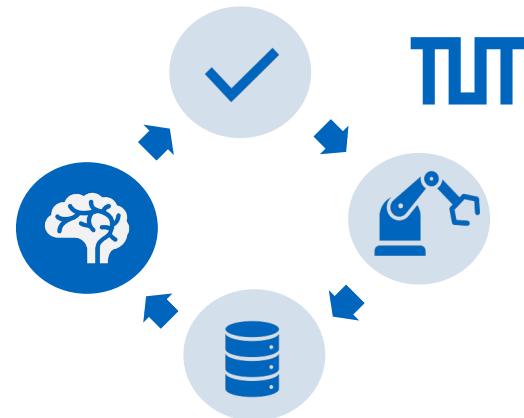
SAC



DQN

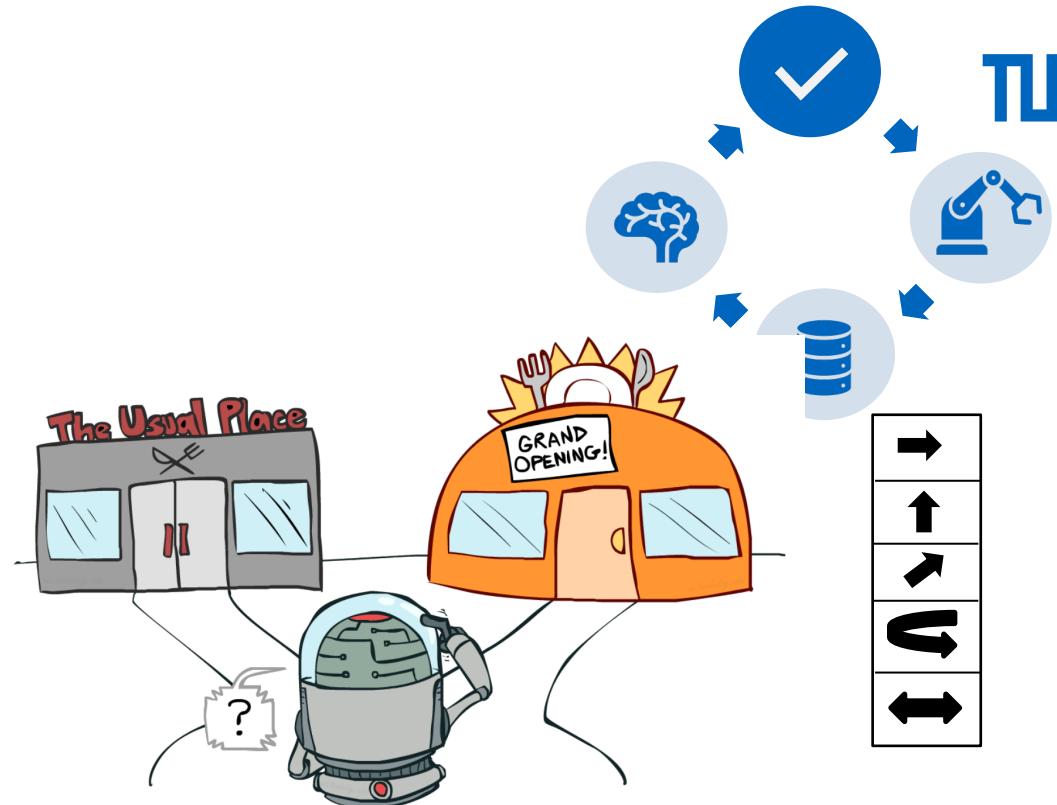


BDQ

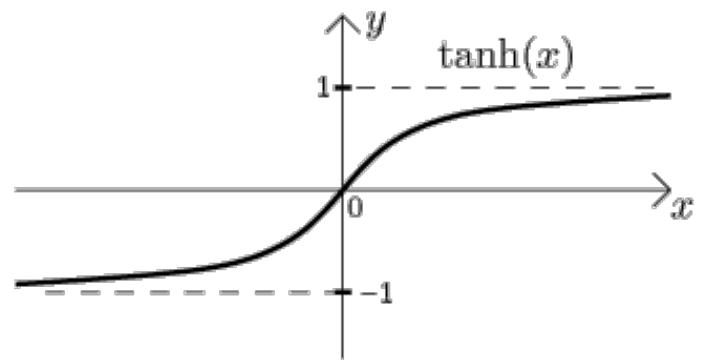
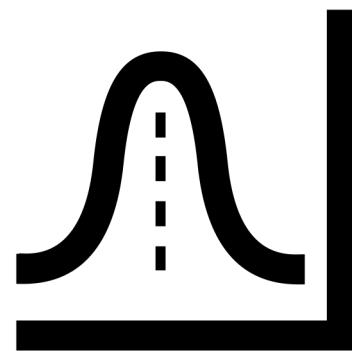


Action Selection

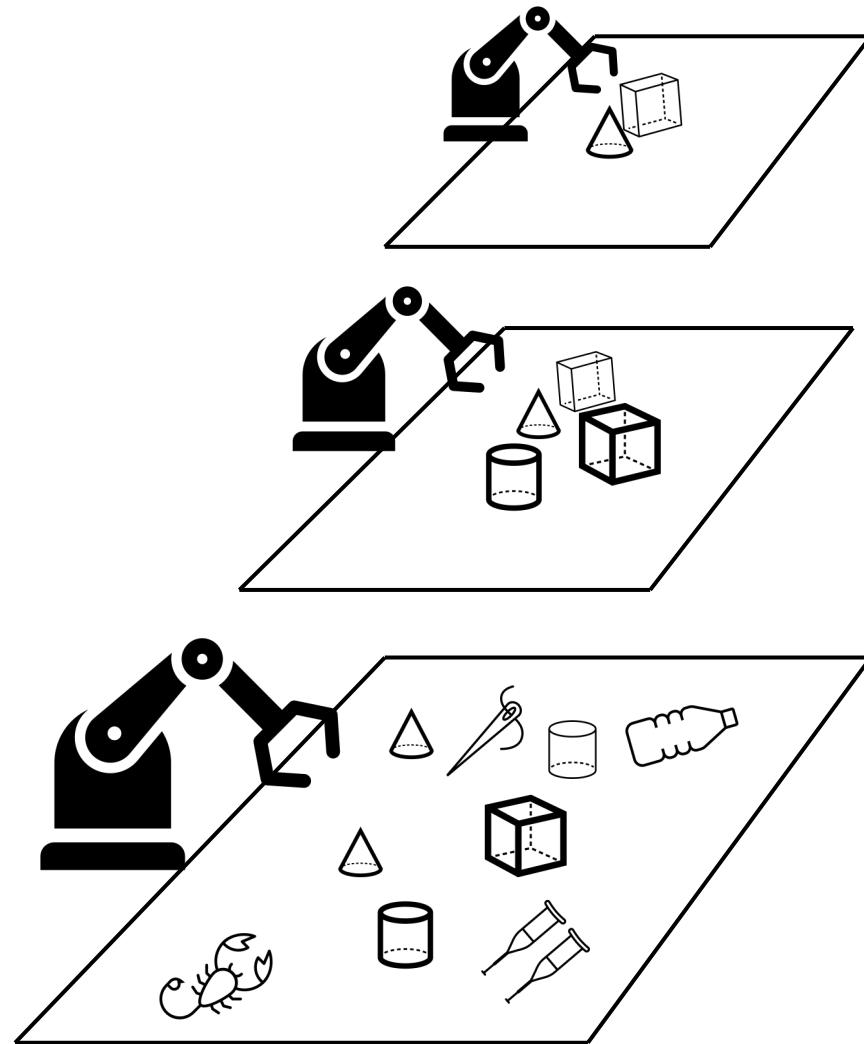
$(1 - \epsilon) \rightarrow \text{argmax} Q(s, a)$
 $\epsilon \rightarrow \text{random action}$



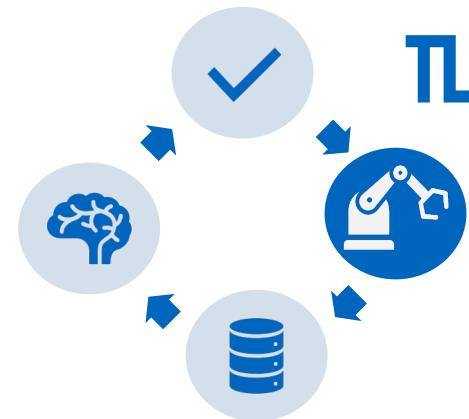
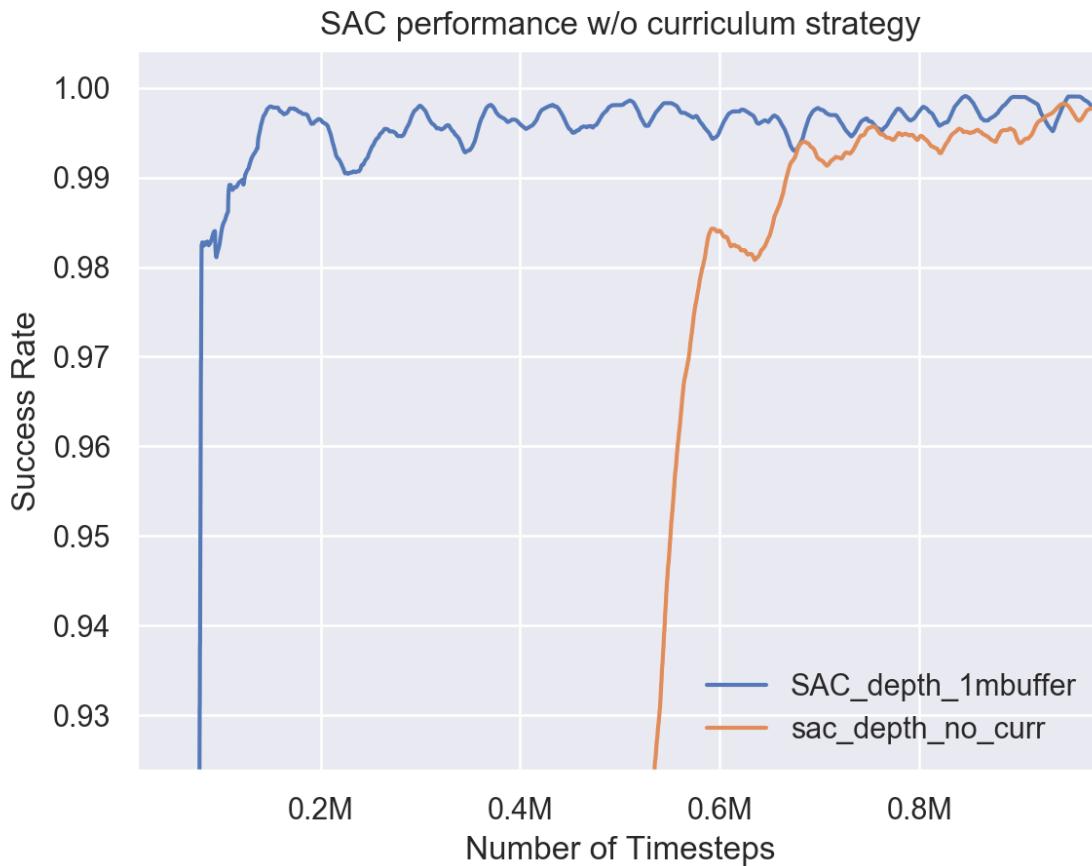
$$a_\theta(s) = \mu_\theta(s) + \sigma_\theta(s)$$



Curriculum Strategy



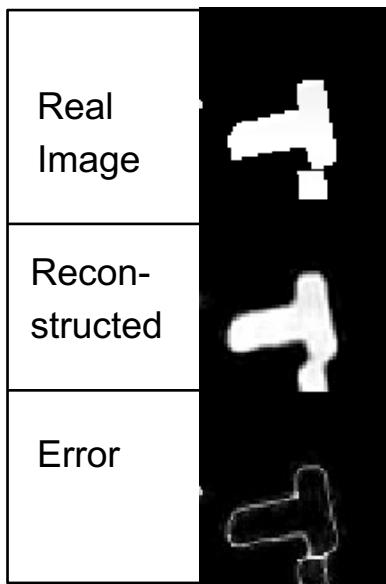
Curriculum Strategy



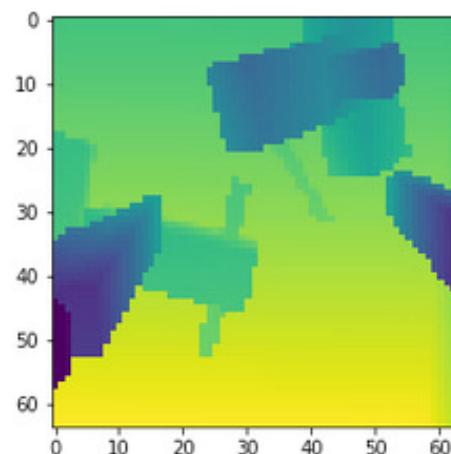
Observation



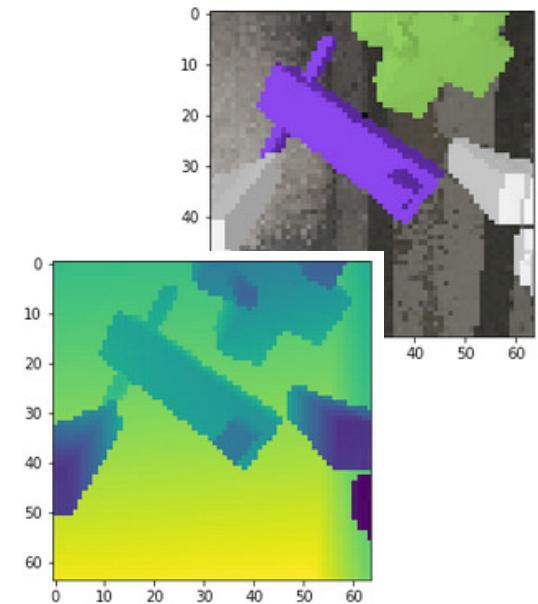
Auto-encoder



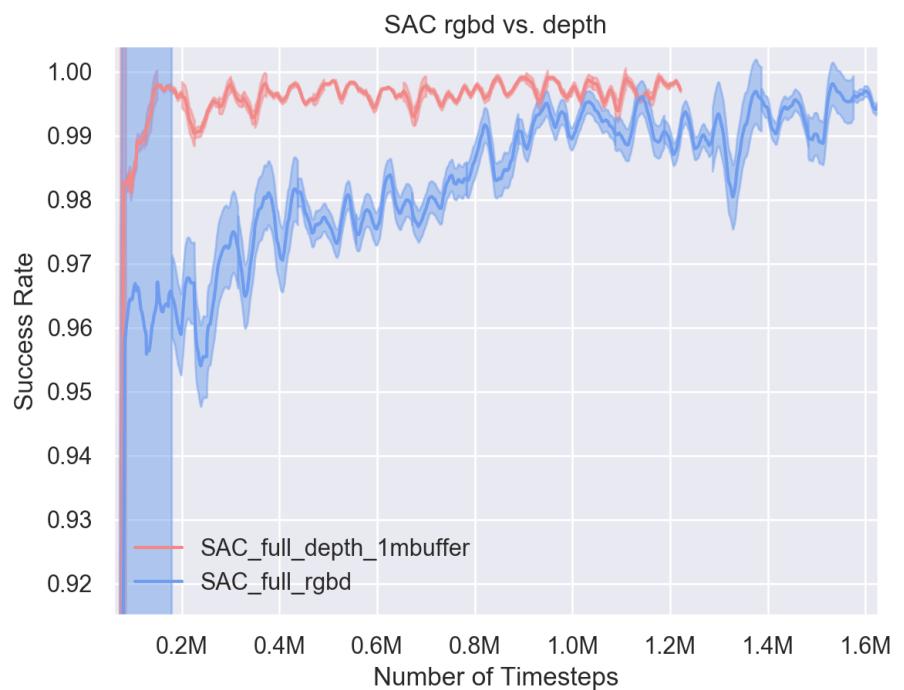
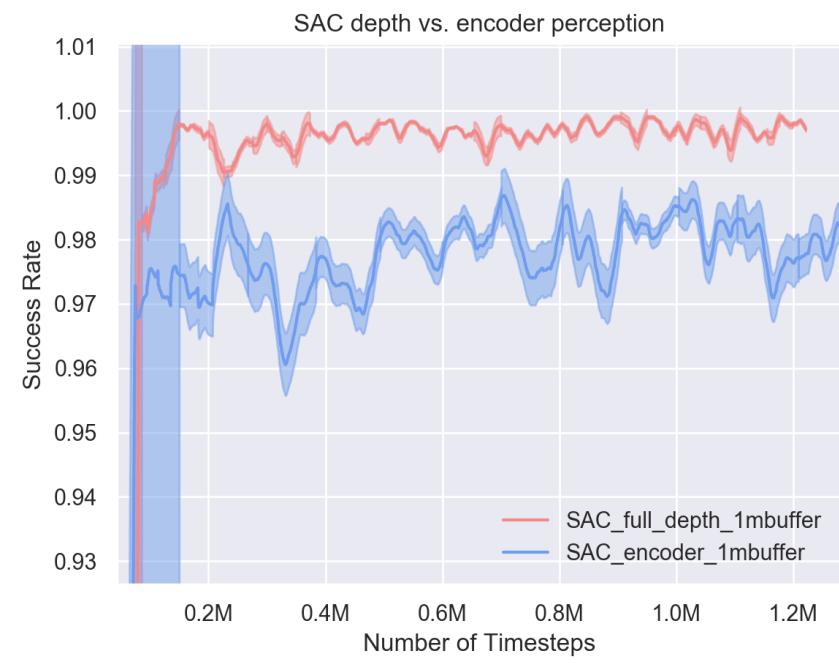
Depth



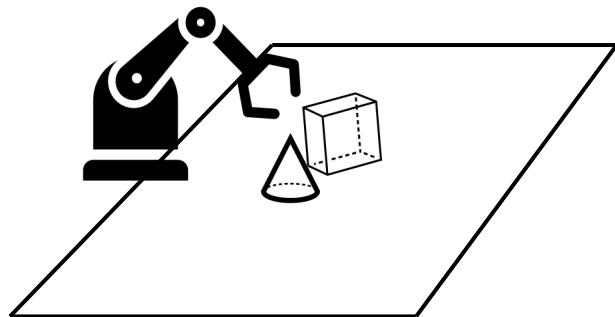
RGB-D



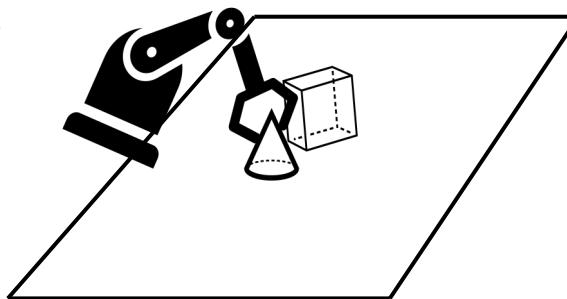
Observation



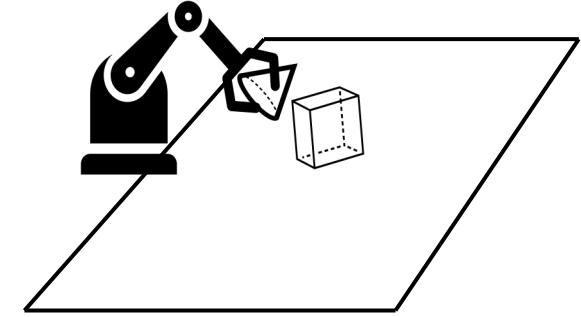
Shaped Reward



€€



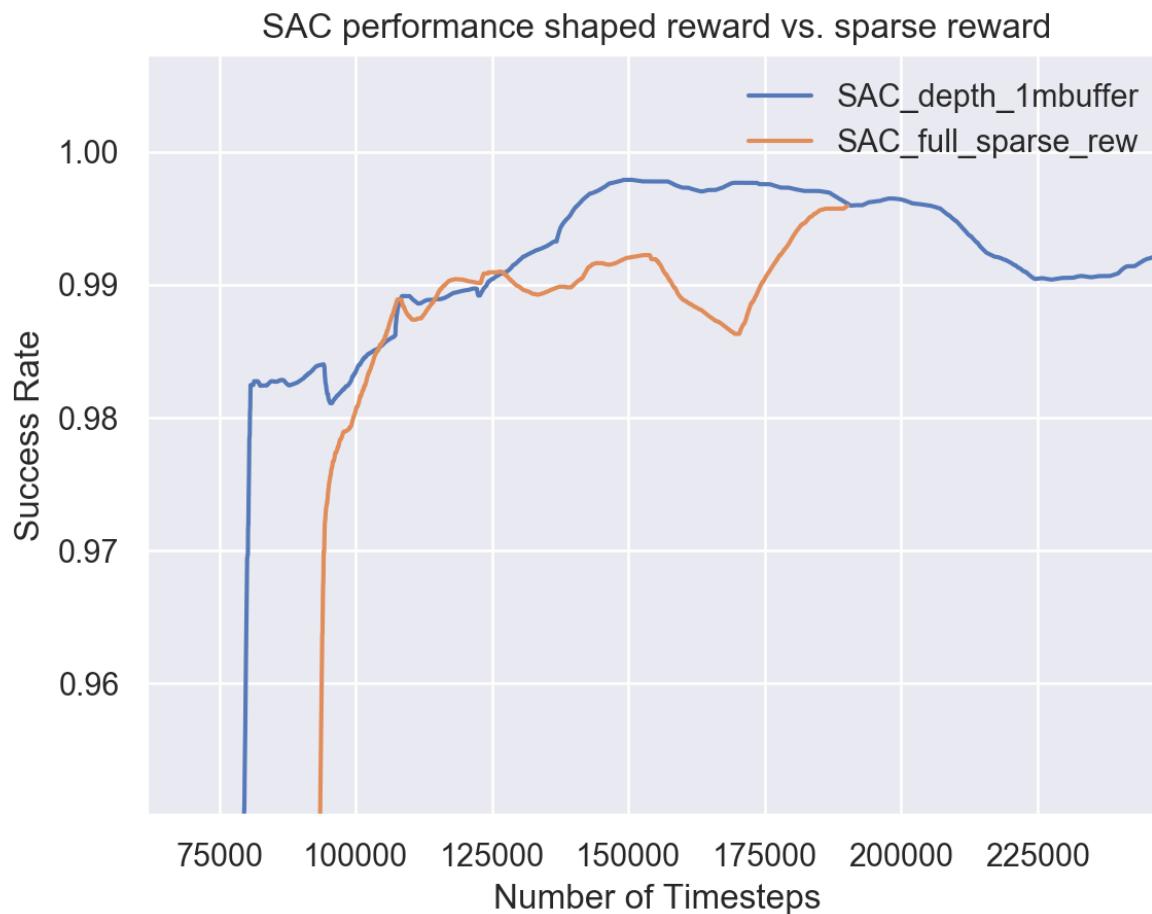
€



€€€

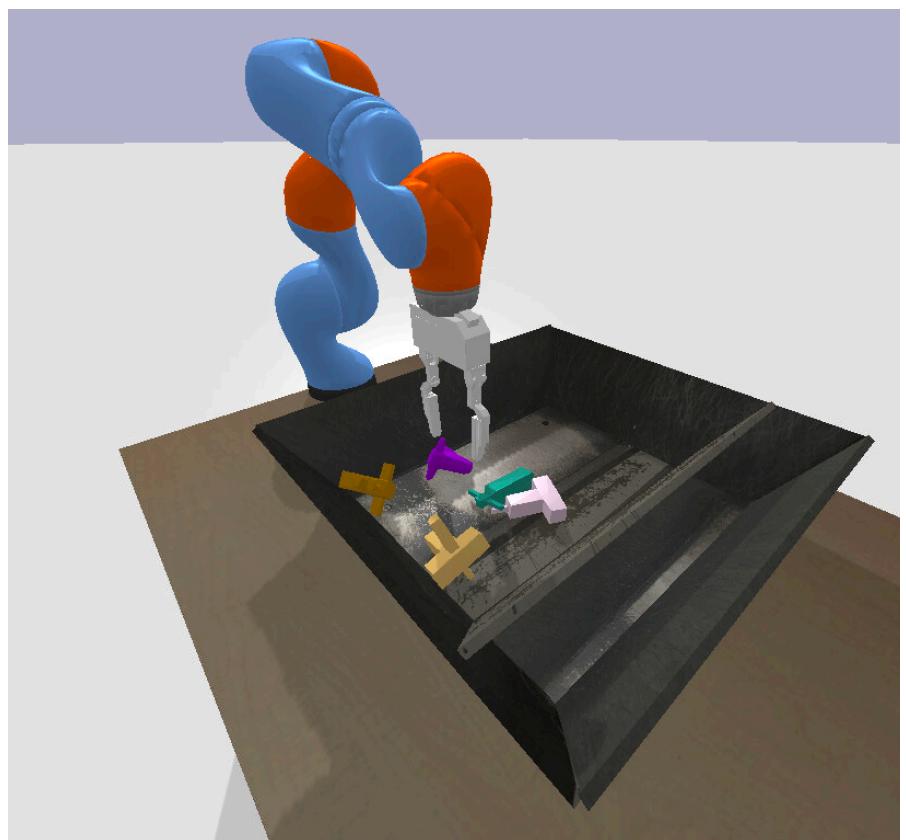


Shaped Reward



Contribution

- Novel grasping model
- Comparison of SAC, BDQ and DQN algorithms
- Comparison of **perception** modules



Future Timeline

- Hardware Experiments
- Multi-agent
- Targeted picking
- Soft Object Extensions



References

- [1] Breyer, M., Furrer, F., Novkovic, T., Siegwart, R., & Nieto, J. (2018). Comparing Task Simplifications to Learn Closed-Loop Object Picking Using Deep Reinforcement Learning. *IEEE Robotics and Automation Letters*, 4(2), 1549–1556. <https://doi.org/10.1109/LRA.2019.2896467>
- [2] Richard S. Sutton and Andrew G. Barto. 2018. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA.
- [3] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., & Hassabis, D. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529–533. <https://doi.org/10.1038/nature14236>
- [4] Tavakoli, A., Pardo, F., & Kormushev, P. (2018). Action branching architectures for deep reinforcement learning. *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*, 4131–4138.
- [5] Mahler, J., Matl, M., Satish, V., Danielczuk, M., DeRose, B., McKinley, S., & Goldberg, K. (2019). Learning ambidextrous robot grasping policies. *Science Robotics*, 4(26).
<https://doi.org/10.1126/scirobotics.aau4984>