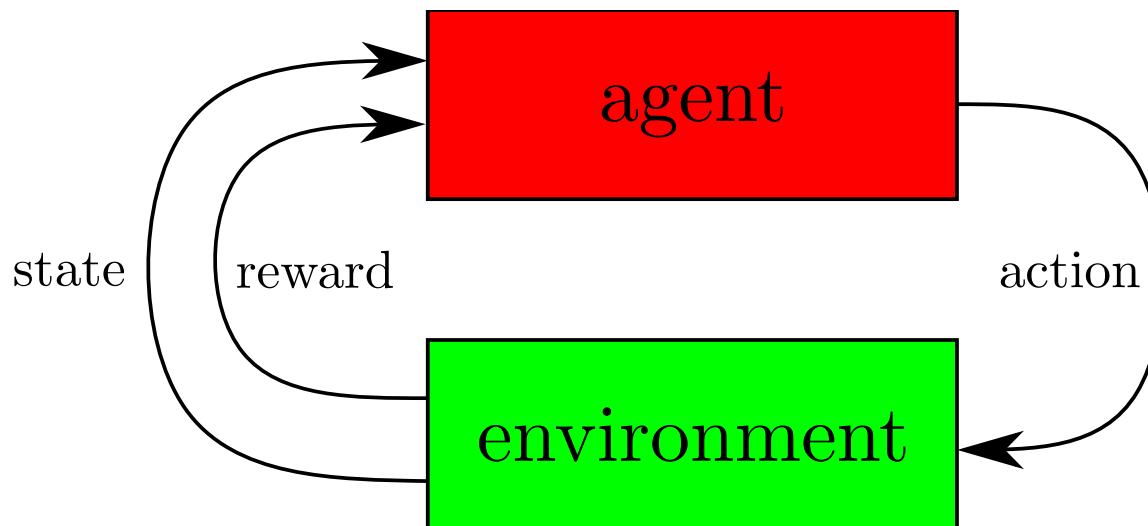


Deep Q-network

Ondřej Podsztavek

Let's Talk ML Prague, 14. 12. 2017



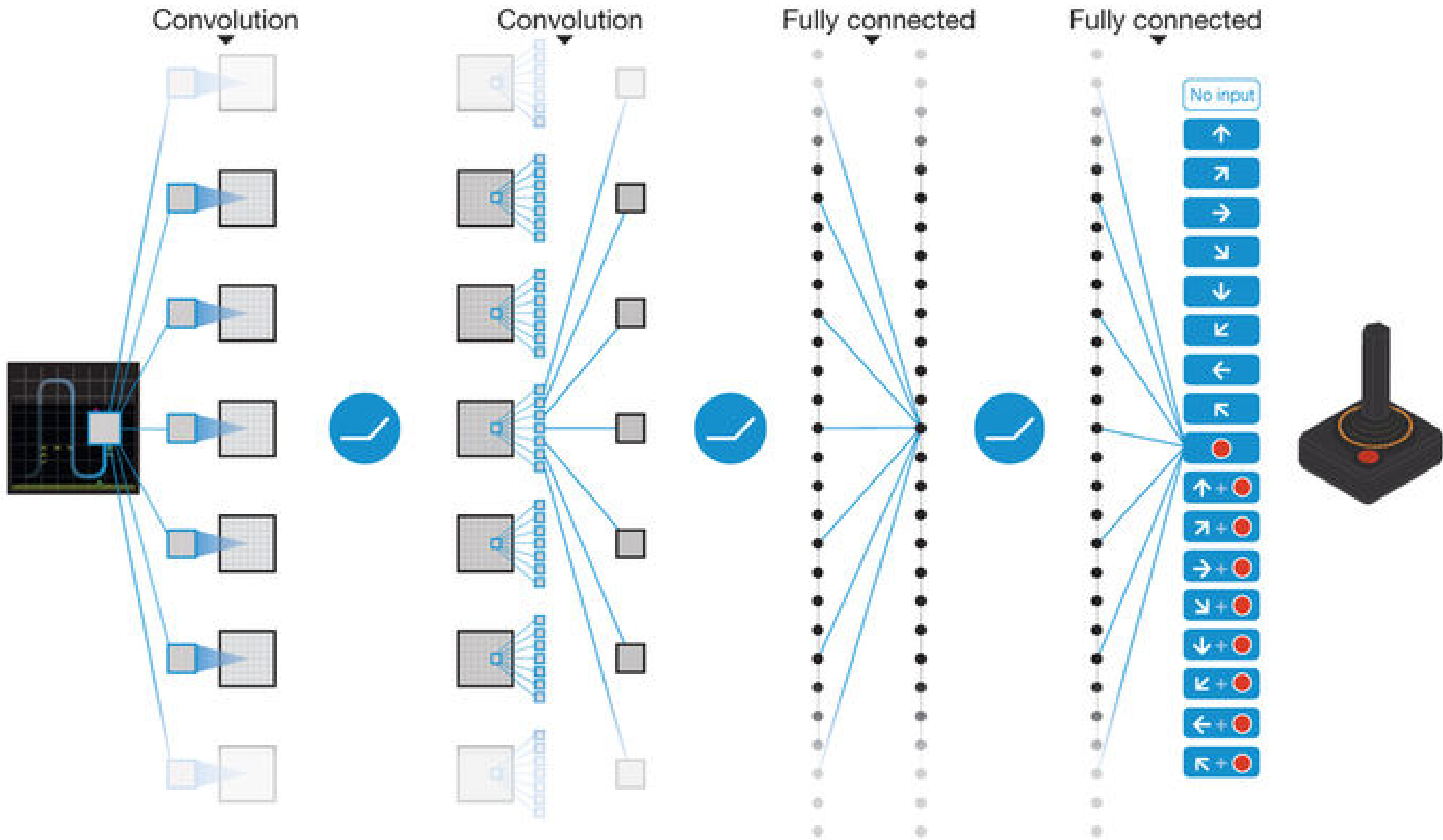
Diverse Array of Tasks

- 49 Atari 2600 games
- the same network architecture and hyperparameters
- pixels and the game score as inputs

Q function

$$Q^*(s, a) = \max_{\pi} \mathbb{E}[r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots | s_t = s, a_t = a, \pi]$$

Deep Convolutional Network

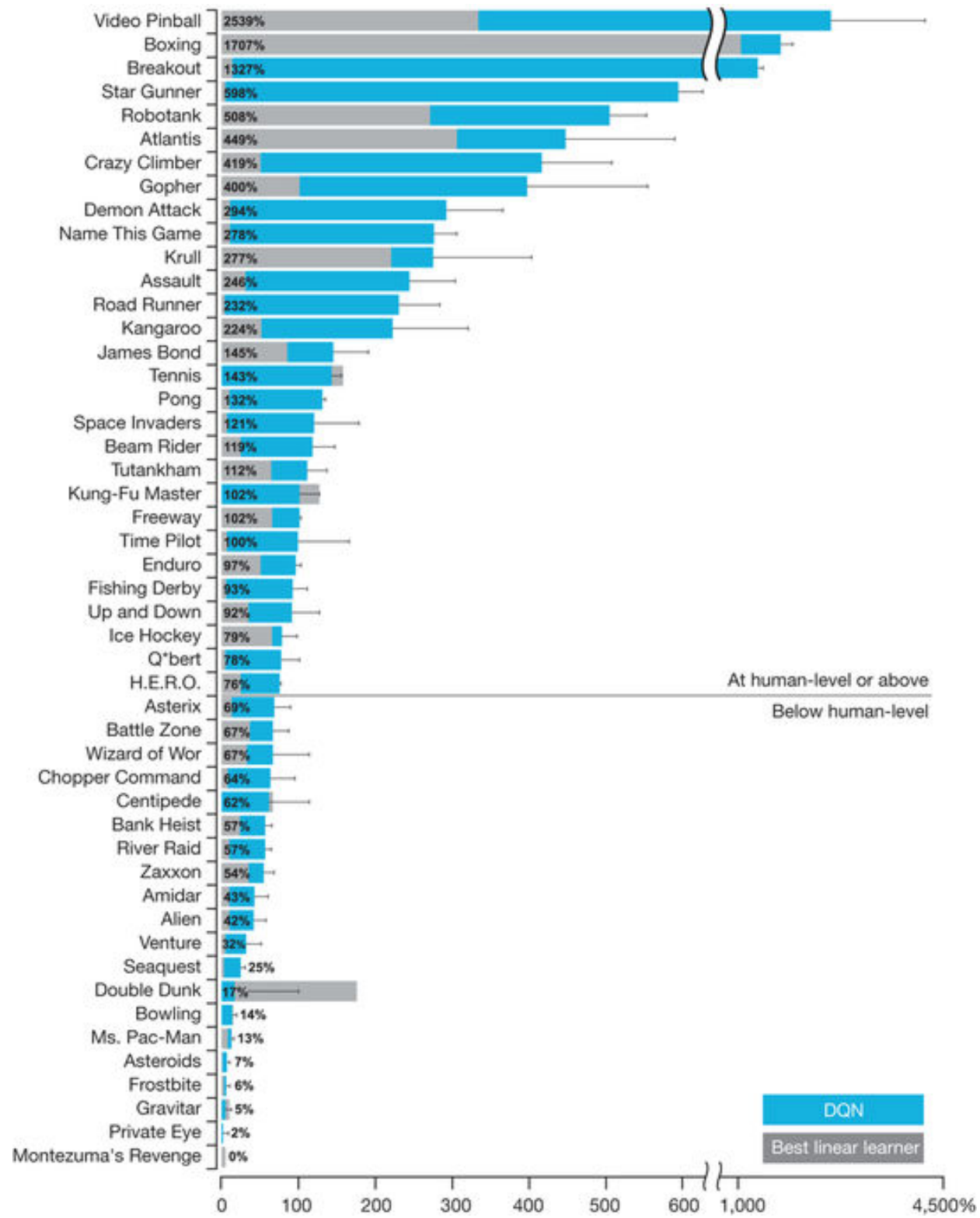


Training Algorithm Modifications

- RL is known to be unstable for a nonlinear function approximator
- separate network for generating the targets for the Q-learning update
- experience replay

$$e_t = (s_t, a_t, r_t, s_{t+1})$$

$$D_t = \{e_1, \dots, e_t\}$$



Variants

- *Double Q-learning*, corrects tendency to overestimate action's values
- *Prioritized Replay*, somehow improves the experience replay
- *Dueling DQN*, one estimates the value at a time step, other calculates advantages of an action

References

- Human-level control through deep reinforcement learning,
<https://www.nature.com/articles/nature14236>
- Deep Reinforcement Learning with Double Q-learning, <https://arxiv.org/abs/1509.06461>
- Learning from Delayed Rewards,
<http://www.cs.rhul.ac.uk/~chrisw/thesis.html>
- OpenAI Baselines: DQN,
<https://blog.openai.com/openai-baselines-dqn/>