



Reinforcement learning

Nando de Freitas

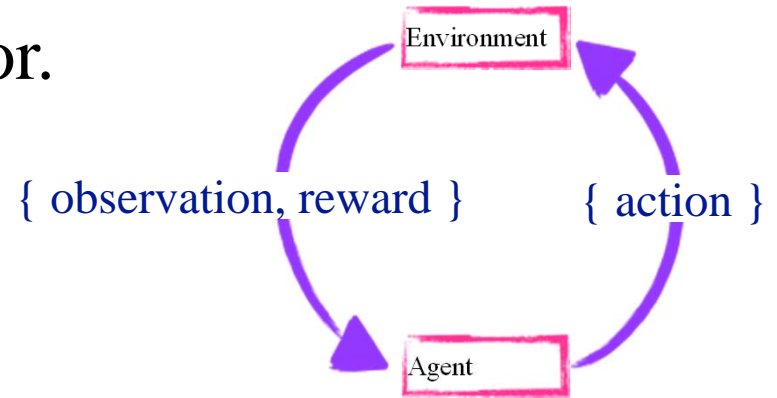


UNIVERSITY OF
OXFORD

The Promise of Reinforcement Learning

Learning to act through trial and error.

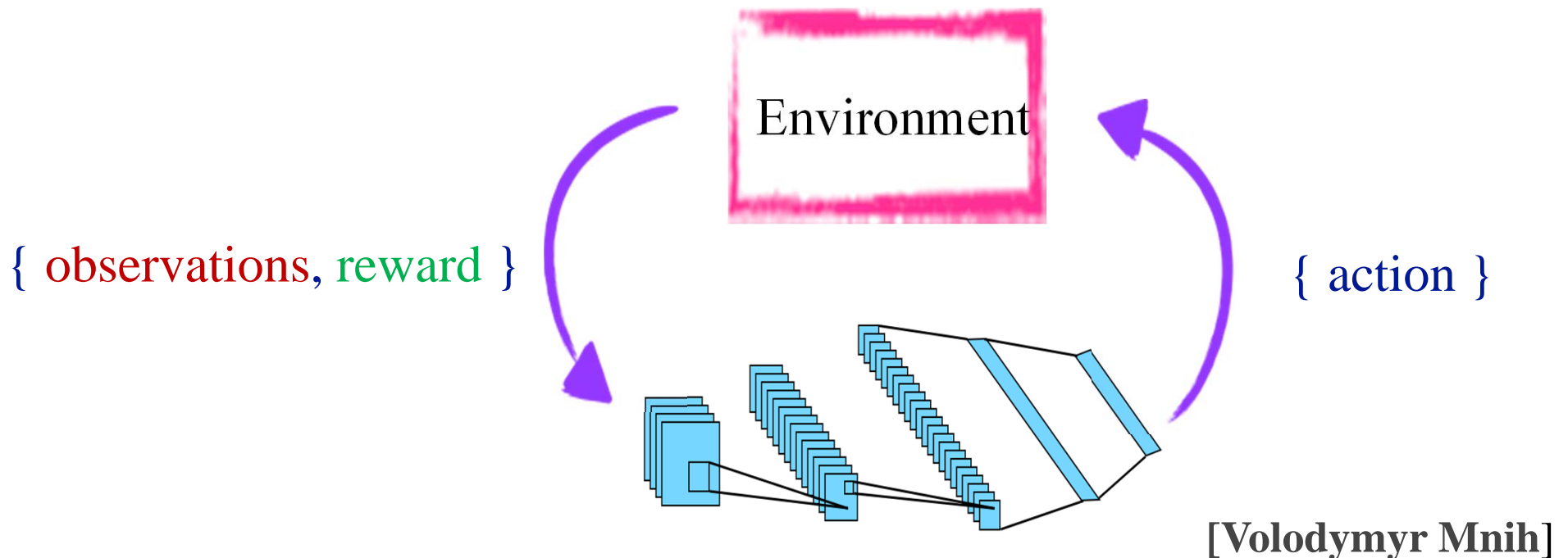
- An agent interacts with an environment and learns by maximizing a scalar reward signal.
- No models, labels, demonstrations, or any other human-provided supervision signal.
- Representation has been a challenge/missing.



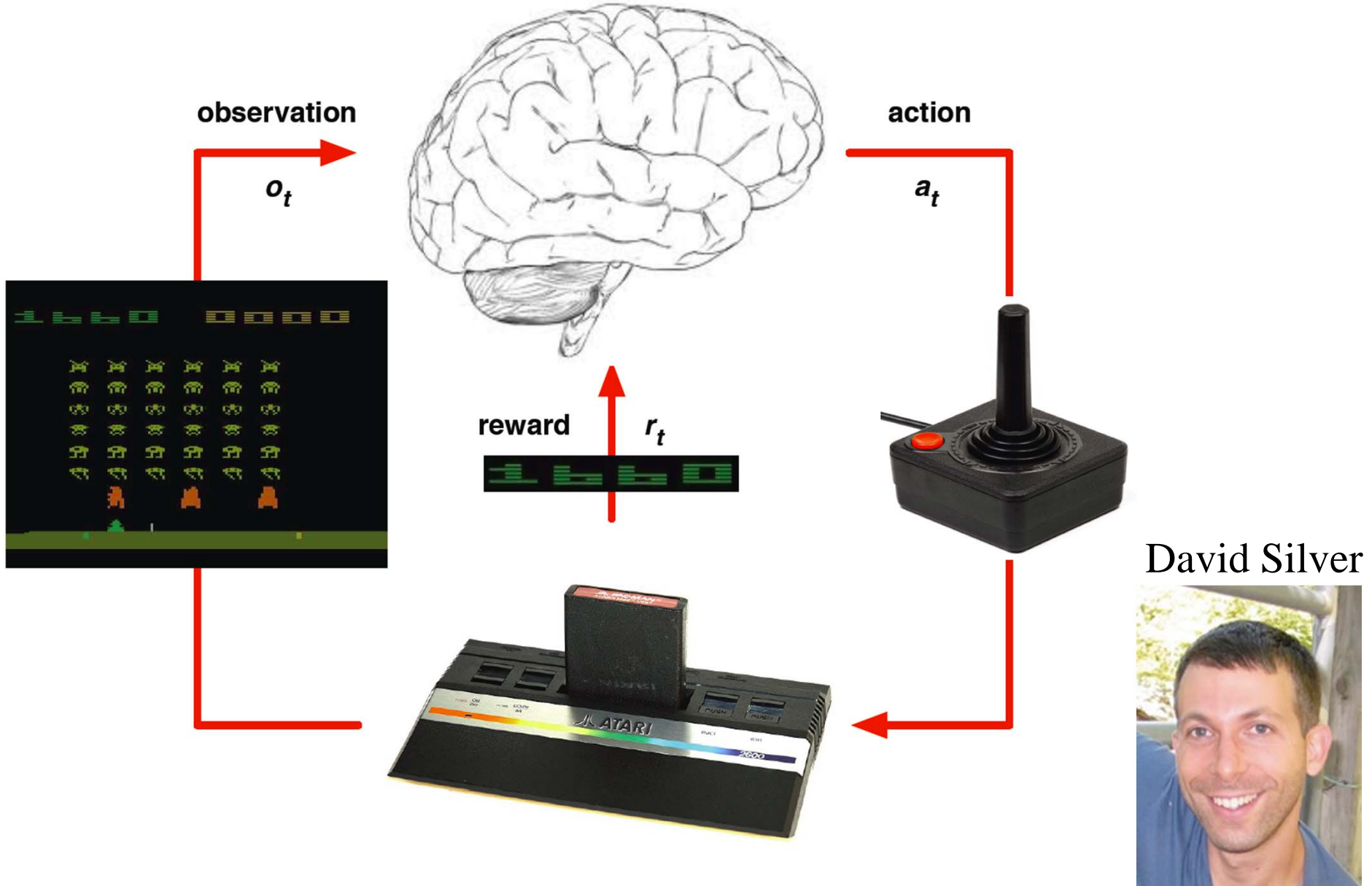
[Volodymyr Mnih]

Deep Reinforcement Learning

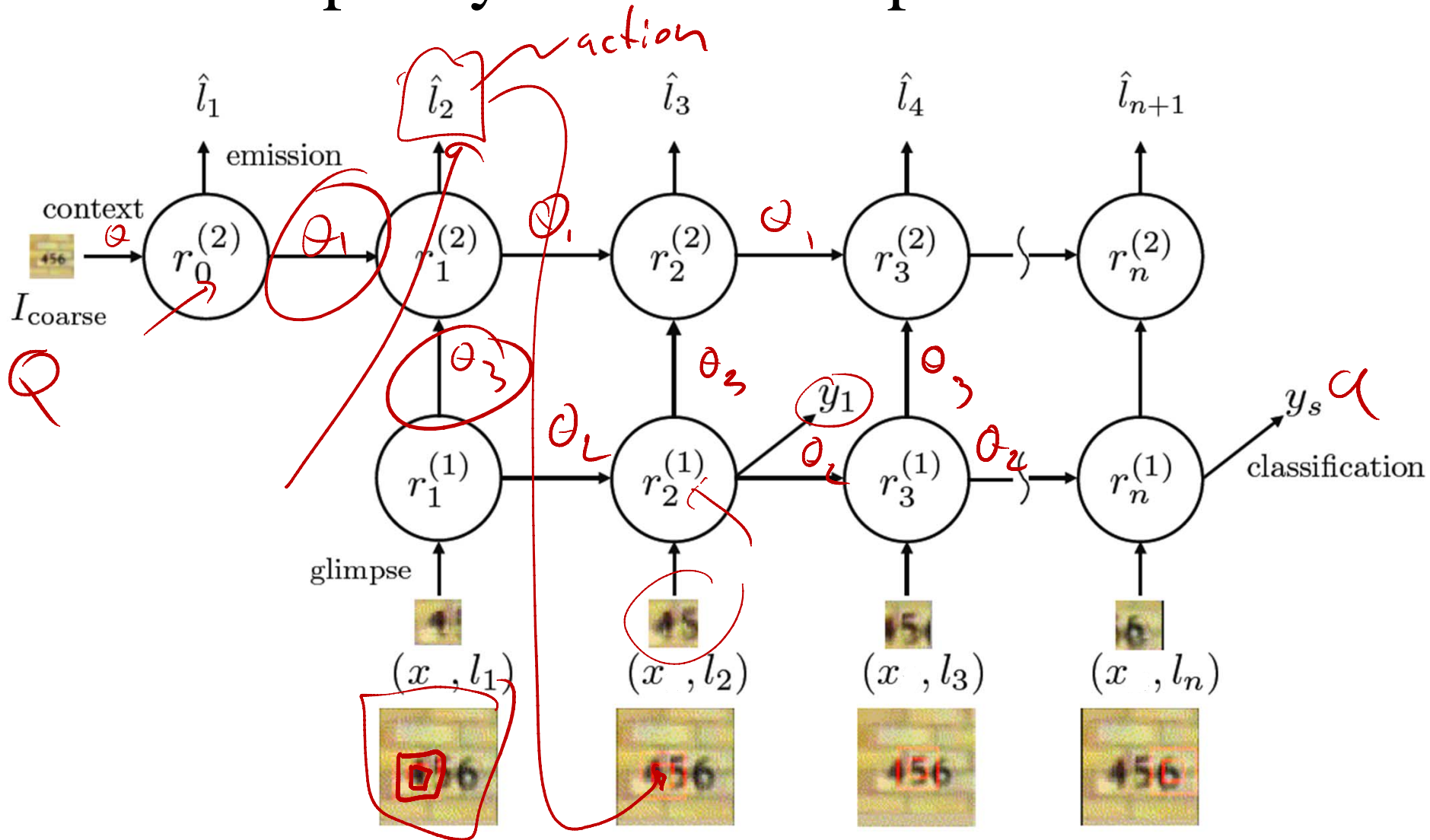
- Combining deep neural networks with RL.
- Learn to act from high-dimensional sensory inputs.
- Is a noisy, sparse, and delayed reward signal sufficient for training deep networks? **Credit assignment** problem.



Example: Learning to play Atari



Direct policy search example: Attention



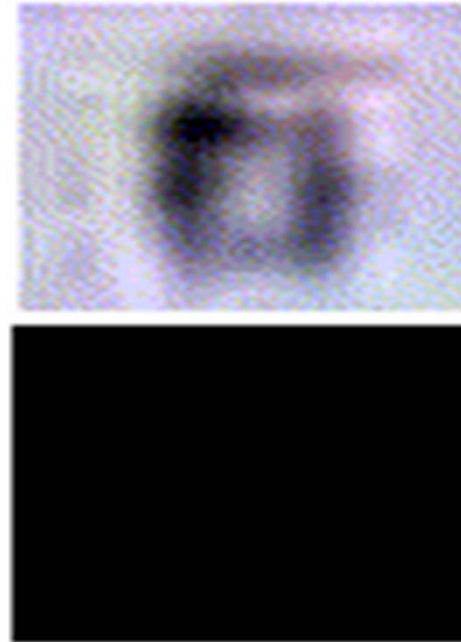
[Ba, Mnih, Kavuckuoglu]

Results

MNIST
sequences



Street View House
Number sequences

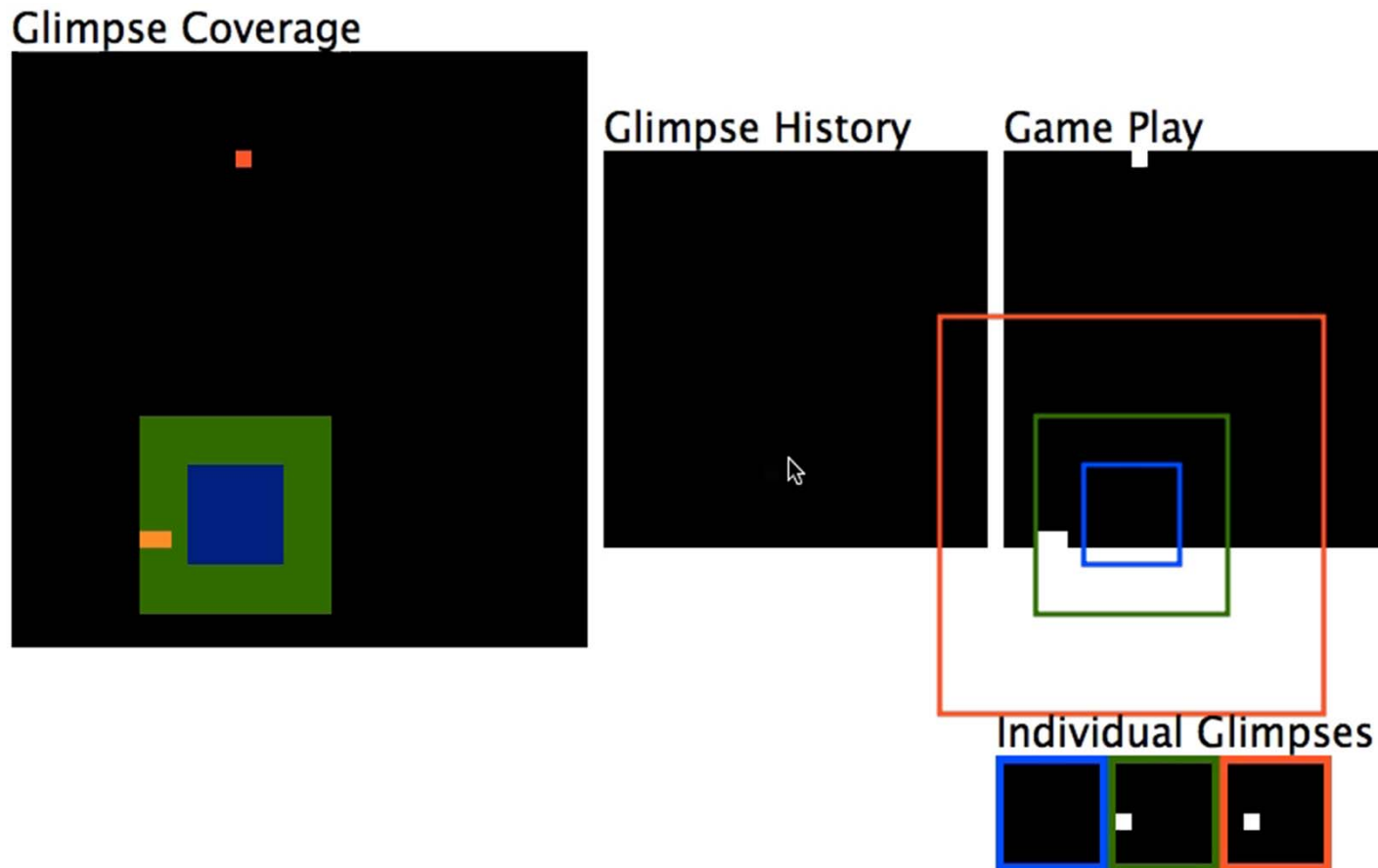


- The attention-based model achieves state-of-the-art accuracy on the SVHN multi-digit task - 3.9% error.
- 4 times fewer floating point operations than the best ConvNet.

[Volodymyr Mnih et al]

Attention-Based Game Agent

- Roughly the same model and training method can be used in a game-playing agent.
- The agent learns to track a ball without being told to do so.



Direct policy search

history ← observations
 $\mathbf{h}_t = \{\mathbf{o}_0, \mathbf{o}_1, \dots, \mathbf{o}_t\}$

Policy

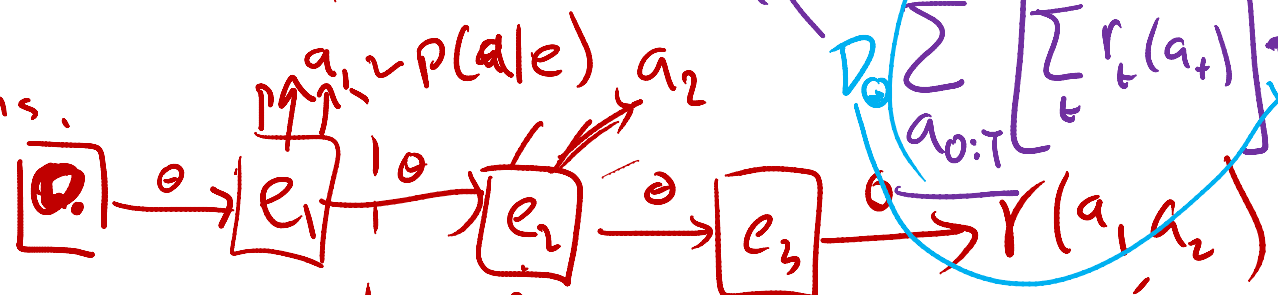
$$\pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) = \prod_{t=0}^T \pi_{\theta}(\mathbf{a}_t | \mathbf{h}_t)$$

actions

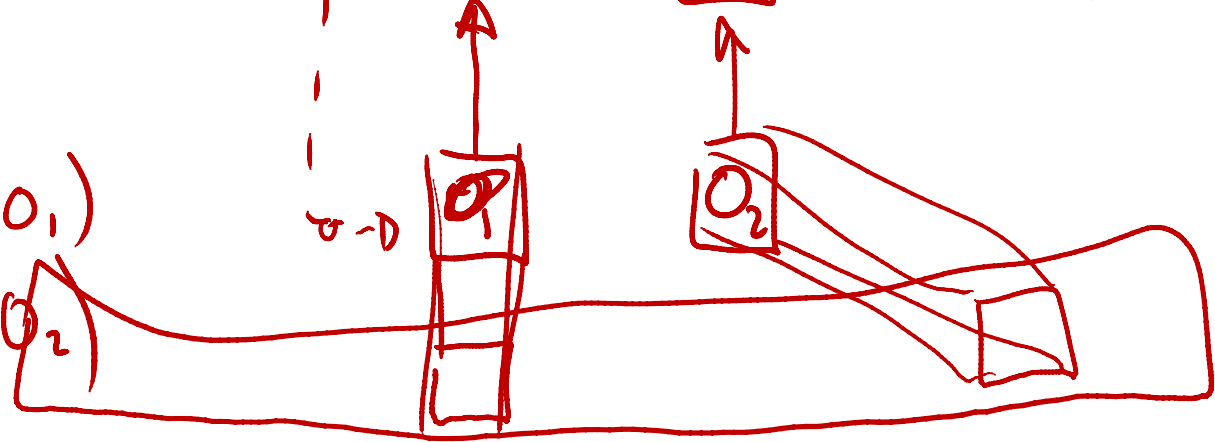
$$J^{\pi}(\theta) = \mathbb{E} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] = \mathbb{E}_{\pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right]$$



↑
 expected returns.



$h_0 = \mathbf{o}_0$
 $h_1 = (\mathbf{o}_0, \mathbf{o}_1)$
 $h_2 = (\mathbf{o}_0, \mathbf{o}_1, \mathbf{o}_2)$



Policy gradients using backprop

$$\begin{aligned}\nabla_{\theta} J^{\pi}(\theta) &= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] \nabla \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) \quad (\log y)' = \frac{y'}{y} \\ &= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] \nabla \log \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) \\ &= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] \left[\sum_{t=0}^T \nabla \log \pi_{\theta}(\mathbf{a}_t | \mathbf{h}_{0:t}) \right] \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) \\ &= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T \nabla \log \pi_{\theta}(\mathbf{a}_t | \mathbf{h}_{0:t}) \sum_{n=t}^T r_n(\mathbf{a}_n) \right] \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})\end{aligned}$$

Policy gradients using backprop

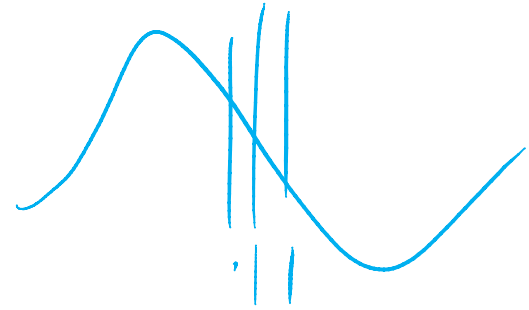
$$\nabla_{\theta} J^{\pi}(\theta) = \sum_{\mathbf{a}_{0:T}} \left[\underbrace{\sum_{t=0}^T \nabla \log \pi_{\theta}(\mathbf{a}_t | \mathbf{h}_{0:t})}_{\text{backprop}} \underbrace{\sum_{n=t}^T r_n(\mathbf{a}_n)}_{\text{backprop}} \right] \underbrace{\pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})}_{\text{backprop}}$$

$$\mathbf{a}_{0:T}^{(i)} \sim \pi_{\theta}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})$$

$$\widehat{\nabla_{\theta} J^{\pi}(\theta)} = \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \underbrace{\nabla \log \pi_{\theta}(\mathbf{a}_t^{(i)} | \mathbf{h}_{0:t})}_{\text{backprop}} \underbrace{\sum_{n=t}^T r_n(\mathbf{a}_n^{(i)})}_{\text{backprop}}$$

↑
Samples

Matt Hoffman



Neuro-dynamic programming

Dynamic programming

Csaba Szepesvári

$$\mathbf{a}_t = \pi(\mathbf{s}_t)$$

\mathbf{s} denotes the state (model abstraction of the environment, e.g. histories)

Value function

$$V^\pi(\mathbf{s}_0) = \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \right]$$

V depends on π
optimal

best value function

discount factor $\gamma \in (0, 1)$

$$V^*(\mathbf{s}_0) = \max_{\pi} \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \right]$$

Dynamic programming

AB + AC
A(B+C)

$$V^*(\mathbf{s}_0) = \max_{\pi} \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \right]$$

$$= \max_{\pi, \mathbf{a}_0} \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots} \left[\underbrace{r_0(\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1)}_{\pi(\mathbf{s}_0)} + \sum_{t=1}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \right]$$

$\mathbf{a}_0 \pm \pi(\mathbf{s}_0)$

$$= \max_{\mathbf{a}_0, \pi} \mathbb{E}_{\mathbf{s}_1, \mathbf{s}_2, \mathbf{s}_3, \dots} \left[\underbrace{r_0(\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1)}_{\pi} + \sum_{t=1}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_0 \right]$$

$$= \max_{\mathbf{a}_0} \mathbb{E}_{\mathbf{s}_1} \left[\underbrace{r_0(\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1)}_{\pi} + \max_{\pi} \mathbb{E}_{\mathbf{s}_2, \mathbf{s}_3, \mathbf{s}_4, \dots} \left\{ \sum_{t=1}^{\infty} \gamma^t r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_1 \right\} \mid \mathbf{s}_0 \right]$$

$$= \max_{\mathbf{a}_0} \mathbb{E}_{\mathbf{s}_1} \left[\underbrace{r_0(\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1)}_{\pi} + \gamma \max_{\pi} \mathbb{E}_{\mathbf{s}_2, \mathbf{s}_3, \mathbf{s}_4, \dots} \left\{ \sum_{t=1}^{\infty} \gamma^{t-1} r_t(\mathbf{s}_t, \pi(\mathbf{s}_t), \mathbf{s}_{t+1}) \mid \mathbf{s}_1 \right\} \mid \mathbf{s}_0 \right]$$

$$= \max_{\mathbf{a}_0} \mathbb{E}_{\mathbf{s}_1} \left[\underbrace{r_0(\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1)}_{\pi} + \gamma \underbrace{V^*(\mathbf{s}_1)} \right]$$

Bellman's equation and TD

$$V^*(\mathbf{s}) = \max_a \mathbb{E}_{\mathbf{s}'} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^*(\mathbf{s}') | \mathbf{s}]$$

current state (circled in red)
future state

$$V^{\pi}(\mathbf{s}) + \eta \underline{V^{\pi}(\mathbf{s})} = \mathbb{E}_{\mathbf{s}'} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^{\pi}(\mathbf{s}')] \eta + V^{\pi}(\mathbf{s})$$

$$V^{\pi}(\mathbf{s}) = V^{\pi}(\mathbf{s}) + \eta \left\{ \mathbb{E}_{\mathbf{s}'} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma V^{\pi}(\mathbf{s}')] - V^{\pi}(\mathbf{s}) \right\}$$

$\tilde{\mathbf{s}}' \sim P(\mathbf{s}' | \dots)$

$$V_{t+1}^{\pi}(\mathbf{s}) = V_t^{\pi}(\mathbf{s}) + \eta \left\{ r(\mathbf{s}, \mathbf{a}, \tilde{\mathbf{s}}') + \gamma V_t^{\pi}(\tilde{\mathbf{s}}') - V_t^{\pi}(\mathbf{s}) \right\}$$

Action-value (Q) functions

$$V^*(\mathbf{s}) = \max_{\mathbf{a}'} Q^*(\mathbf{s}, \mathbf{a}')$$

$$V^*(\mathbf{s}) = \max_{\mathbf{a}} \mathbb{E}_{\mathbf{s}'} [r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \underbrace{V^*(\mathbf{s}')}] | \mathbf{s}] \leftarrow$$

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

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

Q - Learning

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

Neuro-dynamic programming

$Q(\mathbf{s}, \mathbf{a}; \mathbf{w})$, where \mathbf{w} are the parameters

$$L(\mathbf{w}_i) = \mathbb{E}_{\mathbf{s}, \mathbf{a}} \left\{ \left(\underbrace{\mathbb{E}_{\mathbf{s}'} \left[r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}_{i-1}) \right]}_{\substack{\text{target} \\ \text{critic}}} \right) - \underbrace{Q(\mathbf{s}, \mathbf{a}, \mathbf{w}_i)}_{\text{actor}} \right)^2 \right\}$$

$$y_i = \mathbb{E}_{\mathbf{s}'} \left[r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}_{i-1}) \right]$$

target

Neuro-dynamic programming

$$L(\mathbf{w}_i) = \mathbb{E}_{\mathbf{s}, \mathbf{a}} \left\{ \left(\mathbb{E}_{\mathbf{s}'} \left[r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}_{i-1}) \right] - Q(\mathbf{s}, \mathbf{a}, \mathbf{w}_i) \right)^2 \right\}$$

$$\nabla_{\mathbf{w}_i} L(\mathbf{w}_i) = \mathbb{E}_{\mathbf{s}, \mathbf{a}, \mathbf{s}'} \left\{ \left(r(\mathbf{s}, \mathbf{a}, \mathbf{s}') + \gamma \max_{\mathbf{a}'} Q(\mathbf{s}', \mathbf{a}', \mathbf{w}_{i-1}) - Q(\mathbf{s}, \mathbf{a}, \mathbf{w}_i) \right) \nabla_{\mathbf{w}_i} Q(\mathbf{s}, \mathbf{a}, \mathbf{w}_i) \right\}$$

$$\hat{\mathbf{a}} = \arg \max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}; \hat{\mathbf{w}})$$

~~Step~~

① state \mathbf{s}

② Choose \mathbf{a}

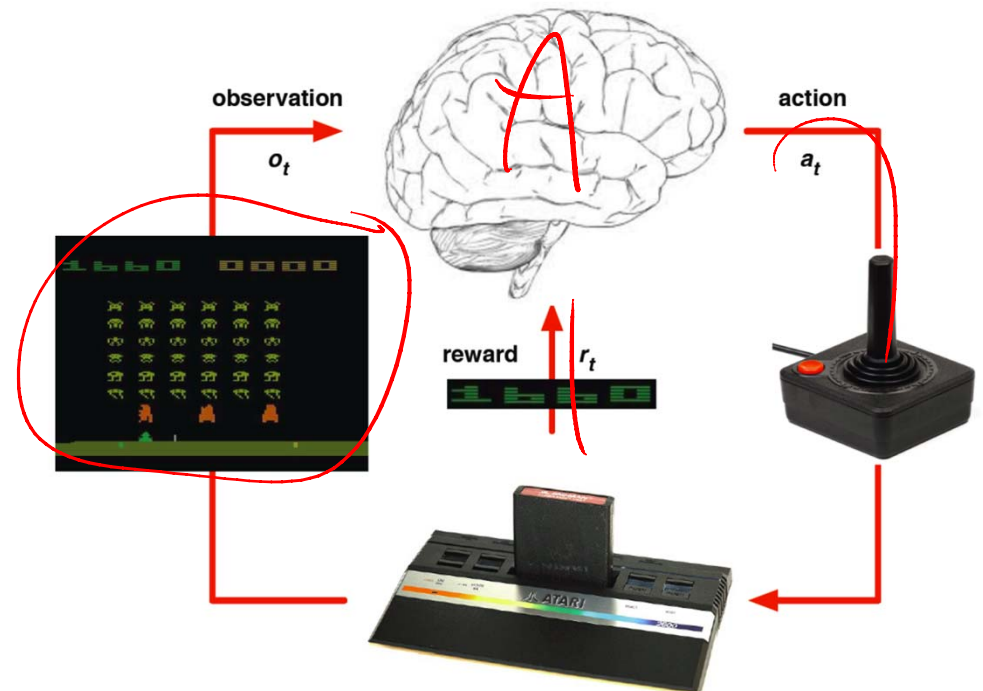
③ Environment $\rightarrow \mathbf{s}'$

④ Update \mathbf{w}

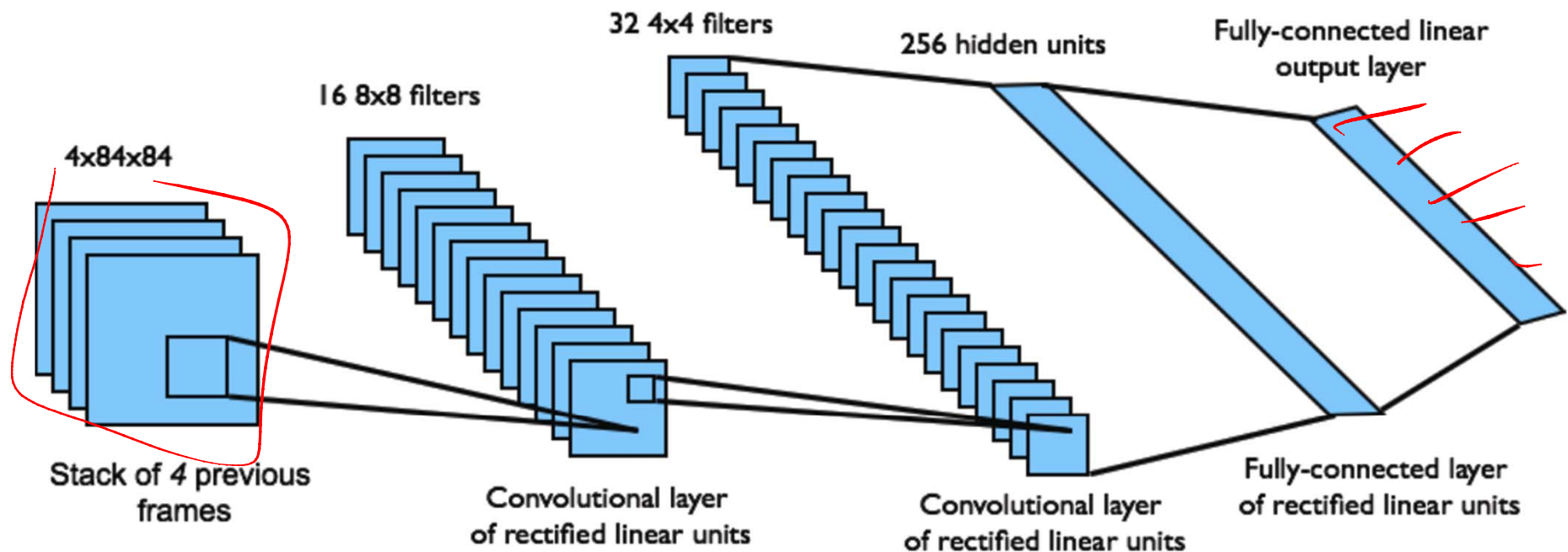
$\left\{ \begin{array}{l} \text{w.p. } \epsilon \quad \mathbf{a} = \arg \max_{\mathbf{a}'} Q(\mathbf{s}, \mathbf{a}') \\ \text{w.p. } 1-\epsilon \quad \mathbf{a} \sim \mathcal{U}[\] \end{array} \right.$

Deepmind's DQN

- Use a deep neural network to represent the action-value function Q .
- Learn Q with end-to-end RL mapping the raw pixels to action values in Atari games.
- New stable online variant of Q-learning:
 - Do updates on samples of past experience.
 - Freeze network used for generating targets and refresh periodically.

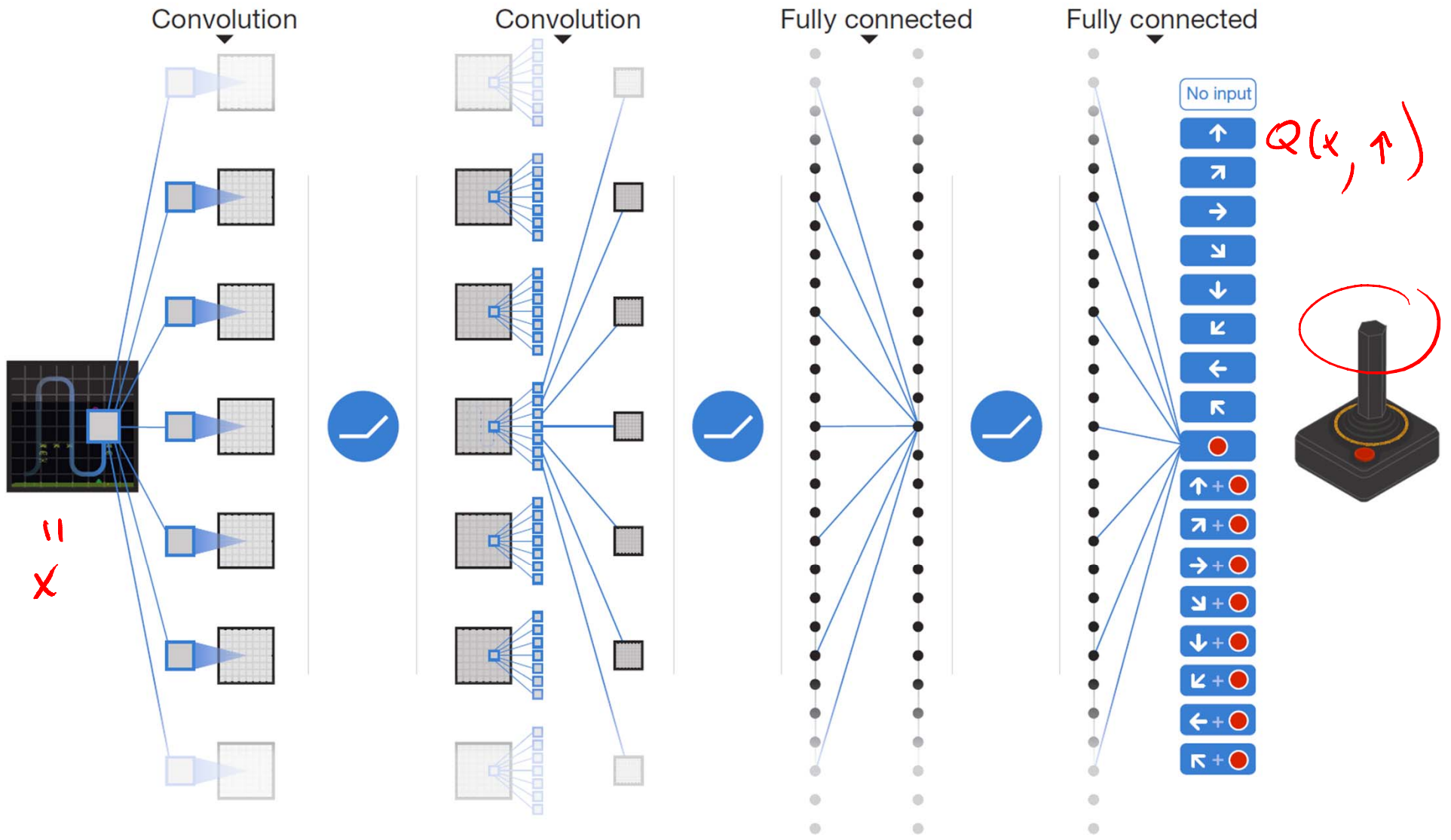


- ▶ End-to-end learning of values $Q(s, a)$ from pixels s
- ▶ Input state s is stack of raw pixels from last 4 frames
- ▶ Output is $Q(s, a)$ for 18 joystick/button positions
- ▶ Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

DQN Convolutional Network



Q-learning with experience re-play

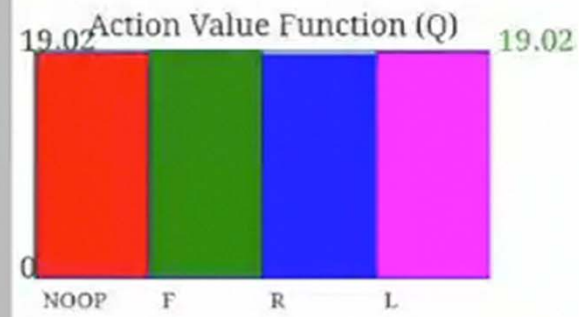
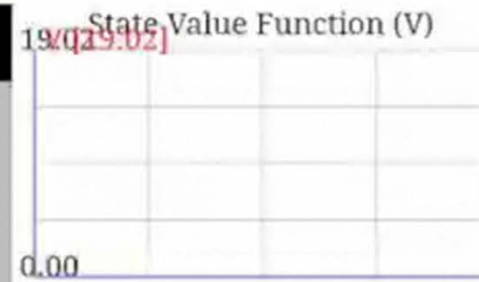
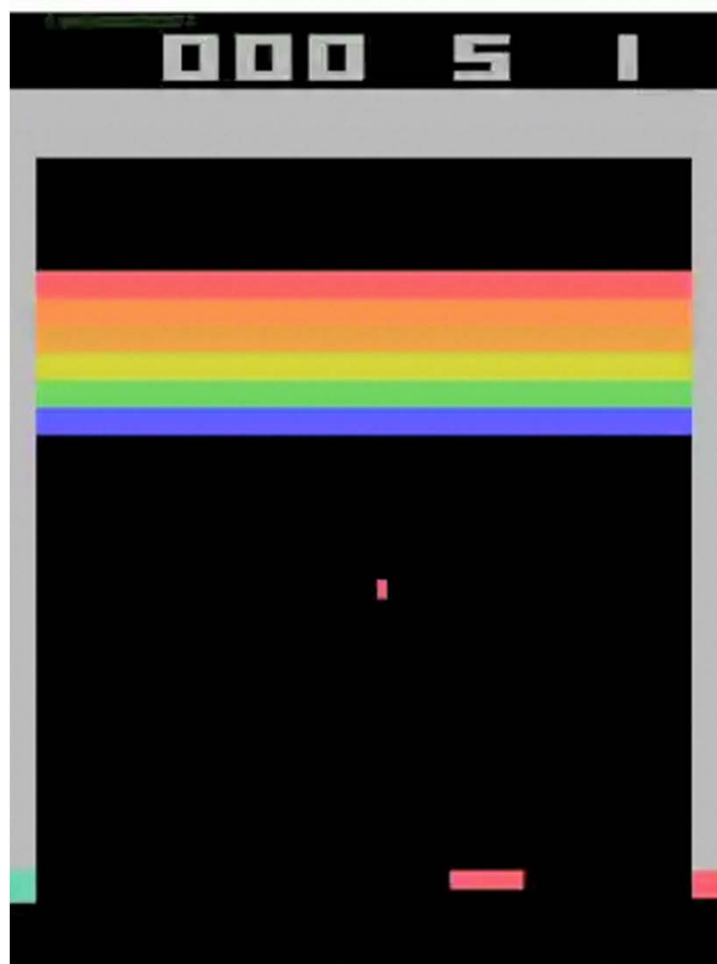
DQN increases stability with **experience replay** and **fixed Q-targets**

- ▶ Take action a_t according to ϵ -greedy policy
- ▶ Store transition $(s_t, a_t, r_{t+1}, s_{t+1})$ in replay memory \mathcal{D}
- ▶ Sample random mini-batch of transitions (s, a, r, s') from \mathcal{D}
- ▶ Compute Q-learning targets w.r.t. old, fixed parameters w^-
- ▶ Optimise MSE between Q-network and Q-learning targets

$$\mathcal{L}_i(w_i) = \mathbb{E}_{s,a,r,s' \sim \mathcal{D}_i} \left[\left(r + \gamma \max_{a'} Q(s', a'; w_i^-) - Q(s, a; w_i) \right)^2 \right]$$

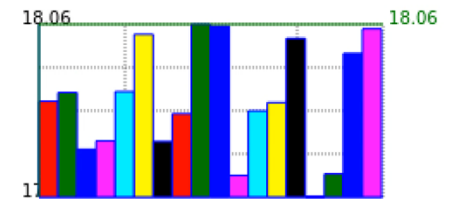
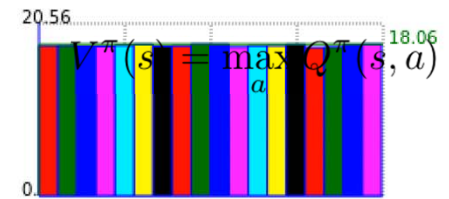
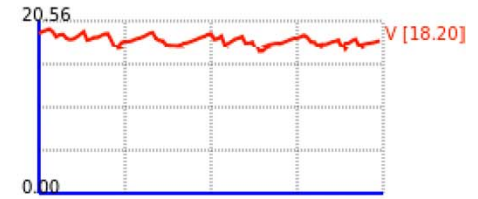
- ▶ Using variant of stochastic gradient descent (Mnih et al.)

Delayed Rewards



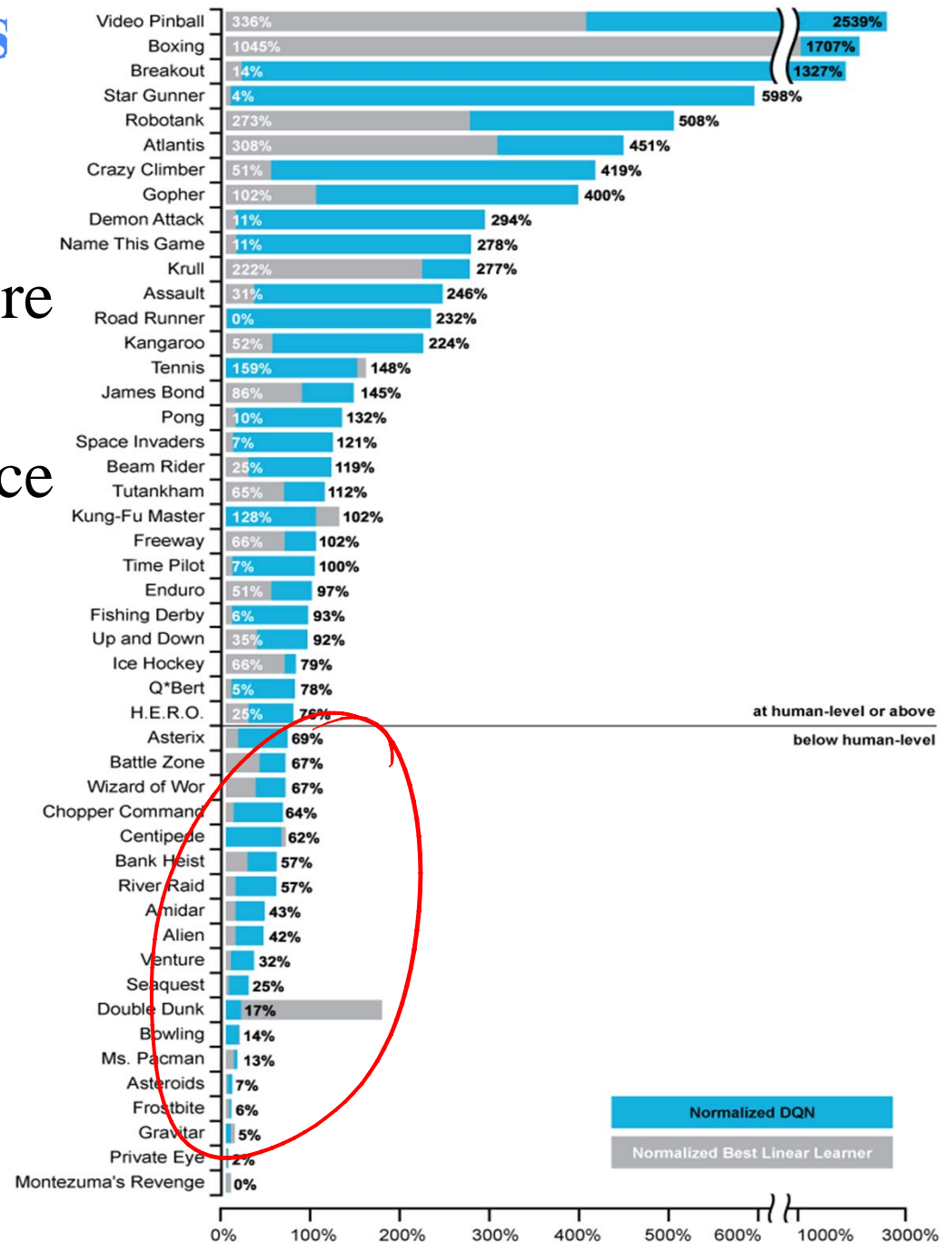
Sacrificing Immediate Rewards

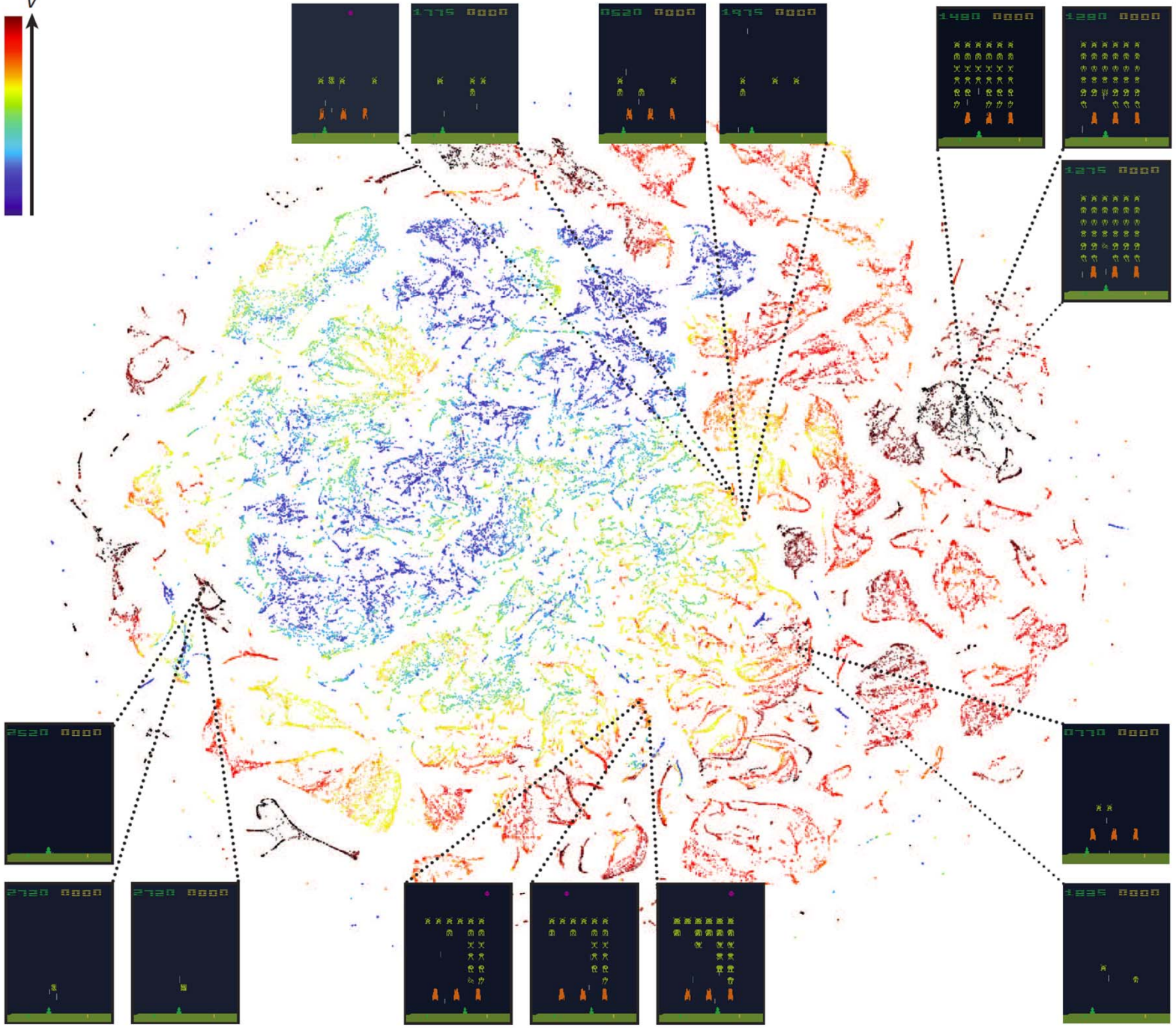
7140 18.198350906372



Results on 49 Games

- The architecture and hyperparameter values were the same for all 49 games.
- DQN achieved performance **comparable to or better than an experienced human on 29 out of 49** games.
- Games with superhuman level play include Pong, Boxing, Breakout, Space Invaders.





INNOVATIONS IN
The microbiome

nature

THE INTERNATIONAL WEEKLY JOURNAL OF SCIENCE

LEARNING CURVE

DQN
AI
And all that

Self-taught AI software
attains human-level
performance in video games

PAGES 486 & 529



Thank you