



Reinforcement learning

Nando de Freitas

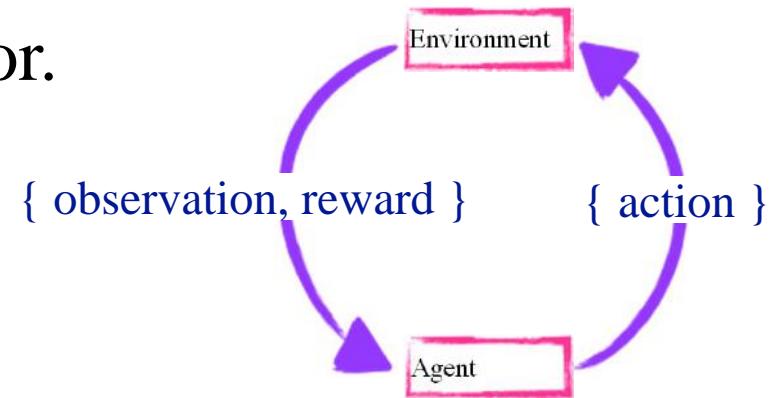


UNIVERSITY OF
OXFORD

The Promise of Reinforcement Learning

Learning to act through trial and error.

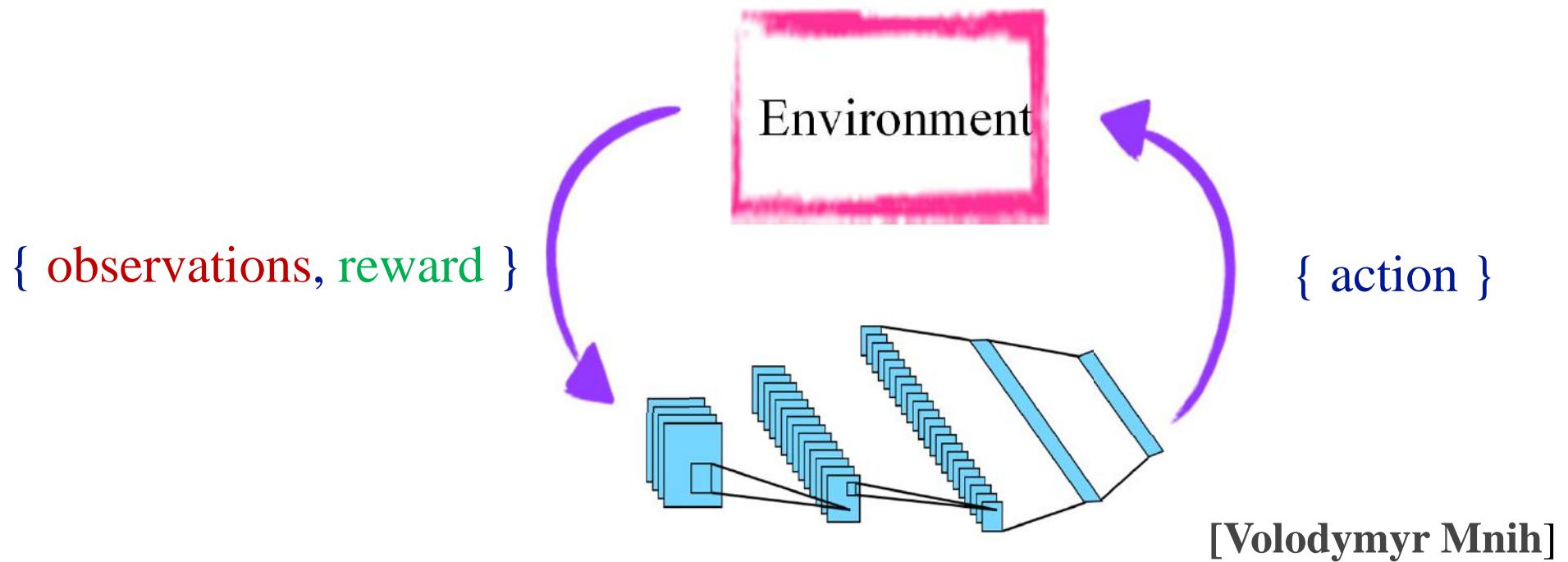
- An agent interacts with an environment and learns by maximizing a scalar reward signal.
- No models, labels, demonstrations, or any other human-provided supervision signal.
- Representation has been a challenge/missing.



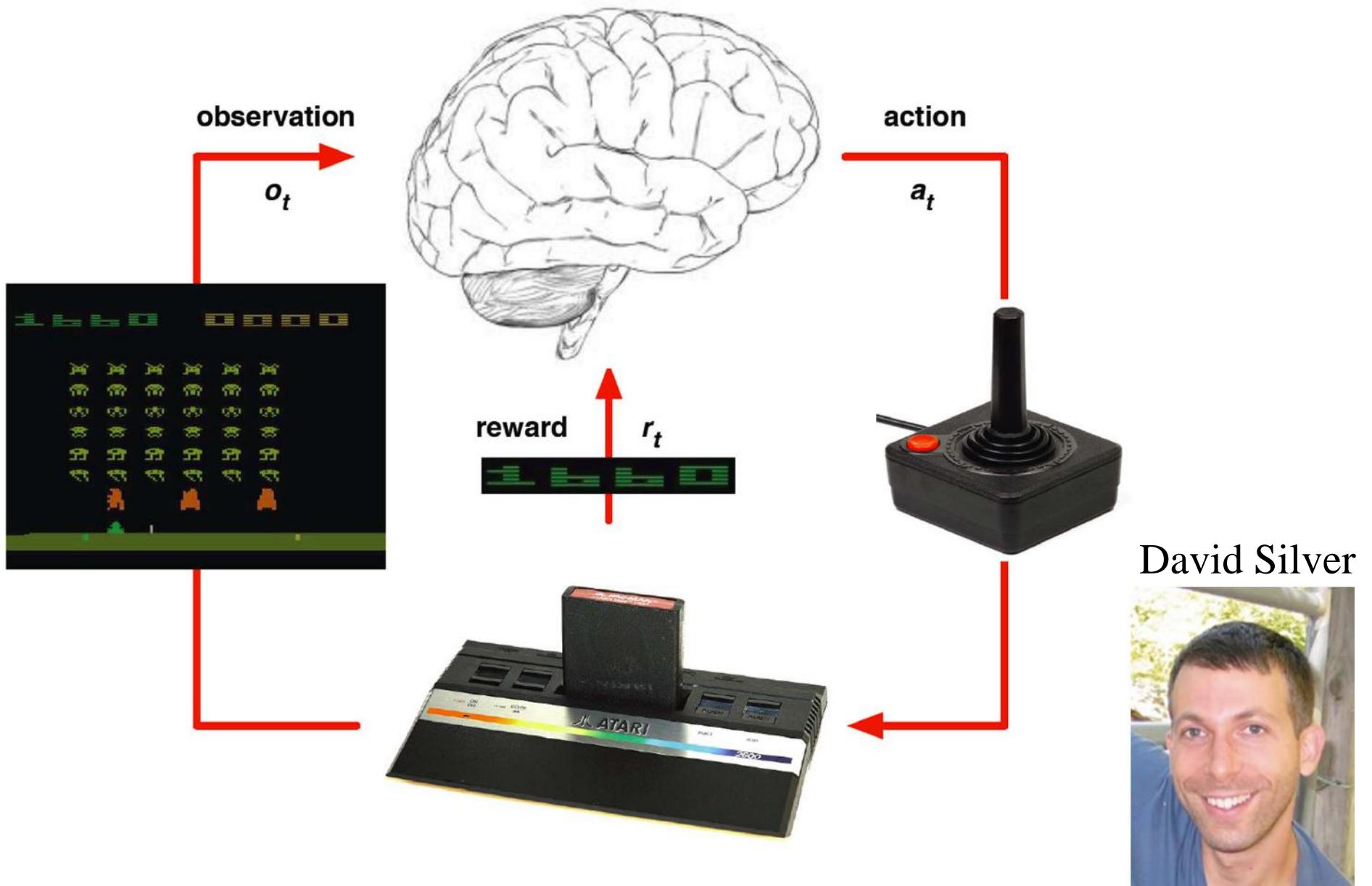
[Volodymyr Mnih]

Deep Reinforcement Learning

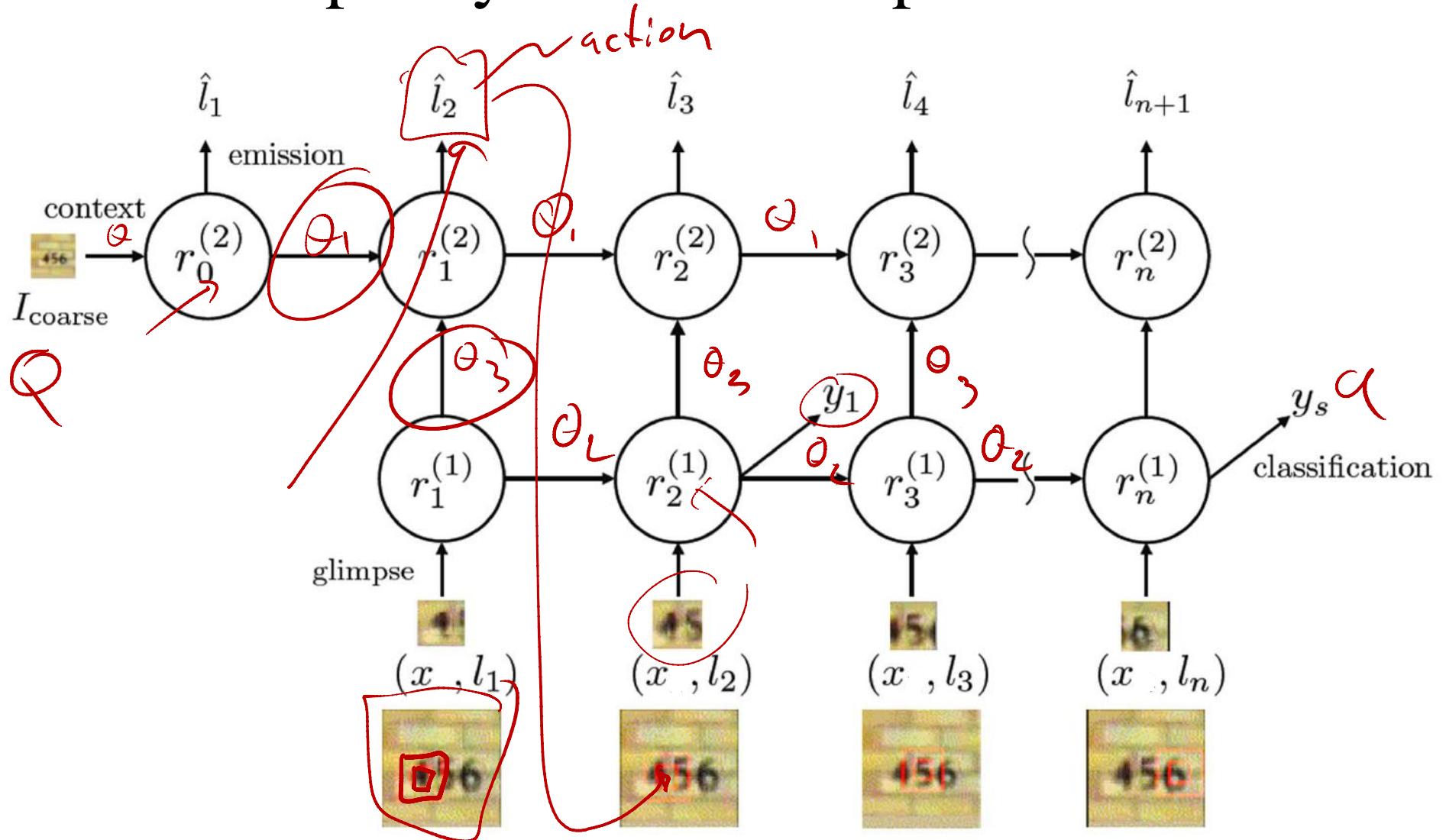
- Combining deep neural networks with RL.
- Learn to act from high-dimensional sensory inputs.
- Is a noisy, sparse, and delayed reward signal sufficient for training deep networks? Credit assignment problem.



Example: Learning to play Atari



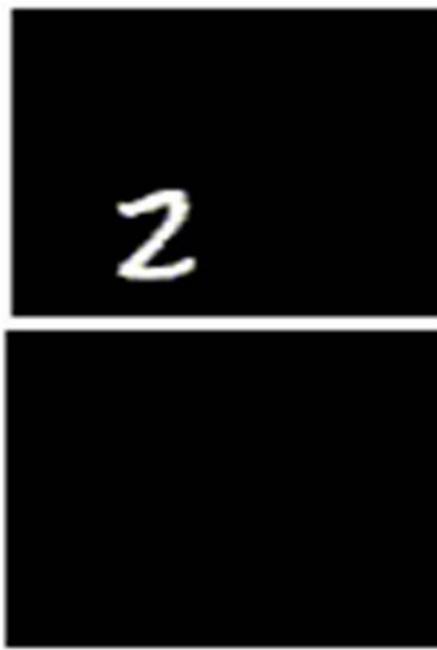
Direct policy search example: Attention



[Ba, Mnih, Kavukcuoglu]

Results

MNIST
sequences



Street View House
Number sequences

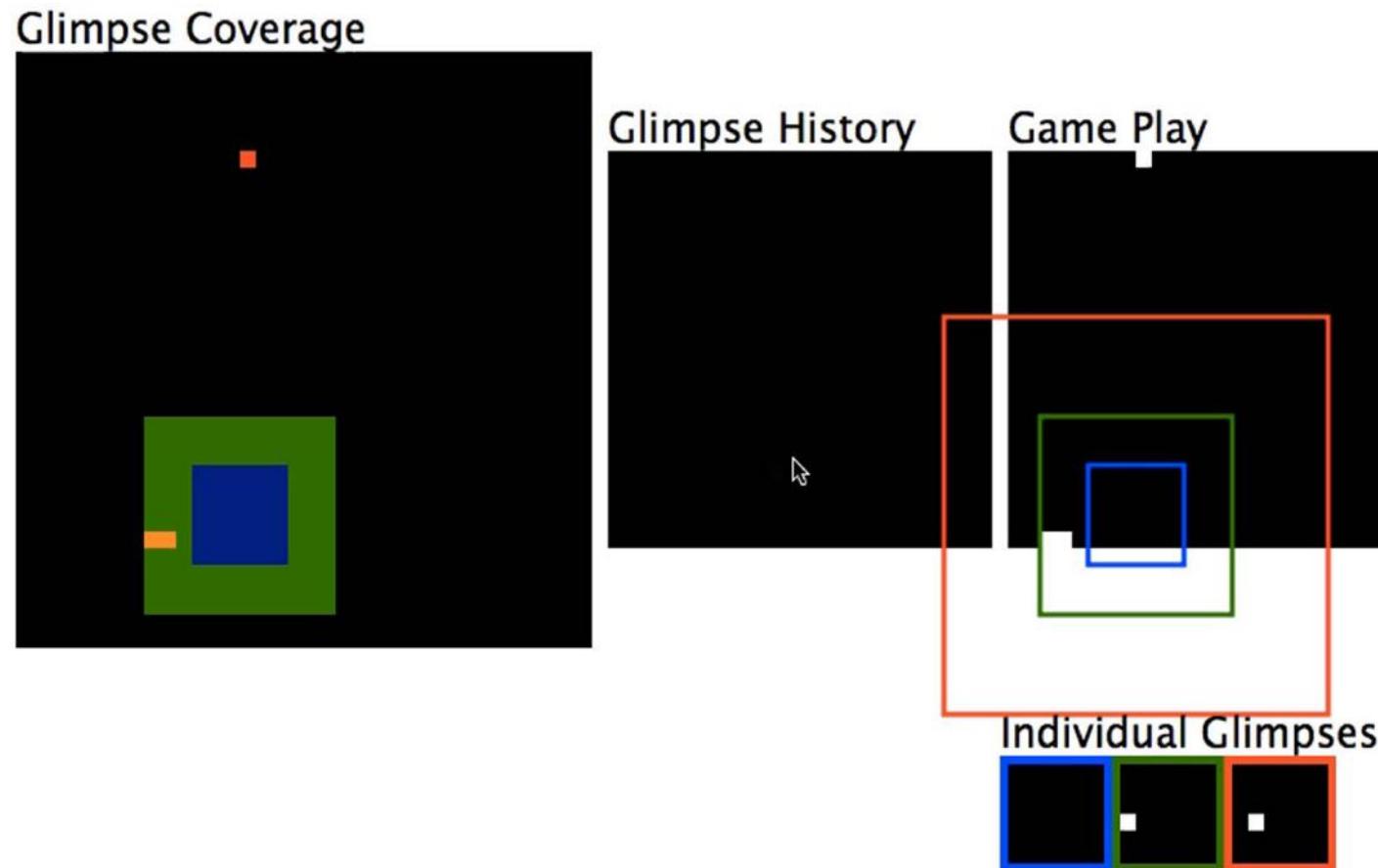


- The attention-based model achieves state-of-the-art accuracy on the SVHN multi-digit task - 3.9% error.
- 4 times fewer floating point operations than the best ConvNet.

[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.



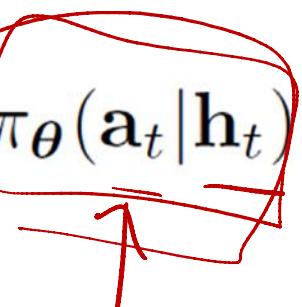
Direct policy search

history \downarrow
 $h_t = \{o_0, o_1, \dots, o_t\}$

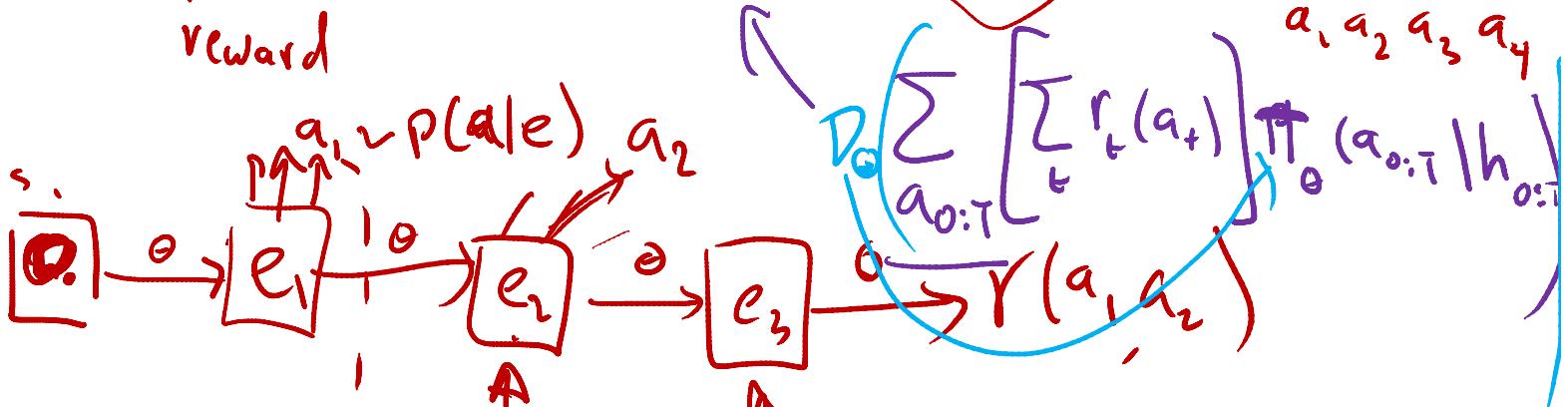
Policy
 $\pi_\theta(a_{0:T} | h_{0:T}) = \prod_{t=0}^T \pi_\theta(a_t | h_t)$

history

$$J^\pi(\theta) = \mathbb{E} \left[\sum_{t=0}^T r_t(a_t) \right] = \mathbb{E}_{\pi_\theta(a_{0:T} | h_{0:T})} \left[\sum_{t=0}^T r_t(a_t) \right]$$



expected returns



$$h_0 = o_0$$

$$h_1 = (o_0, o_1)$$

$$h_2 = (o_0, o_1, o_2)$$

Policy gradients using backprop

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

$(\log \gamma)^I = \frac{\gamma}{\gamma}$

$$= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] \nabla \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T}) \pi_{\boldsymbol{\theta}}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})$$

$$= \sum_{\mathbf{a}_{0:T}} \left[\sum_{t=0}^T r_t(\mathbf{a}_t) \right] \left[\sum_{t=0}^T \nabla \log \pi_{\boldsymbol{\theta}}(\mathbf{a}_t | \mathbf{h}_{0:t}) \right] \pi_{\boldsymbol{\theta}}(\mathbf{a}_{0:T} | \mathbf{h}_{0:T})$$

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

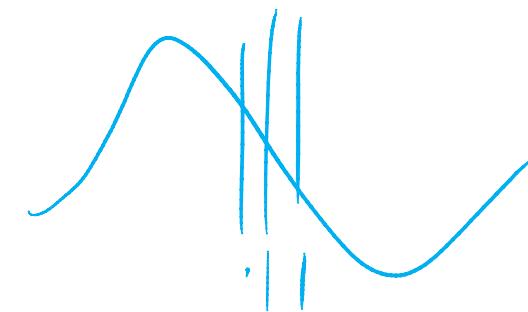
Policy gradients using backprop

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

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

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

Matt Hoffman



Neuro-dynamic programming

Dynamic programming

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

Csaba Szepevari

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

Value function

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

V depends on π
optimal

best value function

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

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

Dynamic programming

$AB + AC$

$$\begin{aligned}
 V^*(\underline{s}_0) &= \max_{\pi} \mathbb{E}_{s_1, s_2, s_3, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_0 \right] && A(B+C) \\
 &= \max_{\pi, a_0} \mathbb{E}_{s_1, s_2, s_3, \dots} \left[r_0(s_0, \underbrace{a_0}_{\pi(s_0)}, s_1) + \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_0 \right] && a_0 \in \Pi(s_0) \\
 &= \max_{a_0, \pi} \mathbb{E}_{s_1, s_2, s_3, \dots} \left[r_0(s_0, \underbrace{a_0}_{\text{a}_0}, s_1) + \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_0 \right] \\
 &= \max_{a_0} \mathbb{E}_{s_1} \left[r_0(s_0, \underbrace{a_0}_{\text{a}_0}, s_1) + \max_{\pi} \mathbb{E}_{s_2, s_3, s_4, \dots} \left\{ \sum_{t=1}^{\infty} \gamma^t r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right\} \middle| s_0 \right] \\
 &= \max_{a_0} \mathbb{E}_{s_1} \left[r_0(s_0, \underbrace{a_0}_{\pi}, s_1) + \gamma \max_{\pi} \mathbb{E}_{s_2, s_3, s_4, \dots} \left\{ \sum_{t=1}^{\infty} \gamma^{t-1} r_t(s_t, \pi(s_t), s_{t+1}) \middle| s_1 \right\} \middle| s_0 \right] \\
 &= \max_{a_0} \mathbb{E}_{s_1} [r_0(\underbrace{s_0, a_0}_{\text{a}_0}, s_1) + \gamma \underbrace{V^*(s_1)}_{\text{a}_0}]
 \end{aligned}$$

Bellman's equation and TD

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

current state *future state*

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

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

$\tilde{s}' \sim P(s'|s, a)$

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

Action-value (Q) functions

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

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

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

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

Q - Learning

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

Neuro-dynamic programming

$Q(s, a; w)$, where w are the parameters

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

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

target

Neuro-dynamic programming

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

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

↑

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

~~s_t~~

① state s

② choose a

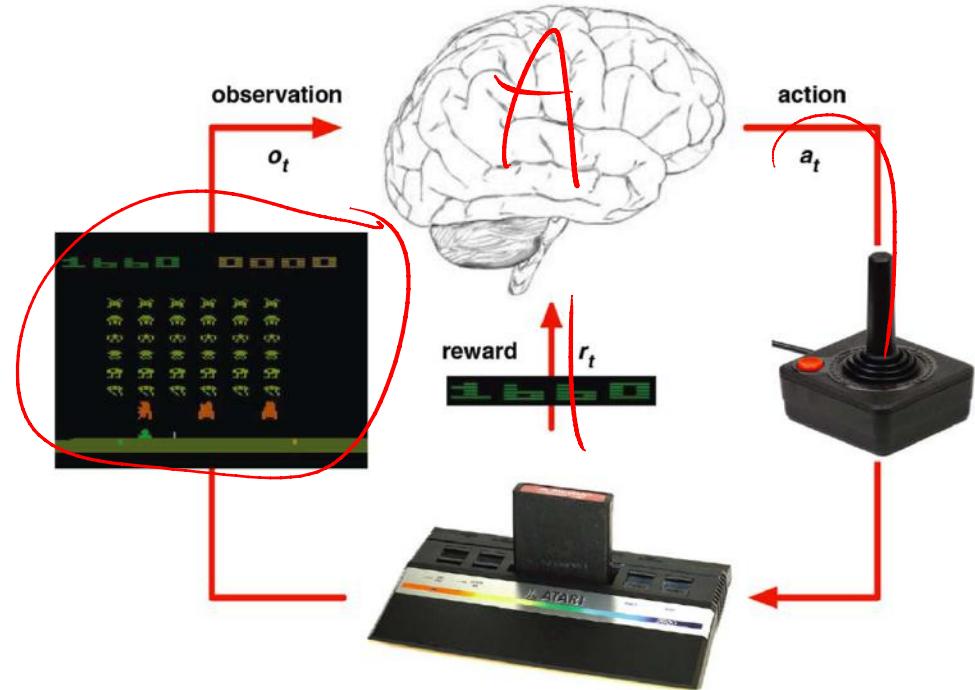
③ Environment $\rightarrow s'$

④ Update ~~\mathbf{w}~~ \mathbf{w}

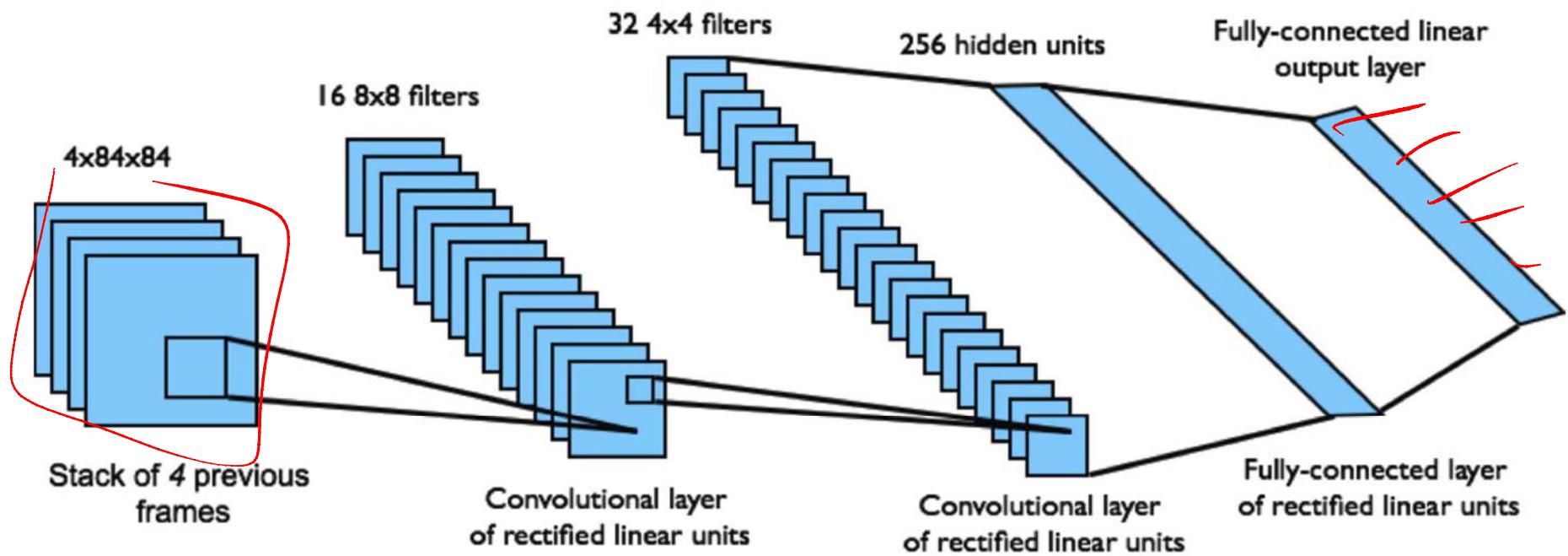
$$\begin{cases} \text{w.p. } \epsilon & a = \arg \max_{\mathbf{a}'} Q(s, \mathbf{a}') \\ & \mathbf{a}' \\ \text{w.p. } 1-\epsilon & a \sim U_{[1]} \end{cases}$$

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.

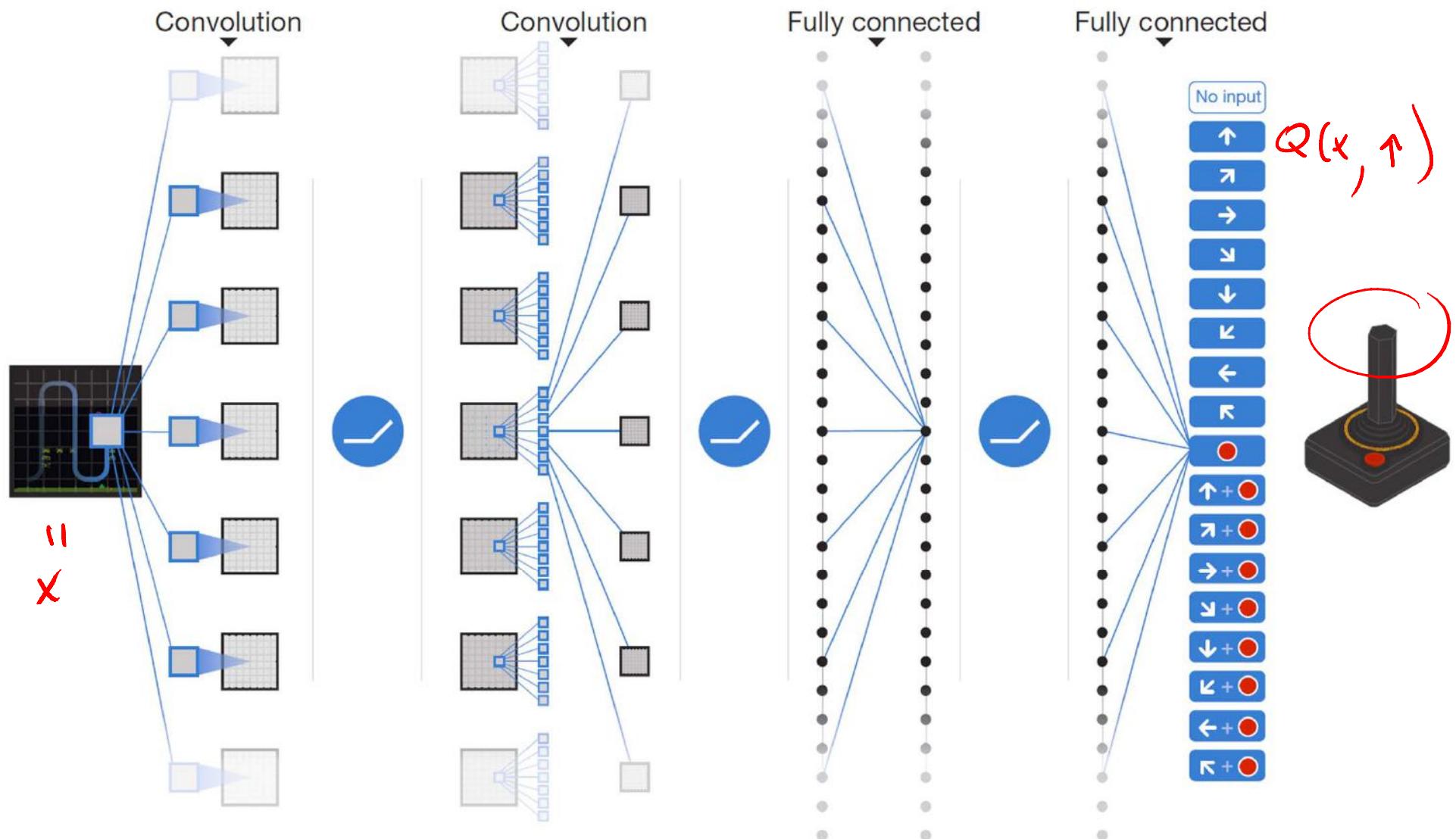


- ▶ End-to-end learning of values $Q(s, a)$ from pixels s
- ▶ Input state s is stack of raw pixels from last 4 frames
- ▶ Output is $Q(s, a)$ for 18 joystick/button positions
- ▶ Reward is change in score for that step



Network architecture and hyperparameters fixed across all games

DQN Convolutional Network



Q-learning with experience re-play

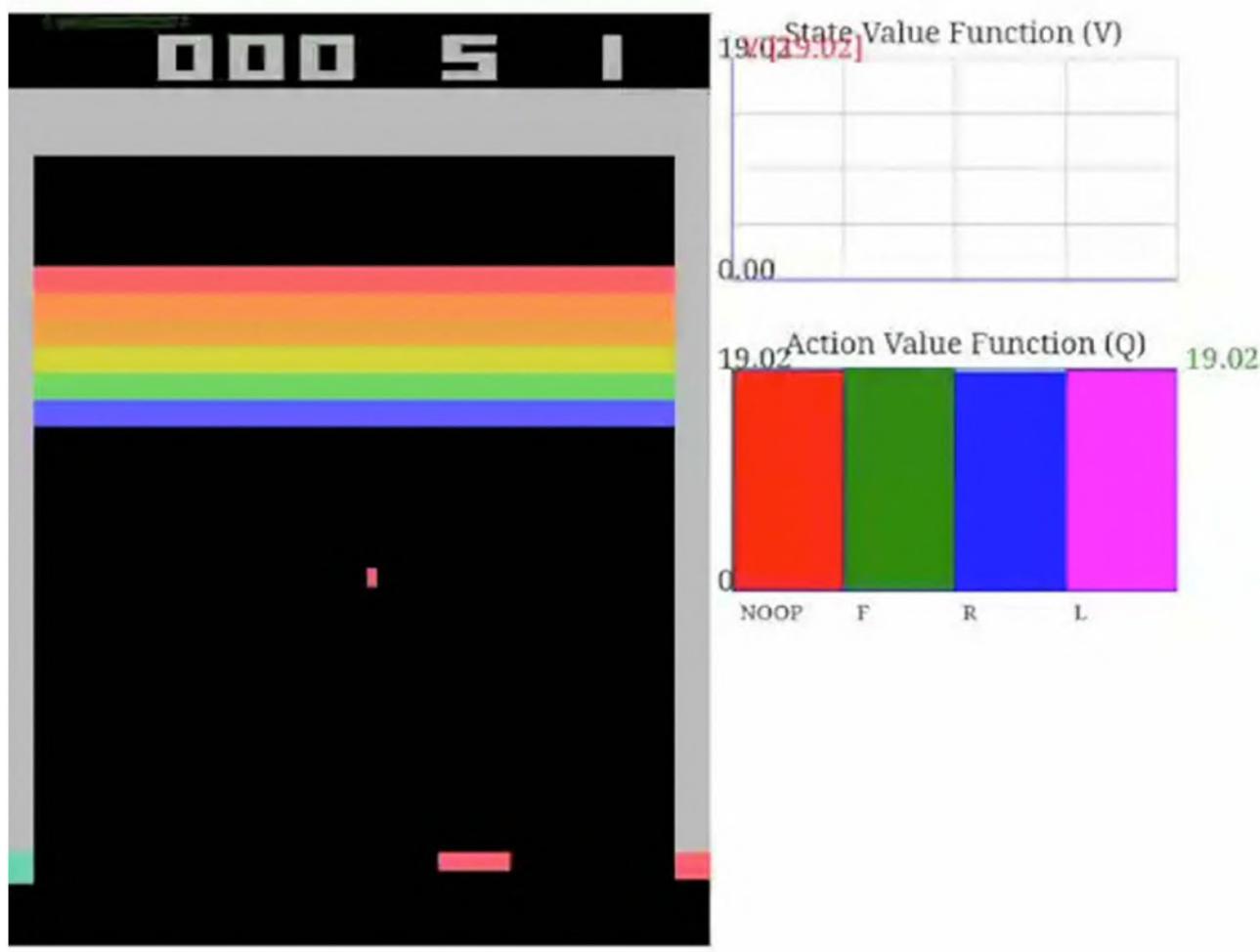
DQN increases stability with **experience replay** and **fixed Q-targets**

- ▶ Take action $\underline{a_t}$ according to ϵ -greedy policy
- ▶ Store transition $(\underline{s_t}, \underline{a_t}, \underline{r_{t+1}}, \underline{s_{t+1}})$ 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

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

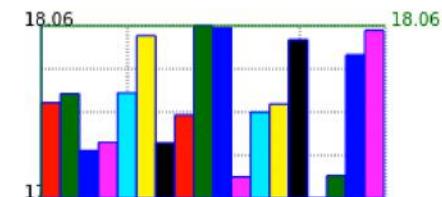
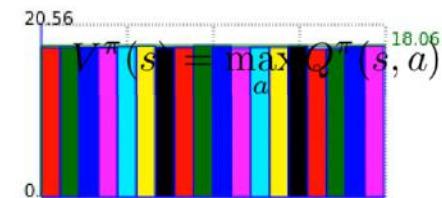
- ▶ Using variant of stochastic gradient descent (Mnih et al.)

Delayed Rewards



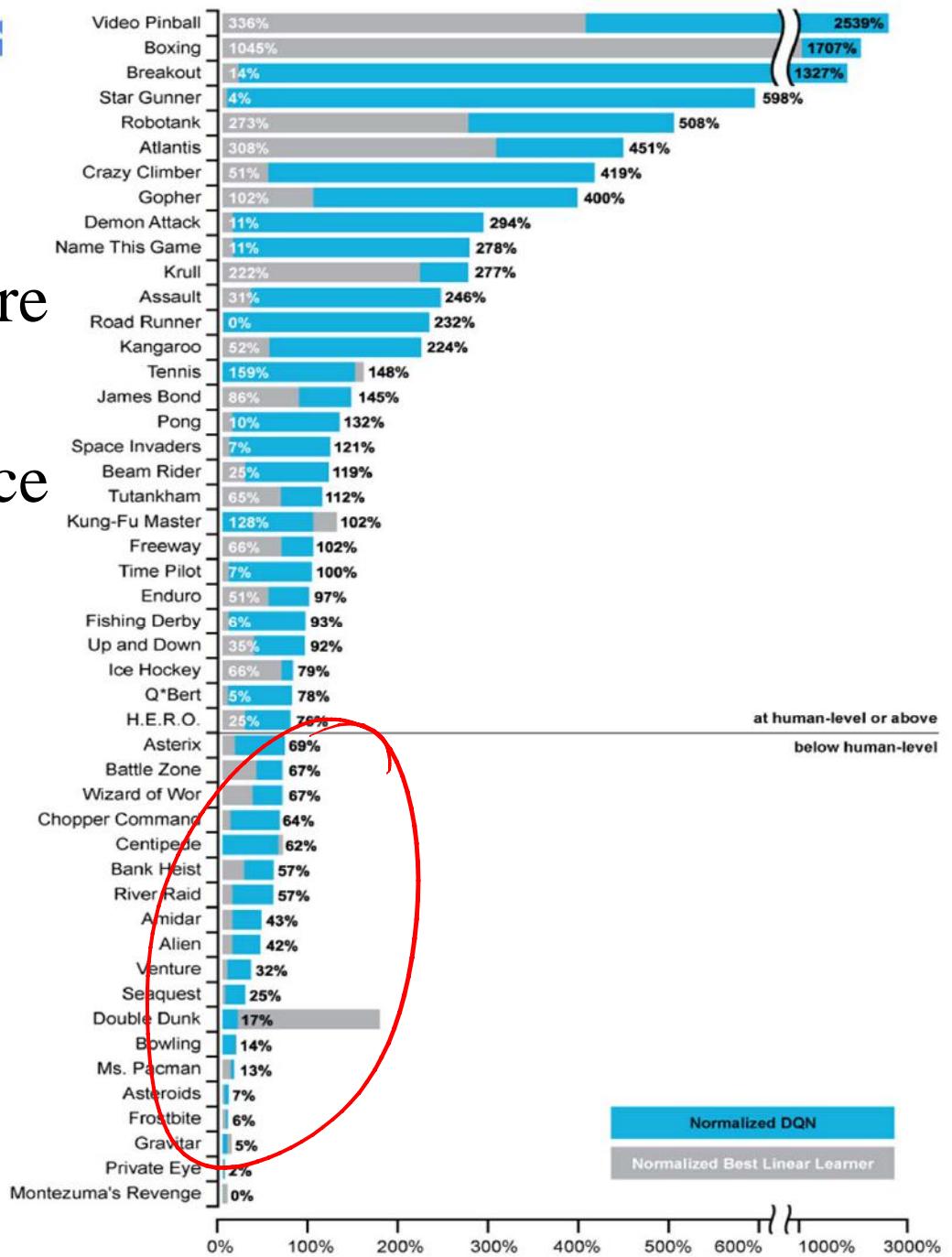
Sacrificing Immediate Rewards

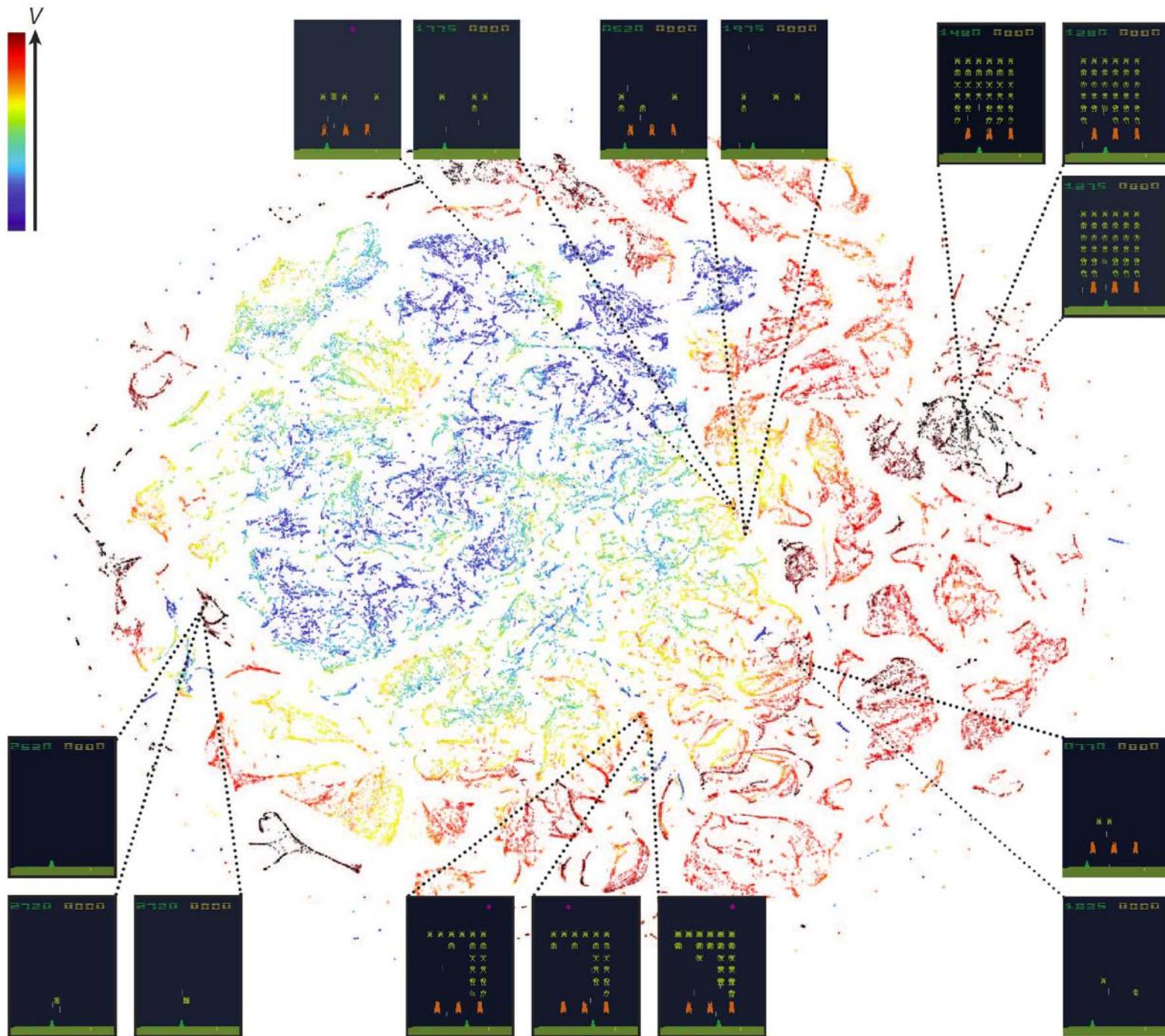
7140 18.198350906372



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.





DQN
AI
And all that



Thank you