

Function Approximation, SGD, DQN

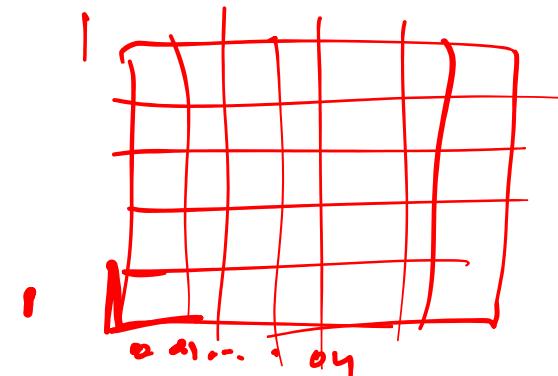
B. Ravindran

Q	v	s_1	$\hat{v}(s_1)$
a_1	a_2	a_3	\vdots
$s_1, \hat{q}(s_1, a_1)$	$\hat{v}(s_1, a_2)$	\dots	s_N
$s_2, \hat{q}(s_2, a_1)$	\dots	\dots	$\hat{v}(s_N)$
s_3	\vdots		
\vdots	\vdots		
			n

Table

Need for Function Approximation

- Issues with large state/action spaces:
 - tabular approaches not memory efficient
 - data sparsity ↗
 - continuous state/action spaces ↗
 - generalization ↗
- Use a parameterized representation
 - Value Functions ✓
 - Policies ✓
 - Models ↗ $p(s', a | s, a)$



linear regression.

$$f(x) = w_1 x_1 + w_2 x_2 + \dots + w_p x_p.$$
$$f(x; \omega)$$
$$f(x; \Theta)$$

Value Function Approximation

- Least squares

$$Q(s_t, a_t) = f(s_t, a_t; \underline{w}_t)$$

$$w_{t+1} = w_t - \frac{1}{2} \alpha \nabla_{w_t} [q_*(s_t, a_t) - Q(s_t, a_t)]^2$$

- But we don't know the target!

- Use the TD target.

$$w_{t+1} = w_t - \frac{1}{2} \alpha \nabla_{w_t} [r_{t+1} + \gamma \max_a Q(s_t, a) - Q(s_t, a_t)]^2$$

Semi Gradient Methods

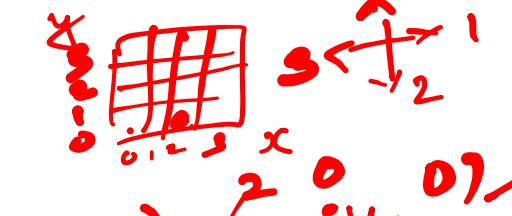
- While computing the gradient of the TD error in Q-learning, we typically ignore the gradient w.r.t the TD target. Hence, it is a **Semi Gradient** method i.e we are computing an approximation of the true gradient.
- Eg: **Linear Functions:**

$$Q(s_t, a_t) = \phi^T(s_t, a_t) \times w_t$$
$$\phi(s_t, a_t) = [\phi_1(s_t, a_t), \phi_2(s_t, a_t), \phi_3(s_t, a_t)]$$

$$\delta_t = r_{t+1} + \gamma \max_a Q(s_t, a) - Q(s_t, a_t)$$

TD Error

$$\nabla_{w_t} [r_{t+1} + \gamma \max_a Q(s_t, a) - Q(s_t, a_t)]^2 = -2\delta_t \phi(s_t, a_t)$$
$$w_{t+1} = w_t + \alpha \delta_t \phi(s_t, a_t)$$


$$\phi(s, a) = [\phi_1(s, a), \phi_2(s, a), \phi_3(s, a)]$$
$$\phi(s, a) = [1, 0, 2, 0, 1, 2, 0, 0]$$

Stochastic Gradient Descent

- In **Stochastic Gradient Descent**, the true gradient of the loss function is approximated by the gradient at a single example.
- In practice, we usually perform **Mini Batch Gradient Descent**, where we compute the gradient using a mini batch of examples. This allows for more efficient computation and smoother convergence.
 - Still SGD

Tile and Coarse coding

Can exploit the possibility that q values of near-by states wouldn't change a lot.

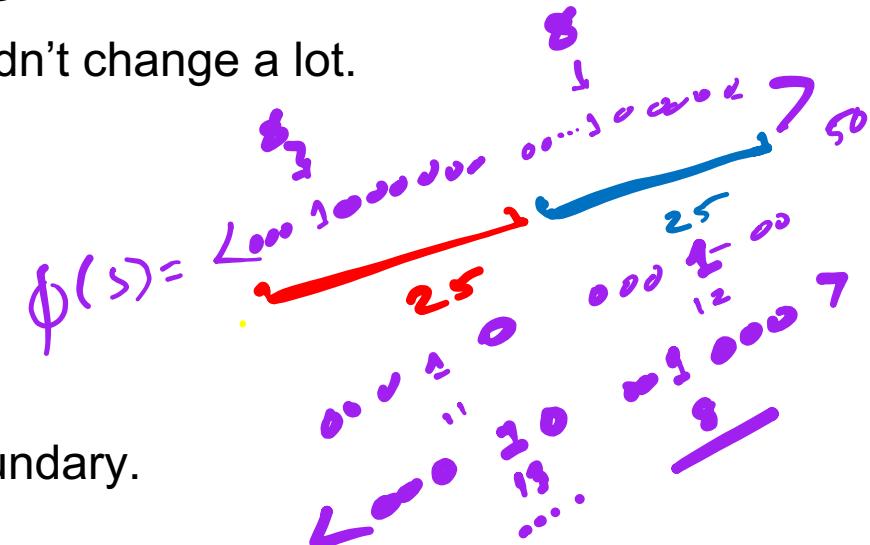
Assume the task is to learn a policy on a big gridworld.

Tabular Lookup!

Naive method:- *State aggregation*

Divide the grid into smaller 10x10 grid.

Abrupt change in Q-value of states that lie on the boundary.



Coarse coding:-

Avoids abrupt changes in the value of the state. Smoothens the qvalues during transition from one cluster to another.

Issue:-

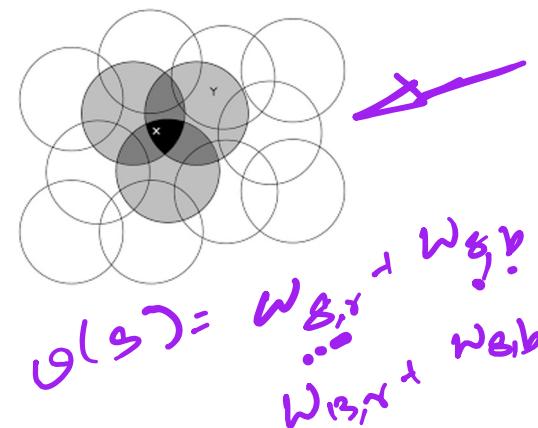
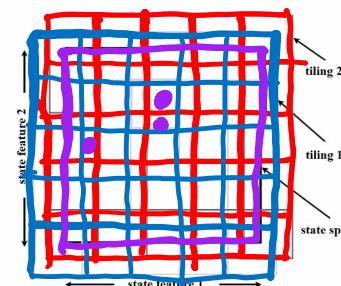
No uniformity in the number of 'ON' bits used to represent a state.

Tile coding:-

Form of coarse coding but systematic.

Number of 'ON' bits == Number of tiles used.

Randomly shift the tiling.

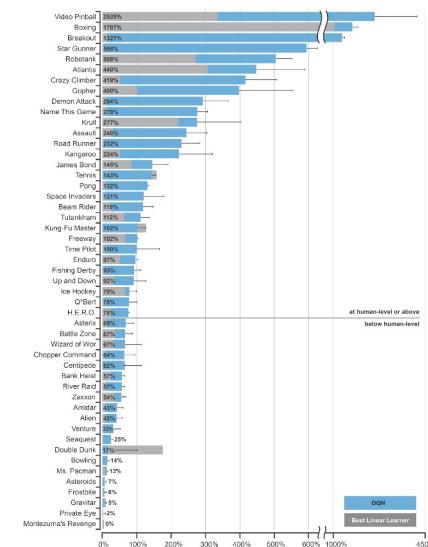
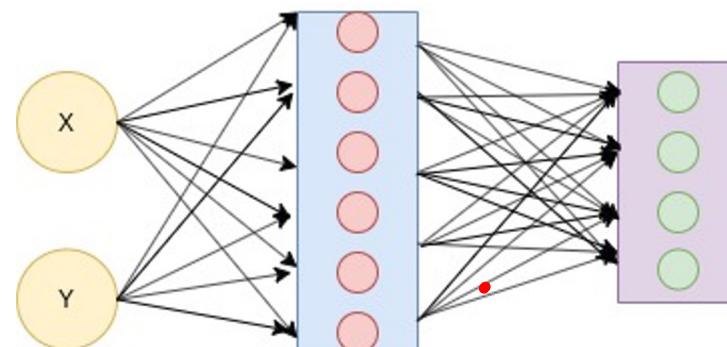


Non-Linear Function Approximator

- Somewhat*
1. Linear function approximators are ~~very~~ restrictive. Can only model linear functions.
Basis expansion does help to generate non-linear functions in the original input space.
 2. Non-linear approximators can model complex functions and are very powerful.
 3. The features are learnt on the fly and are not hardcoded as is the case with tile and sparse coding.
 4. Can generalize to unseen states.

Disadvantage:-

Requires a lot of data and compute.



Human Level Backgammon player

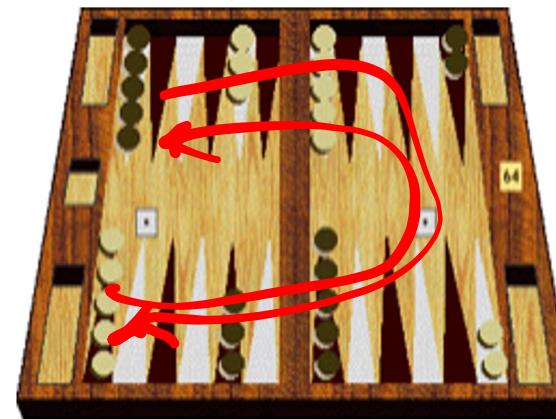
TD-Gammon (Tesauro 92, 94, 95)

Beat the best human
player in 1995



Learnt completely
by *self play*

New moves not recorded by
humans in centuries of play



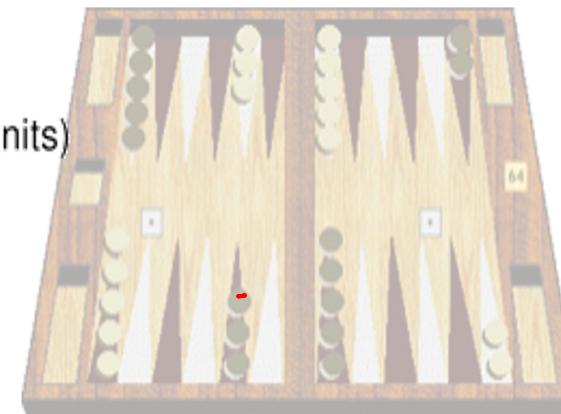
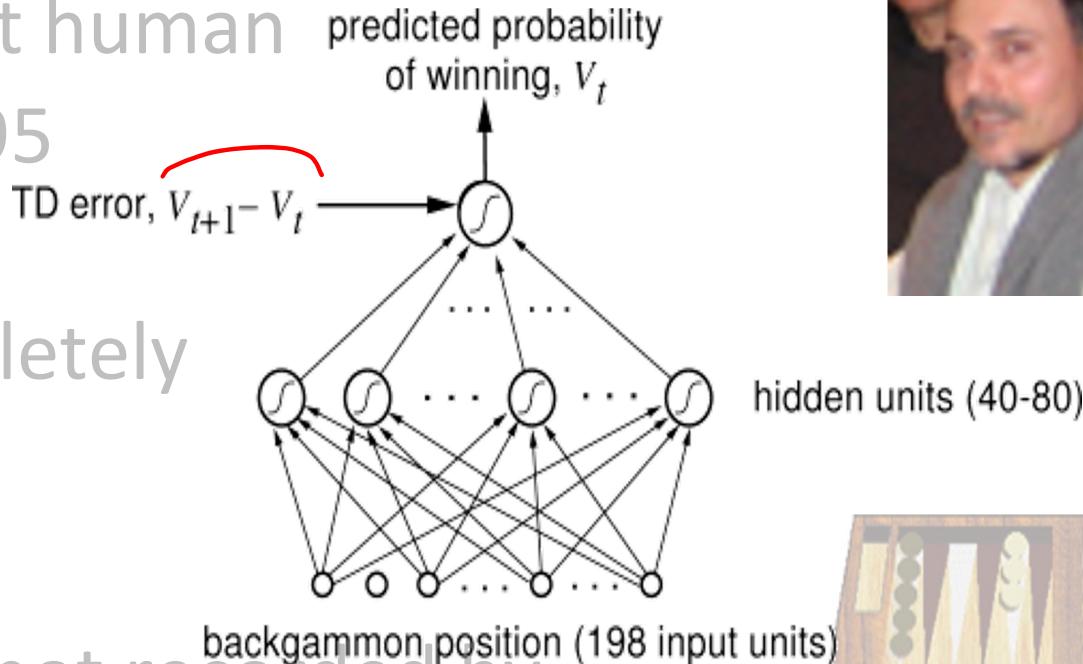
Human Level Backgammon player

TD-Gammon (Tesauro 92, 94, 95)

Beat the best human
player in 1995

Learnt completely
by *self play*

New moves not recorded by
humans in centuries of play



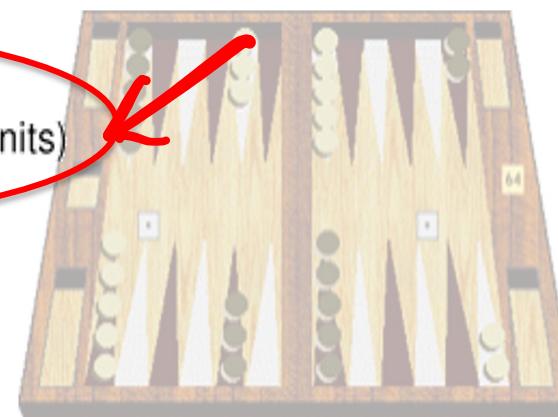
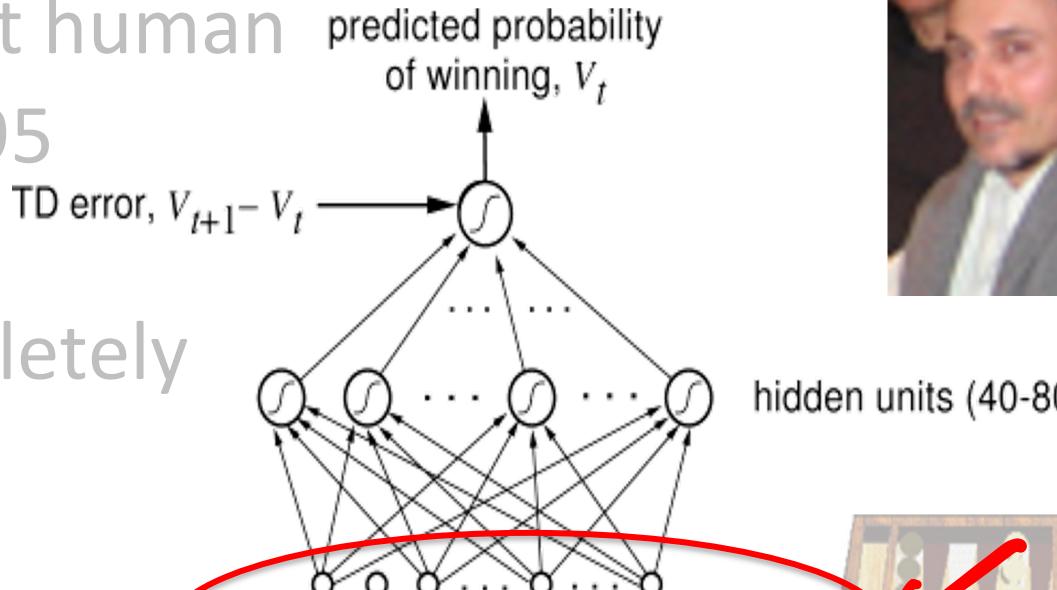
Human Level Backgammon player

TD-Gammon (Tesauro 92, 94, 95)

Beat the best human
player in 1995

Learnt completely
by *self play*

New moves not recorded by
humans in centuries of play



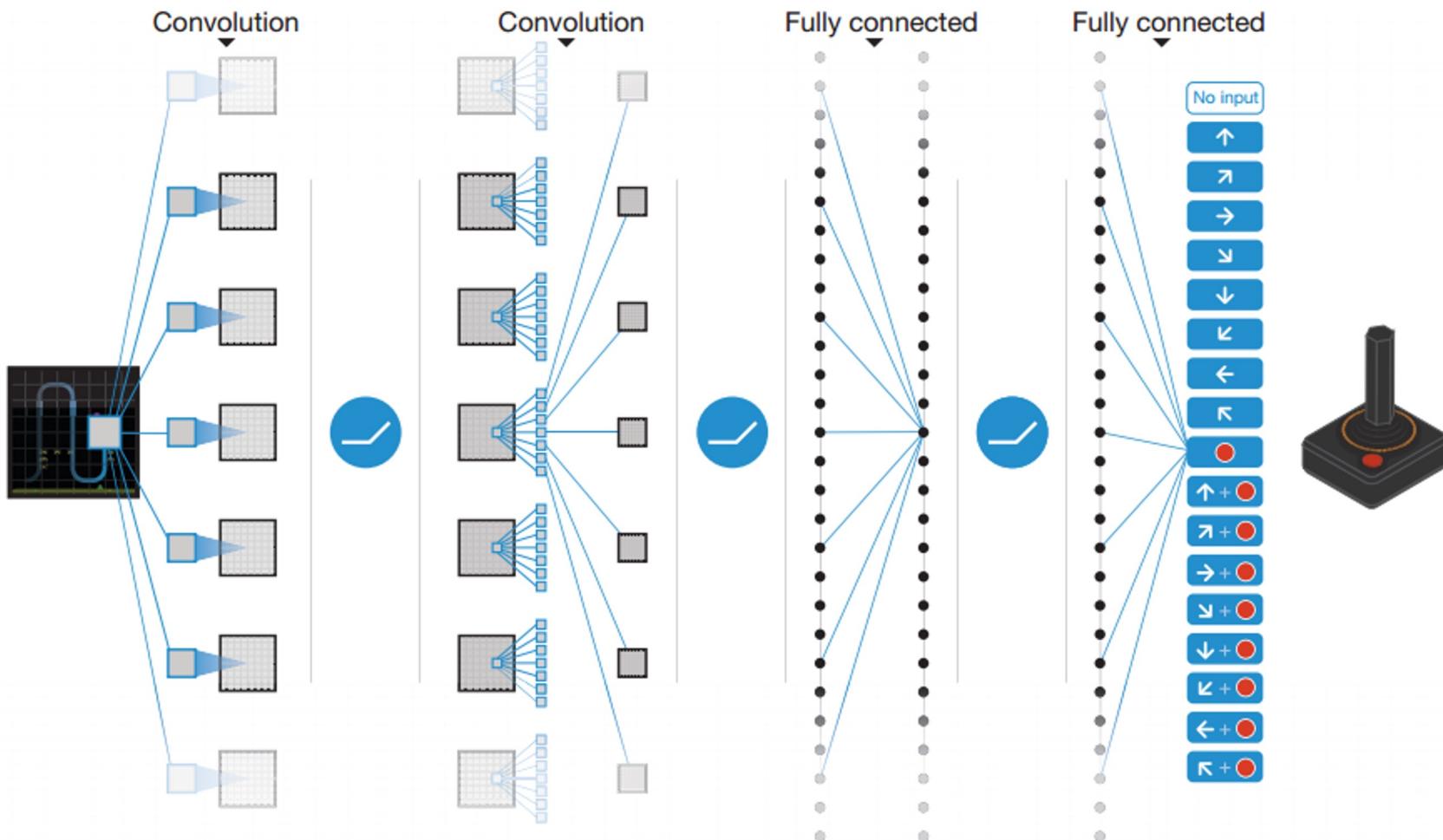
What about the features?



18 actions

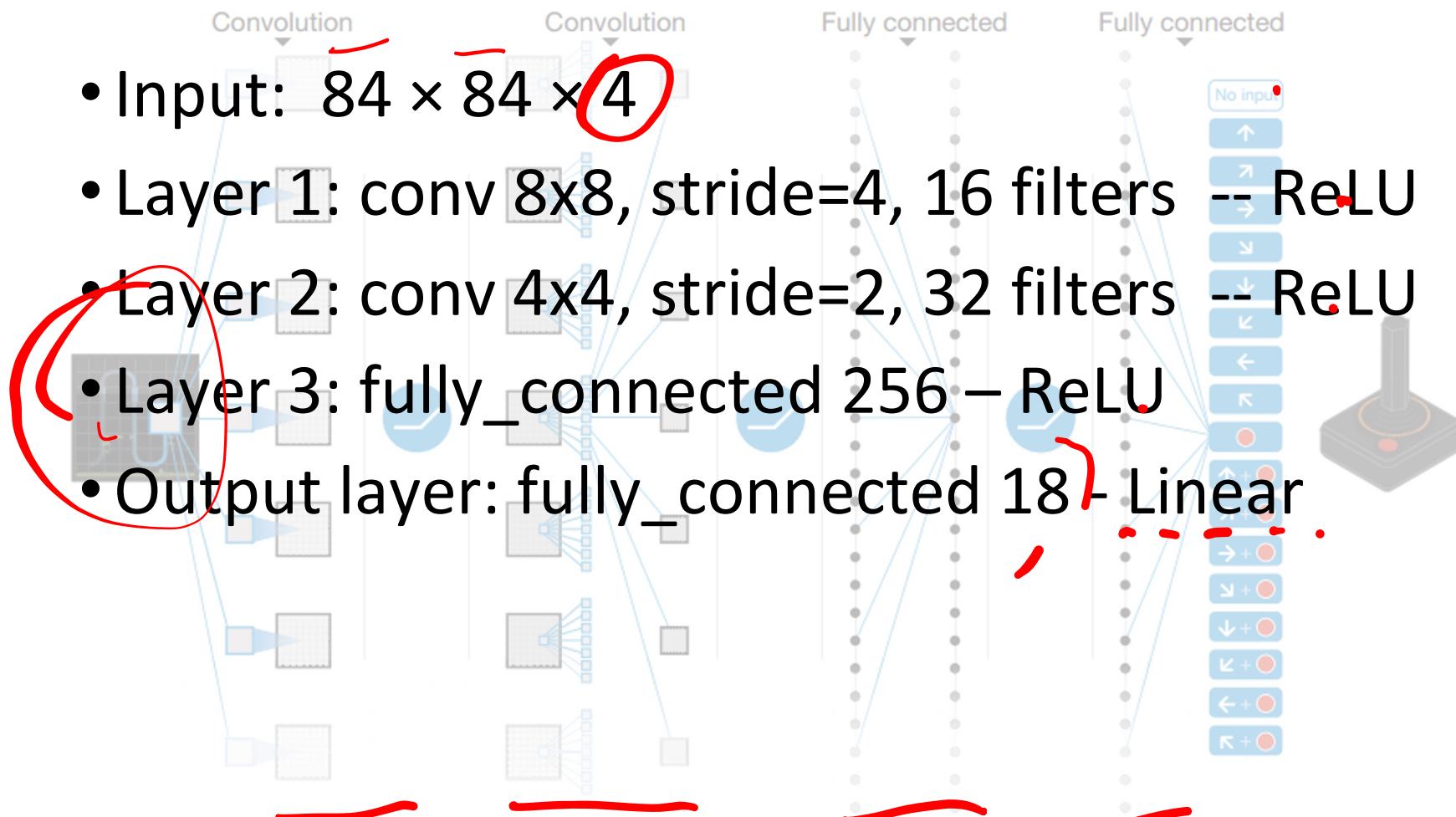
- Learnt to play from video input from scratch!

Deep Q-Learning Network (DQN)



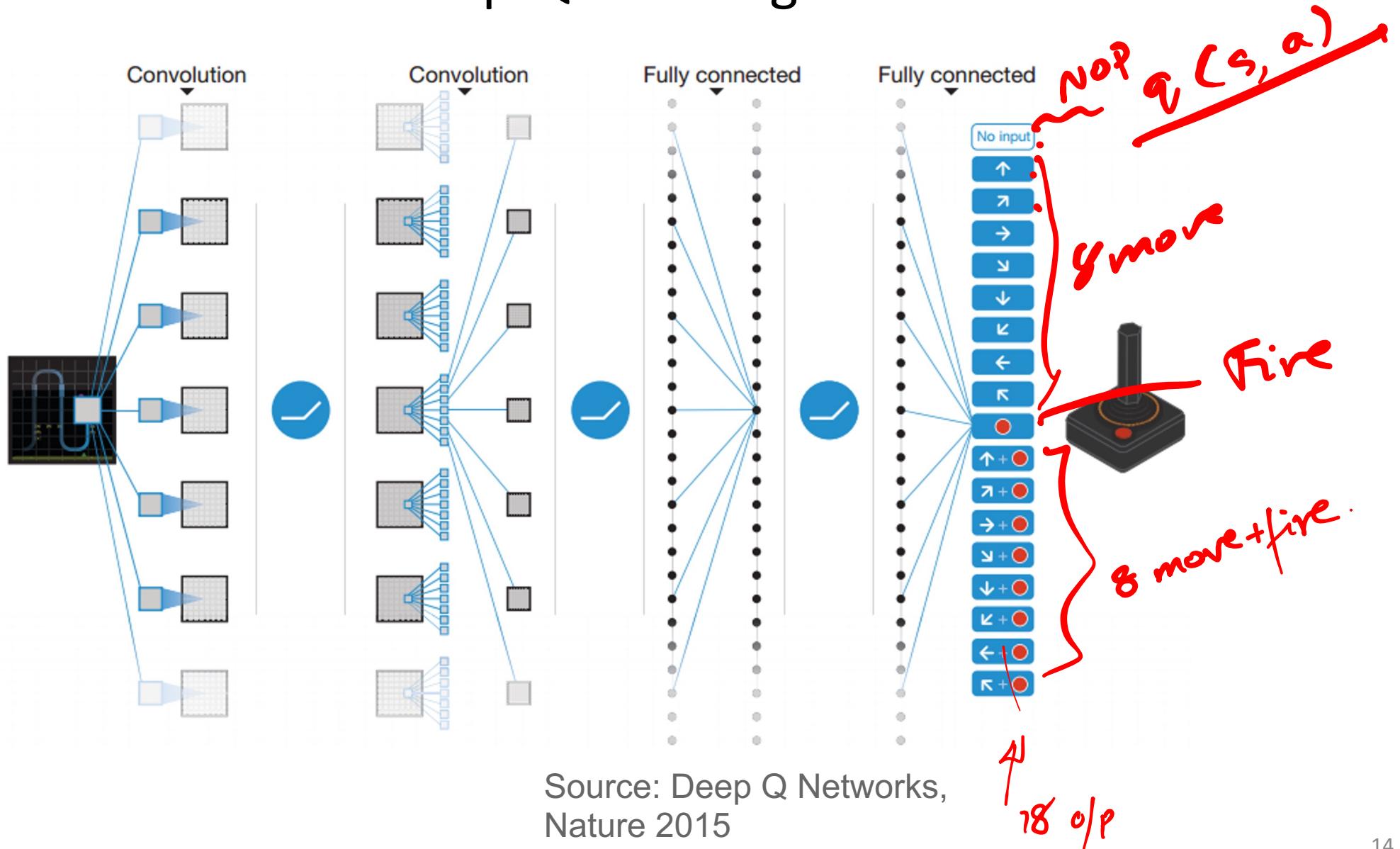
Source: Deep Q Networks,
Nature 2015

Deep Q-Learning



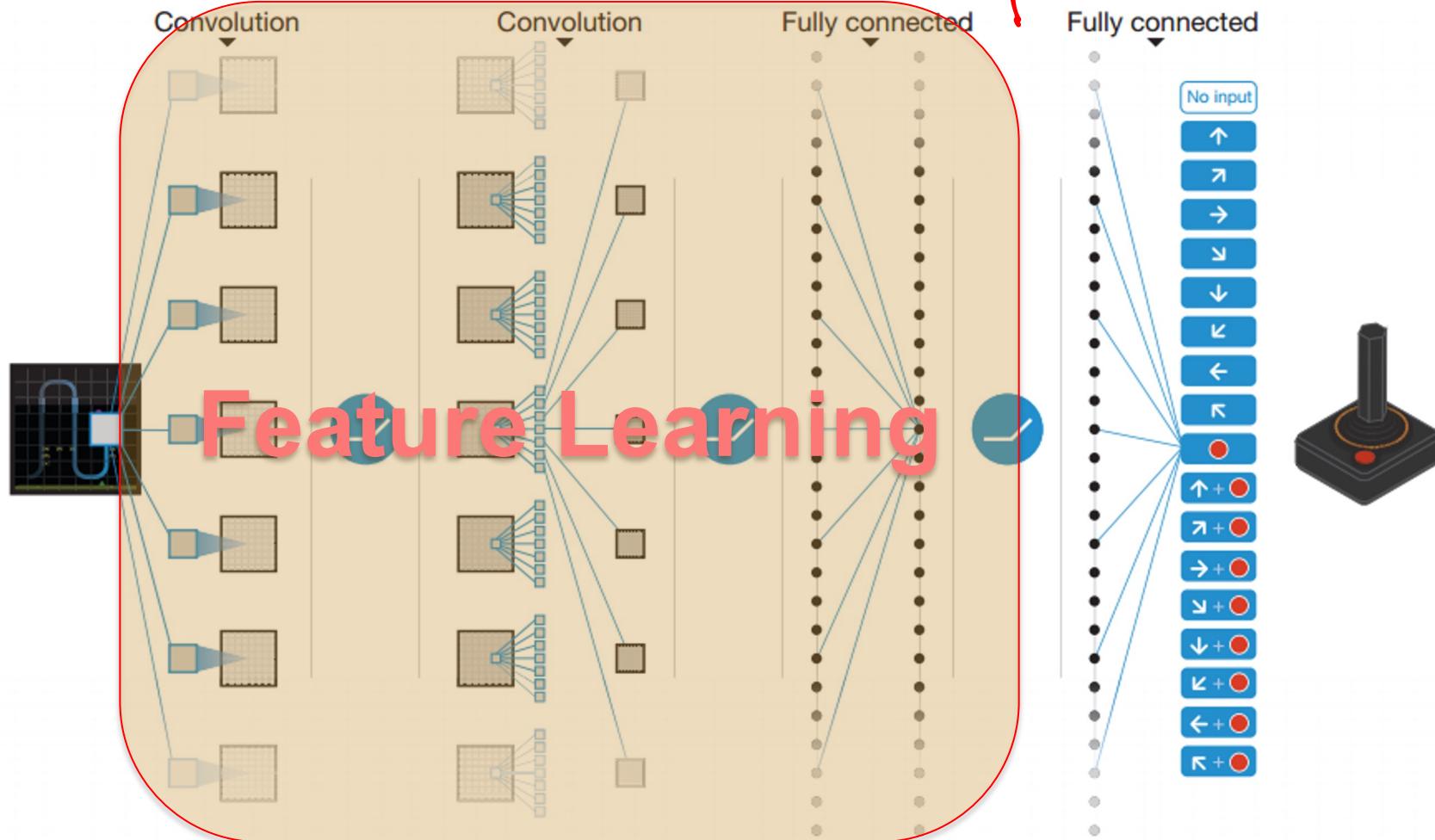
Source: Deep Q Networks,
Nature 2015

Deep Q-Learning



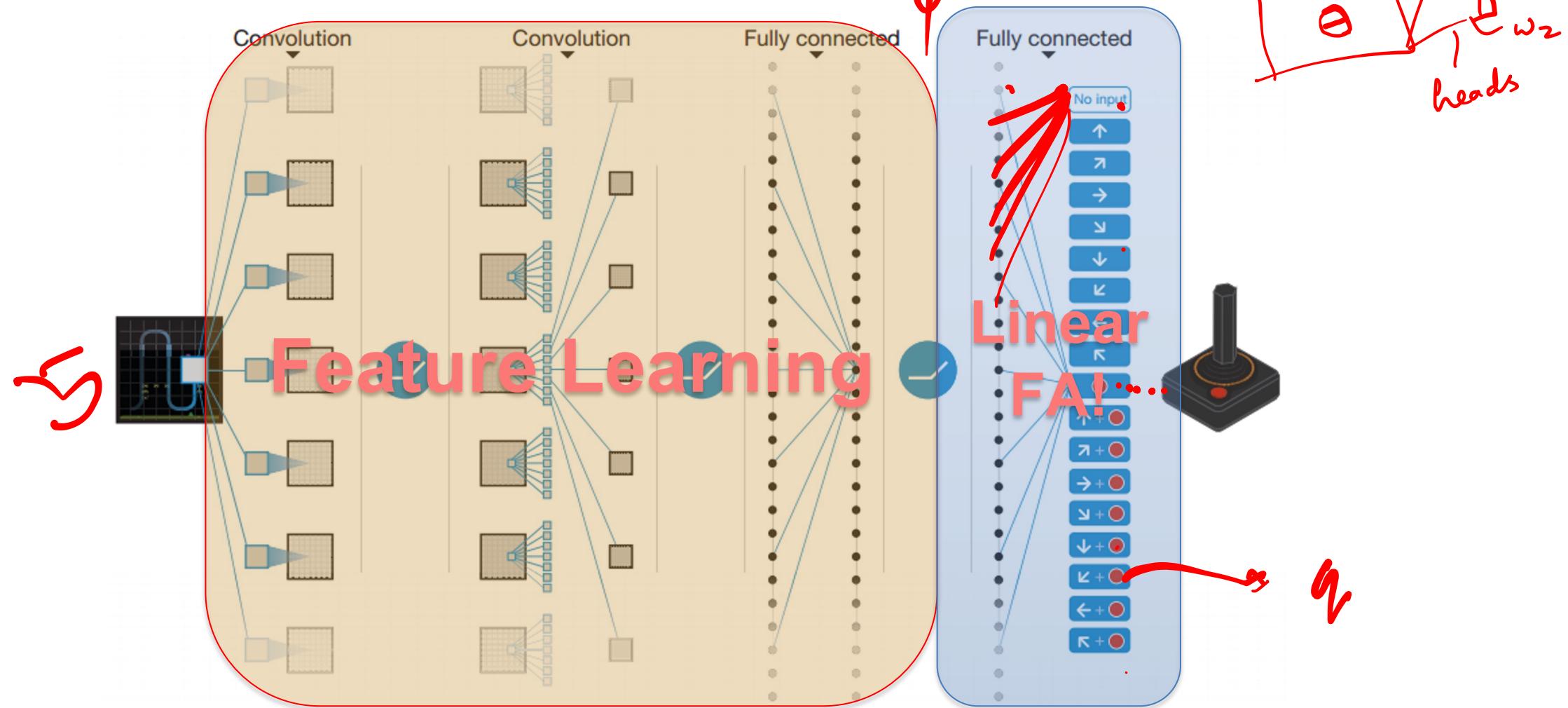
Deep Q-Learning

ϕ



Source: Deep Q Networks,
Nature 2015

Deep Q-Learning



Q-Network Learning

Backprop.

$$w_{t+1} = w_t - \frac{1}{2}\alpha \nabla_{w_t} [r_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a) - \hat{q}(s_t, a_t)]^2$$

Divergence is an issue since the current network is used to decide its own target

- Correlations between samples
- Non-stationarity of the targets
- How do we address these issues?
 - Replay Memory
 - Freeze target network

$$y = \hat{q}(s_t, a_t)$$

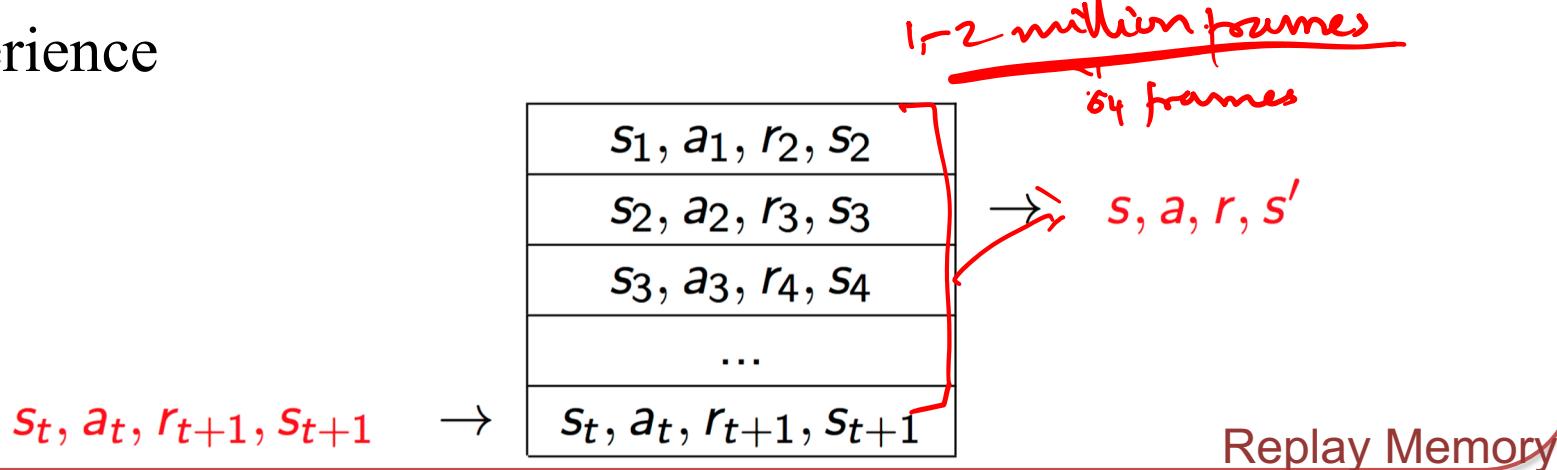
$$\nabla_w (t - y)^2$$

$\langle (s, a) \rangle$

$$\left[r_t + \gamma \max_a \hat{q}(s'_t, a') \right]$$

Q-Network Learning

To remove correlations, we build data-set from the agent's experience



Sample experiences from dataset; w^- frozen (with periodic updates) to address non-stationarity

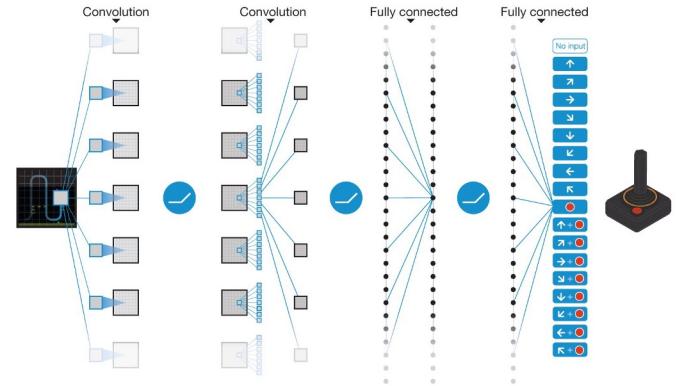
$$\left[\left\{ r_{t+1} + \gamma \max_a \hat{q}(s_{t+1}, a; w^-) \right\} - \hat{q}(s_t, a_t; w) \right]^2$$

Freeze Target Network

Deep-Q Networks

Architecture:-

1. Has a set of convolutional layers which act as feature extractors.
2. These features are then passed through a series of fully connected layers.
3. The output layer has $|A|$ number of nodes which are used to calculate the Q-value for each action.



The above network is updated using huber loss and not regular least squares loss.

Below is the regular expression for huber loss.

$$L_\delta(a) = \begin{cases} \frac{1}{2}a^2 & \text{for } |a| \leq \delta, \\ \delta(|a| - \frac{1}{2}\delta), & \text{otherwise.} \end{cases}$$

