

# ELEC 677: Deep Reinforcement Learning Lecture 11

**Ankit B. Patel, CJ Barberan**

*Baylor College of Medicine (Neuroscience Dept.)*

*Rice University (ECE Dept.)*

11-22-2016

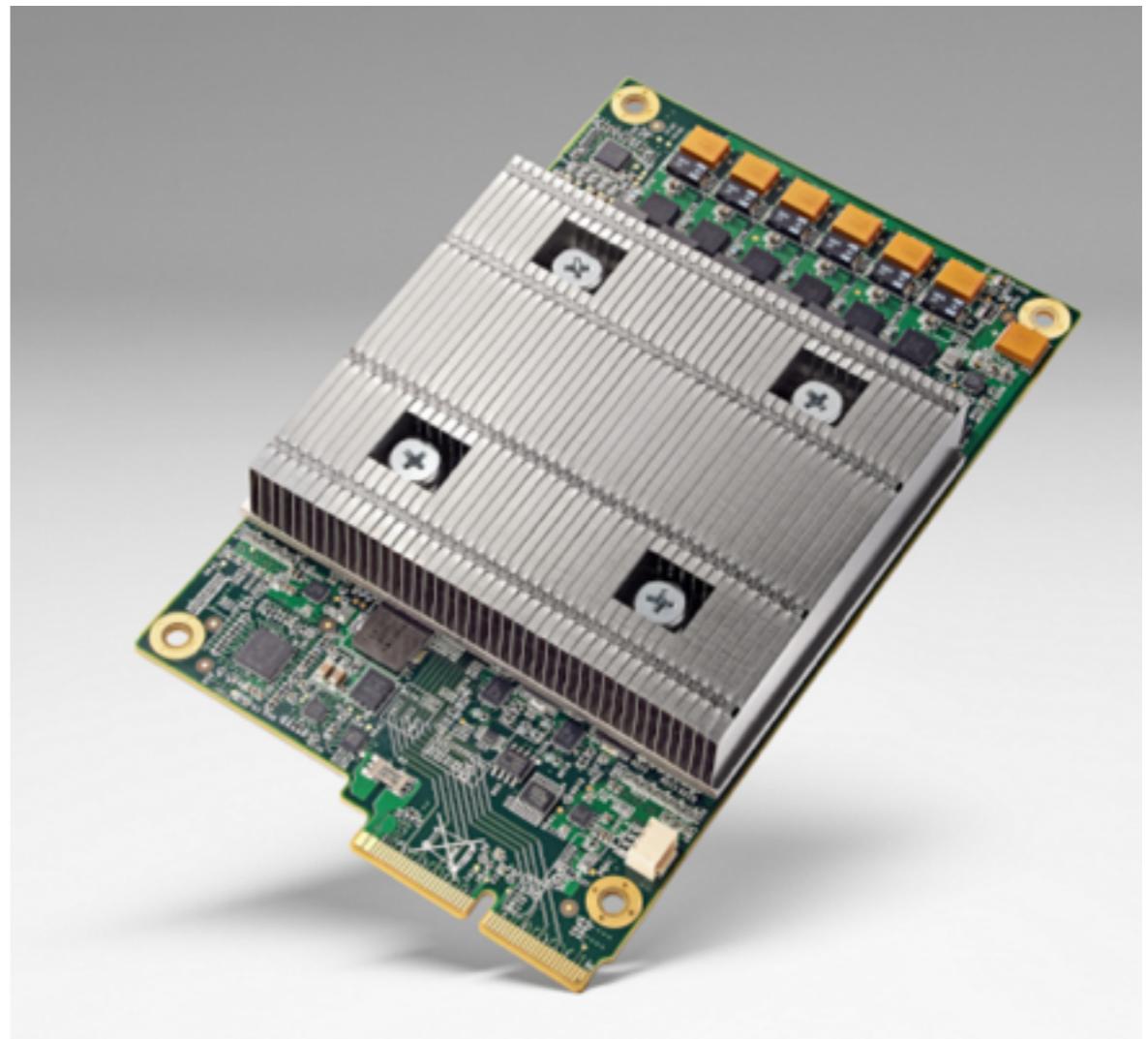
# Latest News

# Google Cloud Machine Learning Family Grows

- Two world-renowned researchers will lead the group
  - Fei-Fei Li
    - Director of Stanford Artificial Intelligence Lab & Stanford Vision Lab
  - Jia Li
    - Head of Research at Snapchat

# Google Cloud Machine Learning New API

- GPUs on the Cloud
  - AMD FirePro S9300 x2
  - Tesla P100
  - Tesla K80s
- Cloud Vision API can use Tensor Processing Unit (TPU) boards
  - Reduce prices



# New AI Research Unit in Montreal

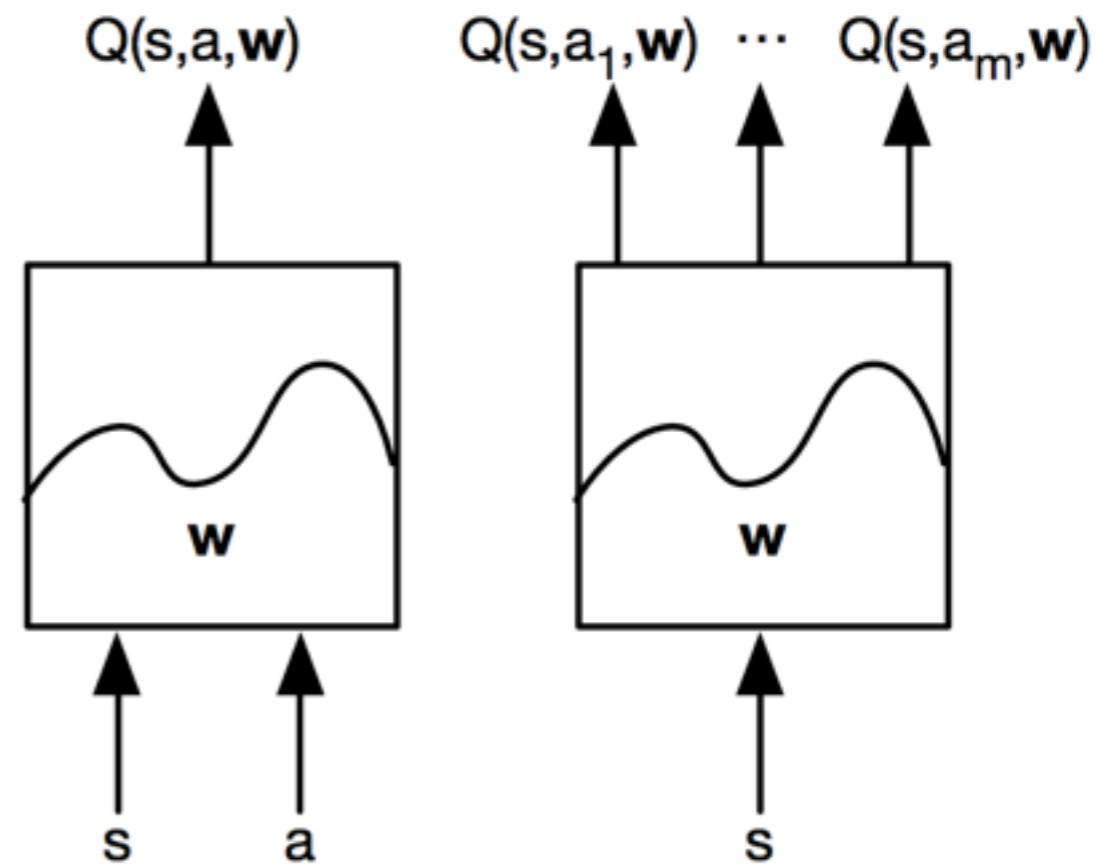
- Google is giving \$3.37 million in grant money to Yoshua Bengio and seven others in the MILA group
- Google AI research group will be supervised by Hugo Larochelle
  - Hugo worked on Twitter's Cortex deep learning division

# Value-Based Deep Reinforcement Learning

# Value-Based Deep RL

- The value function is represented by Q-network with weights  $\mathbf{w}$

$$Q(s, a, \mathbf{w}) \approx Q^*(s, a)$$



# Q-Learning

- Optimal Q-values obey Bellman Equation

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

- Minimise MSE loss via SGD

$$l = \left( r + \gamma \max_a Q(s', a', \mathbf{w}) - Q(s, a, \mathbf{w}) \right)^2$$

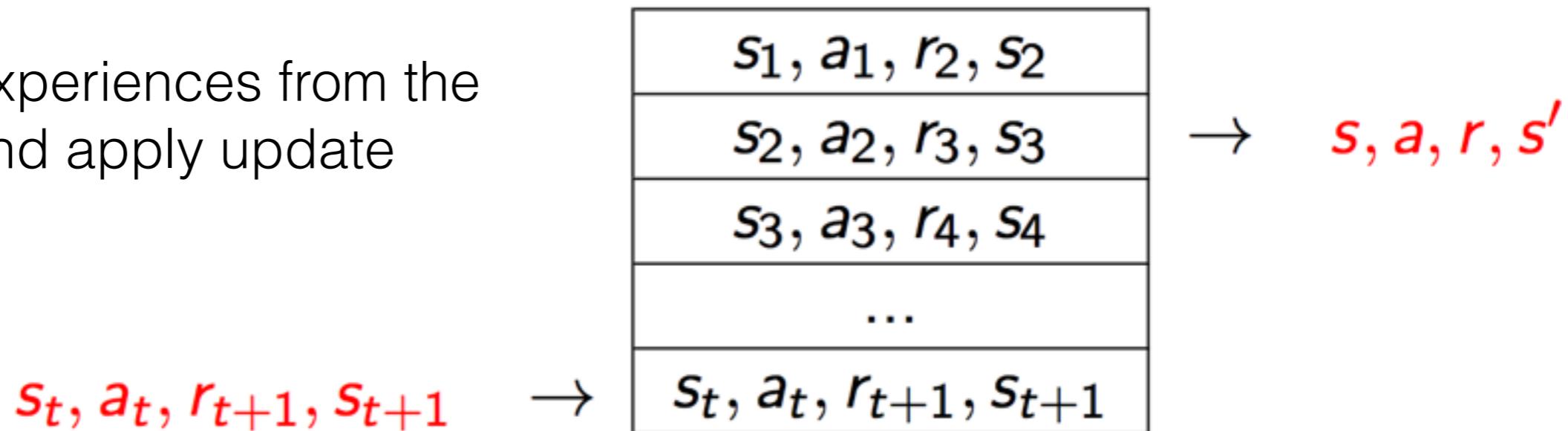
- Diverges due to
  - Correlations between samples
  - Non-stationary targets

# Interesting Question

- What can go wrong in reinforcement learning that doesn't go wrong in standard supervised learning?

# Experience Replay

- Build dataset from agent's own experience
- Sample experiences from the dataset and apply update



# Benefits of Experience Replay

- Greater data efficiency
- Breaks the correlations
- Smoothing out learning and avoiding oscillations or divergence in parameters

# DQN w/ Experience Replay

**Algorithm 1: deep Q-learning with experience replay.**

Initialize replay memory  $D$  to capacity  $N$

Initialize action-value function  $Q$  with random weights  $\theta$

Initialize target action-value function  $\hat{Q}$  with weights  $\theta^- = \theta$

**For** episode = 1,  $M$  **do**

    Initialize sequence  $s_1 = \{x_1\}$  and preprocessed sequence  $\phi_1 = \phi(s_1)$

**For**  $t = 1, T$  **do**

        With probability  $\varepsilon$  select a random action  $a_t$

        otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $D$

        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $D$

        Set  $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

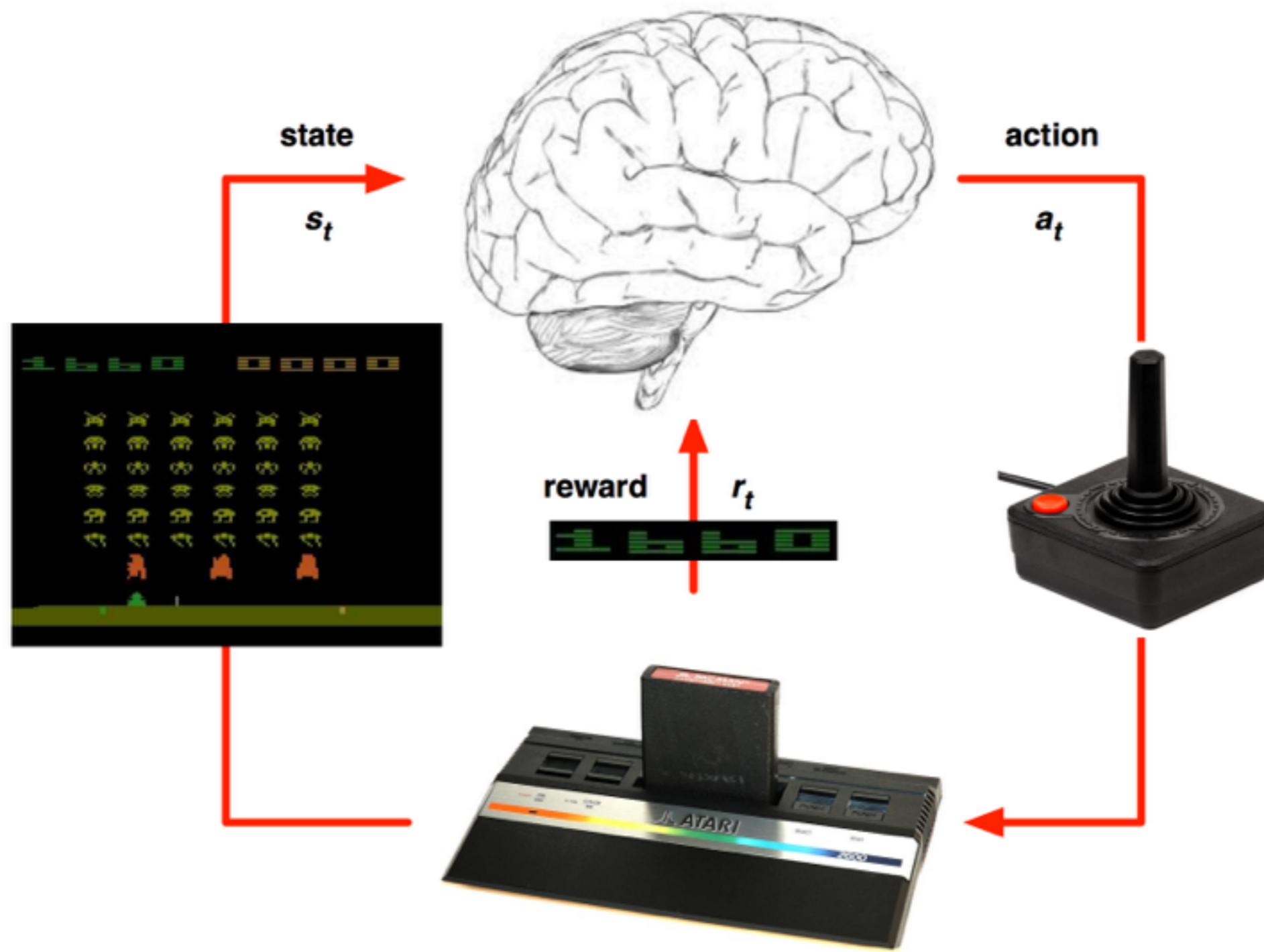
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  with respect to the network parameters  $\theta$

        Every  $C$  steps reset  $\hat{Q} = Q$

**End For**

**End For**

# DQN: Atari 2600



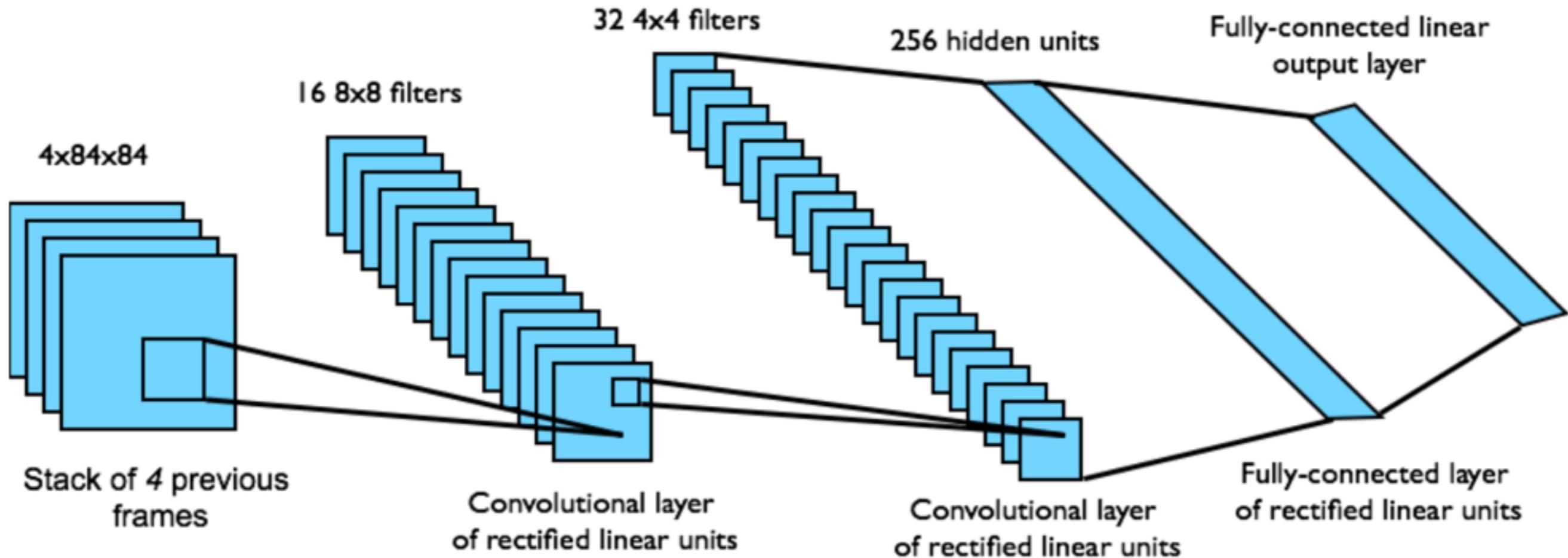
[David Silver ICML 2016]

# DQN: Atari 2600

- End-to-end learning of  $Q(s,a)$  from the frames
- Input state is stack of pixels from last 4 frames
- Output is  $Q(s,a)$  for the joystick/button positions
  - Varies with different games (~3-18)
- Reward is change in score for that step

# DQN: Atari 2600

- Network architecture and hyper parameters are fixed across all the games



# Improvements on the Basic DQN Algorithm

# Double DQN

- Double DQN
  - Current Q-network  $w$  is used to select actions
  - Older Q-network is  $w^-$  is used to evaluate actions
- Prioritized reply
  - Store experience in priority queue by the DQN error

# Dueling Network

- Duelling Network
  - Action-independent value function  $V(s,v)$
  - Action-dependent advantage function  $A(s,a,\mathbf{w})$
  - $Q(s,a) = V(s,v) + A(s,a,\mathbf{w})$

# Prioritized Experience Replay

---

**Algorithm 1** Double DQN with proportional prioritization

---

- 1: **Input:** minibatch  $k$ , step-size  $\eta$ , replay period  $K$  and size  $N$ , exponents  $\alpha$  and  $\beta$ , budget  $T$ .
- 2: Initialize replay memory  $\mathcal{H} = \emptyset$ ,  $\Delta = 0$ ,  $p_1 = 1$
- 3: Observe  $S_0$  and choose  $A_0 \sim \pi_\theta(S_0)$
- 4: **for**  $t = 1$  to  $T$  **do**
- 5:   Observe  $S_t, R_t, \gamma_t$
- 6:   Store transition  $(S_{t-1}, A_{t-1}, R_t, \gamma_t, S_t)$  in  $\mathcal{H}$  with maximal priority  $p_t = \max_{i < t} p_i$
- 7:   **if**  $t \equiv 0 \pmod K$  **then**
- 8:     **for**  $j = 1$  to  $k$  **do**
- 9:       Sample transition  $j \sim P(j) = p_j^\alpha / \sum_i p_i^\alpha$
- 10:       Compute importance-sampling weight  $w_j = (N \cdot P(j))^{-\beta} / \max_i w_i$
- 11:       Compute TD-error  $\delta_j = R_j + \gamma_j Q_{\text{target}}(S_j, \arg \max_a Q(S_j, a)) - Q(S_{j-1}, A_{j-1})$
- 12:       Update transition priority  $p_j \leftarrow |\delta_j|$
- 13:       Accumulate weight-change  $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_\theta Q(S_{j-1}, A_{j-1})$
- 14:     **end for**
- 15:     Update weights  $\theta \leftarrow \theta + \eta \cdot \Delta$ , reset  $\Delta = 0$
- 16:     From time to time copy weights into target network  $\theta_{\text{target}} \leftarrow \theta$
- 17:   **end if**
- 18:   Choose action  $A_t \sim \pi_\theta(S_t)$
- 19: **end for**

---

# Prioritized Experience Replay

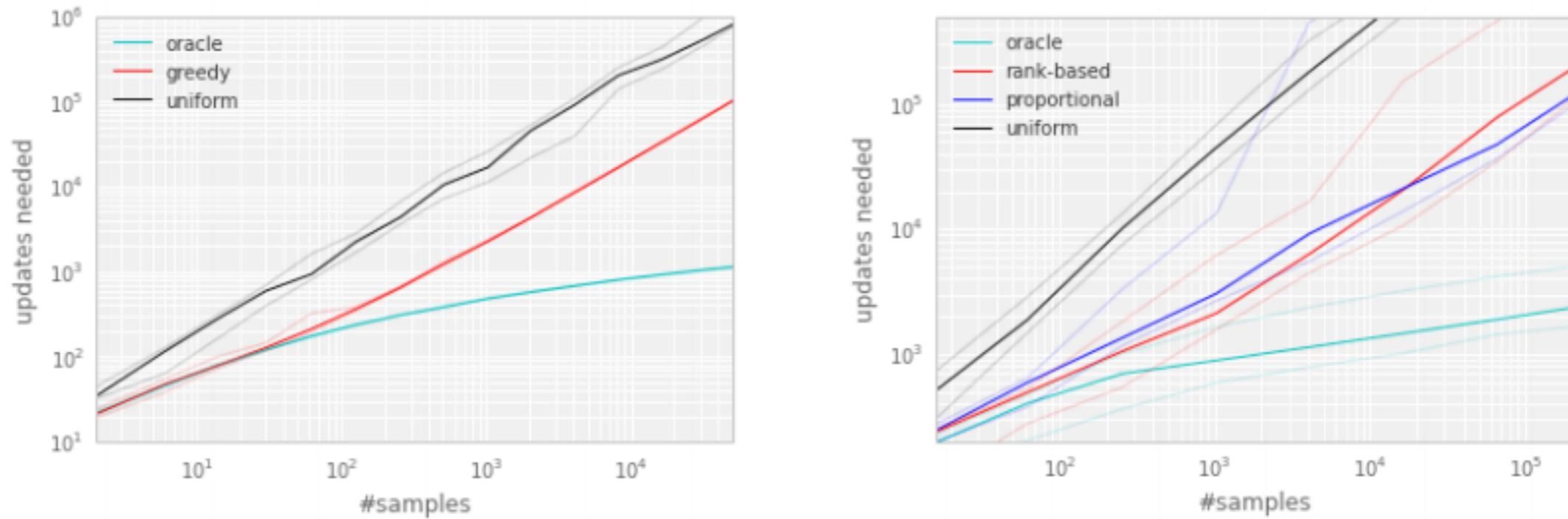
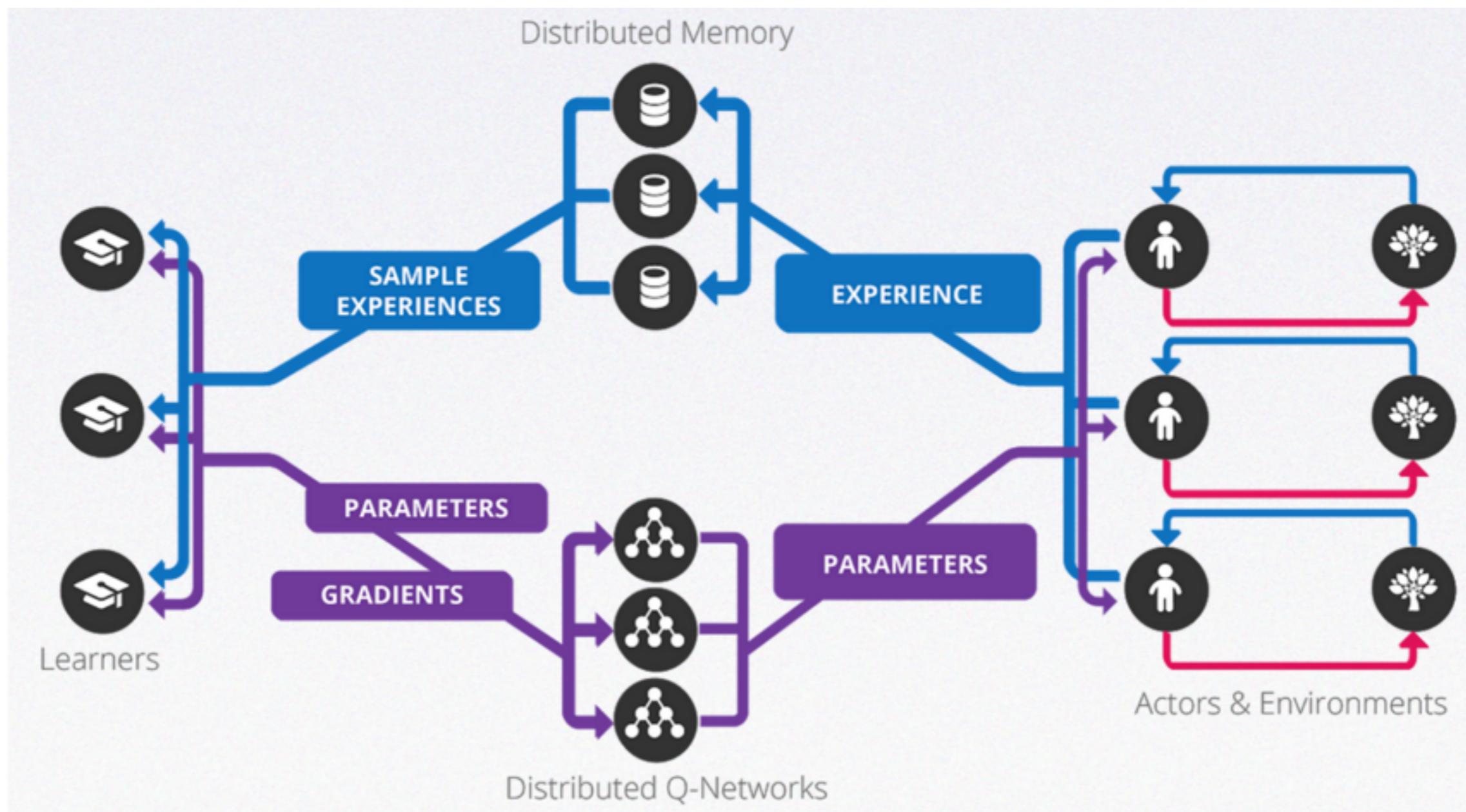


Figure 2: Median number of updates required for Q-learning to learn the value function on the Blind Cliffwalk example, as a function of the total number of transitions (only a single one of which was successful and saw the non-zero reward). Faint lines are min/max values from 10 random initializations. Black is uniform random replay, cyan uses the hindsight-oracle to select transitions, red and blue use prioritized replay (rank-based and proportional respectively). The results differ by multiple orders of magnitude, thus the need for a log-log plot. In both subplots it is evident that replaying experience in the right order makes an enormous difference to the number of updates required. See Appendix B.1 for details. **Left:** Tabular representation, greedy prioritization. **Right:** Linear function approximation, both variants of stochastic prioritization.

# General Reinforcement Learning Architecture



# Policy-based Deep Reinforcement Learning

# Deep Policy Network

- The policy is represented by a network with weights  $\mathbf{u}$

$$a = \pi(a|s, \mathbf{u}) \text{ or } a = \pi(s, \mathbf{u})$$

- Define objective function as total discounted reward

$$L(\mathbf{u}) = \mathbb{E} [r_1 + \gamma r_2 + \gamma^2 r_3 + \dots | \pi(\cdot, \mathbf{u})]$$

- Optimise objective via SGD
- Adjust policy parameters  $\mathbf{u}$  to achieve higher reward

# Policy Gradients

- Gradient of a stochastic policy

$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E} \left[ \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q^\pi(s, a) \right]$$

- Gradient of a deterministic policy

$$\frac{\partial L(\mathbf{u})}{\partial \mathbf{u}} = \mathbb{E} \left[ \frac{\partial Q^\pi(s, a)}{\partial a} \frac{\partial a}{\partial \mathbf{u}} \right]$$

# Actor Critic Algorithm

- Estimate value function  $Q(s, a, \mathbf{w}) \approx Q^\pi(s, a)$
- Update policy parameters  $\mathbf{u}$  via SGD

$$\frac{\partial I}{\partial \mathbf{u}} = \frac{\partial \log \pi(a|s, \mathbf{u})}{\partial \mathbf{u}} Q(s, a, \mathbf{w})$$

$$\frac{\partial I}{\partial \mathbf{u}} = \frac{\partial Q(s, a, \mathbf{w})}{\partial a} \frac{\partial a}{\partial \mathbf{u}}$$

# A3C: Labyrinth

- End-to-end learning of softmax policy from raw pixels
- Observations are pixels of the current frame
- State is an LSTM  $f(o_1, \dots, o_t)$
- Outputs both value  $V(s)$  & softmax over actions  $\pi(a|s)$
- Task is to collect apples (+1) and escape (+10)

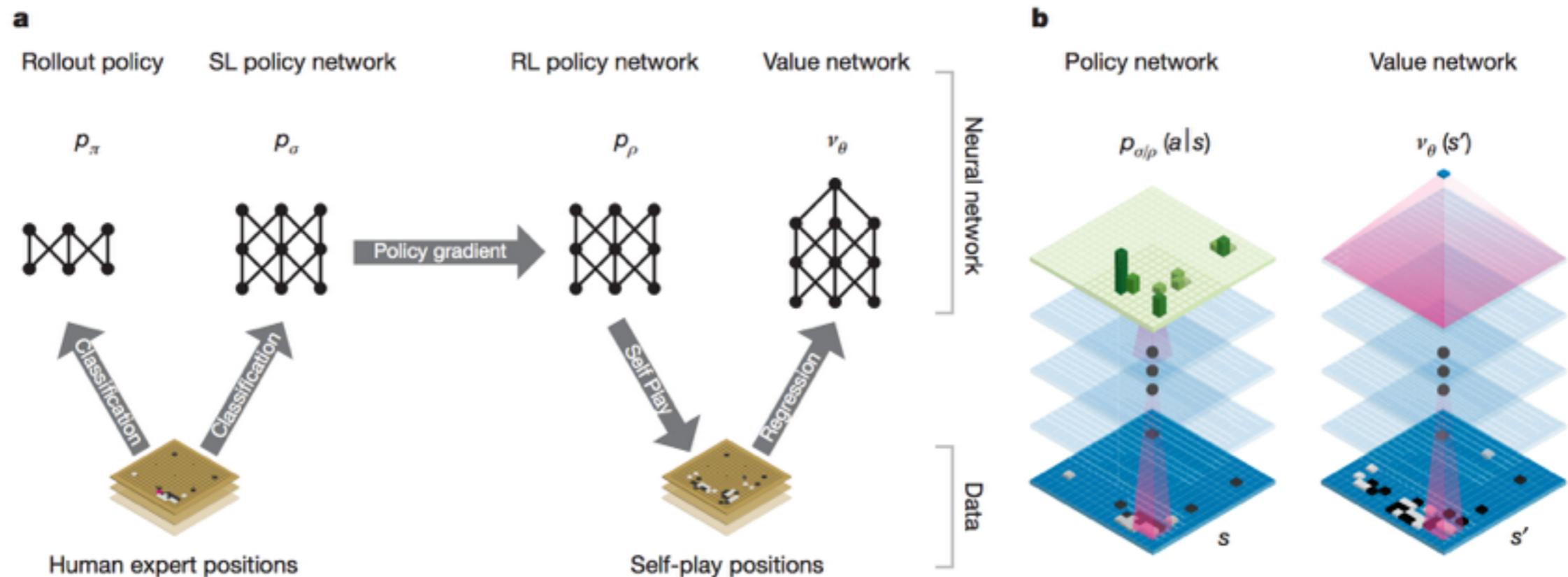
# Model-based Deep Reinforcement Learning

# Learning Models of the Environment

- Generative model of Atari 2600
- Issues
  - Errors in transition model compound over the trajectory
  - Planning trajectory differ from executed trajectories
  - Long, unusual trajectory rewards are totally wrong

# Newer Implementations

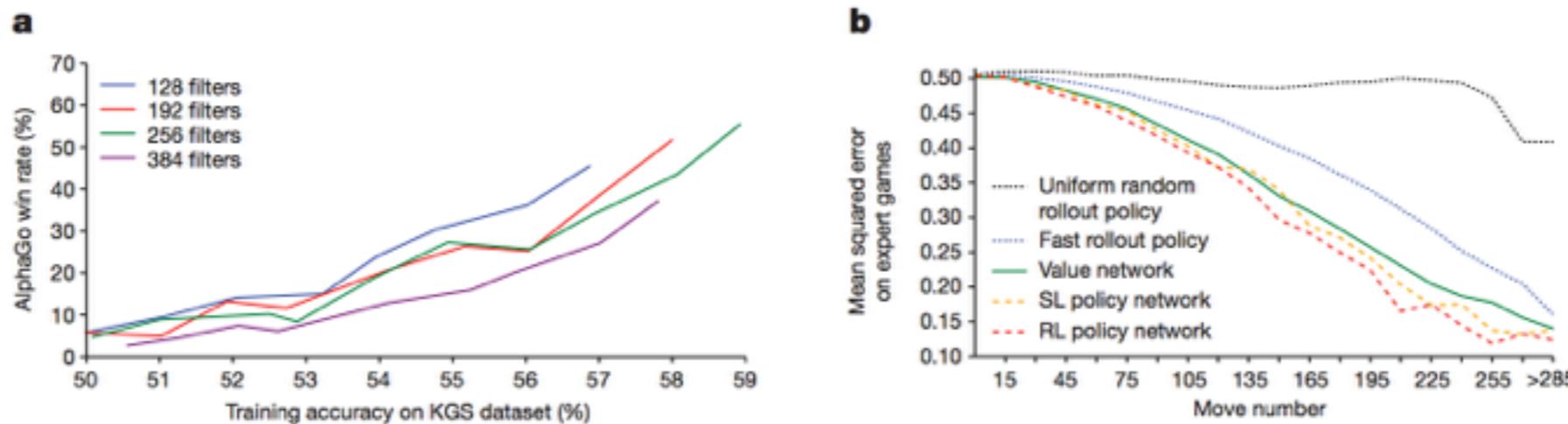
# AlphaGo



**Figure 1 | Neural network training pipeline and architecture.** **a**, A fast rollout policy  $p_\pi$  and supervised learning (SL) policy network  $p_\sigma$  are trained to predict human expert moves in a data set of positions. A reinforcement learning (RL) policy network  $p_\rho$  is initialized to the SL policy network, and is then improved by policy gradient learning to maximize the outcome (that is, winning more games) against previous versions of the policy network. A new data set is generated by playing games of self-play with the RL policy network. Finally, a value network  $v_\theta$  is trained by regression to predict the expected outcome (that is, whether

the current player wins) in positions from the self-play data set. **b**, Schematic representation of the neural network architecture used in AlphaGo. The policy network takes a representation of the board position  $s$  as its input, passes it through many convolutional layers with parameters  $\sigma$  (SL policy network) or  $\rho$  (RL policy network), and outputs a probability distribution  $p_\sigma(a|s)$  or  $p_\rho(a|s)$  over legal moves  $a$ , represented by a probability map over the board. The value network similarly uses many convolutional layers with parameters  $\theta$ , but outputs a scalar value  $v_\theta(s')$  that predicts the expected outcome in position  $s'$ .

# Policy and Value Networks

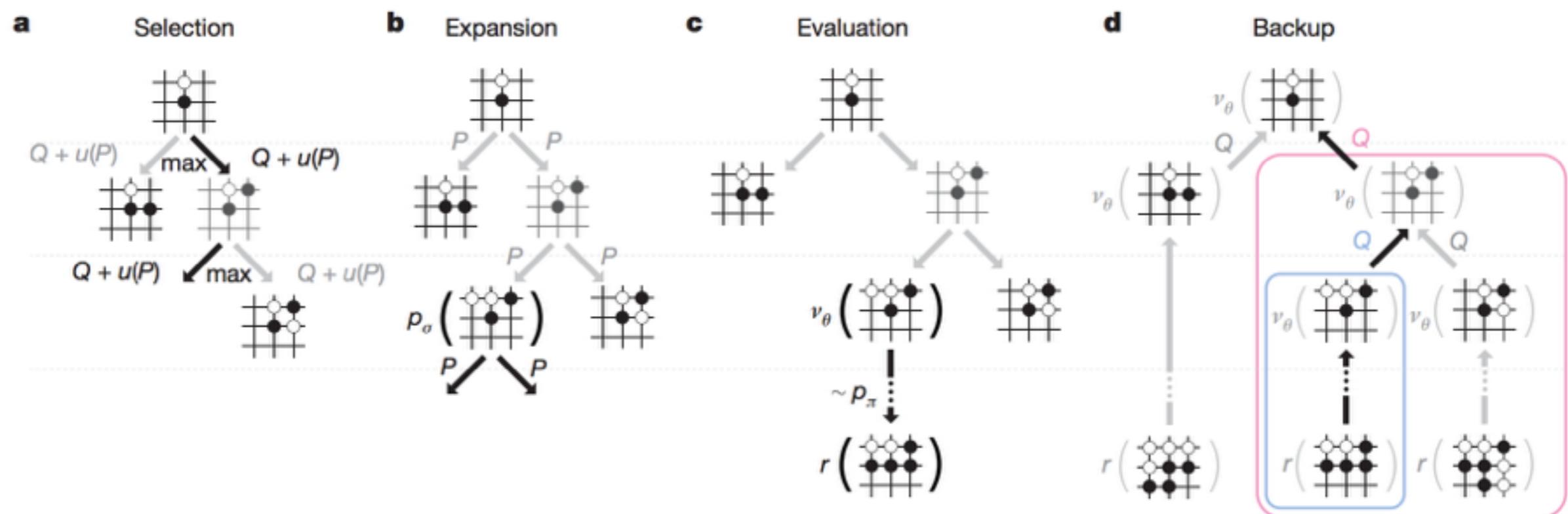


**Figure 2 | Strength and accuracy of policy and value networks.**

a, Plot showing the playing strength of policy networks as a function of their training accuracy. Policy networks with 128, 192, 256 and 384 convolutional filters per layer were evaluated periodically during training; the plot shows the winning rate of AlphaGo using that policy network against the match version of AlphaGo. b, Comparison of evaluation accuracy between the value network and rollouts with different policies.

Positions and outcomes were sampled from human expert games. Each position was evaluated by a single forward pass of the value network  $v_\theta$ , or by the mean outcome of 100 rollouts, played out using either uniform random rollouts, the fast rollout policy  $p_\pi$ , the SL policy network  $p_\sigma$  or the RL policy network  $p_\rho$ . The mean squared error between the predicted value and the actual game outcome is plotted against the stage of the game (how many moves had been played in the given position).

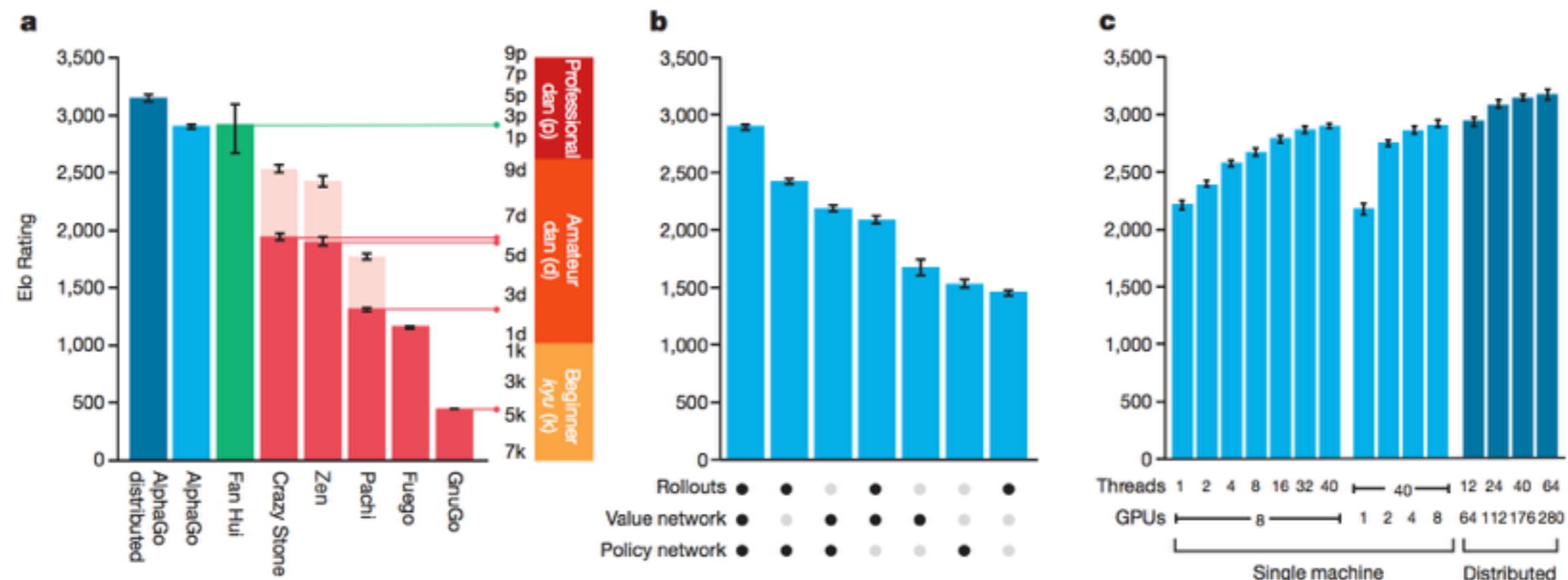
# Monte Carlo Tree Search



**Figure 3 | Monte Carlo tree search in AlphaGo.** **a**, Each simulation traverses the tree by selecting the edge with maximum action value  $Q$ , plus a bonus  $u(P)$  that depends on a stored prior probability  $P$  for that edge. **b**, The leaf node may be expanded; the new node is processed once by the policy network  $p_\sigma$  and the output probabilities are stored as prior probabilities  $P$  for each action. **c**, At the end of a simulation, the leaf node

is evaluated in two ways: using the value network  $v_\theta$ ; and by running a rollout to the end of the game with the fast rollout policy  $p_\pi$ , then computing the winner with function  $r$ . **d**, Action values  $Q$  are updated to track the mean value of all evaluations  $r(\cdot)$  and  $v_\theta(\cdot)$  in the subtree below that action.

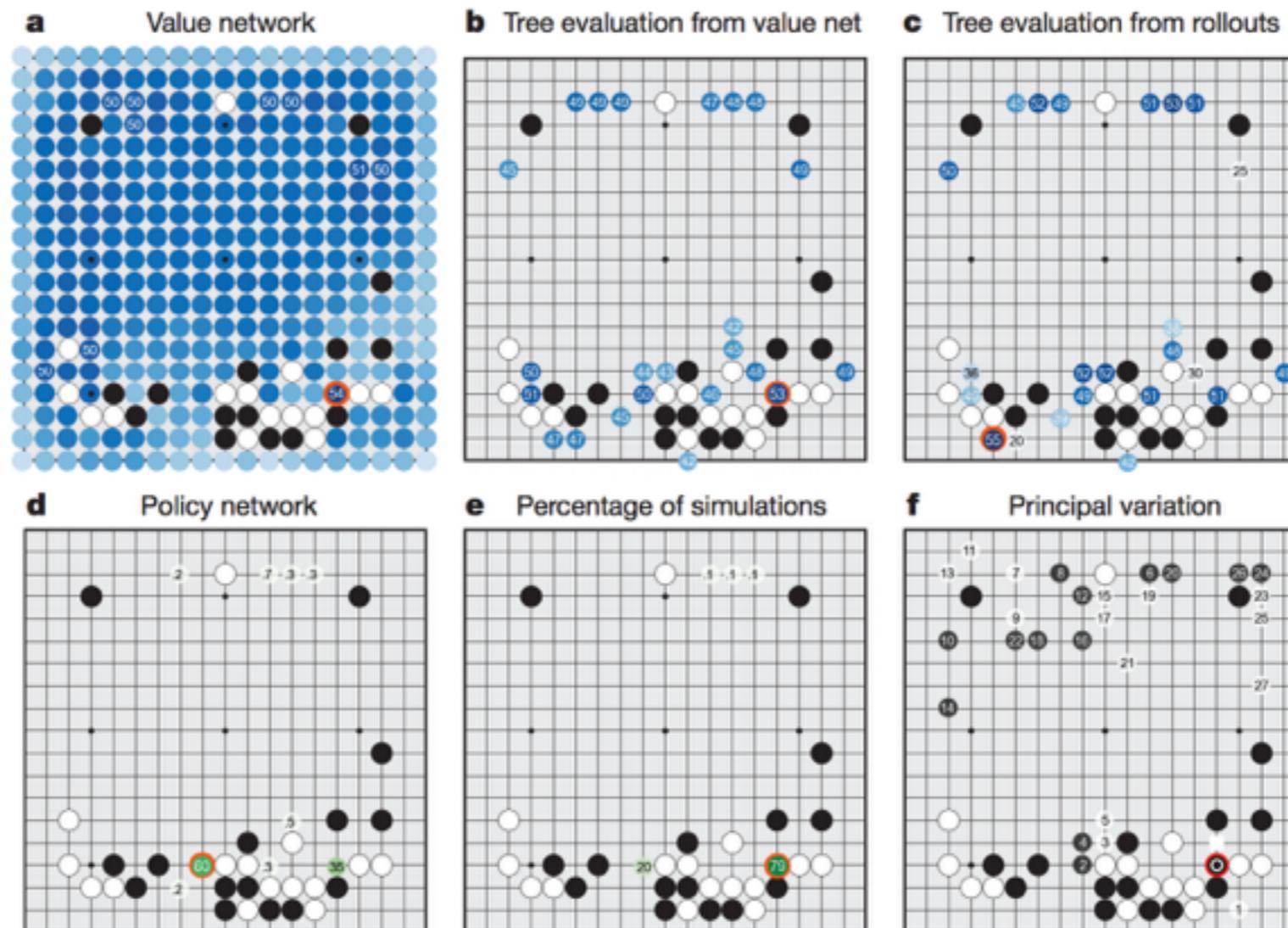
# AlphaGo Results



**Figure 4 | Tournament evaluation of AlphaGo. a,** Results of a tournament between different Go programs (see Extended Data Tables 6–11). Each program used approximately 5 s computation time per move. To provide a greater challenge to AlphaGo, some programs (pale upper bars) were given four handicap stones (that is, free moves at the start of every game) against all opponents. Programs were evaluated on an Elo scale<sup>37</sup>: a 230 point gap corresponds to a 79% probability of winning, which roughly corresponds to one amateur *dan* rank advantage on KGS<sup>38</sup>; an approximate correspondence to human ranks is also shown,

horizontal lines show KGS ranks achieved online by that program. Games against the human European champion Fan Hui were also included; these games used longer time controls. 95% confidence intervals are shown. **b**, Performance of AlphaGo, on a single machine, for different combinations of components. The version solely using the policy network does not perform any search. **c**, Scalability study of MCTS in AlphaGo with search threads and GPUs, using asynchronous search (light blue) or distributed search (dark blue), for 2 s per move.

# Which Move to Make

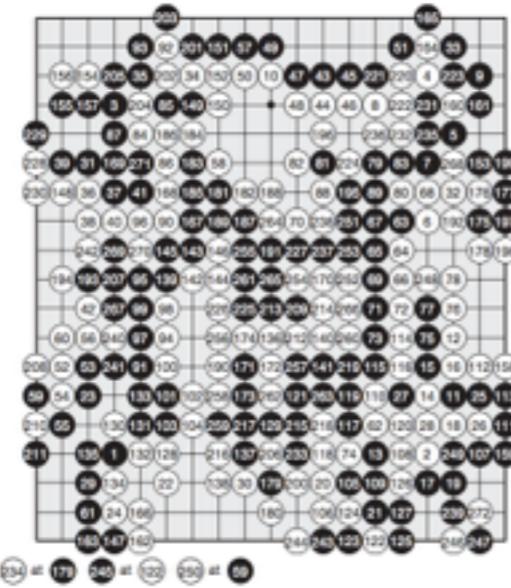


**Figure 5 | How AlphaGo (black, to play) selected its move in an informal game against Fan Hui.** For each of the following statistics, the location of the maximum value is indicated by an orange circle. **a**, Evaluation of all successors  $s'$  of the root position  $s$ , using the value network  $v_\theta(s')$ ; estimated winning percentages are shown for the top evaluations. **b**, Action values  $Q(s, a)$  for each edge  $(s, a)$  in the tree from root position  $s$ ; averaged over value network evaluations only ( $\lambda = 0$ ). **c**, Action values  $Q(s, a)$ , averaged over rollout evaluations only ( $\lambda = 1$ ).

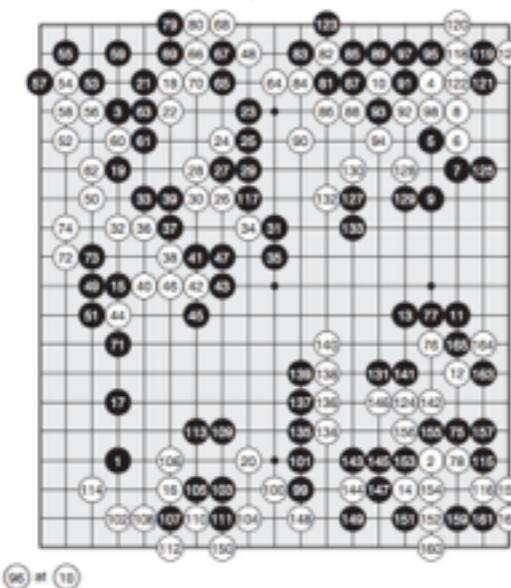
**d**, Move probabilities directly from the SL policy network,  $p_\sigma(a|s)$ ; reported as a percentage (if above 0.1%). **e**, Percentage frequency with which actions were selected from the root during simulations. **f**, The principal variation (path with maximum visit count) from AlphaGo's search tree. The moves are presented in a numbered sequence. AlphaGo selected the move indicated by the red circle; Fan Hui responded with the move indicated by the white square; in his post-game commentary he preferred the move (labelled 1) predicted by AlphaGo.

# Five Matches Against Fan Hui

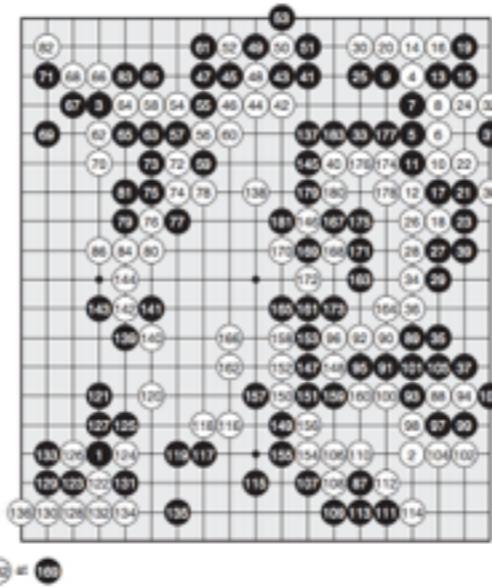
Game 1  
Fan Hui (Black), AlphaGo (White)  
AlphaGo wins by 2.5 points



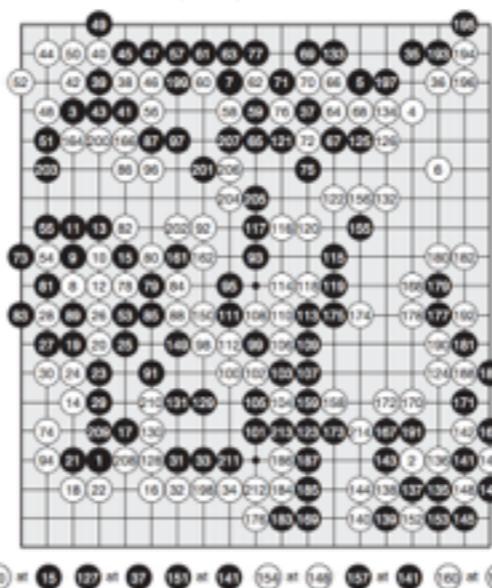
Game 4  
AlphaGo (Black), Fan Hui (White)  
AlphaGo wins by resignation



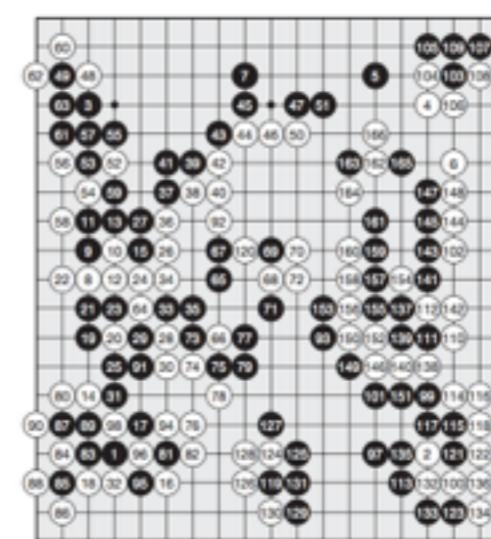
Game 2  
AlphaGo (Black), Fan Hui (White)  
AlphaGo wins by resignation



Game 5  
Fan Hui (Black), AlphaGo (White)  
AlphaGo wins by resignation



Game 3  
Fan Hui (Black), AlphaGo (White)  
AlphaGo wins by resignation



# Neural Architecture Search with Reinforcement Learning

- Uses an RNN to generate model descriptions of the NNs
- Trains the RNN with RL to maximize expected accuracy of generated architectures on validation set

# Neural Architecture Search

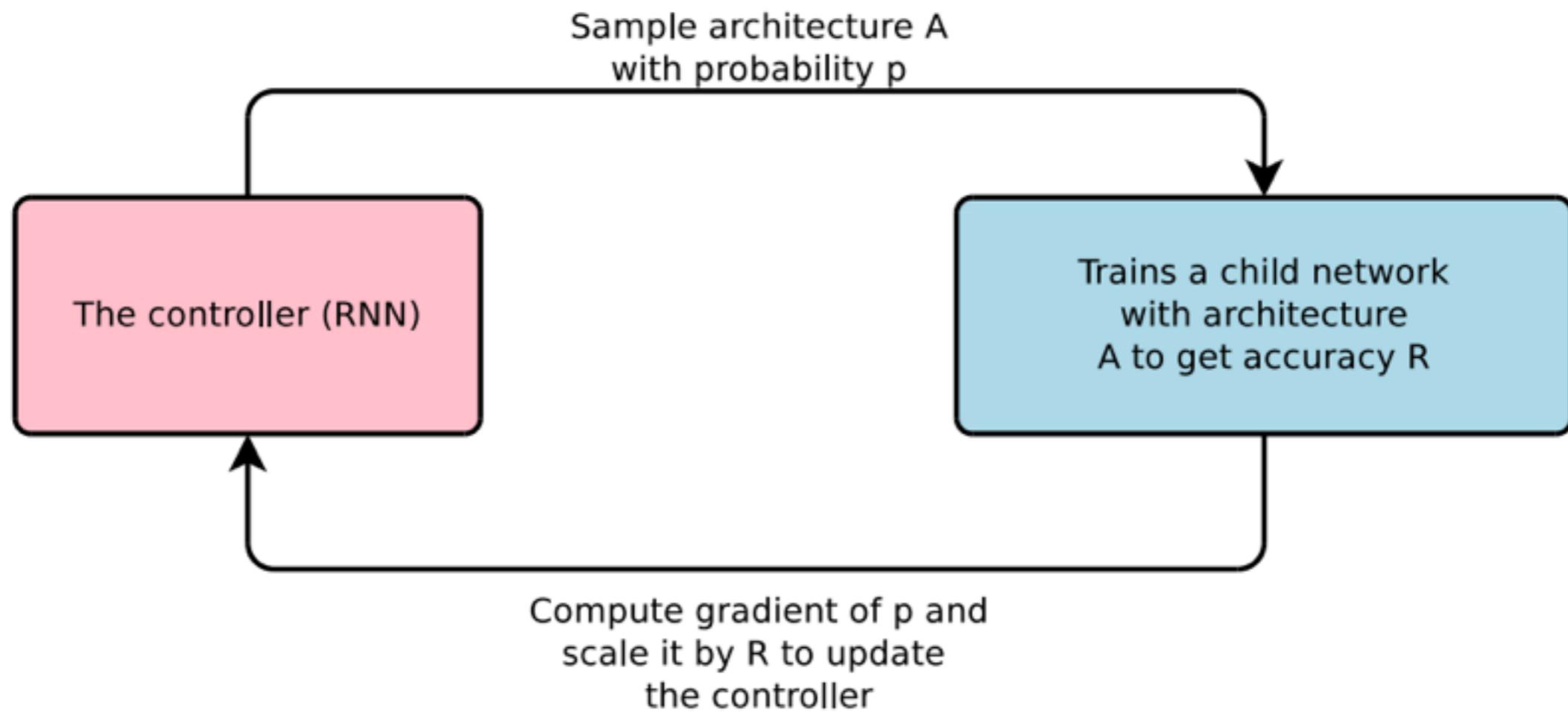


Figure 1: An overview of Neural Architecture Search.

# RNN Controller

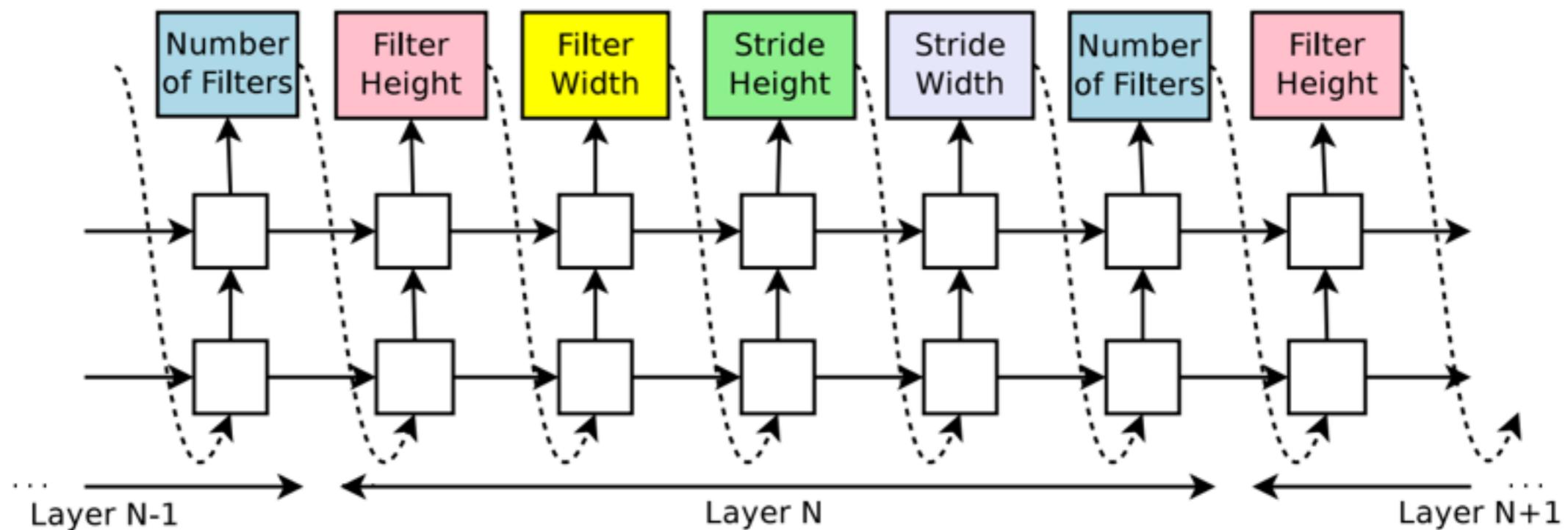


Figure 2: How our controller recurrent neural network samples a simple convolutional network. It predicts filter height, filter width, stride height, stride width, and number of filters for one layer and repeats. Every prediction is carried out by a softmax classifier and then fed into the next time step as input.

# Training with REINFORCE

$$J(\theta_c) = E_{P(a_{1:T}; \theta_c)}[R]$$

$$\nabla_{\theta_c} J(\theta_c) = \sum_{t=1}^T E_{P(a_{1:T}; \theta_c)} \left[ \nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R \right]$$

$$\frac{1}{m} \sum_{k=1}^m \sum_{t=1}^T \nabla_{\theta_c} \log P(a_t | a_{(t-1):1}; \theta_c) R_k$$

Where  $m$  is the number of different architectures that the controller samples in one batch and  $T$  is the number of hyperparameters our controller has to predict to design a neural network architecture.

# Parallel Training

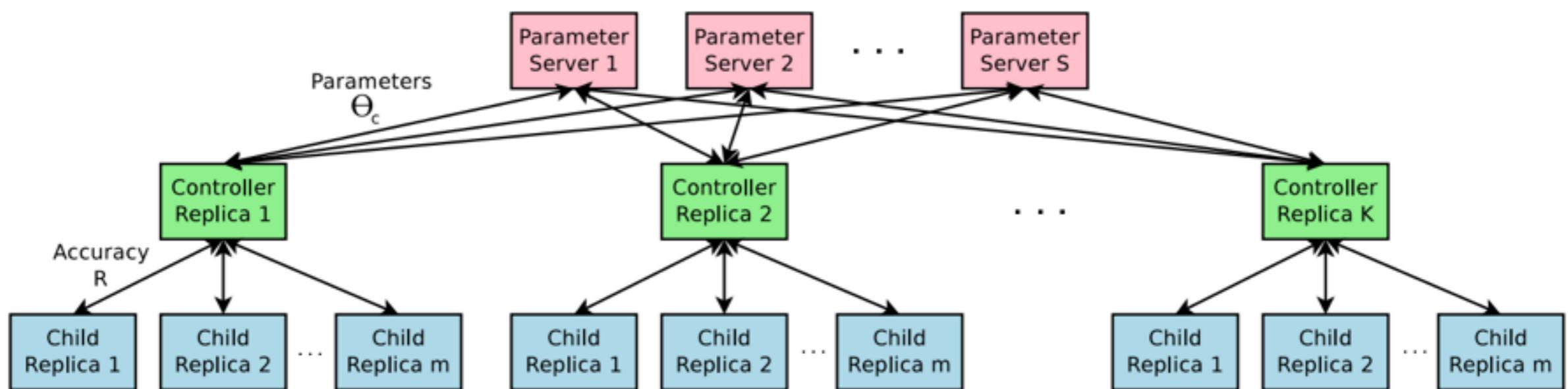


Figure 3: Distributed training for Neural Architecture Search. We use a set of  $S$  parameter servers to store and send parameters to  $K$  controller replicas. Each controller replica then samples  $m$  architectures and run the multiple child models in parallel. The accuracy of each child model is recorded to compute the gradients with respect to  $\theta_c$ , which are then sent back to the parameter servers.

# Increase Architecture Complexity

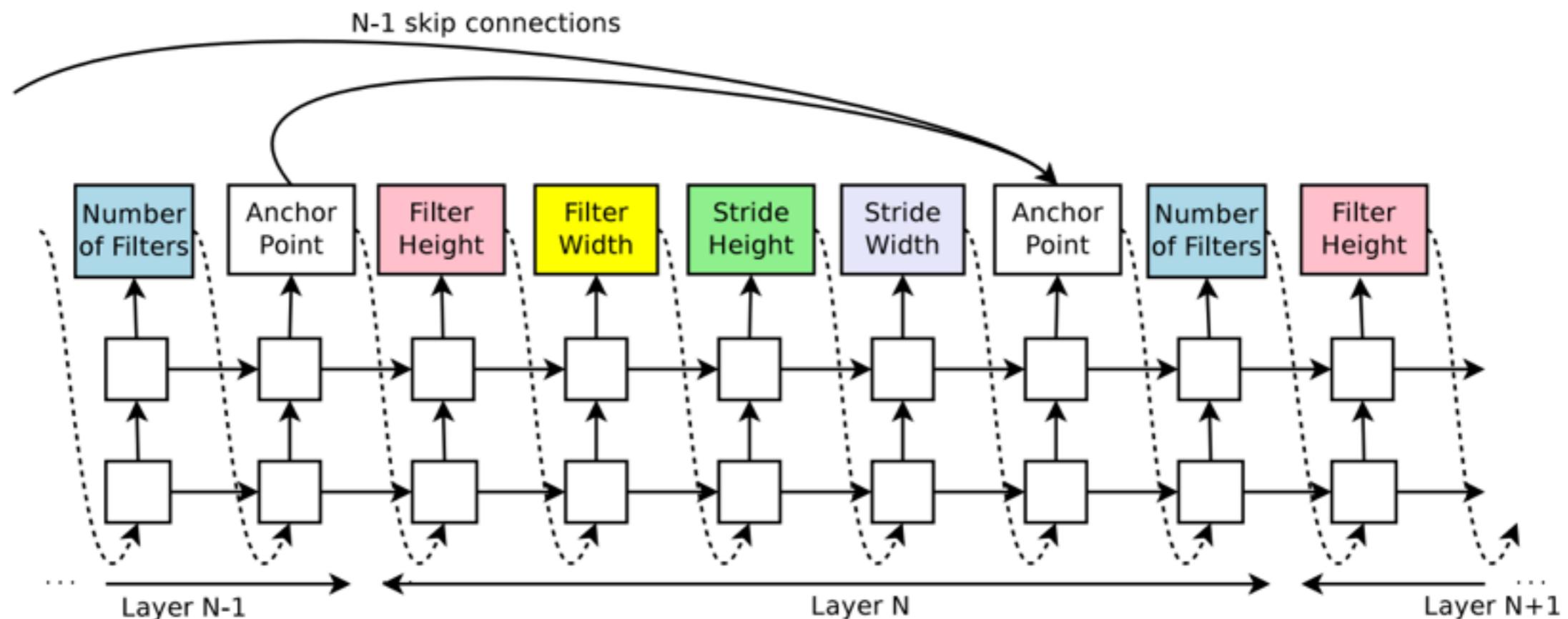


Figure 4: The controller uses anchor points, and set-selection attention to form skip connections.

# Constructing the Network

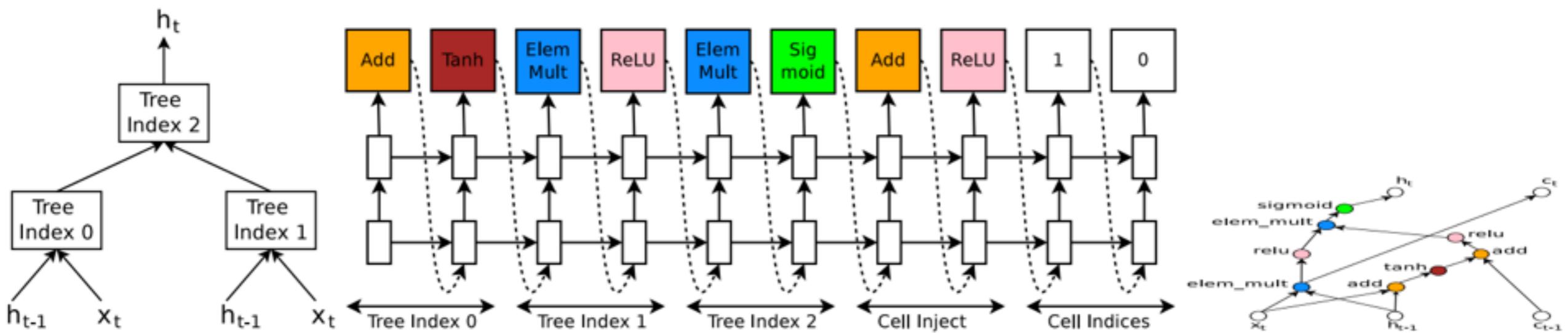


Figure 5: An example of a recurrent cell constructed from a tree that has two leaf nodes (base 2) and one internal node. Left: the tree that defines the computation steps to be predicted by controller. Center: an example set of predictions made by the controller for each computation step in the tree. Right: the computation graph of the recurrent cell constructed from example predictions of the controller.

# Results

Model	Depth	Parameters	Error rate (%)
Network in Network (Lin et al., 2013)	-	-	8.81
All-CNN (Springenberg et al., 2014)	-	-	7.25
Deeply Supervised Net (Lee et al., 2015)	-	-	7.97
Highway Network (Srivastava et al., 2015)	-	-	7.72
Scalable Bayesian Optimization (Snoek et al., 2015)	-	-	6.37
FractalNet (Larsson et al., 2016) with Dropout/Drop-path	21 21	38.6M 38.6M	5.22 4.60
ResNet (He et al., 2016a)	110	1.7M	6.61
ResNet (reported by Huang et al. (2016b))	110	1.7M	6.41
ResNet with Stochastic Depth (Huang et al., 2016b)	110 1202	1.7M 10.2M	5.23 4.91
Wide ResNet (Zagoruyko & Komodakis, 2016)	16 28	11.0M 36.5M	4.81 4.17
ResNet (pre-activation) (He et al., 2016b)	164 1001	1.7M 10.2M	5.46 4.62
DenseNet ( $L = 40, k = 12$ ) Huang et al. (2016a)	40	1.0M	5.24
DenseNet( $L = 100, k = 12$ ) Huang et al. (2016a)	100	7.0M	4.10
DenseNet ( $L = 100, k = 24$ ) Huang et al. (2016a)	100	27.2M	3.74
Neural Architecture Search v1 no stride or pooling	15	4.2M	5.50
Neural Architecture Search v2 predicting strides	20	2.5M	6.01
Neural Architecture Search v3 max pooling	39	7.1M	4.47
Neural Architecture Search v3 max pooling + more filters	39	32.0M	3.84

Table 1: Performance of Neural Architecture Search and other state-of-the-art models on CIFAR-10.

# Reinforcement Learning with Unsupervised Auxiliary Tasks

- Train an agent that maximises other pseudo-reward functions simultaneously by RL
  - Those tasks have to develop in absence of extrinsic rewards
- Outperforms previous state-of-the-art on atari
  - Averaging 880% expert human performance

# UNREAL Agent

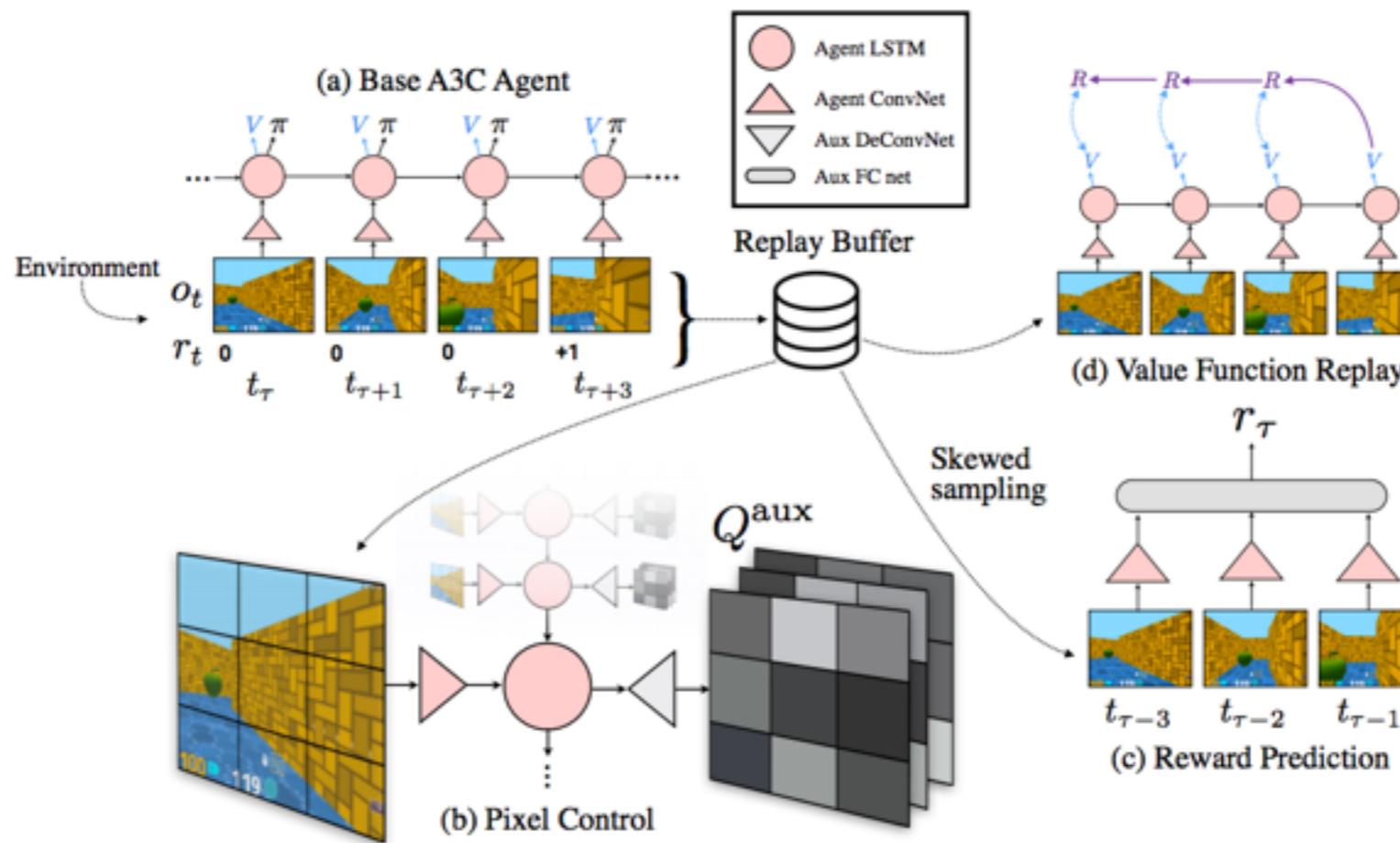


Figure 1: Overview of the *UNREAL* agent. (a) The base agent is a CNN-LSTM agent trained on-policy with the A3C loss (Mnih et al., 2016). Observations, rewards, and actions are stored in a small replay buffer which encapsulates a short history of agent experience. This experience is used by auxiliary learning tasks. (b) Pixel Control – auxiliary policies  $Q^{\text{aux}}$  are trained to maximise change in pixel intensity of different regions of the input. The agent CNN and LSTM are used for this task along with an auxiliary deconvolution network. This auxiliary control task requires the agent to learn how to control the environment. (c) Reward Prediction – given three recent frames, the network must predict the reward that will be obtained in the next unobserved timestep. This task network uses instances of the agent CNN, and is trained on reward biased sequences to remove the perceptual sparsity of rewards. (d) Value Function Replay – further training of the value function using the agent network is performed to promote faster value iteration. Further visualisation of the agent can be found in <https://youtu.be/Uz-zGYrYEjA>

# Auxiliary Control Tasks

Given a set of auxiliary control tasks  $\mathcal{C}$ , let  $\pi^{(c)}$  be the agent's policy for each auxiliary task  $c \in \mathcal{C}$  and let  $\pi$  be the agent's policy on the base task. The overall objective is to maximise total performance across all these auxiliary tasks,

$$\arg \max_{\theta} \mathbb{E}_{\pi}[R_{1:\infty}] + \lambda_c \sum_{c \in \mathcal{C}} \mathbb{E}_{\pi_c}[R_{1:\infty}^{(c)}], \quad (1)$$

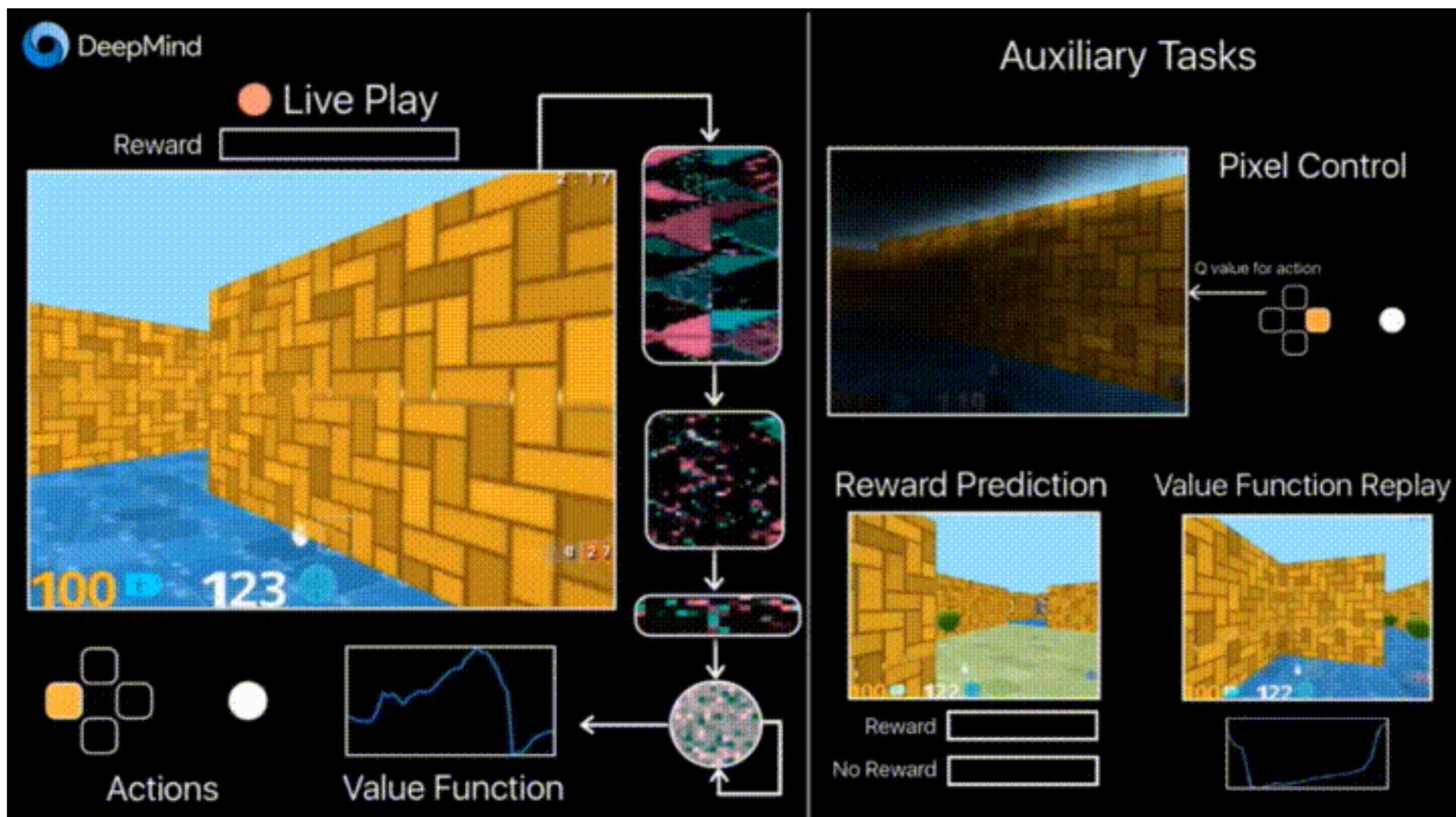
$R_{t:t+n}^{(c)} = \sum_{k=1}^n \gamma^k r_t^{(c)}$  is the discounted return for auxiliary reward  $r^{(c)}$

$\theta$  is the set of parameters of  $\pi$  and all  $\pi^{(c)}$ 's

$n$ -step Q-learning loss  $\mathcal{L}_Q^{(c)} =$

$$\mathbb{E} \left[ (R_{t:t+n} + \gamma^n \max_{a'} Q^{(c)}(s', a', \theta^-) - Q^{(c)}(s, a, \theta))^2 \right]$$

# Auxiliary Tasks



# UNREAL Algorithm

$$\mathcal{L}_{UNREAL}(\theta) = \mathcal{L}_{A3C} + \lambda_{VR}\mathcal{L}_{VR} + \lambda_{PC} \sum_c \mathcal{L}_Q^{(c)} + \lambda_{RP}\mathcal{L}_{RP}$$

auxiliary control loss  $\mathcal{L}_{PC}$

auxiliary reward prediction loss  $\mathcal{L}_{RP}$

replayed value loss  $\mathcal{L}_{VR}$

A3C Loss is minimised on policy

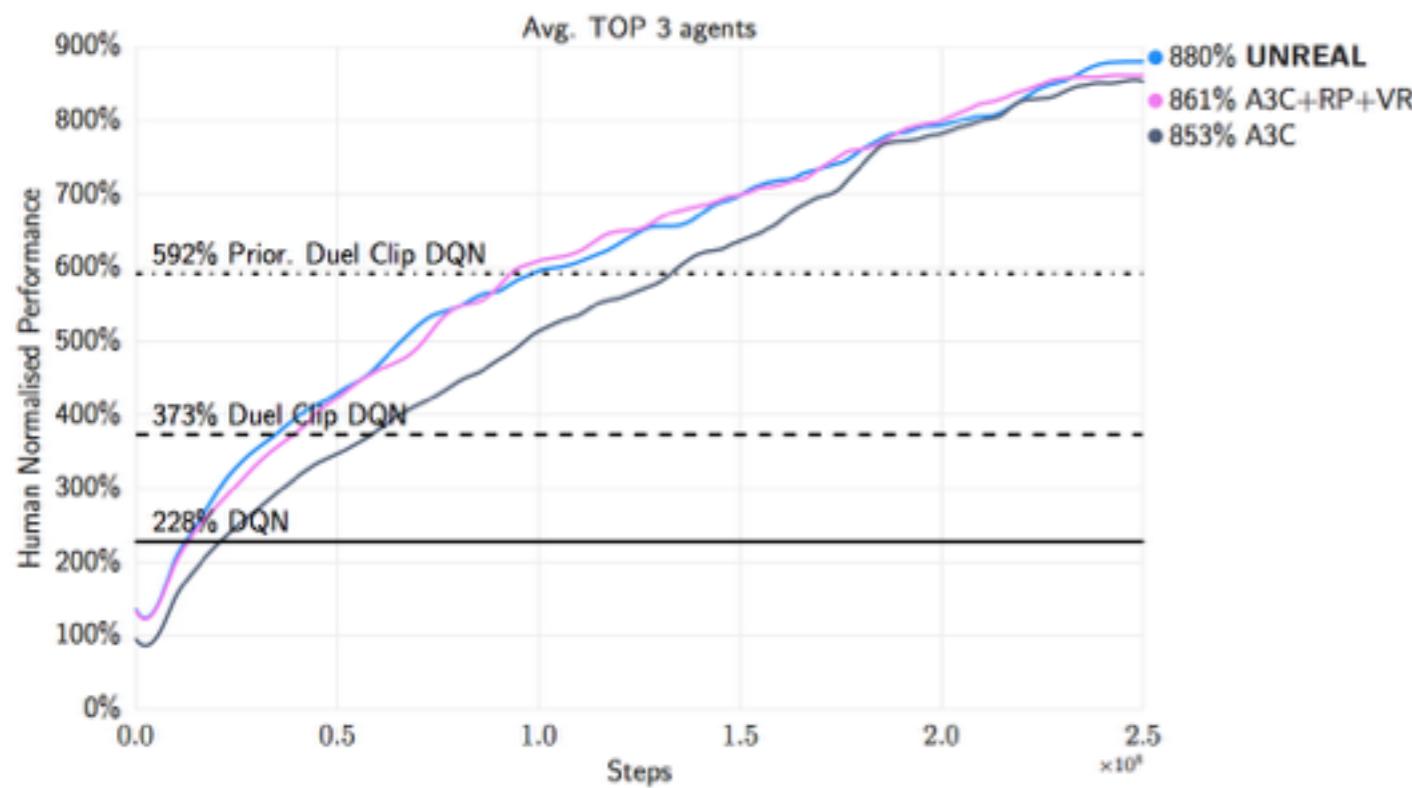
Value function is optimised from replayed data

Auxiliary control loss is optimised off-policy from replied data

Reward loss is optimised from rebalanced replay data

# Atari Results

Atari Performance



Atari Robustness

