# Deep Reinforcement Learning at Scale
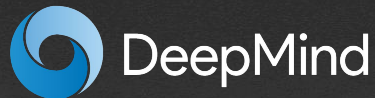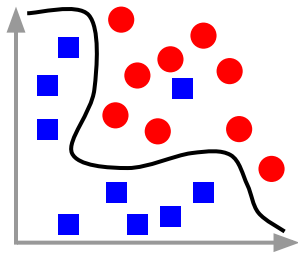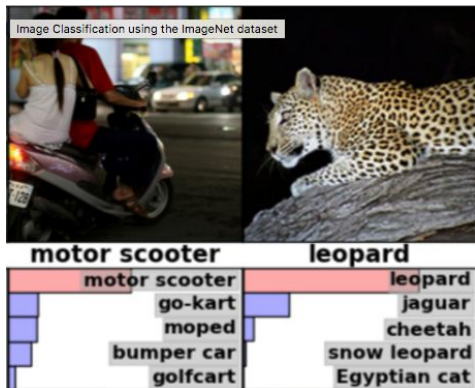
Timothy Lillicrap
Research Scientist, DeepMind & UCL

Deep Learning at Supercomputer Scale | NIPS Workshop

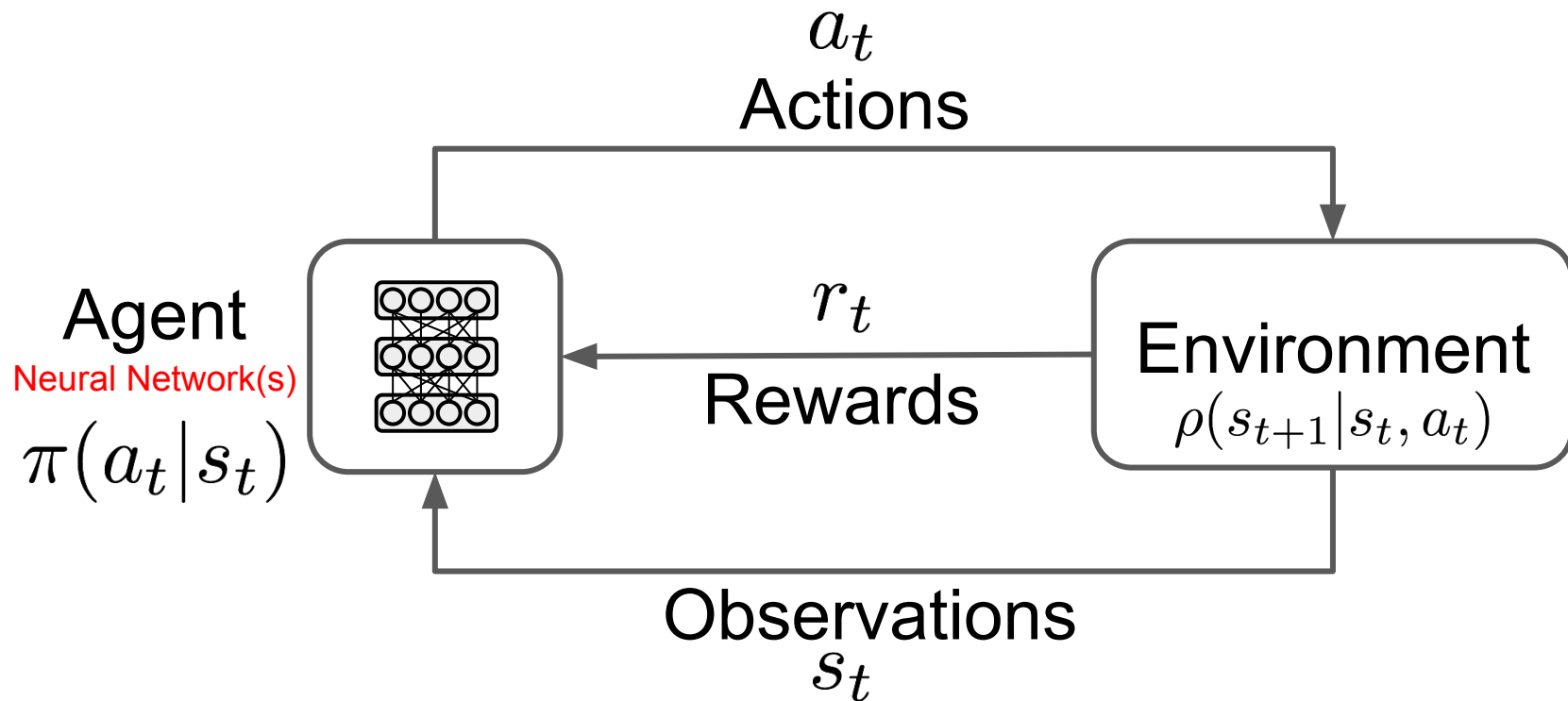DeepMind

# What is Reinforcement Learning?

## Supervised Learning



**Fixed dataset**
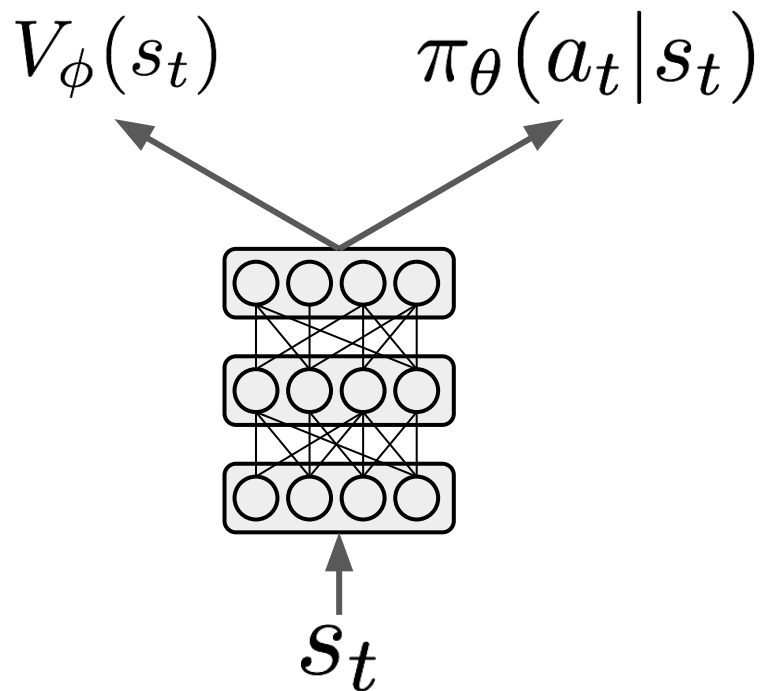
## Reinforcement Learning
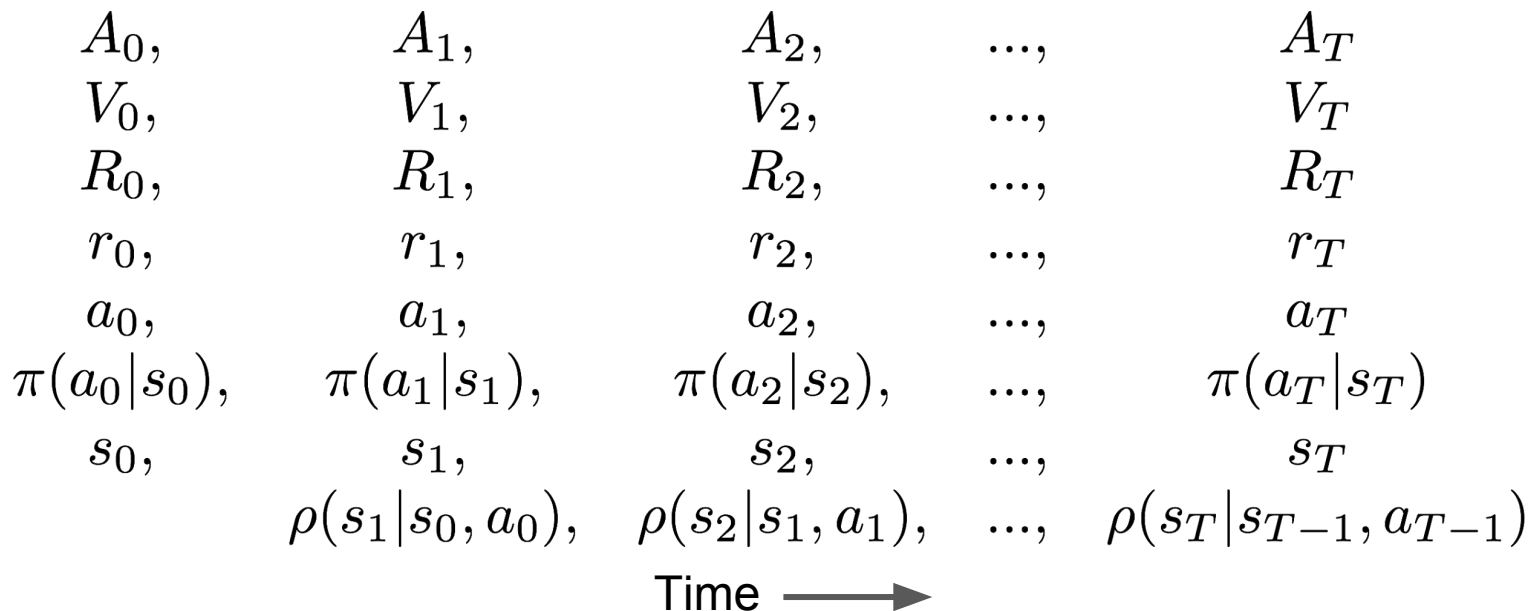


**Data depends on actions taken in environment**

# Formalizing the Agent-Environment Loop



$a_t$
Actions

Agent
Neural Network(s)
$\pi(a_t|s_t)$

$r_t$
Rewards

Environment
$\rho(s_{t+1}|s_t, a_t)$

Observations
$s_t$

# Advantage Actor-Critic (A3C)

$$V_\phi(s_t) \qquad \pi_\theta(a_t|s_t)$$

$$s_t$$

Mnih et al., *ICML* 2016

# A Single Trial (with Advantage Actor-Critic)

$$
\begin{array}{lllll}
A_0, & A_1, & A_2, & ..., & A_T \\
V_0, & V_1, & V_2, & ..., & V_T \\
R_0, & R_1, & R_2, & ..., & R_T \\
r_0, & r_1, & r_2, & ..., & r_T \\
a_0, & a_1, & a_2, & ..., & a_T \\
\pi(a_0|s_0), & \pi(a_1|s_1), & \pi(a_2|s_2), & ..., & \pi(a_T|s_T) \\
s_0, & s_1, & s_2, & ..., & s_T \\
& \rho(s_1|s_0,a_0), & \rho(s_2|s_1,a_1), & ..., & \rho(s_T|s_{T-1},a_{T-1})
\end{array}
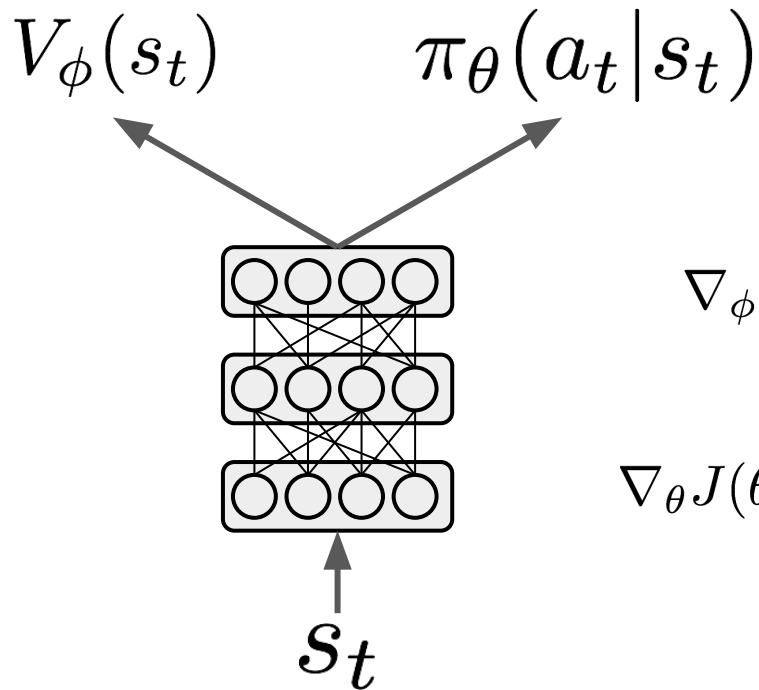$$

Time $\longrightarrow$

$$
R_t = \sum_{k=t}^{T} \gamma^{t-k} r_t \qquad V_t = V_\phi(s_t) \qquad A_t = R_t - V_\phi(s_t)
$$

# Combating Variance: Advantage Actor-Critic

$$V_\phi(s_t) \qquad \pi_\theta(a_t|s_t)$$
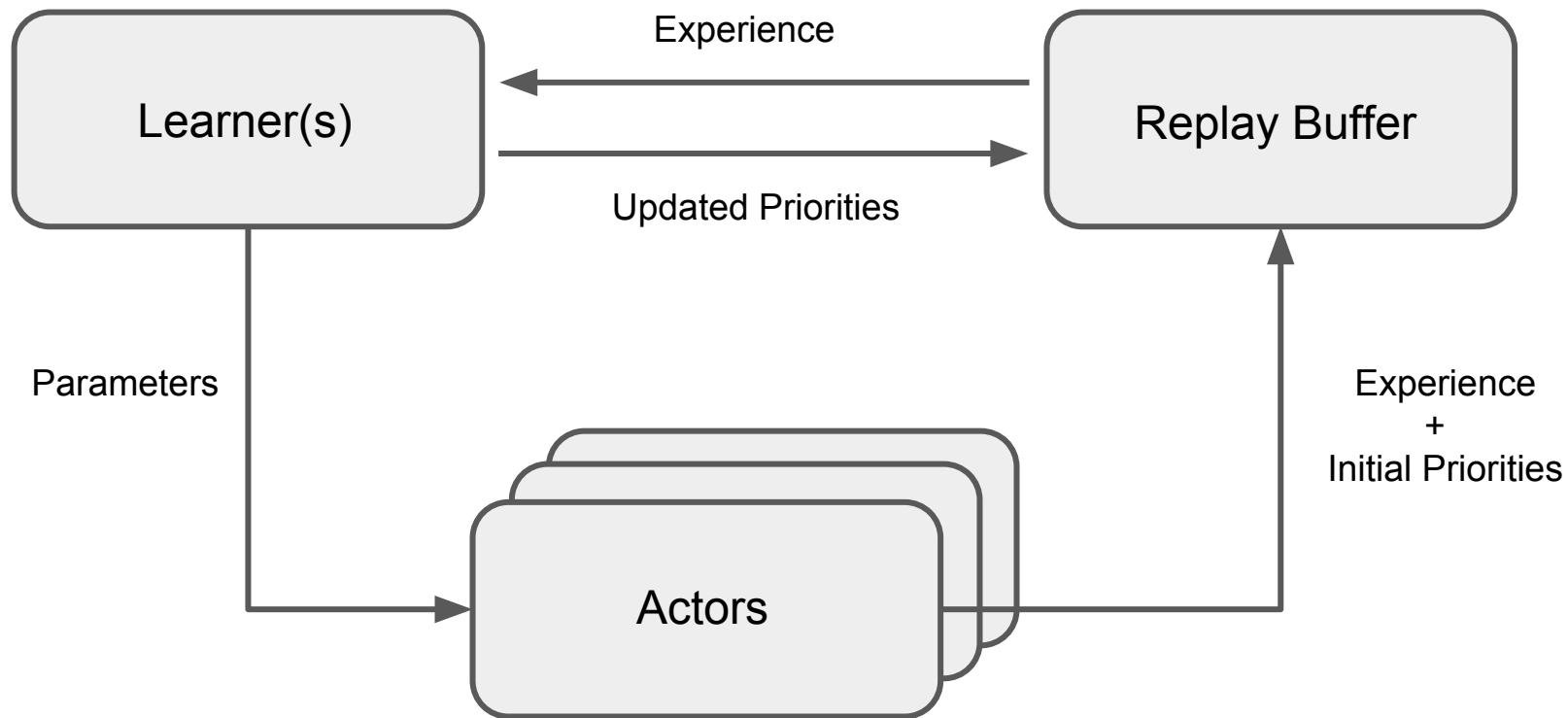


$$\nabla_\phi \mathcal{L} = \sum_{t=0}^{T} \nabla_\phi (R_t - V_\phi(s_t))^2$$

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_t|s_t)(R_t - V_\phi(s_t)) \right]$$
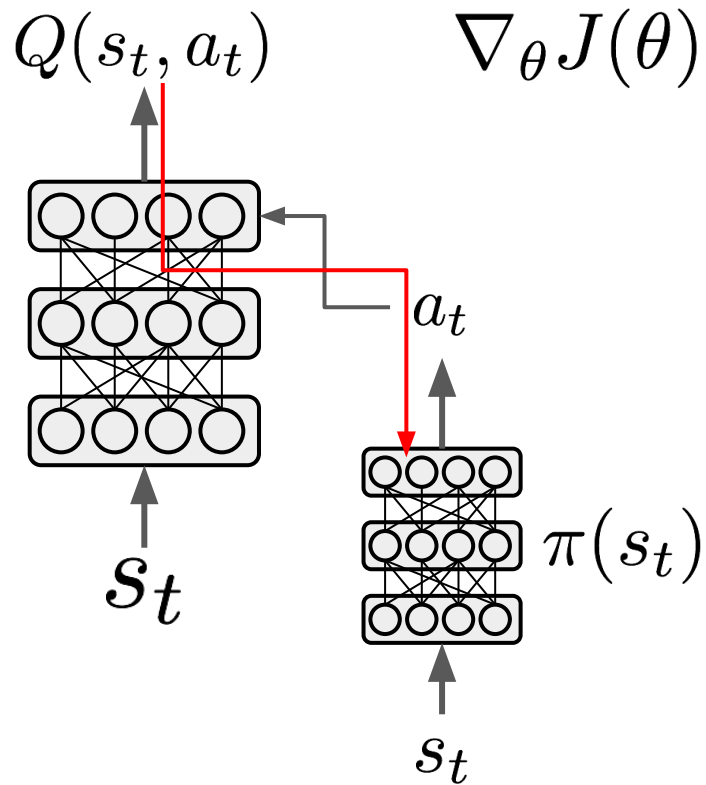
$$s_t$$

Mnih et al., *ICML* 2016

# Scaling Reinforcement Learning (A3C)



Mnih et al., *ICML* 2016

# Scaling Reinforcement Learning



Horgan et al., 2017 & Schaul et al. 2015

# Off-policy Actor-Critic for Continuous Actions



$$\nabla_\theta J(\theta) \approx \mathbb{E}_\mathcal{D} \left[ \nabla_a Q(s, a; \phi)|_{a=\pi(s)} \nabla_\theta \pi(s) \right]$$

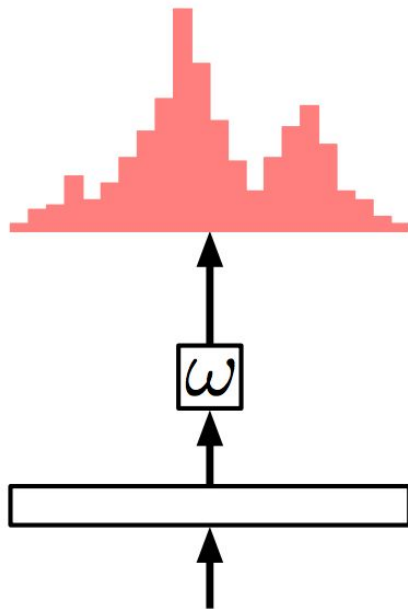$$(s, a, r, s') \sim U(\mathcal{D})$$
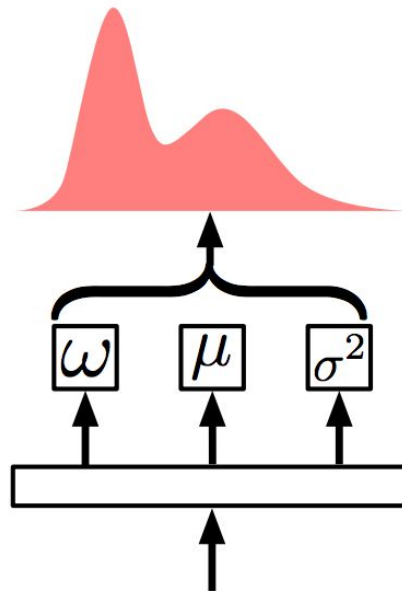
$$y = r + \gamma Q(s', \pi(s'); \phi^{\text{target}})$$
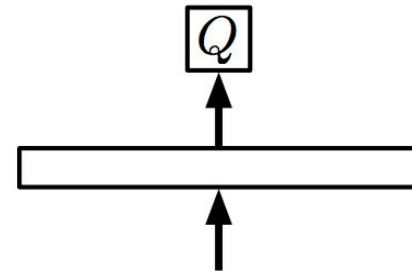
$$L(\phi) = (y - Q(s, a; \phi))^2$$

$Q(s_t, a_t)$

$a_t$

$s_t$

$\pi(s_t)$

$s_t$

Lillicrap et al., *ICLR* 2016

# Distributional Distributed DDPG (D4PG)



Hoffman, Barth-Maron et al., 2017
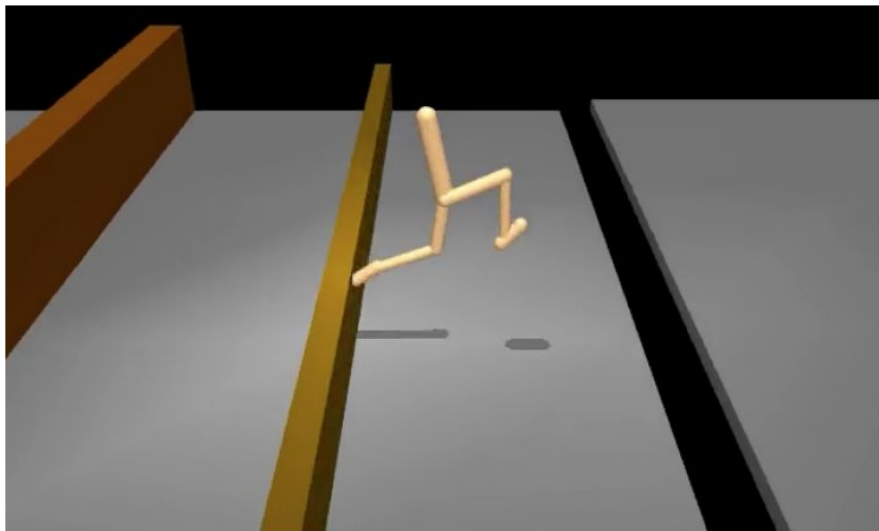
Acrobot(Swingup) · Acrobot(Swingup Sparse) · Cartpole(Swingup) · Cartpole(Swingup Sparse) · Cheetah(Walk) · Finger(Turn Easy) · Finger(Turn Hard) · Fish(Swim) · Fish(Upright) · Hopper(Stand) · Humanoid(Run) · Humanoid(Stand) · Humanoid(Walk) · Manipulator(Bring Ball) · Swimmer(Swimmer15) · Swimmer(Swimmer6)

Episode Returns vs Training Time (Hours)

D3PG, Non-Prioritized, N = 1
D3PG, Non-Prioritized, N = 5
D3PG, Prioritized, N = 1
D3PG, Prioritized, N = 5
D4PG, Non-Prioritized, N = 1
D4PG, Non-Prioritized, N = 5
D4PG, Prioritized, N = 1
D4PG, Prioritized, N = 5
DDPG

Hoffman, Barth-Maron et al., 2017

Acrobot(Swingup) · Acrobot(Swingup Sparse) · Finger(Turn Hard) · Fish(Swim) · Humanoid(Run) · Manipulator(Bring Ball) · Swimmer(Swimmer15) · Swimmer(Swimmer6)

D4PG Categorical, 5e-05
D4PG MoG, 5e-05
D4PG MoG, 0.0001

Hoffman, Barth-Maron et al., 2017

Hoffman, Barth-Maron et al., 2017

Catch, N = 1 · Match Moving Target In Hand, N = 1 · Pickup And Orient, N = 1

Catch, N = 5 · Match Moving Target In Hand, N = 5 · Pickup And Orient, N = 5

Episode Returns vs Training Time (Hours)

D3PG, Non-Prioritized
D3PG, Prioritized
D4PG, Non-Prioritized
D4PG, Prioritized

Hoffman, Barth-Maron et al., 2017

# Distributional Distributed DDPG (D4PG)

Standard Networks

Parkour Networks



Hoffman, Barth-Maron et al., 2017

N = 1

N = 5

N = 1

N = 5

- ——— D3PG, Non-Prioritized
- ——— D3PG, Prioritized
- ——— D4PG, Non-Prioritized
- ——— D4PG, Prioritized
- – – – PPO

Hoffman, Barth-Maron et al., 2017

Hoffman, Barth-Maron et al., 2017

# Playing Go with Deep Networks and Planning



Policy network

$p_{\sigma/\rho}(a|s)$

Value network

$v_\theta(s')$

$$\rho(s_{t+1} | s_t, a_t)$$

Use environment model
in order to plan!

Silver, Huang et al., *Nature*, 2016

# Training Policy and Value Networks



Silver, Huang et al., *Nature*, 2016

# Planning with an Environment Model & MCTS



Silver, Huang et al., *Nature*, 2016

# Planning with an Environment Model



Silver, Huang et al., *Nature*, 2016

# Playing Go with Without Human Knowledge



Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

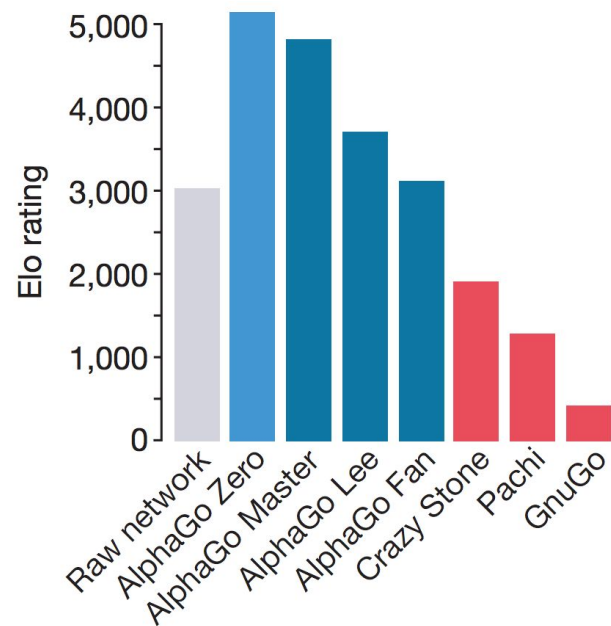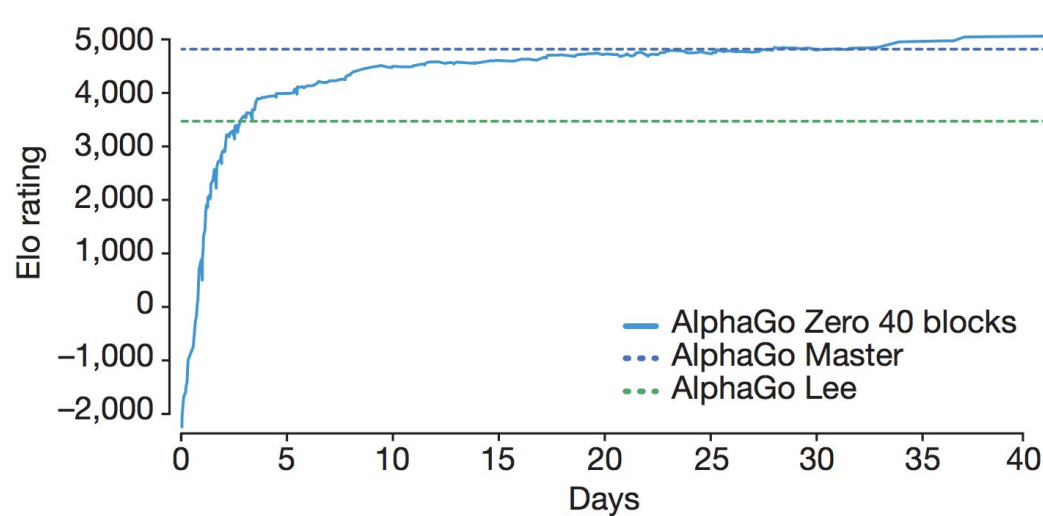# Playing Go with ~~Without~~ Without Human Knowledge

Neural network training



$$(\boldsymbol{p}, v) = f_\theta(s) \ \text{ and } \ l = (z - v)^2 - \boldsymbol{\pi}^{\mathrm{T}} \log \boldsymbol{p} + c\|\theta\|^2$$
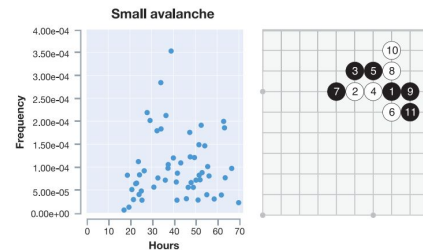
Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

# Playing Go with Without Human Knowledge



Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

# Playing Go with Without Human Knowledge



Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

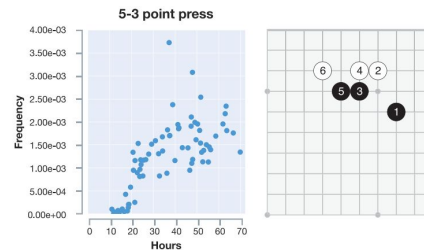# Playing Go with Without Human Knowledge

# Playing Go with Without Human Knowledge



Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

# Playing Go with Without Human Knowledge



Silver, Schrittwieser, Simonyan, et al. *Nature*, 2017

# Questions?