# **Deep Reinforcement Learning**
# Introduction and State-of-the-art

Arjun Chandra
Research Scientist
Telenor Research / Telenor-NTNU AI Lab
arjun.chandra@telenor.com
@boelger

24 October 2017

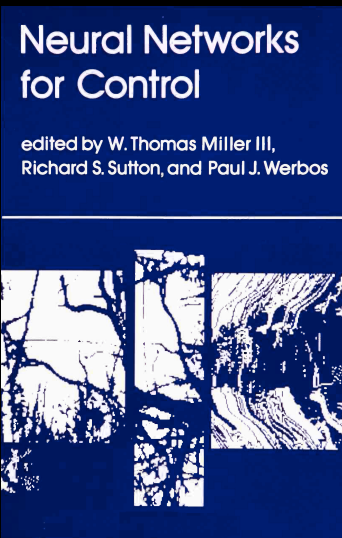https://join.slack.com/t/deep-rl-tutorial/signup

# The Plan

- Some history
- RL and Deep RL in a nutshell
- Deep RL Toolbox
- Challenges and State-of-the-art
  - Data Efficiency
  - Exploration
  - Temporal Abstractions
  - Generalisation

# Robot Motor Skill Coordination with EM-based Reinforcement Learning

Petar Kormushev, Sylvain Calinon, and Darwin G. Caldwell
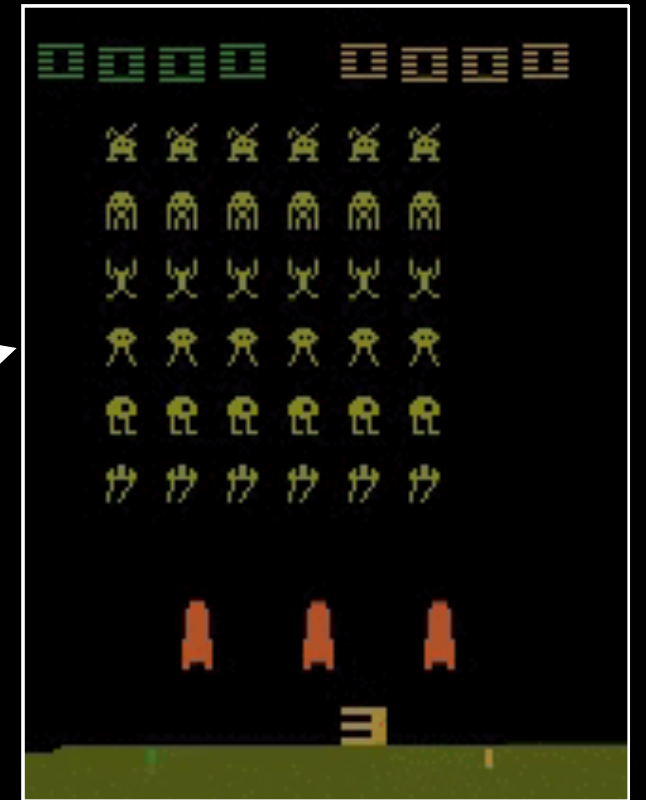
Italian Institute of Technology

# Brief History

Rich Sutton et al.

**Neural Networks for Control**

edited by W. Thomas Miller III, Richard S. Sutton, and Paul J. Werbos

late
1980s

RL for robots using NNs, L-J Lin. **PhD 1993, CMU**

Gerald Tesauro

1995

Stanford

http://heli.stanford.edu/

2004

Google DeepMind

Vlad Mnih et. al.

2013 —

David Silver et. al.

2015 —

# Problem Characteristics

**dynamic**

**uncertainty**/volatility

uncharted/**unimagined**/ exception laden

**delayed** consequences

requires **strategy**

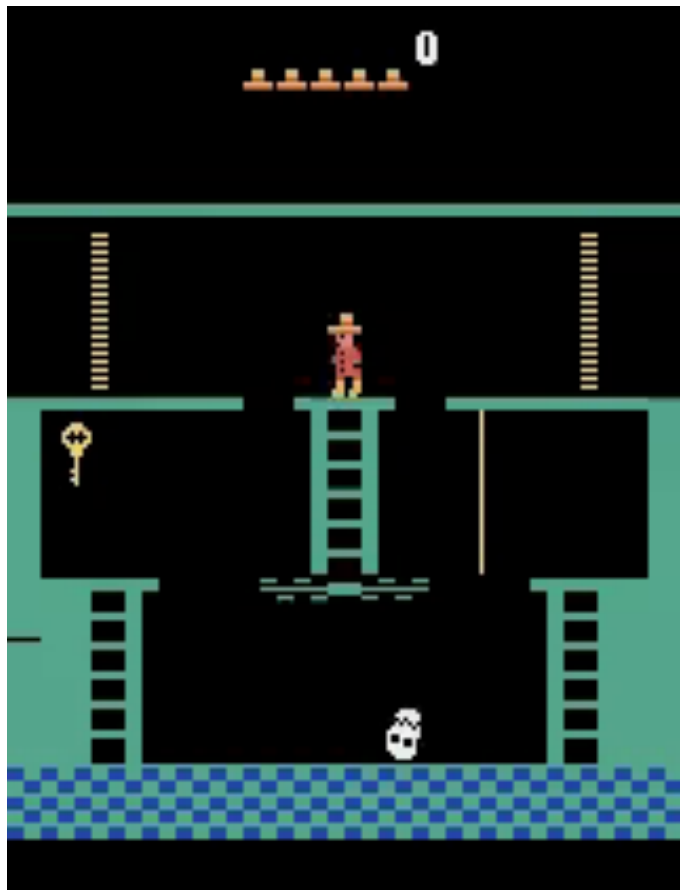# Solution

machine with **agency** which <span style="color:#a9c5d6">learn</span>, <span style="color:#c9a9c5">plan</span>, and <span style="color:#b5c08a">act</span> to find a strategy for solving the problem
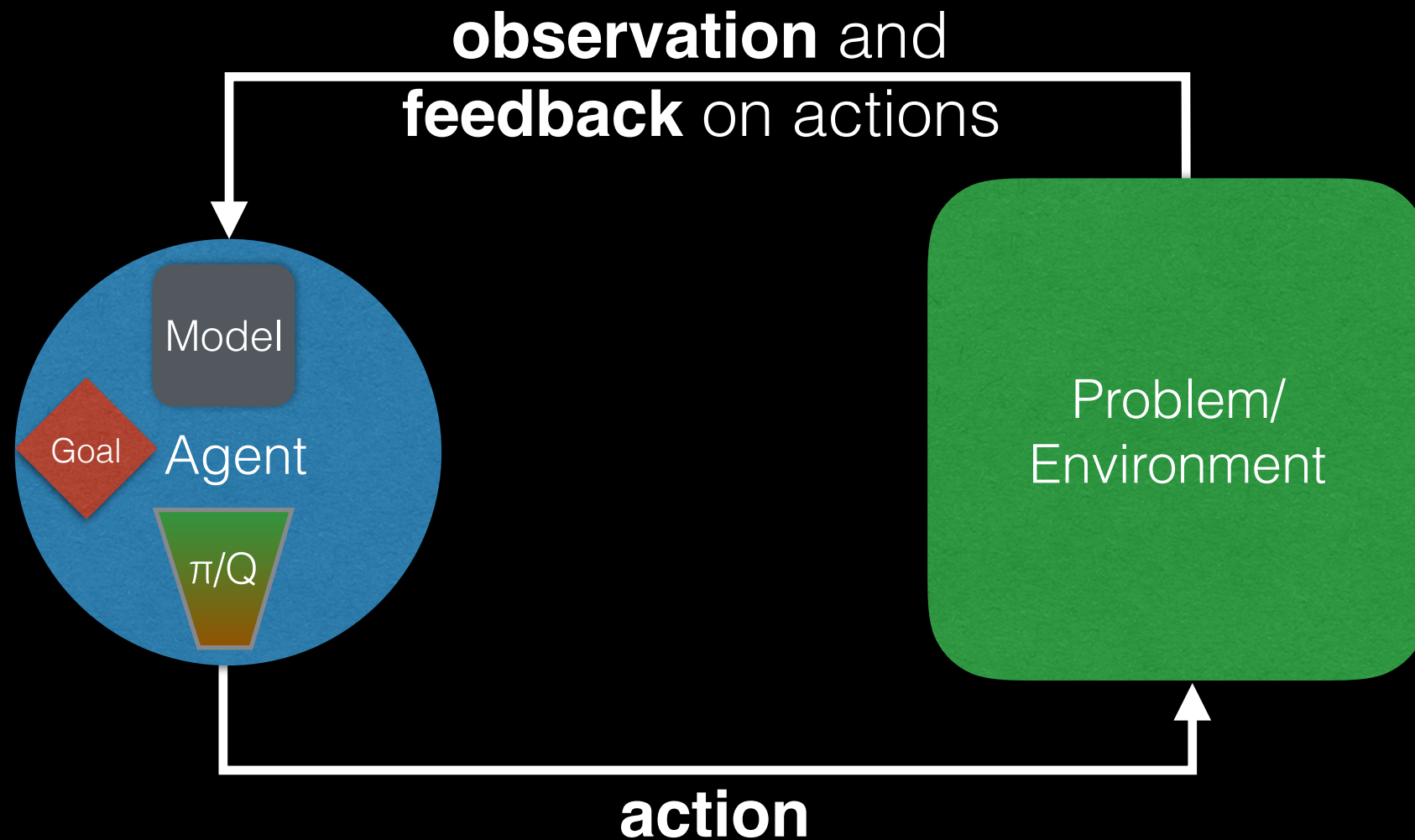


**autonomous** to some extent

**probe** and **learn from feedback**

focus on the **long-term objective**
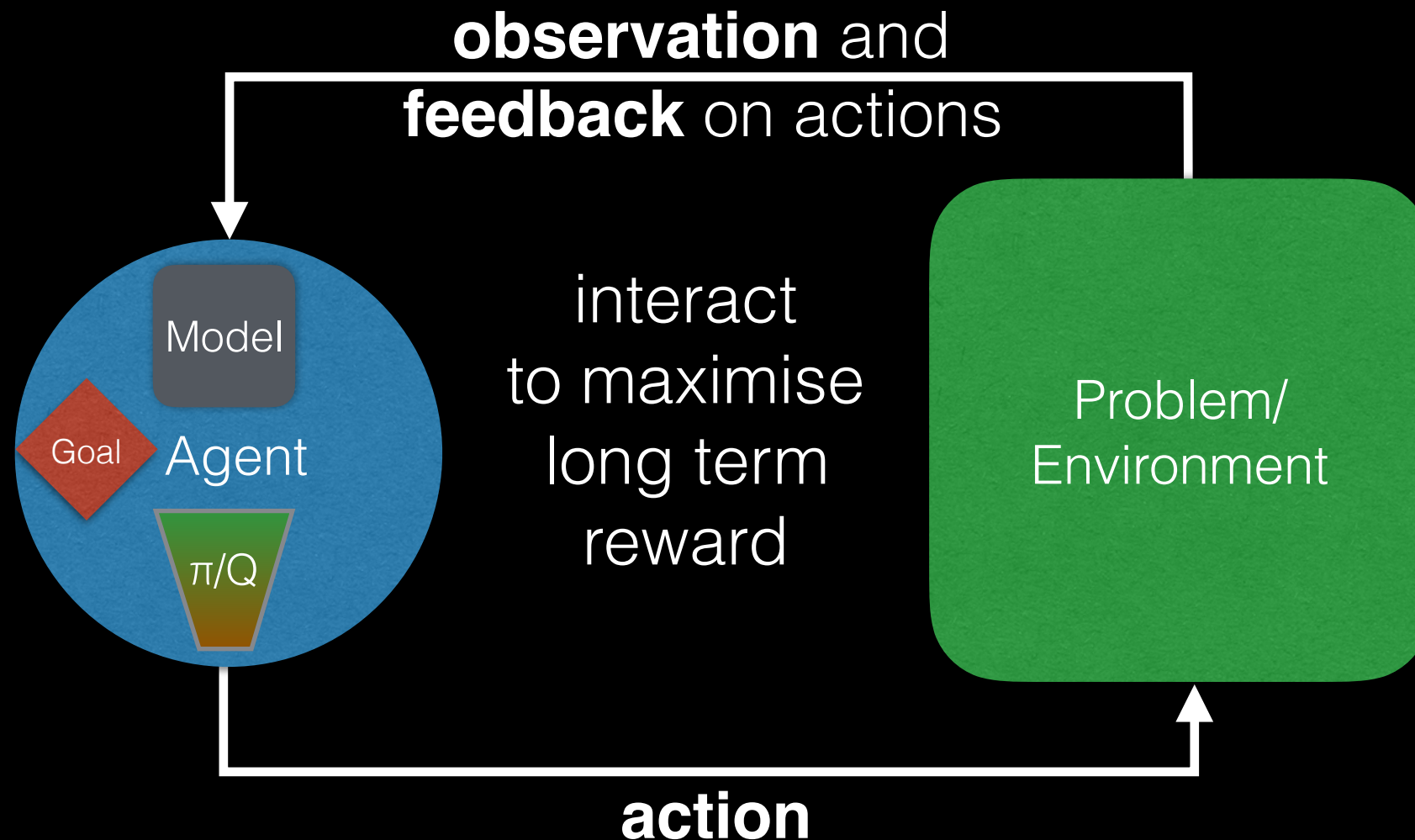
**explore** and **exploit**

# Reinforcement Learning



**observation** and
**feedback** on actions

Model

Goal    Agent

π/Q

Problem/
Environment

**action**

Goal  maximise return **E{R}**    Model  dynamics model    π/Q  policy/value function

# The MDP game!



**observation** and
**feedback** on actions

Model

Goal  Agent

π/Q

interact
to maximise
long term
reward

Problem/
Environment

**action**

Goal  maximise return **E{R}**

Inspired by Prof. Rich Sutton's tutorial:
https://www.youtube.com/watch?v=ggqnxyjaKe4

# The MDP (S,A,P,R,ϒ)

R: immediate reward function R(s, a)
P: state transition probability P(s'|s, a)

R=-10±3
P=0.99

R=10±3
P=1.00

R=40±3
P=0.01

R=-10±3
P=0.01

R=20±3
P=0.01

R=40±3
P=0.99

R=20±3
P=0.99

https://github.com/traai/basic-rl

# Terminology

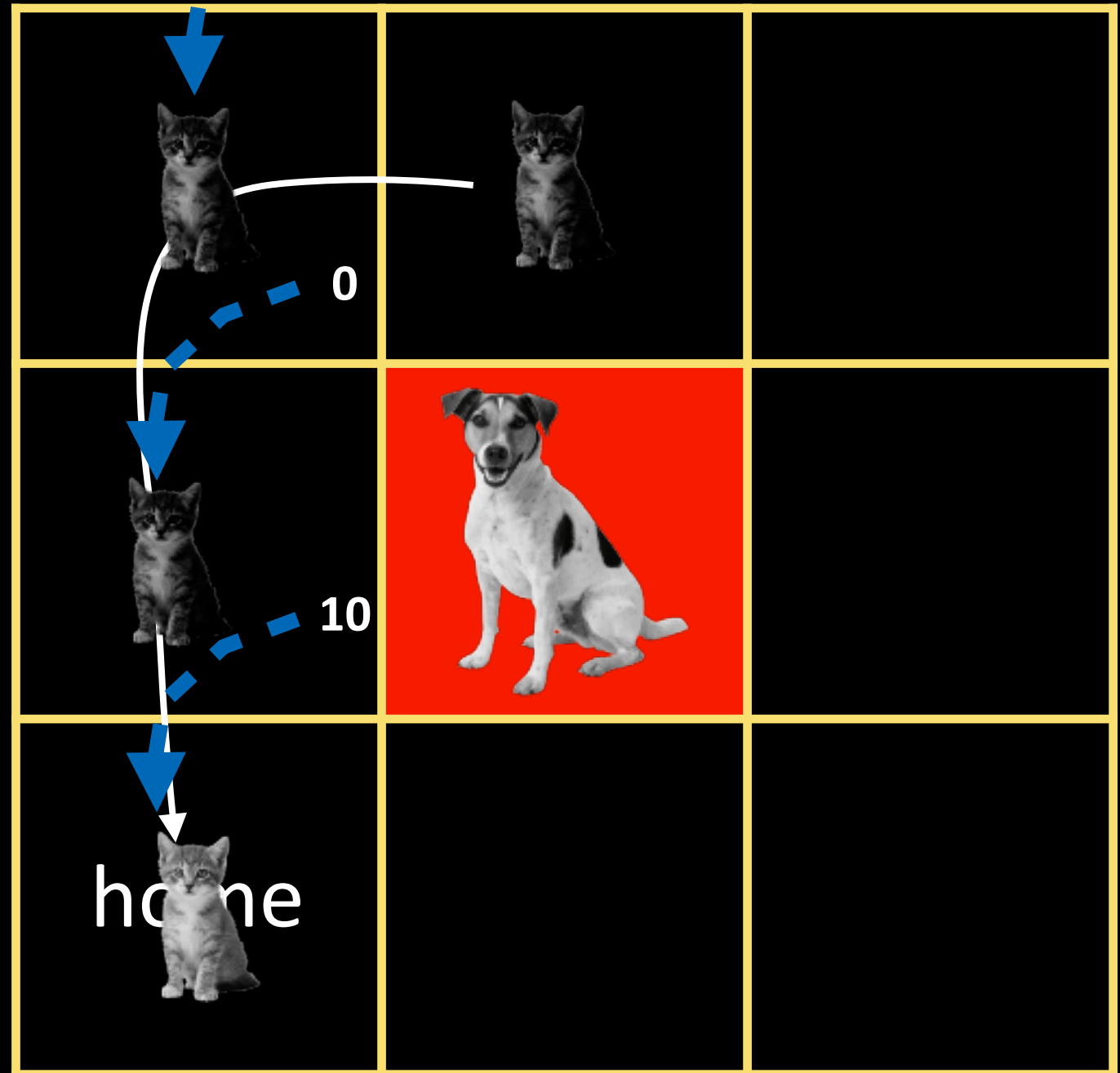state or action
value function

policy

dynamics model

reward

goal

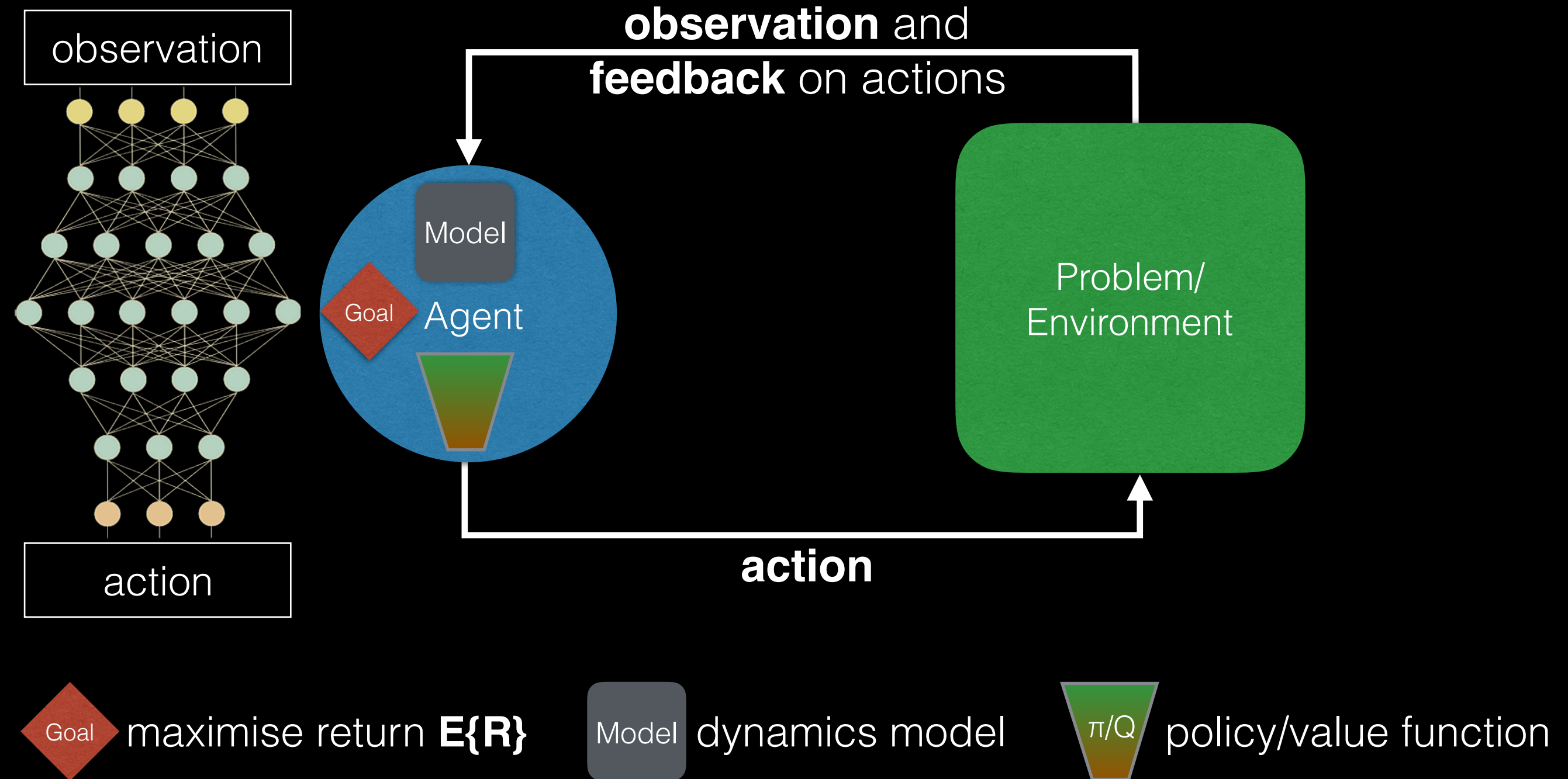

home

# Terminology

**state or action**
**value function**
Q(s,a) V(s)

policy

dynamics model

reward

goal

Q

Q  V  Q

Q

home

goal

# Terminology

state or action
value function

**policy** $\pi(s|a)$
$\pi(s)$

dynamics model

reward

home

goal

# Terminology

state or action
value function

policy

dynamics model

**reward**

goal

home

# Terminology

state or action
value function
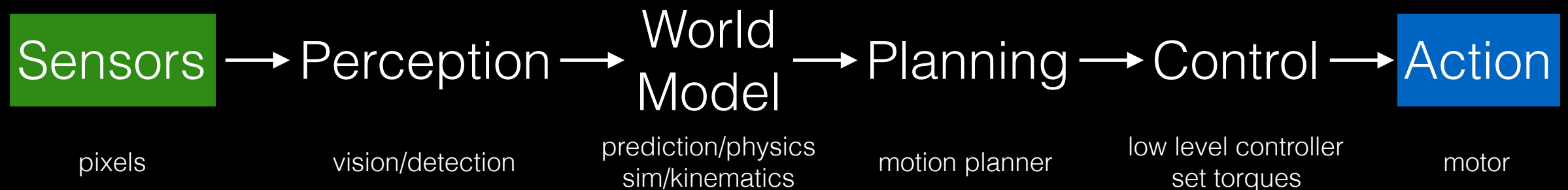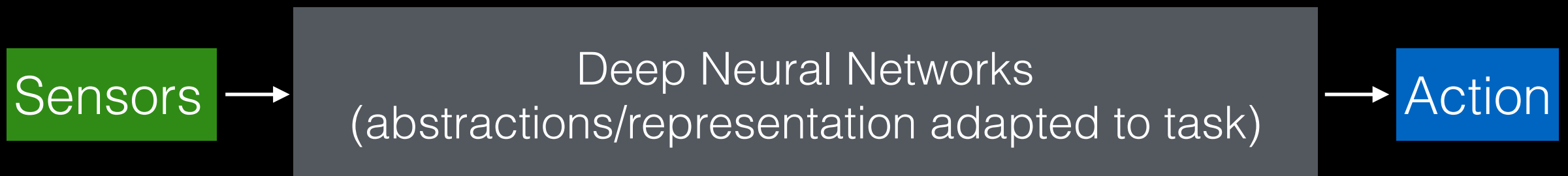
policy

dynamics model

reward

**goal**

# Deep Reinforcement Learning



observation

action

**observation** and
**feedback** on actions

Model

Goal  Agent

Problem/
Environment

**action**

Goal  maximise return **E{R}**   Model  dynamics model   π/Q  policy/value function

# Deep Reinforcement Learning

Sensors $\rightarrow$ Perception $\rightarrow$ World Model $\rightarrow$ Planning $\rightarrow$ Control $\rightarrow$ Action

pixels     vision/detection     prediction/physics sim/kinematics     motion planner     low level controller set torques     motor

## abstractions ~ info loss (manual craft)

Sensors $\rightarrow$ Deep Neural Networks (abstractions/representation adapted to task) $\rightarrow$ Action

**Explaining How a Deep Neural Network Trained with End-to-End Learning Steers a Car**, Bojarski et. al., https://arxiv.org/pdf/1704.07911.pdf
2017

# SL + RL



https://www.youtube.com/watch?v=NJU9ULQUwng



https://www.youtube.com/watch?v=KnPiP9PkLAs



data mismatch

# Toolbox

Standard algorithms to give you a
**flavour of the norm**!

# DQN



**image**
**score change**
on action

Buffer

Goal Agent

NN

**action**

**Human-level control through deep reinforcement learning**,
Mnih et. al., Nature 518, Feb 2015

# experience replay buffer



**save** transition in
memory

randomly **sample**
from memory
for training
= i.i.d

# freeze
# target

freeze

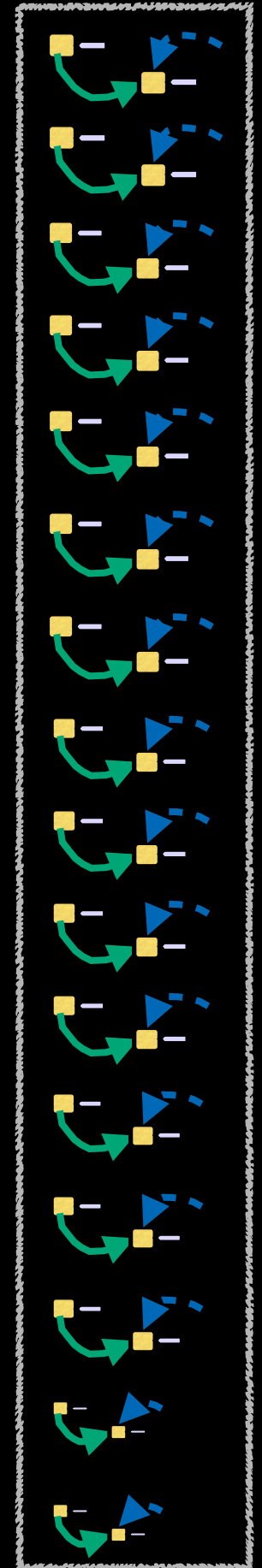$$\left( r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, \mathbf{w}) \right)^2$$

https://storage.googleapis.com/deepmind-media/dqn/
DQNNaturePaper.pdf

**Human-level control through deep reinforcement learning**, Mnih et. al., Nature 518, Feb 2015

# prioritised experience replay

sample
from memory
based on surprise

$$\left| r + \gamma \max_{a'} Q(s', a', \mathbf{w}^-) - Q(s, a, w) \right|$$

**Prioritised Experience Replay,** Schaul et. al., **ICLR 2016**

# dueling architecture



$$Q(s, a) = V(s) + A(s, a)$$

**Dueling Network Architectures for Deep RL** Wang et. al., **ICML 2016**

however
training is

SLOOooₒₒ....W

# Parallel Asynchronous Training

## value and policy based methods



https://youtu.be/0xo1Ldx3L5Q

https://youtu.be/Ajjc08-iPx8

https://youtu.be/nMR5mjCFZCw

parallel
agents

shared
parameters

lock-free
updates

**Asynchronous Methods for Deep Reinforcement Learning**, Mnih et. al., **ICML 2016**

parallel learners

Agent Copy

**Hogwild!** updates

shared params

Agent

**Hogwild!** updates

https://github.com/traai/async-deep-rl

# So 2016...

# Can we train even faster?

# PAAC
# (**P**arallel **A**dvantage **A**ctor-**C**ritic)



**1 GPU/CPU**

**Reduced** training time

**SOTA** performance

https://github.com/alfredvc/paac

**Efficient Parallel Methods for Deep Reinforcement Learning,**
A. V. Clemente, H. N. Castejón, and A. Chandra, **RLDM 2017**

Alfredo Clemente

# Challenges and SOTA

Data Efficiency

Exploration

Temporal Abstractions

Generalisation

# Data Efficiency

# Demonstrations



**observation** and
**feedback** on action

Buffer

Goal  Agent

NN

past
**observations,
action,
feedback**

**action**

# Deep RL with Unsupervised Auxiliary Tasks

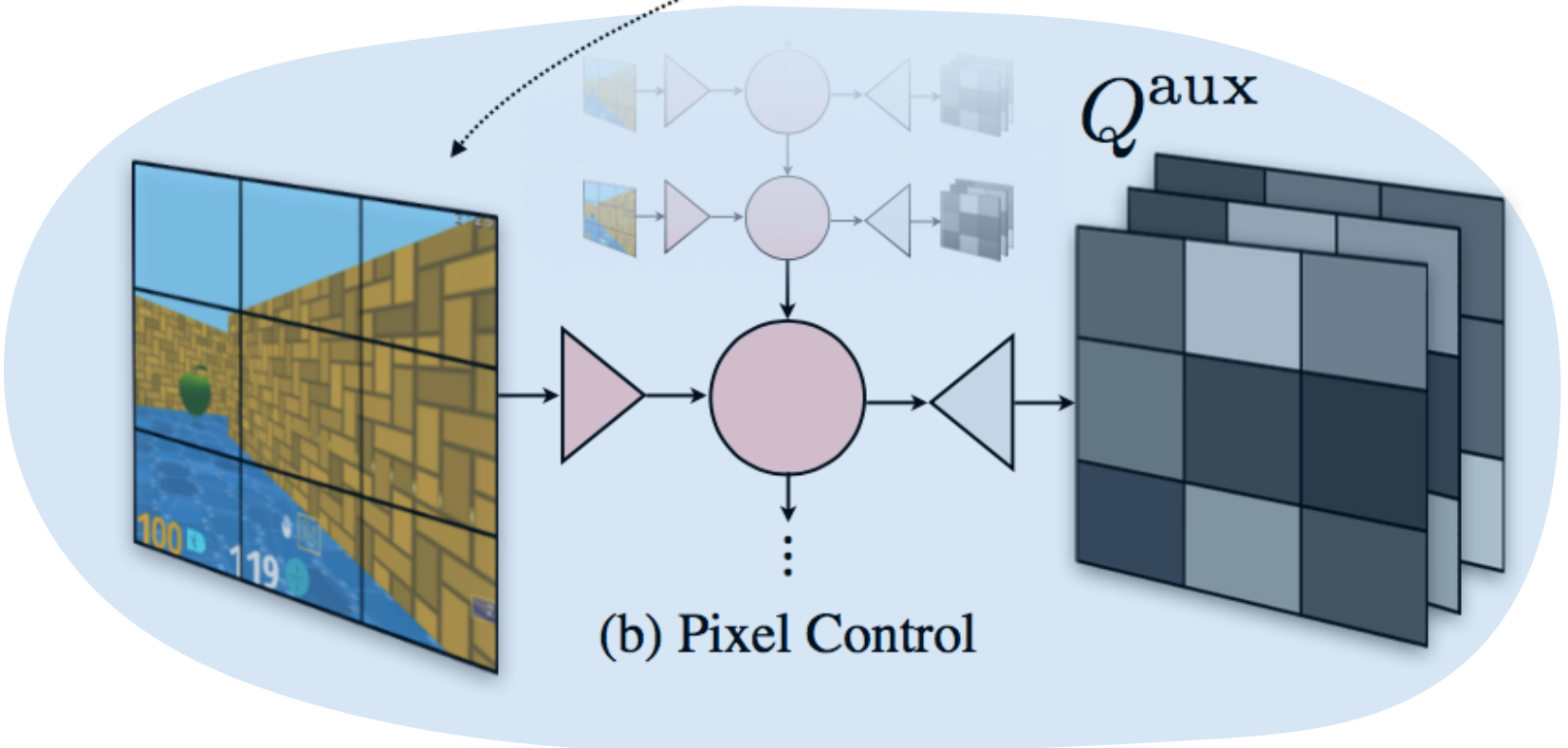Use replay buffer wisely

**observation** and
**feedback** on actions

Buffer

Goal     Agent

Problem/
Environment

**action**

**Reinforcement Learning with Unsupervised Auxiliary Tasks**,
Jaderberg et. al. ICML 2017

(a) Base A3C Agent

$V$ $\pi$ $V$ $\pi$ $V$ $\pi$ $V$ $\pi$

Agent LSTM

Agent ConvNet

Aux DeConvNet

Aux FC net

Environment

$o_t$

$r_t$   0   0   0   +1

$t_\tau$   $t_{\tau+1}$   $t_{\tau+2}$   $t_{\tau+3}$

Replay Buffer

$R \leftarrow R \leftarrow R \leftarrow$

$V$ $V$ $V$ $V$

(d) Value Function Replay

Skewed sampling

$r_\tau$

$Q^{\text{aux}}$

(b) Pixel Control

$t_{\tau-3}$   $t_{\tau-2}$   $t_{\tau-1}$

(c) Reward Prediction

**Reinforcement Learning with Unsupervised Auxiliary Tasks**, Jaderberg et. al. ICML 2017

(a) Base A3C Agent

Agent LSTM
Agent ConvNet
Aux DeConvNet
Aux FC net

Environment

$o_t$

$r_t$ 0 0 0 +1

$t_\tau$ $t_{\tau+1}$ $t_{\tau+2}$ $t_{\tau+3}$

$V$ $\pi$ $V$ $\pi$ $V$ $\pi$ $V$ $\pi$

Replay Buffer

$Q^{aux}$

(b) Pixel Control

learn to **act to affect pixels**

e.g. if grabbing fruit makes it disappear, agent would do it

(a) Base A3C Agent

Agent LSTM

Agent ConvNet

Aux DeConvNet

Aux FC net

Replay Buffer

Environment

$o_t$

$r_t$ 0 0 0 +1

$t_\tau$ $t_{\tau+1}$ $t_{\tau+2}$ $t_{\tau+3}$

predict

**short term reward**

e.g. replay pick key
series of frames

Skewed
sampling

$r_\tau$

$t_{\tau-3}$ $t_{\tau-2}$ $t_{\tau-1}$

(c) Reward Prediction

(a) Base A3C Agent

Agent LSTM

Agent ConvNet

Aux DeConvNet

Aux FC net

Replay Buffer

(d) Value Function Replay

predict

**long term reward**

**Reinforcement Learning with Unsupervised Auxiliary Tasks**, Jaderberg et. al. ICML 2017

https://deepmind.com/blog/reinforcement-learning-unsupervised-auxiliary-tasks/

# Distributional RL



**A Distributional Perspective on Reinforcement Learning,**

Bellemare et. al., ICML 2017

# Normal DQN **target**:
[sample **reward** after step + **discounted** previous **return** estimate from then on]

# **BUT this:**
[fuse **R** with **discounted** previous **return** **distribution**]

*Figure 4.* Learned value distribution during an episode of SPACE INVADERS. Different actions are shaded different colours. Returns below 0 (which do not occur in SPACE INVADERS) are not shown here as the agent assigns virtually no probability to them.

"If I shoot now, it is game over for me"

A Distributional Perspective on Reinforcement Learning, Bellemare et. al., ICML 2017

wrong/fatal actions

bimodal

under pressure

**A Distributional Perspective on Reinforcement Learning,** Bellemare et. al., ICML 2017

# Exploration

# Curiosity Driven Exploration

# Curiosity Driven Exploration



curiosity as next state prediction error
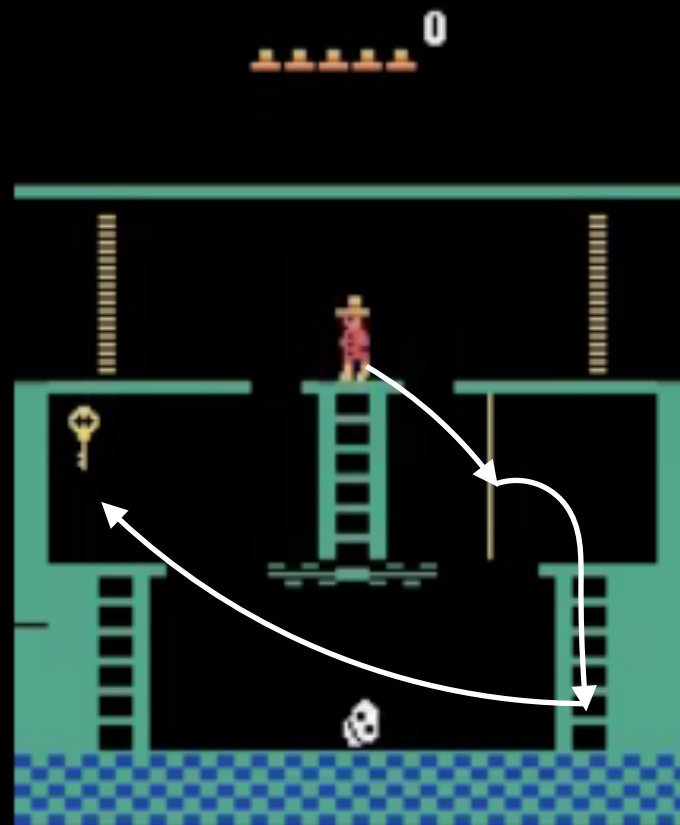
… only focus on **relevant parts of state**

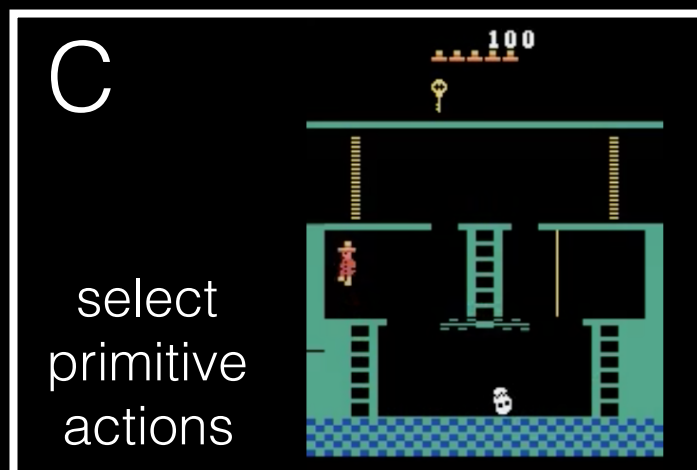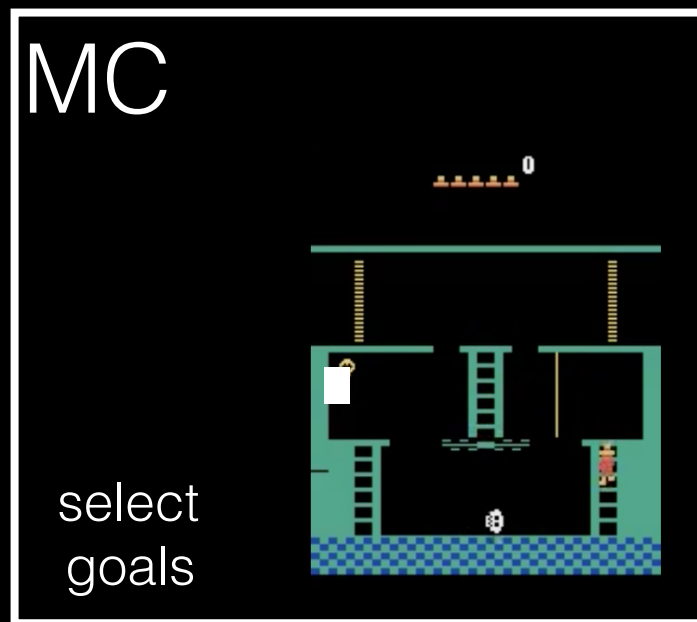**Curiosity-driven Exploration by Self-supervised Prediction**,
Pathak, Agrawal et al., ICML 2017.

Curiosity Driven Exploration
by Self-Supervised
Prediction

ICML 2017

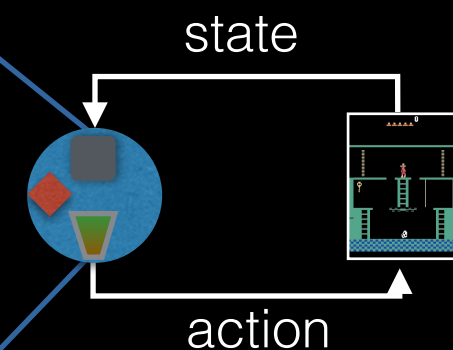Deepak Pathak, Pulkit Agrawal, Alexei Efros, Trevor Darrell
UC Berkeley

https://github.com/pathak22/noreward-rl
https://pathak22.github.io/noreward-rl/

# Temporal Abstractions

# HRL with pre-set Goals



MC

select goals

C

select primitive actions

**meta-controller** chooses **goals**

state

action

**controller** chooses **actions**

**Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation,** T. D. Kulkarni, K. R. Narasimhan et. al. NIPS 2016
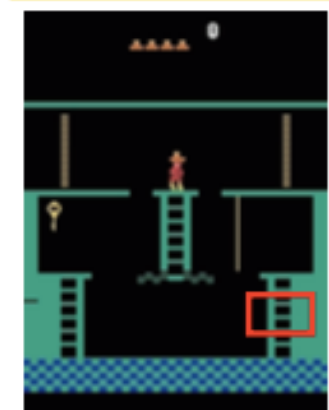
Meta Controller

Controller

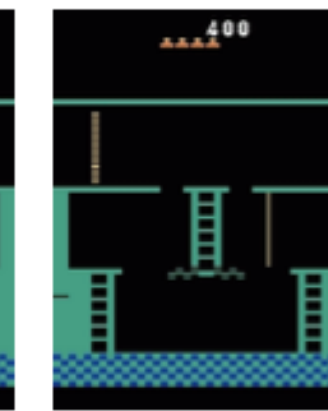termination (death)

goal reached

1     2     3     4     5     6

Meta Controller

Controller

goal reached

7     8     9     10     11     12

**Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation,** T. D. Kulkarni, K. R. Narasimhan et. al. NIPS 2016
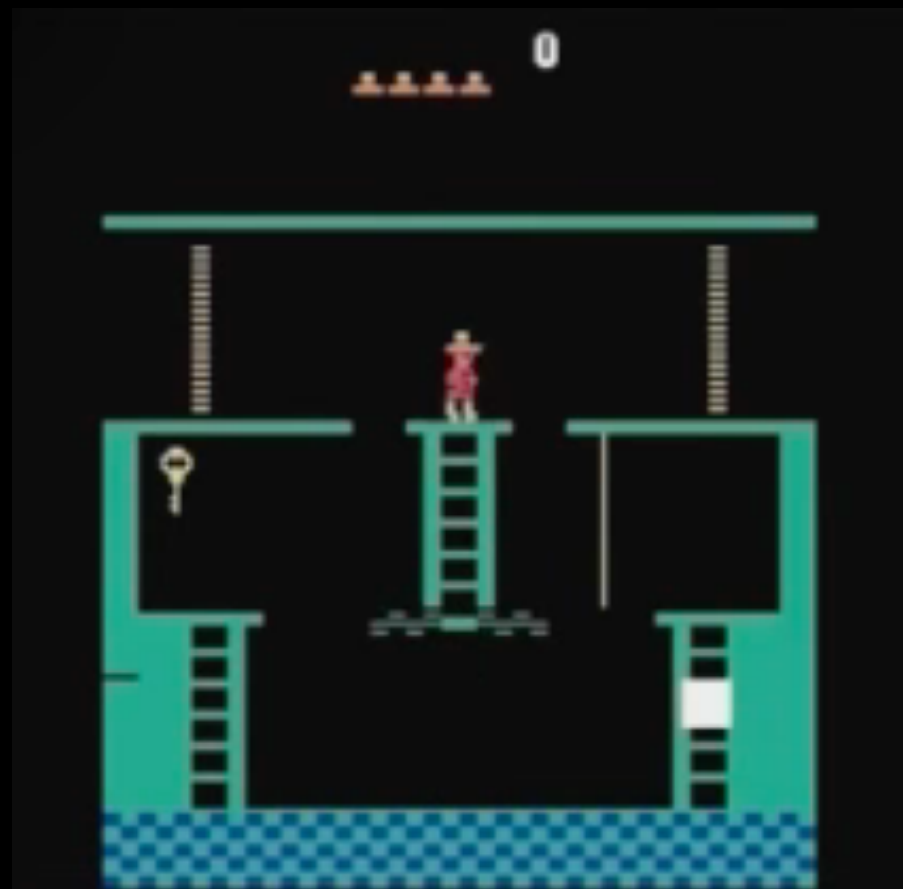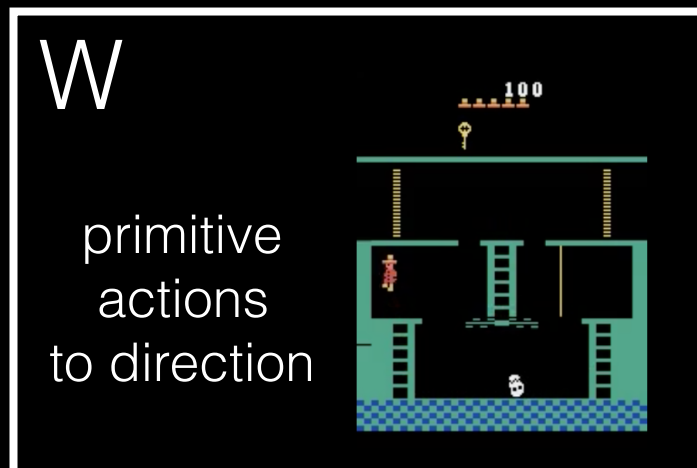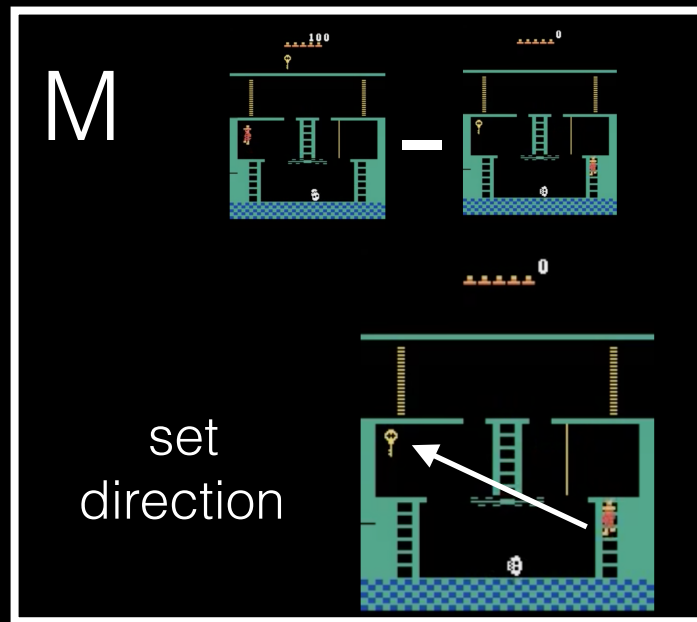
pre-defined goal
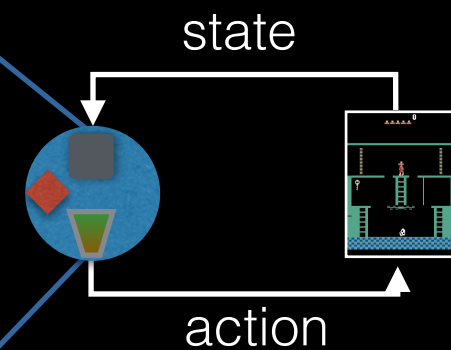selected by
meta-controller



**Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation,** T. D. Kulkarni, K. R. Narasimhan et. al. NIPS 2016
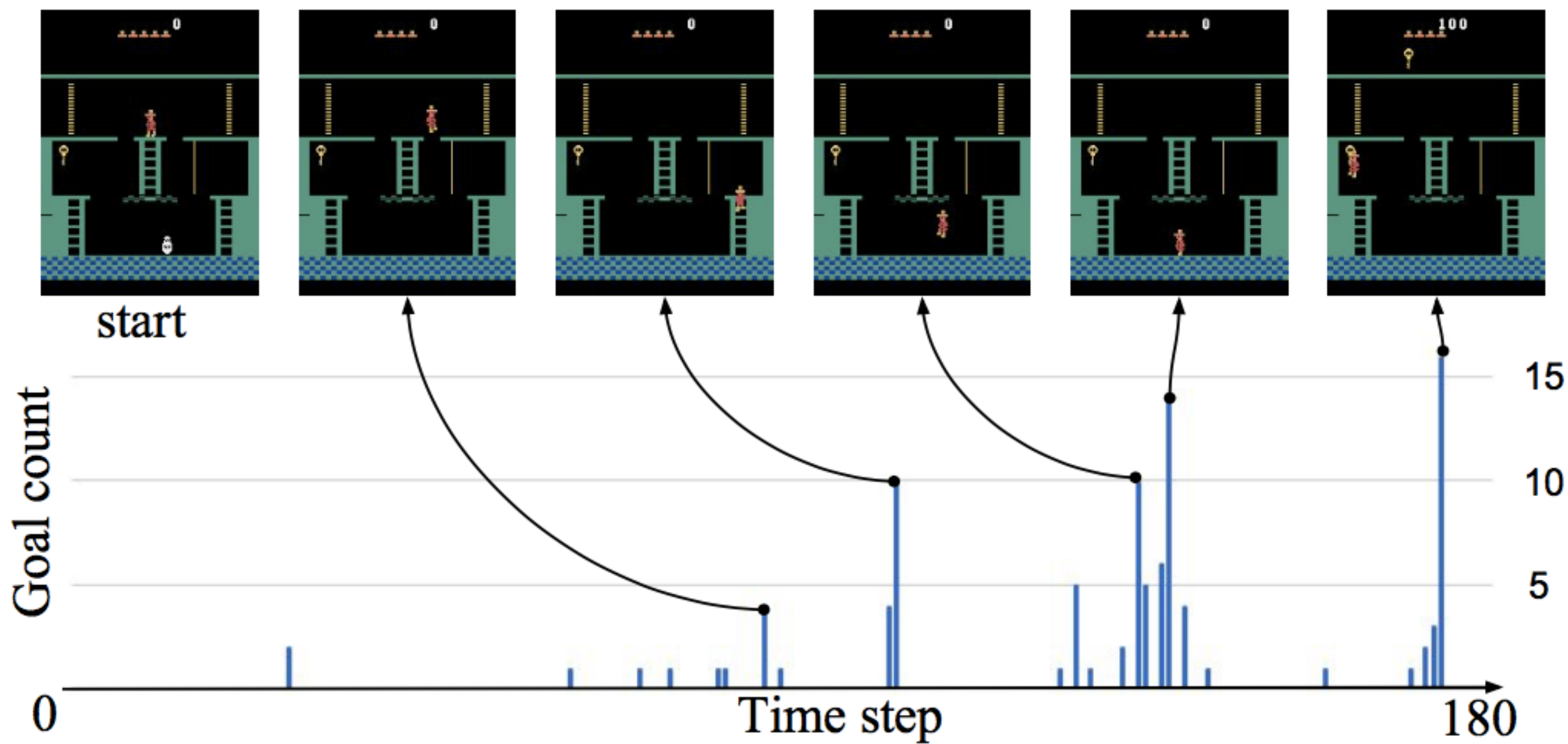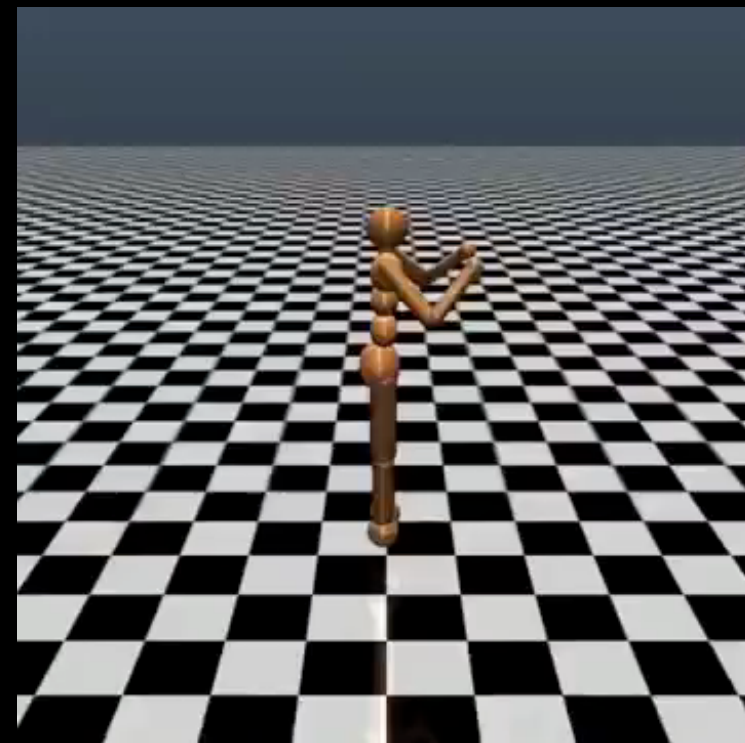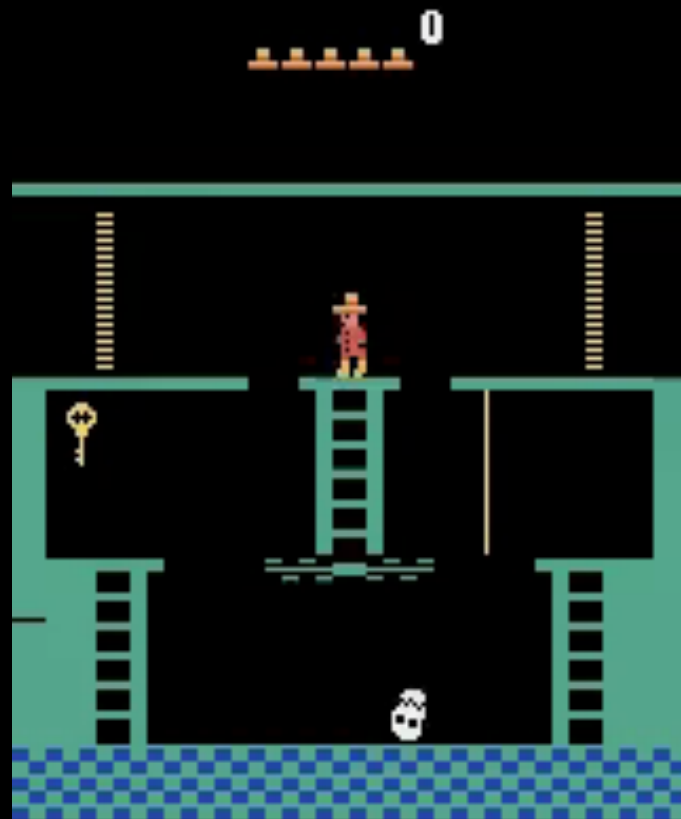
# FeUdal Networks for HRL



**manager** tries to finds **good** directions

**worker** tries to **achieve** them

state

action

M

set direction

W

primitive actions to direction

start

Goal count

Time step

0                                                                180

15

10

5

**FeUdal Networks for Hierarchical Reinforcement Learning**, Vezhnevets et. al. ICML 2017

# Generalisation
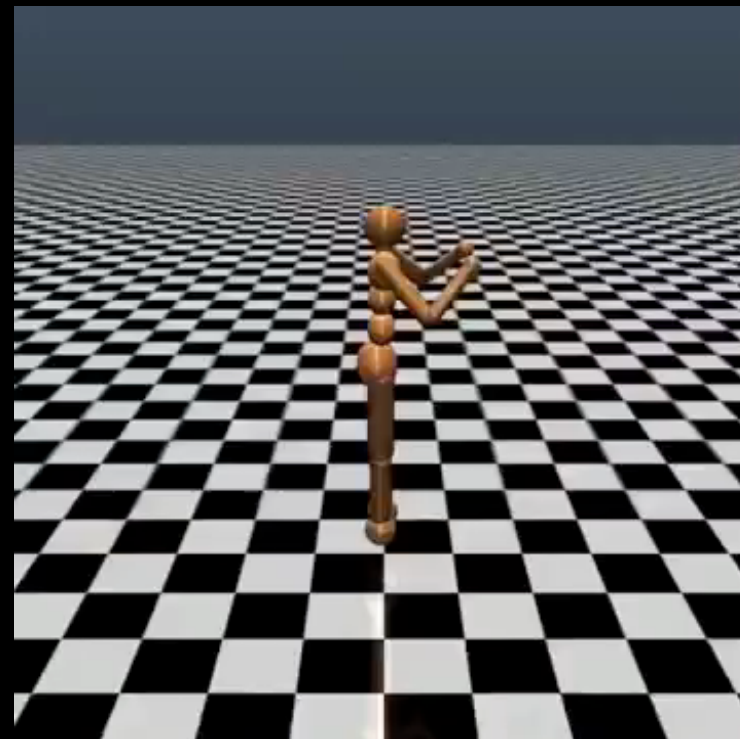
# Meta-learning
# (Learn to Learn)
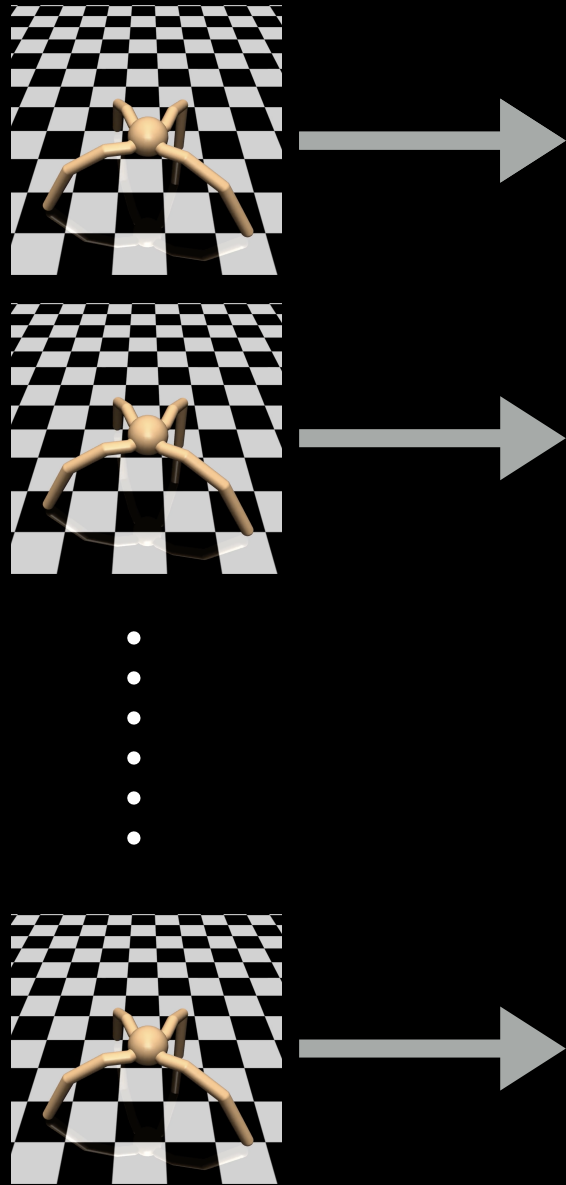## **Versatile** agents!

**Transfer** learning works with images



http://www.derinogrenme.com/2015/07/29/makale-imagenet-large-scale-visual-recognition-challenge/
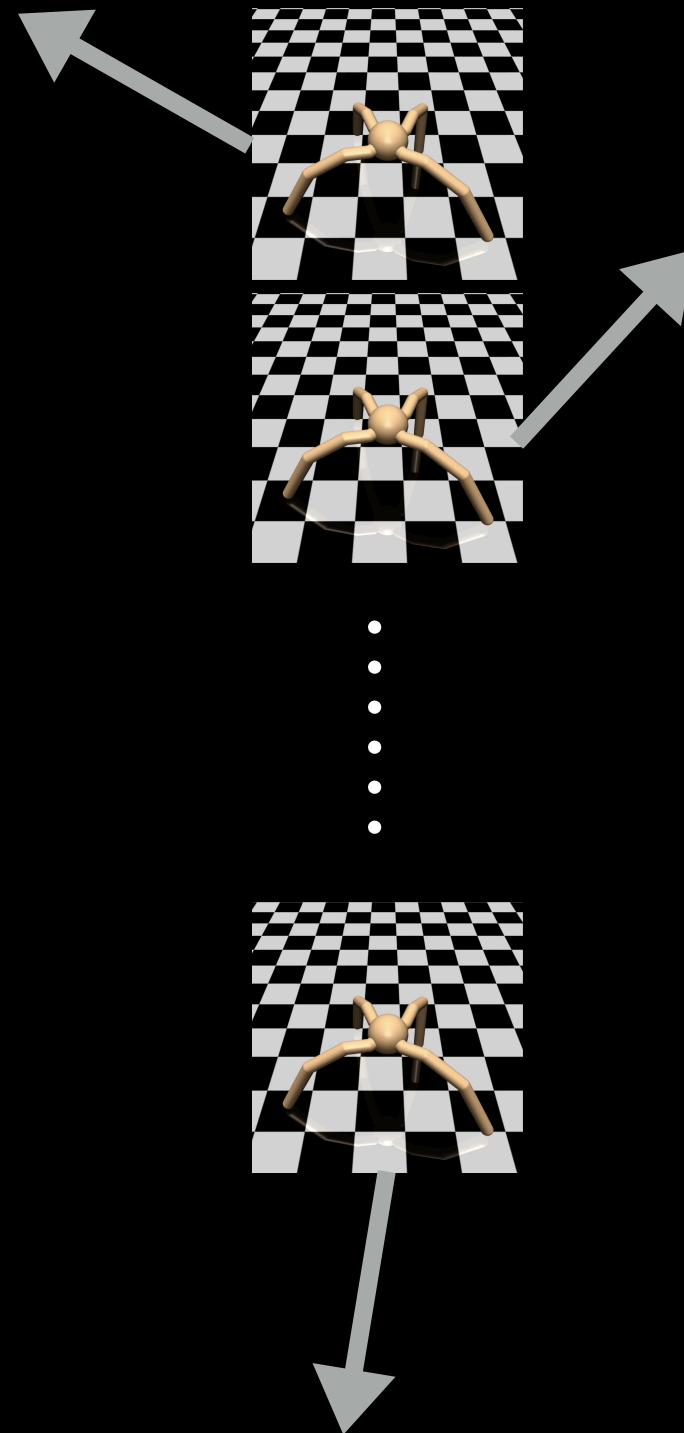
Good **features for decision making**?

learn
to go East

learn to
reduce learning
time to go to X

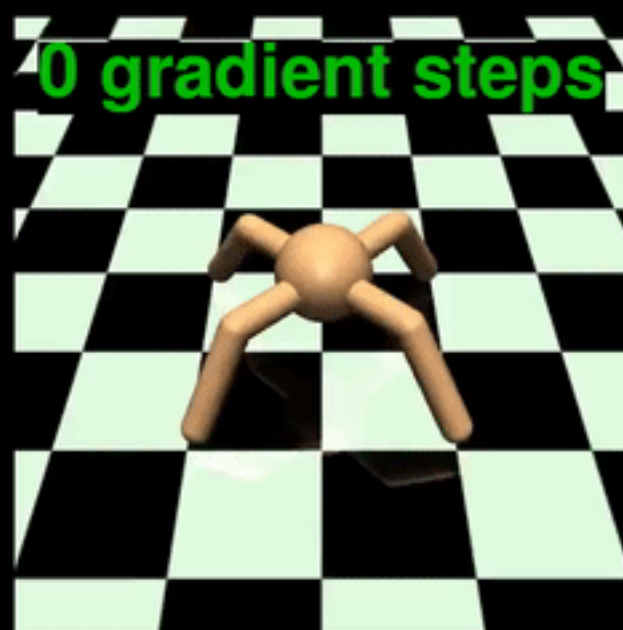**Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks.**
C. Finn, P. Abbeel, S. Levine. ICML 2017.

MAML — 0 gradient steps

0 grad/opt step:
policy ready
to learn

MAML — 0 gradient steps

1 grad/opt step:
learnt to
achieve goal

http://bair.berkeley.edu/blog/2017/07/18/learning-to-learn/
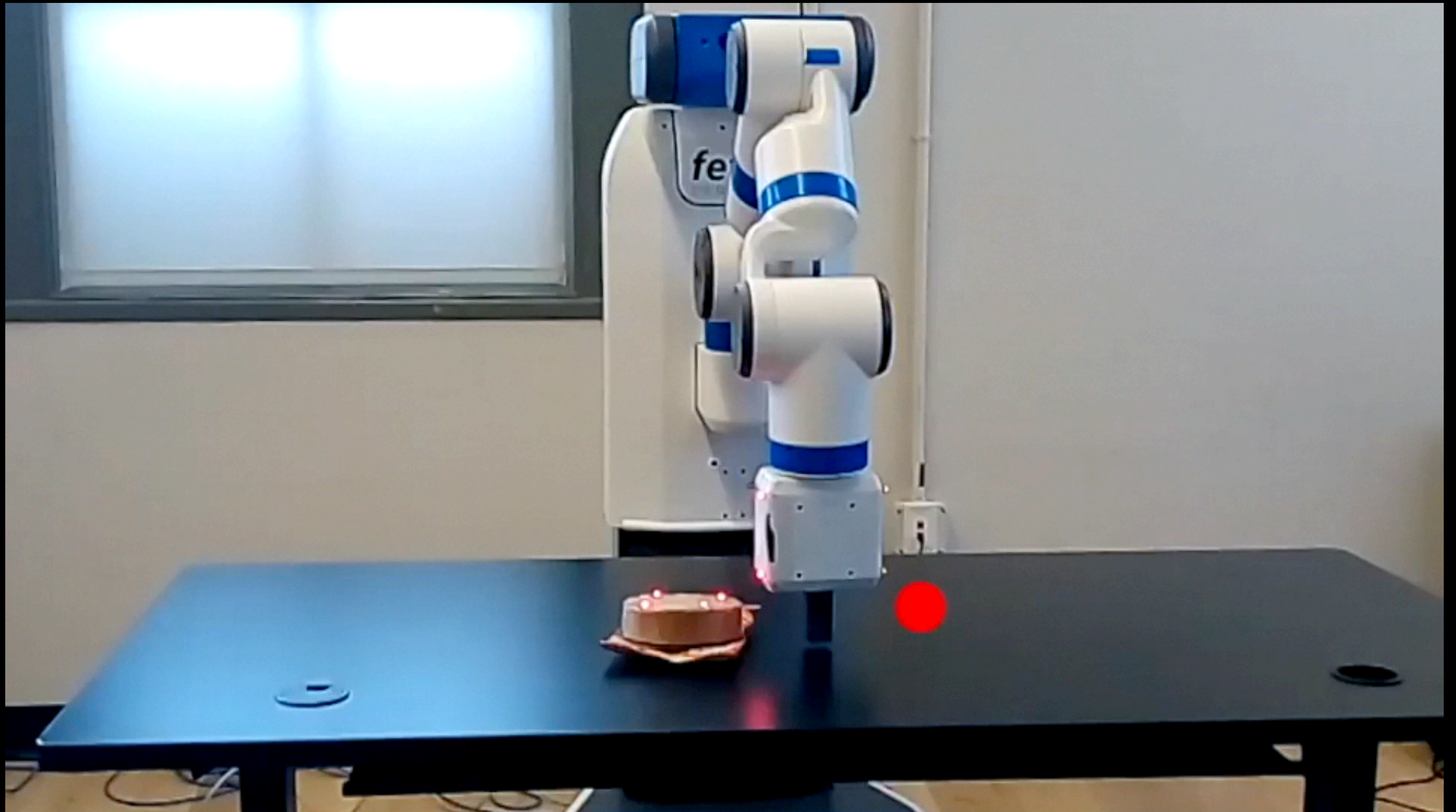Code: https://github.com/cbfinn/maml_rl
Videos: https://sites.google.com/view/maml

# Domain Randomisation

# Generalising
from Simulation

# Sim-to-Real Transfer of Robotic Control with Dynamics Randomization, Peng et al. arXiv preprint, **18 Oct 2017**
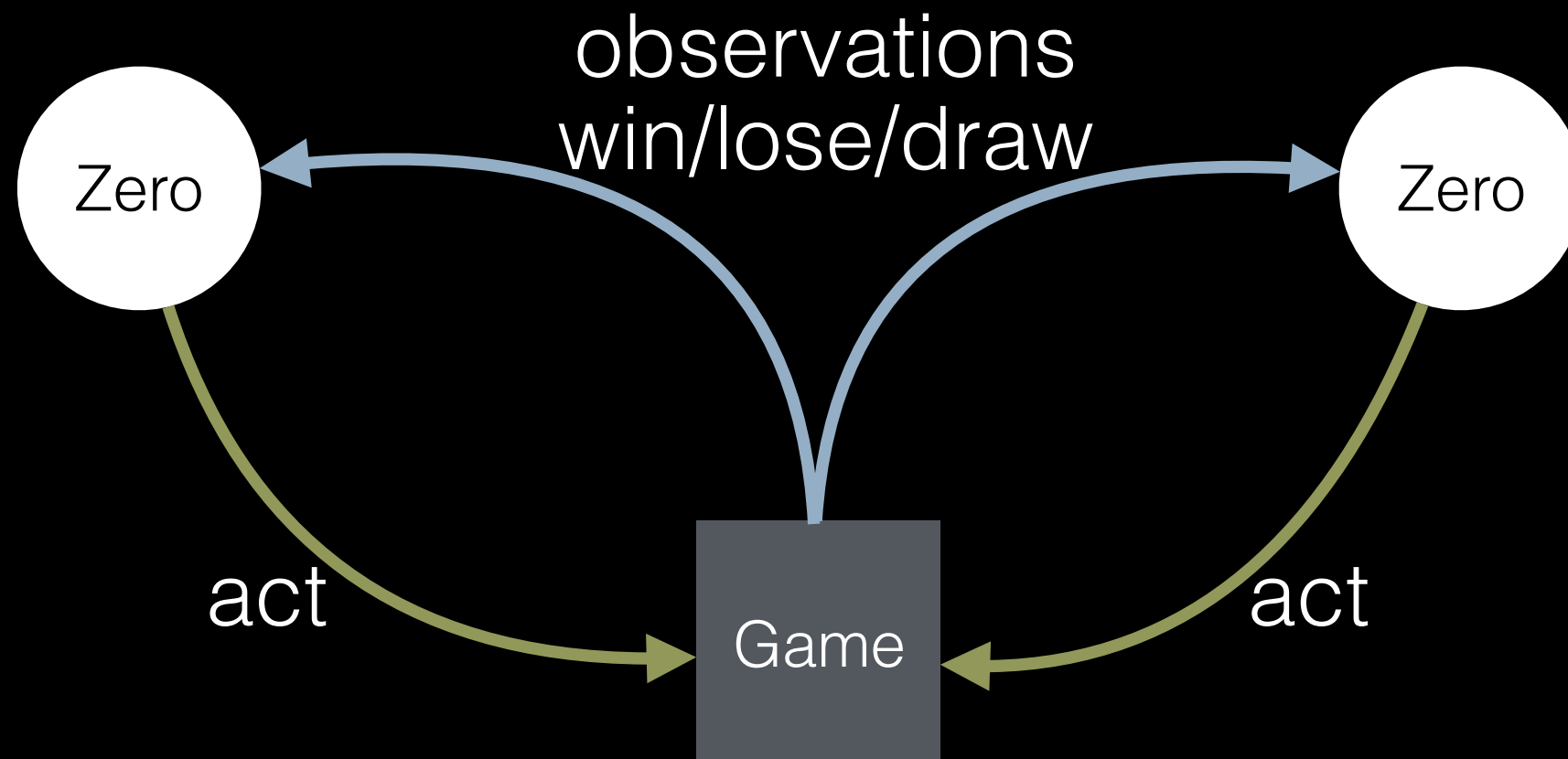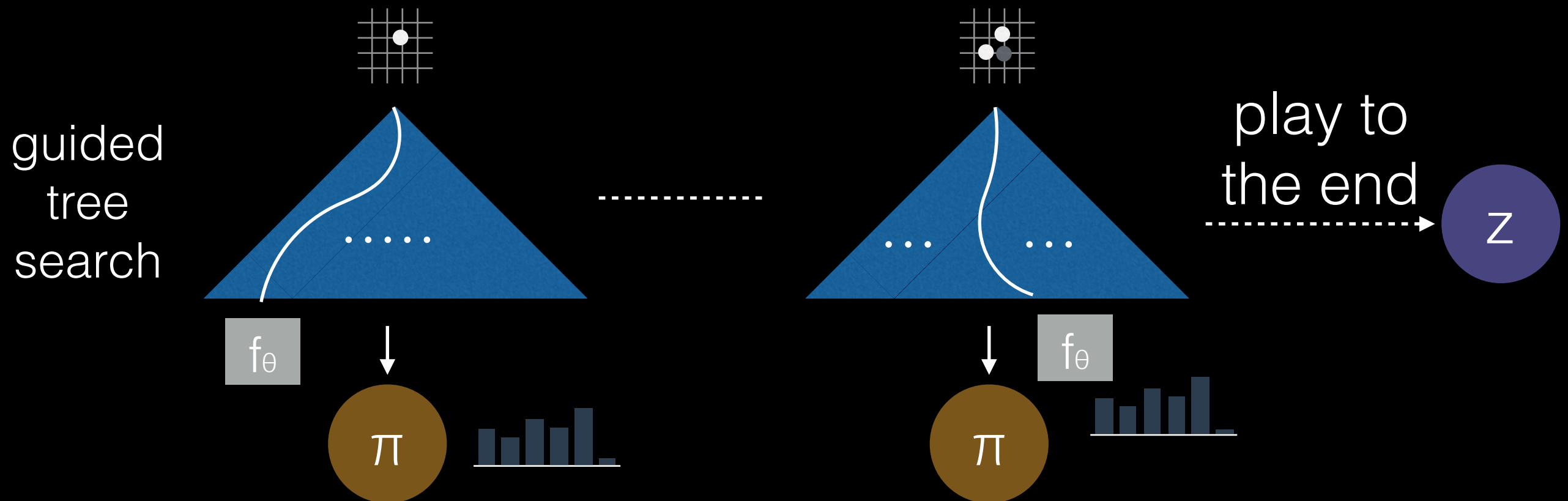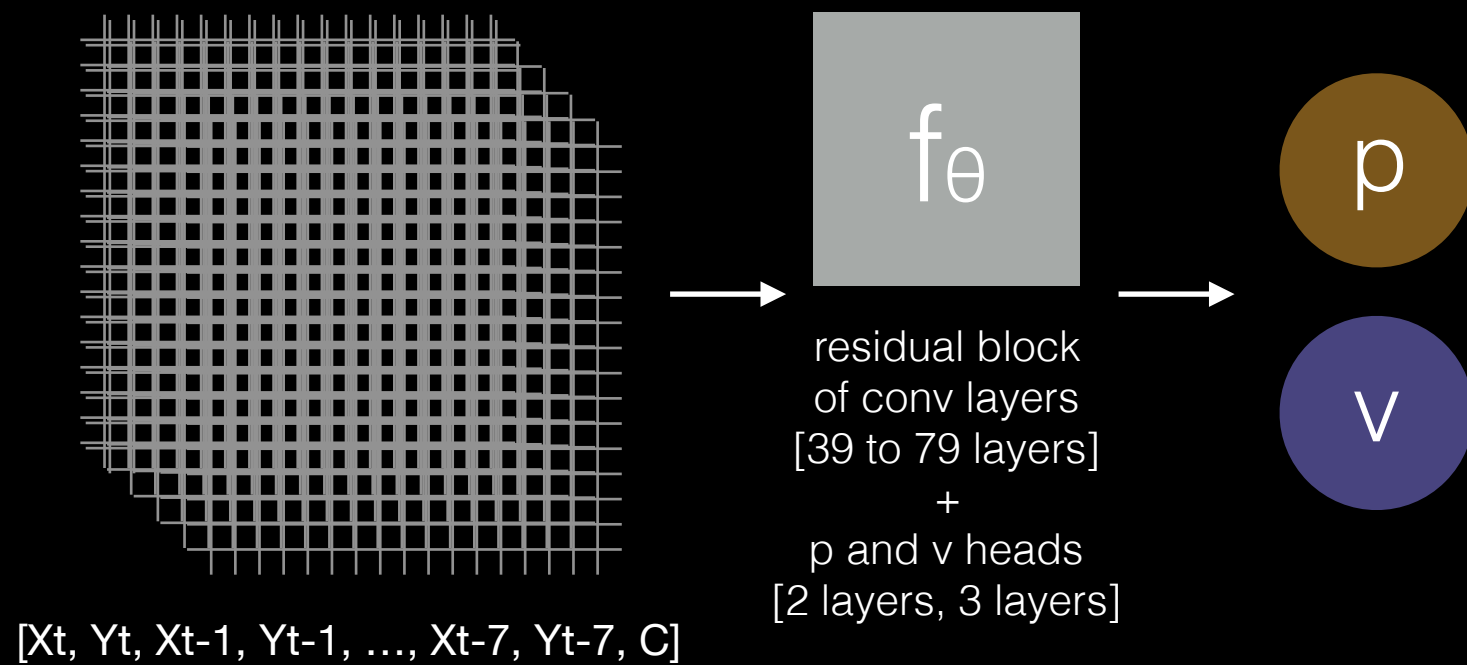


https://blog.openai.com/generalizing-from-simulation/

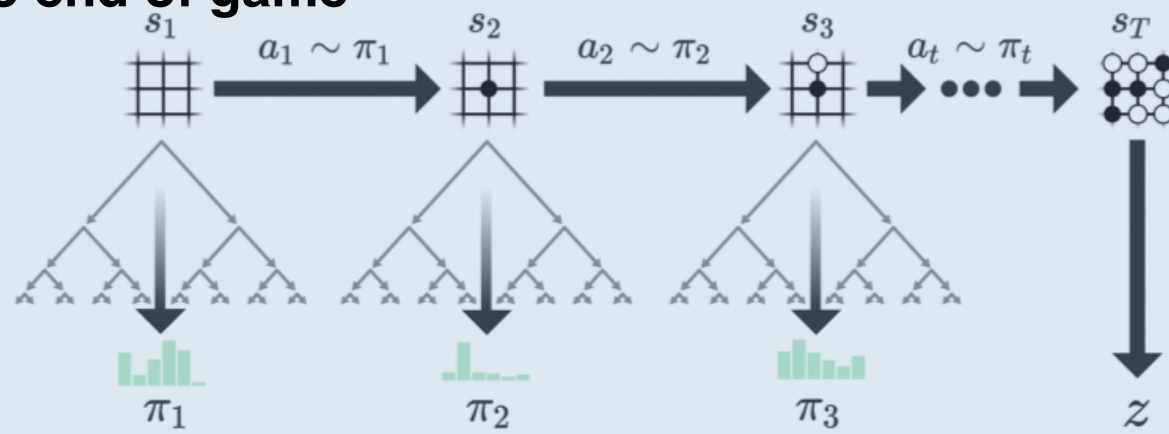# Generalisation via Self-play

# Deep RL in **AlphaGo Zero**

Improve
**thinking** and **intuition**
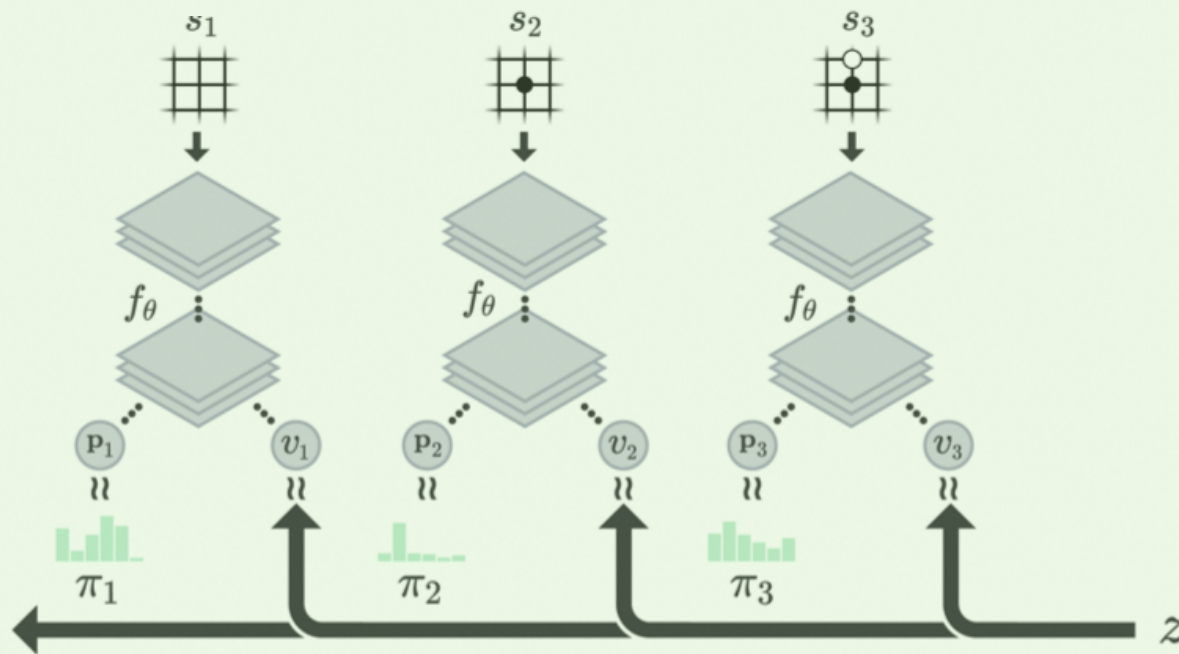with **feedback** **from** **self-play**
[**zero** human game data]

# Very High Level Mechanics



residual block
of conv layers
[39 to 79 layers]
+
p and v heads
[2 layers, 3 layers]

[Xt, Yt, Xt-1, Yt-1, ..., Xt-7, Yt-7, C]

guided
tree
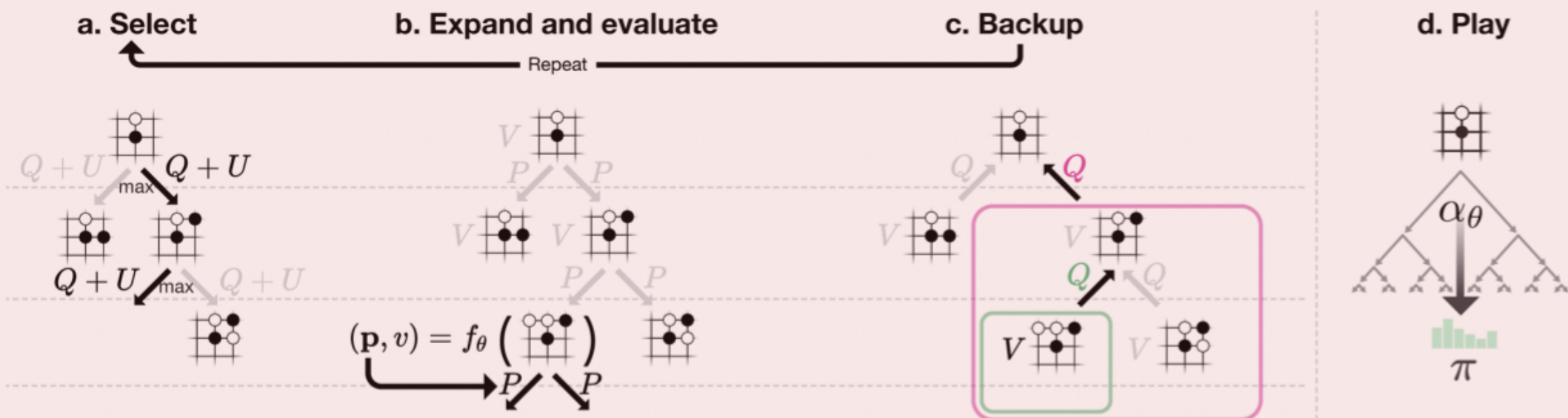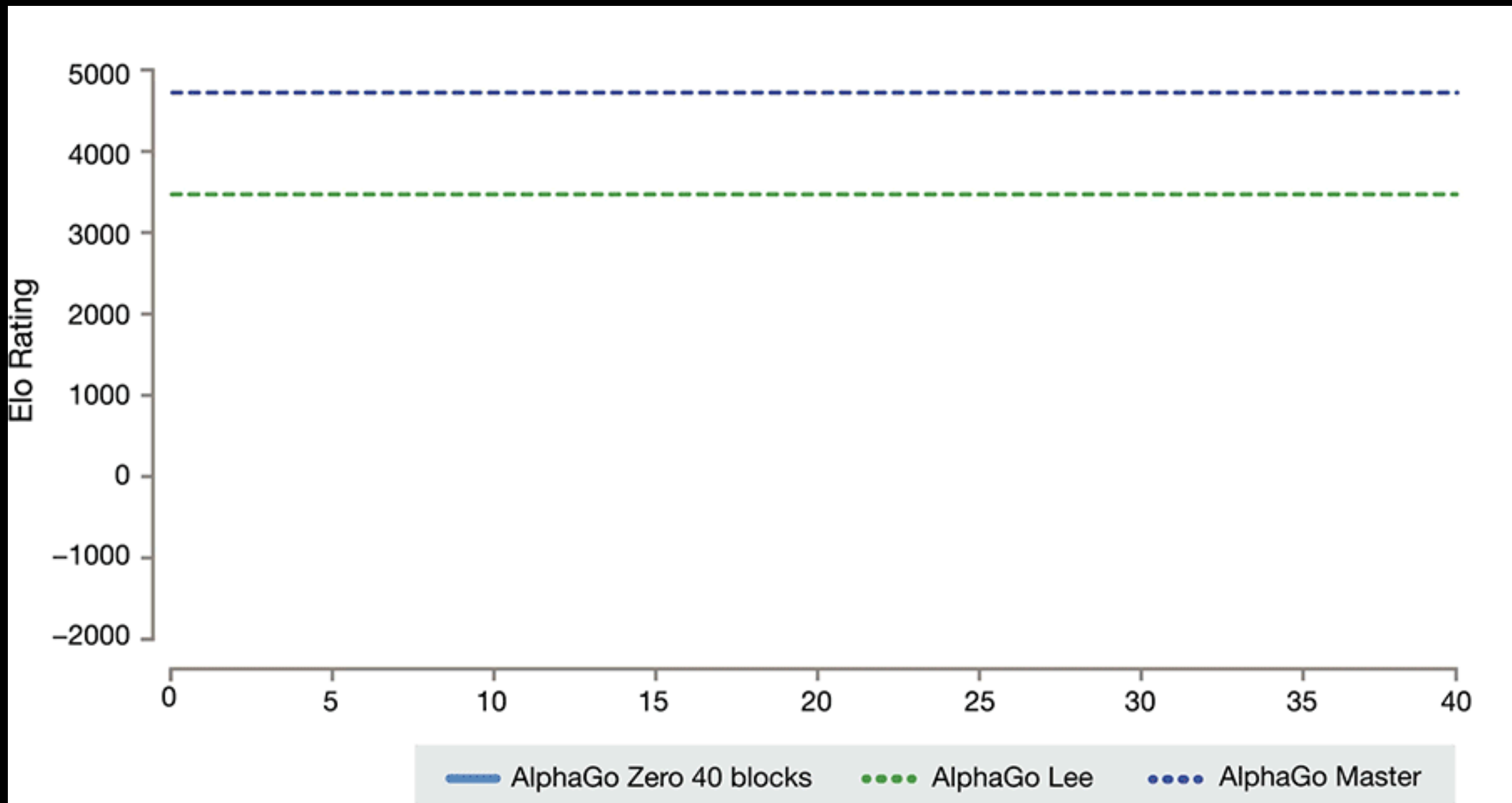search

play to
the end

**Self-play to end of game**

**NN training: learn to evaluate**

$$l = (z - v)^2 - \pi^{\mathrm{T}} \log p + c\|\theta\|^2$$

**Self-play step: select move by simulation + evaluation**

a. Select    b. Expand and evaluate    c. Backup    d. Play

**Mastering the game of Go without human knowledge**, Silver et.al., Nature, Vol. 550, **October 19, 2017**

https://deepmind.com/blog/alphago-zero-learning-scratch/

AlphaGo Zero
Discovering new knowledge

https://deepmind.com/blog/alphago-zero-learning-scratch/
https://www.youtube.com/watch?v=WXHFqTvfFSw

# Inspired to
# study RL much?

Next lecture:
Building Blocks of (Deep) RL
November 8, 2017

https://join.slack.com/t/deep-rl-tutorial/signup