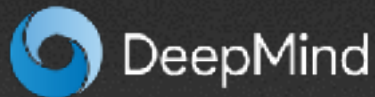# Deep Reinforcement Learning and Complex Environments

Raia Hadsell

DeepMind

# End-to-end Deep Learning for robots?

# End-to-end Deep Learning for robots?

**2010**: Speech Recognition

Audio → **Deep Net** → Text

# End-to-end Deep Learning for robots?

**2010**: Speech Recognition

Audio → **Deep Net** → Text

**2012:** Computer Vision

Pixels → **Deep Net** l → Labels

# End-to-end Deep Learning for robots?

**2010**: Speech Recognition

Audio → **Deep Net** → Text

**2012:** Computer Vision

Pixels → **Deep Net** → Labels

**2014:** Machine Translation

Text → **Deep Net** → Text

# End-to-end Deep Learning for robots?

**2010**: Speech Recognition

Audio → | **Deep Net** | Text

**2012:** Computer Vision

Pixels → | **Deep Net** | → Labels

**2014:** Machine Translation

Text → | **Deep Net** | → Text

**2017: Robotics?**

Sensors → Perception → World Model → Planning → Control → Action

# Robotics is different



LABELS

Google DeepMind

# Robotics is different

SENSORS $\longrightarrow$ ACTIONS
$\longleftarrow$

Google DeepMind

# Deep Reinforcement Learning



GOAL

OBSERVATIONS

REWARD?

Agent

Environment

neural
network

ACTIONS

Google DeepMind

General Artificial Intelligence

# General Atari Player



[Mnih et al, *Playing Atari with Deep Reinforcement Learning,* 2014]

# 9DOF Random reacher
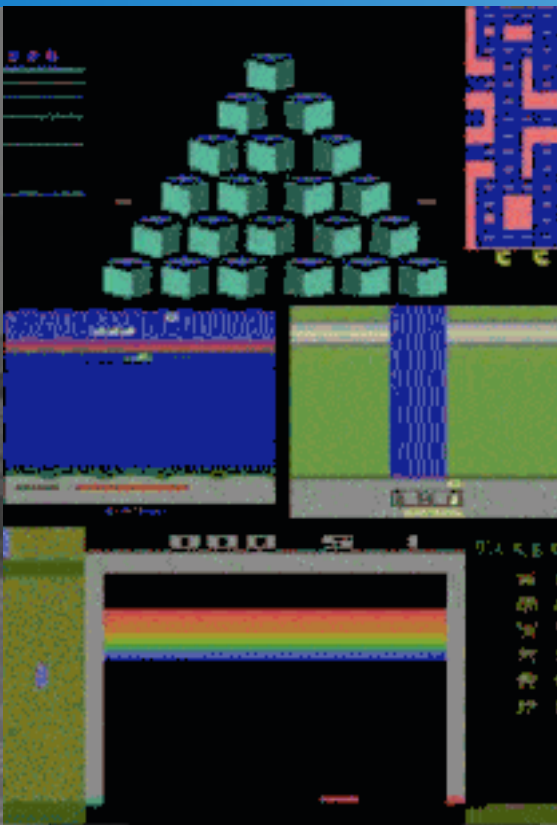
- Can deep RL agents learn multiple tasks?

- Can deep RL agents learn efficiently?

- Can deep RL agents learn from real data?

- Can deep RL agents learn continuous control?

Lab Mazes

Multiple Tasks
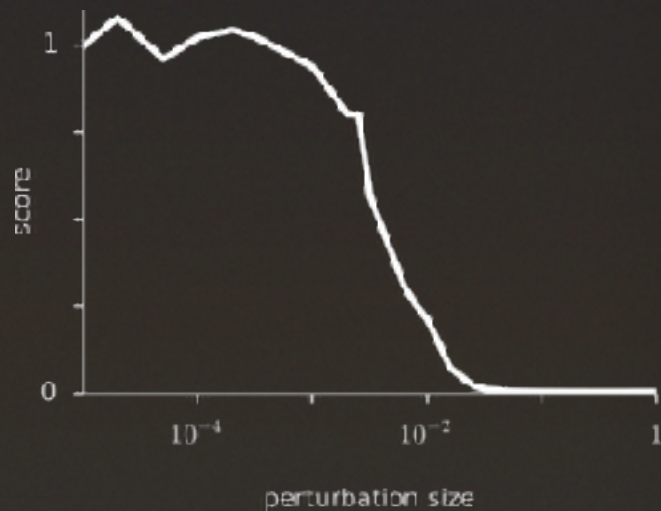&
Lifelong learning

StreetLearn

Parkour

# Lifelong Learning - 3 challenges

1. Catastrophic forgetting
2. Positive transfer
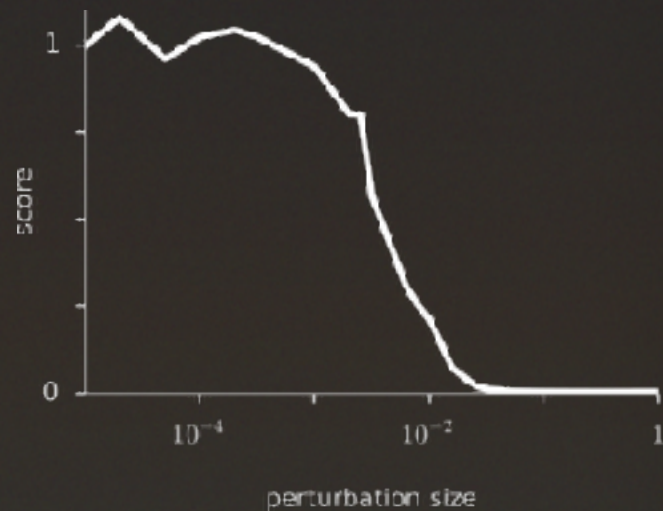3. Specialization and generalization

# Catastrophic forgetting

- Well-known phenomenon
- Especially severe in Deep RL

DeepMind

# Catastrophic forgetting

- Well-known phenomenon
- Especially severe in Deep RL

Raia Hadsell 2017

# Catastrophic forgetting

# Catastrophic forgetting

Raia Hadsell 2017
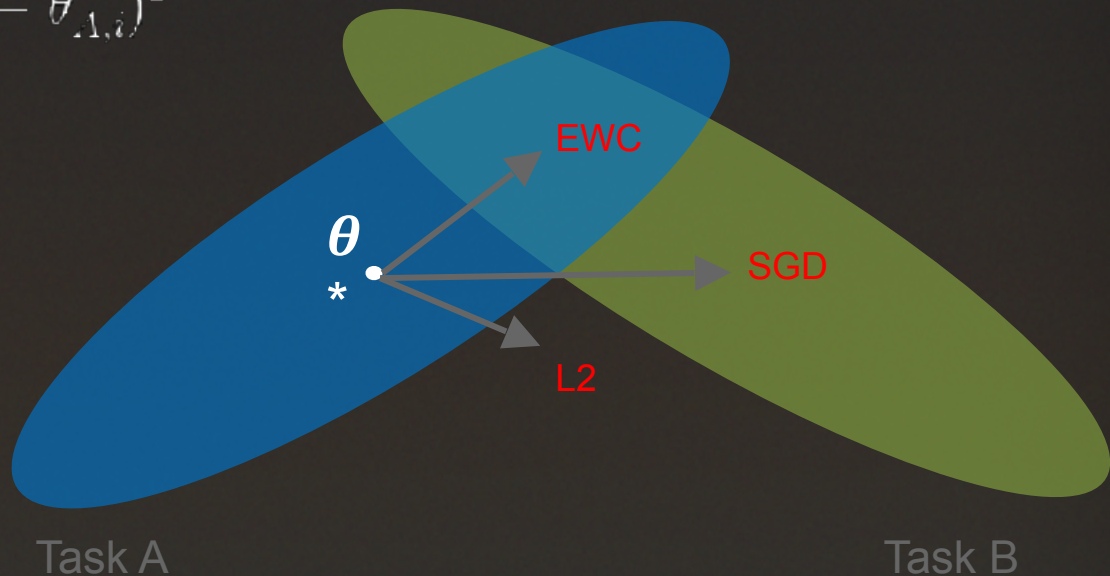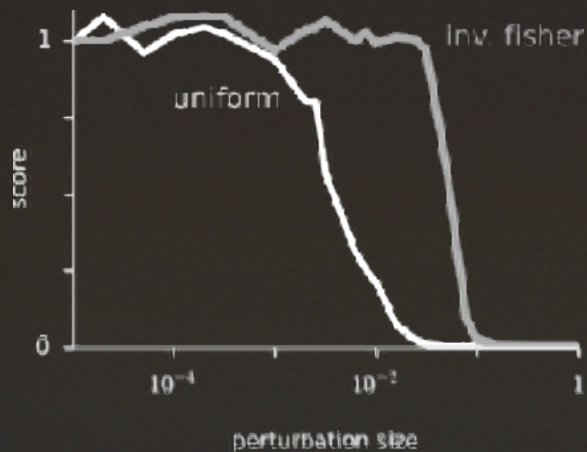
# Elastic Weight Consolidation



$$\mathcal{L}(\theta) = \mathcal{L}_B(\theta) + \sum_i \frac{\lambda}{2} F_i (\theta_i - \theta_{A,i}^*)^2$$

EWC

SGD

$\theta$
*

L2

Task A

Task B

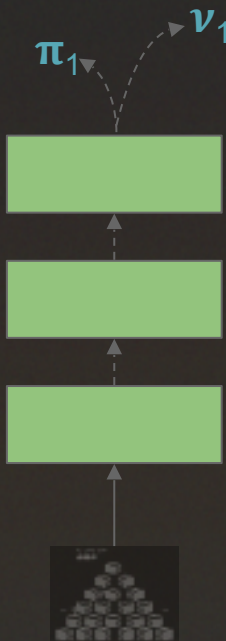Raia Hadsell 2017

DeepMind

# What if my tasks really don't get along?

# What if my tasks really don't get along?

**Progressive Nets**

- add columns for new tasks
- freeze params of learnt columns
- layer-wise neural connections

→ capacity for task-specific features
→ enables deep compositionality
→ precludes forgetting



$\pi_1$ $\nu_1$

# What if my tasks really don't get along?

**Progressive Nets**

- add columns for new tasks
- freeze params of learnt columns
- layer-wise neural connections

→     capacity for task-specific features
→     enables deep compositionality
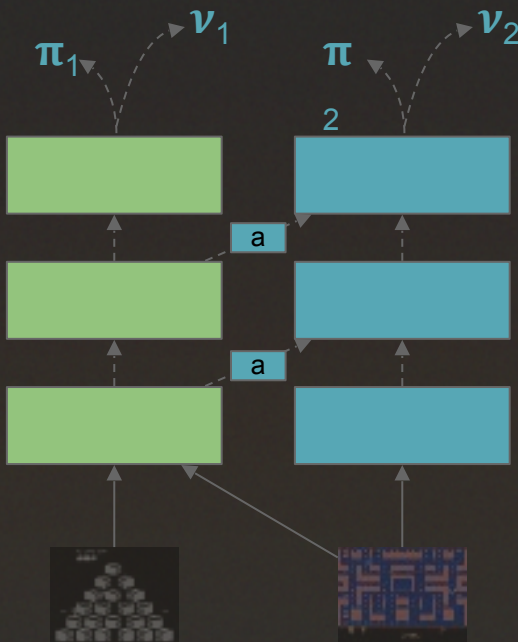→     precludes forgetting

DeepMind

# What if my tasks really don't get along?

**Progressive Nets**

- add columns for new tasks
- freeze params of learnt columns
- layer-wise neural connections

→ capacity for task-specific features
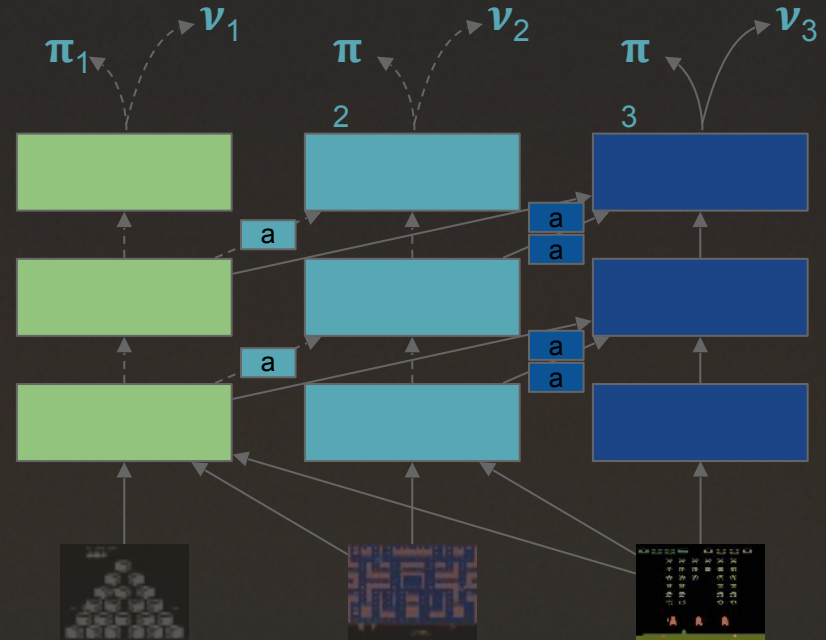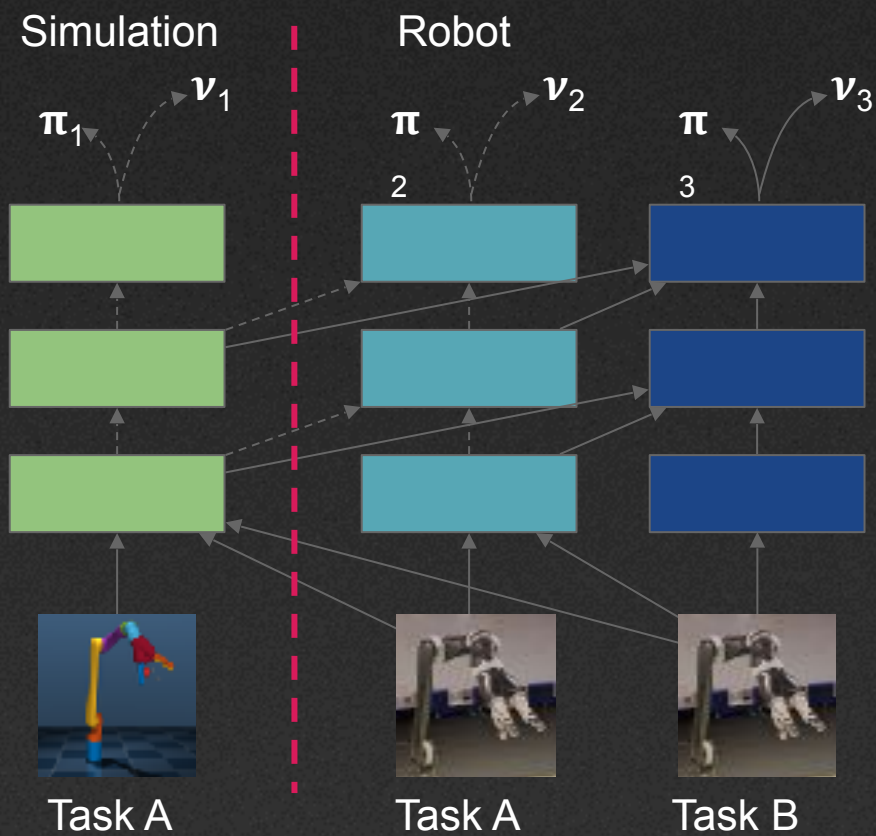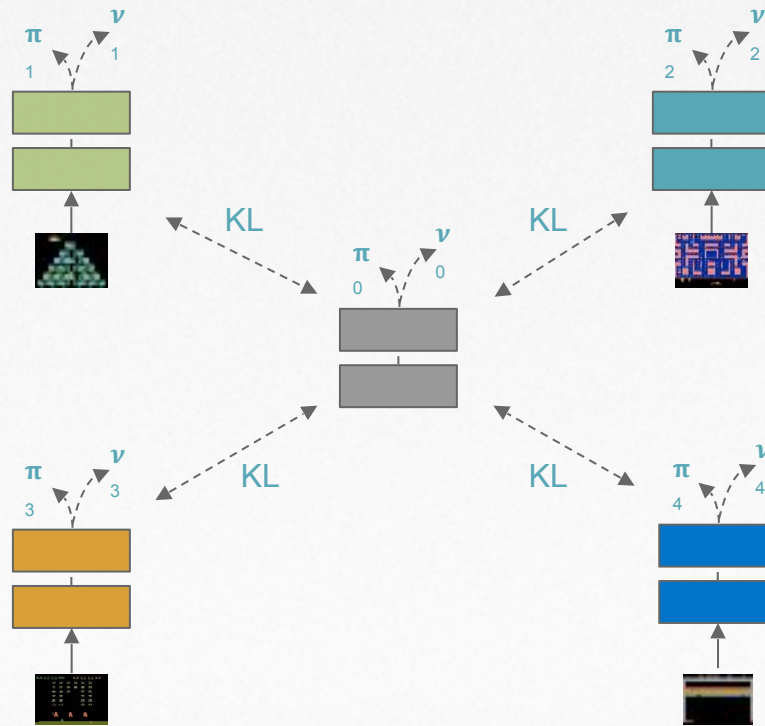→ enables deep compositionality
→ precludes forgetting



Andrei Rusu et al (2016), "Progressive Neural Networks"

Raia Hadsell 2017

Sim-to-Real

# What if my tasks really don't get along?

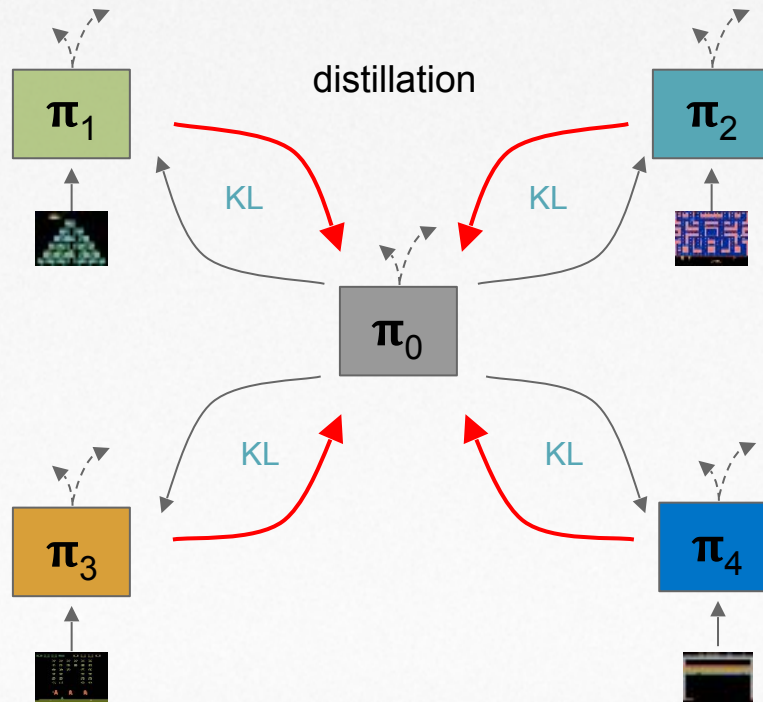DeepMind

# **Distral** (Distill and Transfer Learning)

- Task-specific networks plus shared network

- KL Divergence constraint

- Regularisation in policy space rather than parameter space

- Shared policy as a communication channel between tasks

Raia Hadsell

DeepMind

# **Distral** (Distill and Transfer Learning)

- Task-specific networks plus shared network

- Regularisation in policy space rather than parameter space

- Shared policy as a communication channel between tasks

→ *Distillation* of knowledge into shared model enables *transfer* to tasks



distillation

$\pi_1$

$\pi_2$

KL

KL

$\pi_0$

KL

KL

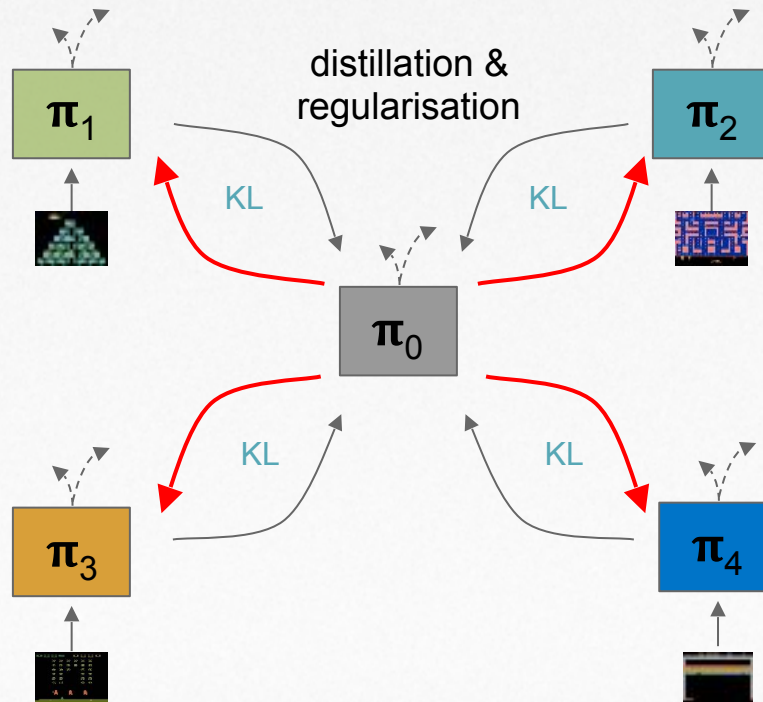$\pi_3$

$\pi_4$

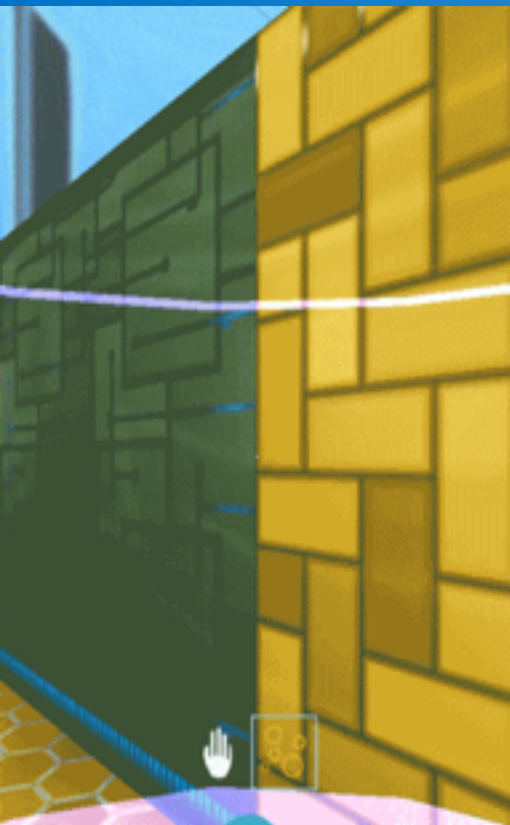Raia Hadsell

# **Distral** (Distill and Transfer Learning)

- Task-specific networks plus shared network

- Regularisation in policy space rather than parameter space

- Shared policy as a communication channel between tasks

→ *Distillation* of knowledge into shared model enables *transfer* to tasks

→ *Regularisation* of shared model gives stability and robustness

Yee Whye Teh et al (2017), "Distral: Robust Multitask Reinforcement Learning"

Raia Hadsell

Lab Mazes
&
Auxiliary Learning

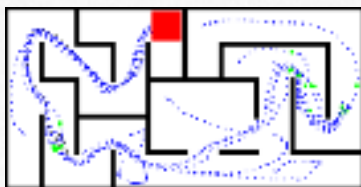Multiple Tasks
&
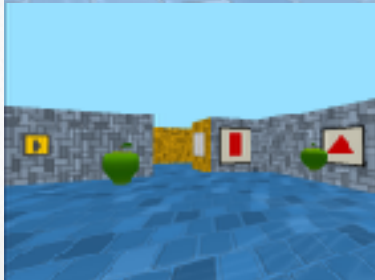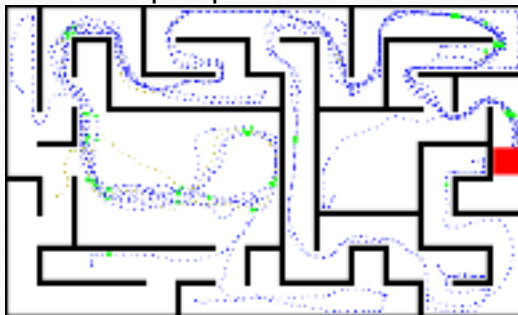Lifelong learning

StreetLearn

Parkour

# Navigation mazes



3600 steps/episode



10800 steps/episode



Game **episode:**

1. Random start
2. Find the goal (+10)
3. Teleport randomly
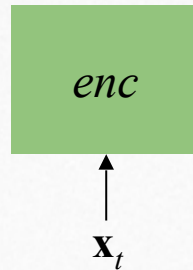4. Re-find the goal (+10)
5. Repeat (limited time)

Variants:

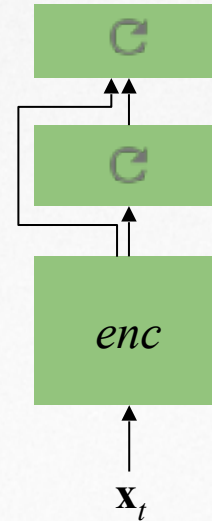Static maze, static goal
Static maze, random goal
Random maze

# Nav agent architecture

1. Convolutional encoder and RGB inputs

Piotr Mirowski, Razvan Pascanu et al (2017) "Learning to navigate in complex environments"
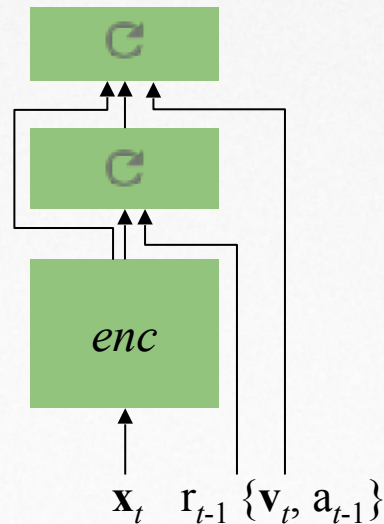
# Nav agent architecture

1. Convolutional encoder and RGB inputs

2. Single or stacked LSTM with skip connection

# Nav agent architecture

1. Convolutional encoder and RGB inputs

2. Stacked LSTM

3. Additional inputs (reward, action, and velocity)



$$\mathbf{x}_t \quad \mathrm{r}_{t-1} \; \{\mathbf{v}_t, \mathrm{a}_{t-1}\}$$

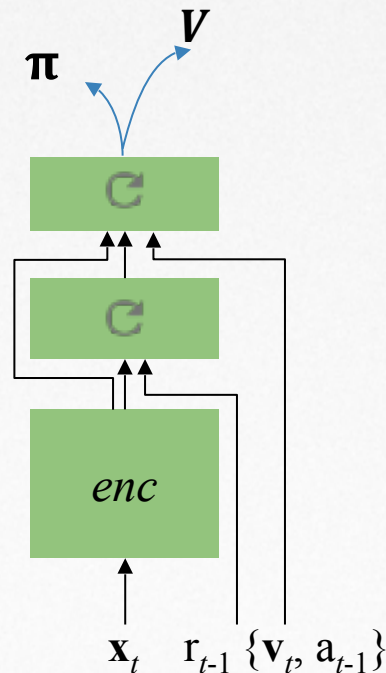# Nav agent architecture

1. Convolutional encoder and RGB inputs

2. Stacked LSTM

3. Additional inputs (reward, action, and velocity)

4. RL: Asynchronous advantage actor critic (A3C)

$V$

$\pi$

$enc$

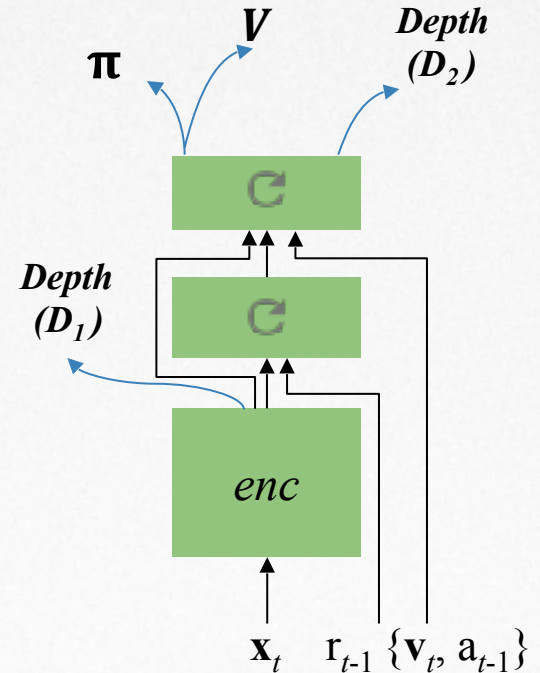$\mathbf{x}_t$    $\mathrm{r}_{t-1}$  $\{\mathbf{v}_t, \mathrm{a}_{t-1}\}$

# Nav agent architecture

1. Convolutional encoder and RGB inputs

2. Stacked LSTM

3. Additional inputs (reward, action, and velocity)

4. RL: Asynchronous advantage actor critic (A3C)

5. Aux task 1: Depth predictors



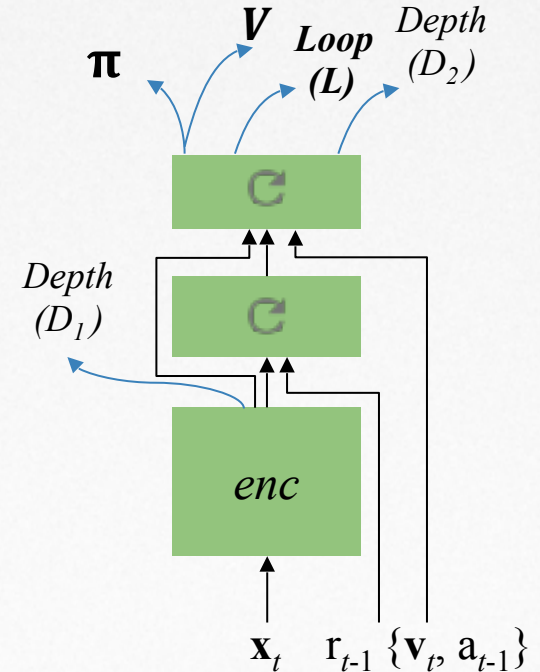Piotr Mirowski, Razvan Pascanu et al (2017) "Learning to navigate in complex environments"
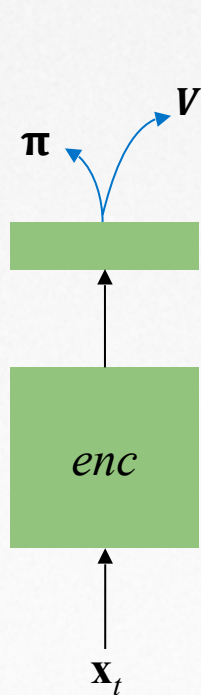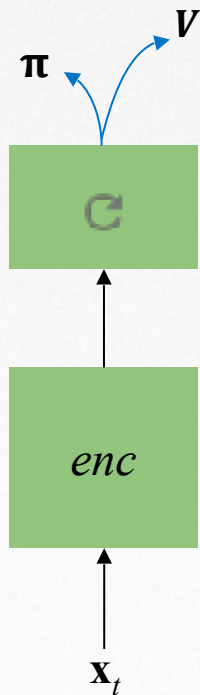
# Nav agent architecture

1. Convolutional encoder and RGB inputs

2. Stacked LSTM

3. Additional inputs (reward, action, and velocity)

4. RL: Asynchronous advantage actor critic (A3C)

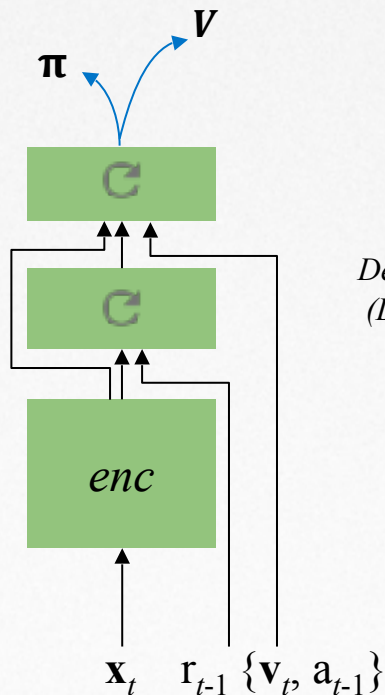5. Aux task 1: Depth predictor

6. Aux task 2: Loop closure predictor

# Variations in architecture



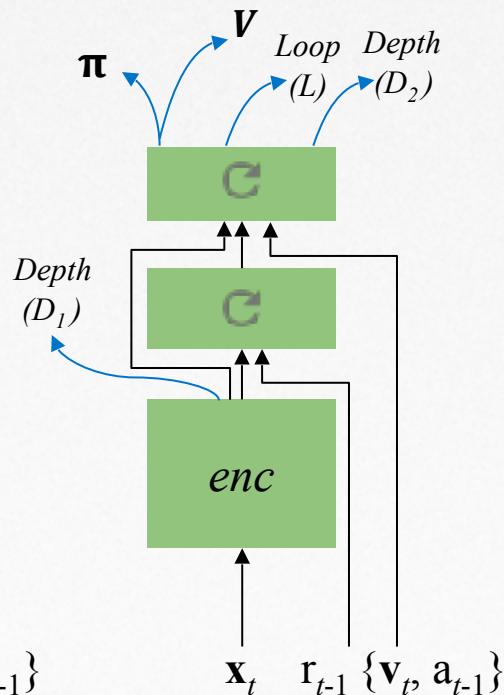a. FF A3C      b. LSTM A3C      c. Nav A3C      d. Nav A3C +$D_1D_2$L

# Results on large maze with static goal



+10          +1



FF A3C* (81)
LSTM A3C* (154)
Nav A3C* (157)
Nav A3C+L (181)
Nav A3C+D1 (192)
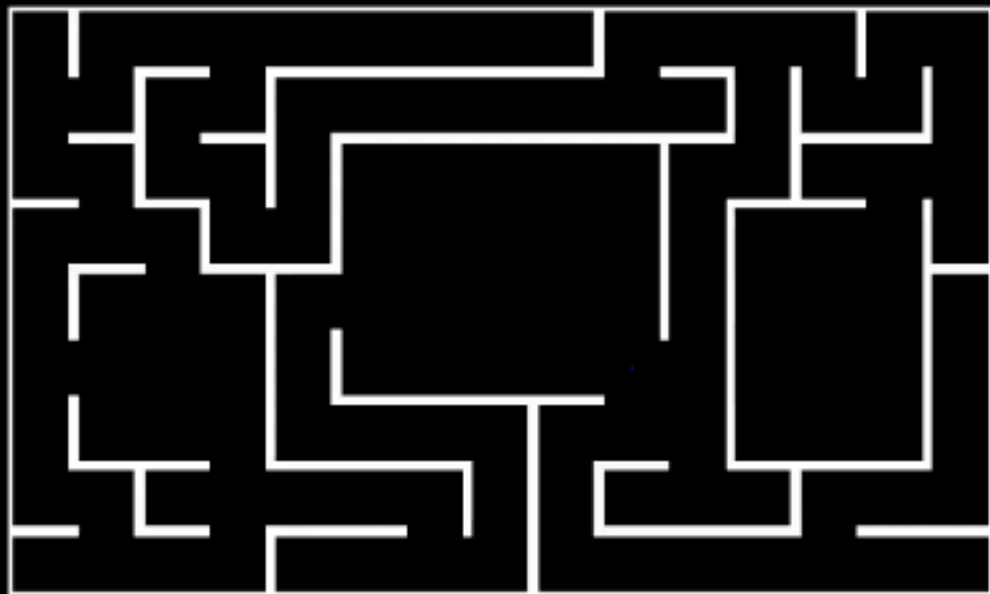Nav A3C+D2 (199)
Nav A3C+D1D2L (191)
Human expert (172)

Lab Mazes
&
Auxiliary Learning

Multiple Tasks
&
Lifelong learning

StreetLearn
&
Real woRld RL
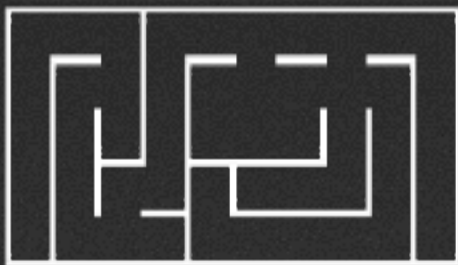
Parkour

# Navigation mazes in the real world?



observation

observation

structure

structure

DeepMind

# StreetView as an RL environment: StreetLearn


observation


observation



- RGB image cropped from panorama (84x84)
- Goal location

Actions: move to next node, rotate view 20° or 60°


structure


structure

# StreetView as an RL environment: StreetLearn



left or right?

DeepMind

# StreetView as an RL environment: StreetLearn



Looks like a road, but it's a park entrance

DeepMind

# StreetView as an RL environment: StreetLearn



west side highway

DeepMind

# StreetView as an RL environment: StreetLearn



curved roads and tunnels

DeepMind

# StreetView as an RL environment: StreetLearn



really, tunnels!

# StreetLearn: The Courier Task



1. Spawn randomly and navigate to a random target location.

2. Start receiving reward when close to target (within 400m).

3. If target is reached (100m), navigate to a new random target.

Raia Hadsell

# Agent architecture



Absolute heading prediction

Policy (π, V)

Local graph neighbour prediction

Global pathway

LSTM

LSTM

LSTM

Relative pathway

CNN

target

image

$r_{t-1}$, $a_{t-1}$

DeepMind

Raia Hadsell

# Agent architecture

Absolute heading prediction

Policy (π, V)

Local graph neighbour prediction

Global pathway

LSTM

LSTM

LSTM

Relative pathway

CNN

target


image

$r_{t-1}$, $a_{t-1}$

DeepMind

Raia Hadsell

# Agent architecture

Global pathway

Relative pathway

Absolute heading prediction

Policy (π, V)

Local graph neighbour prediction

LSTM

LSTM

LSTM

CNN

target

image

$r_{t-1}$, $a_{t-1}$

DeepMind

Raia Hadsell

# Agent architecture



Absolute heading prediction

Policy (π, V)

Local graph neighbour prediction

Global pathway

LSTM

LSTM

LSTM

Relative pathway

CNN

target

image

$r_{t-1}$, $a_{t-1}$

Raia Hadsell

Lab Mazes
&
Auxiliary Learning

Multiple Tasks
&
Lifelong learning

StreetLearn
&
Real woRld RL

Parkour
&
Continuous control

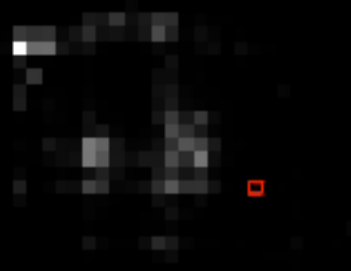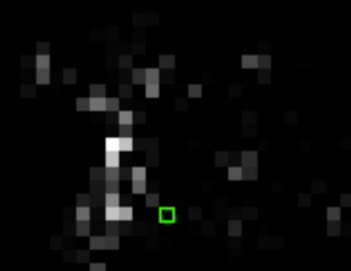# Proprioceptive and exteroceptive observations

**Proprioceptive** --
"near the body":

- Joint angles & velocities
- Touch sensors
- Positions and velocities of limbs in body coordinate frame
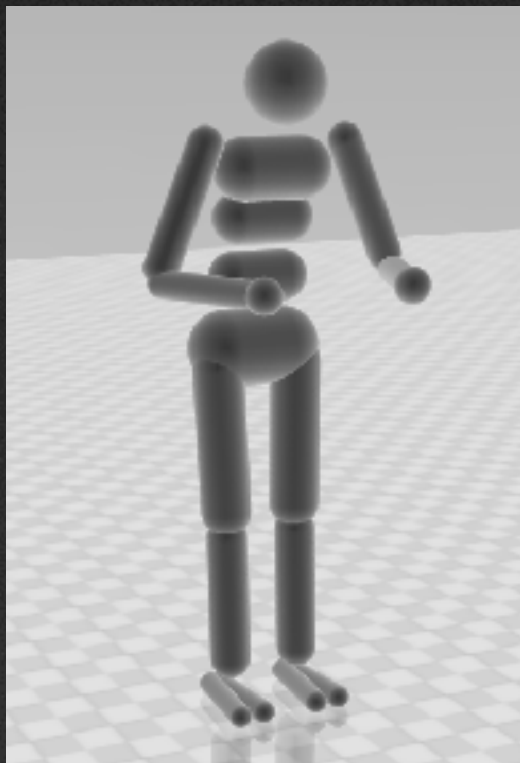
# Proprioceptive and exteroceptive observations

**Proprioceptive** --
"near the body":

- Joint angles & velocities
- Touch sensors
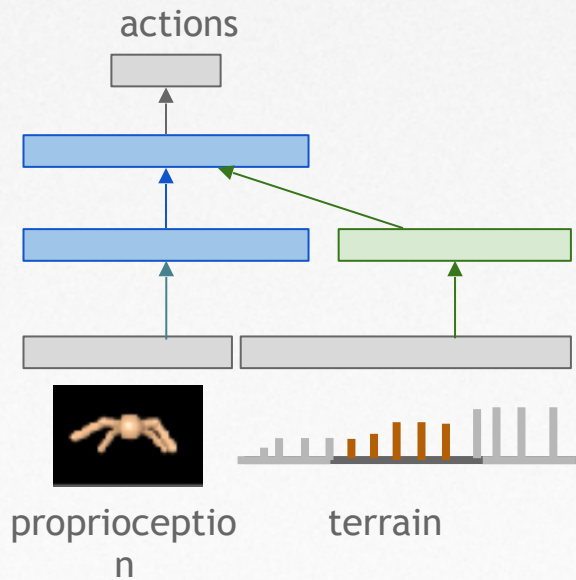- Positions and velocities of limbs in body coordinate frame



**Exteroceptive** --
"away from the body":

- Position / velocity in global coordinate frame
- Task-related (e.g. goal position)
- Vision

DeepMind

# Rich environments for skill discovery: setup



actions

proprioception

terrain

## Training

- Proximal policy optimization
  [Schulman et al.]
- Batched policy gradient
- Trust region
  ("gradient-based TRPO")
- High-performance
  implementation:
  - Distributed (multiple workers)
  - Synchronous gradient updates

Nicolas Heess, et al. 2016: "Learning and transfer of modulated locomotor controllers"    Raia Hadsell

# Single uniform reward, based on forward progress

Nicolas Heess, et al. 2017:
"Emergence of Locomotion Behaviours in Rich Environments"

DeepMind

Raia Hadsell

# Humanoid: learned behaviors



- 27 DoFs
- 21 actuators

Nicolas Heess, et al. 2017:
"Emergence of Locomotion Behaviours in Rich Environments"

Deep RL — **Raia Hadsell**

- Can deep RL agents learn multiple tasks?

- Can deep RL agents learn efficiently?

- Can deep RL agents learn from real data?

- Can deep RL agents learn continuous control?

DeepMind

*Overcoming catastrophic forgetting in NNs,* 2016
James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A. Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, Demis Hassabis, Claudia Clopath, Dharshan Kumaran, Raia Hadsell

*Progressive Neural Networks,* 2016
Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray Kavukcuoglu, Razvan Pascanu, Raia Hadsell

*Distral: Robust Multitask RL,* 2017
Yee Whye Teh, Victor Bapst, Wojciech Marian Czarnecki, John Quan, James Kirkpatrick, Raia Hadsell, Nicolas Heess, Razvan Pascanu

*Learning to navigate in complex environments,* 2017
Piotr Mirowski*, Razvan Pascanu*, Fabio Viola, Hubert Soyer, Andrew J. Ballard, Andrea Banino, Misha Denil, Ross Goroshin, Laurent Sifre, Koray Kavukcuoglu, Dharshan Kumaran, Raia Hadsell

*Learning and transfer of modulated locomotor controllers,* 2016
Nicolas Heess, Greg Wayne, Yuval Tassa, Timothy Lillicrap, Martin Riedmiller, David Silver

*Emergence of Locomotion Behaviours in Rich Environments,* 2017
Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, S. M. Ali Eslami, Martin Riedmiller, David Silver

# Thank you!

DeepMind