

REINFORCEMENT LEARNING

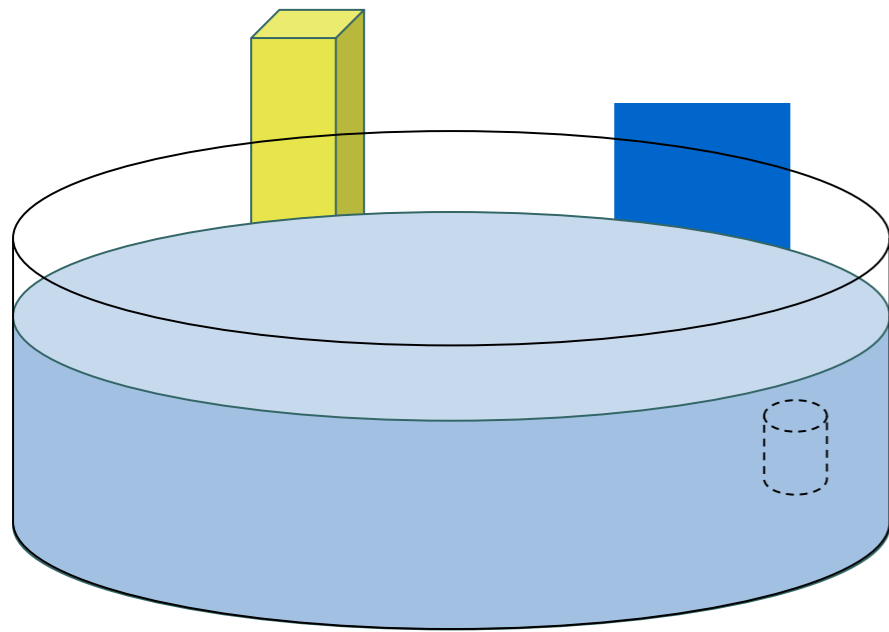
Eleni Vasilaki, University of Sheffield, UK



REINFORCEMENT LEARNING EXAMPLES



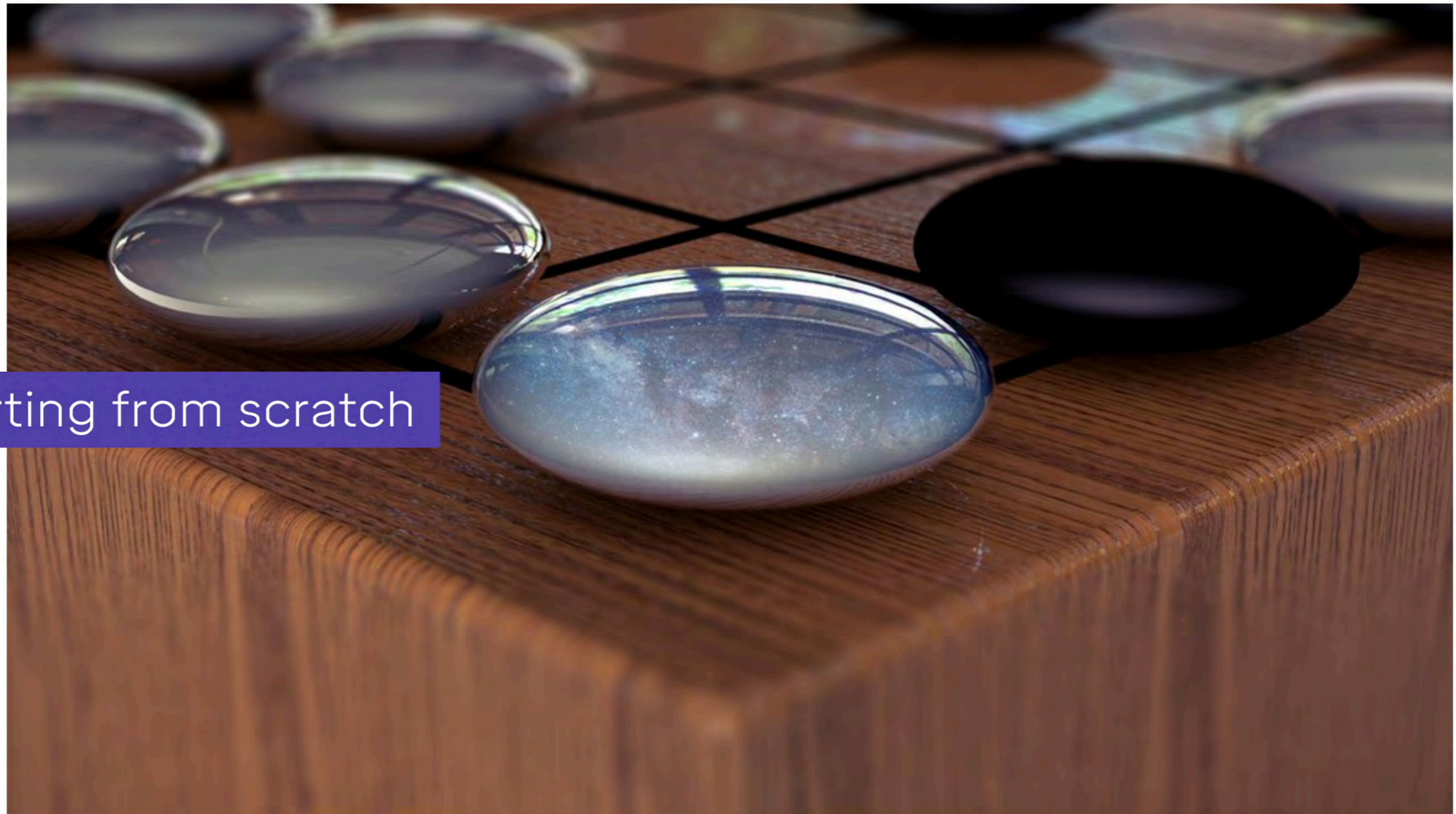
- Reward/Punishment
- Exploration/Exploitation



Morris, 1981



GOOGLE DEEPMIND ALPHAGO ZERO



▶ Starting from scratch

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

Silver et al, 2017 <https://www.nature.com/articles/nature24270>

EPICURUS & REINFORCEMENT LEARNING



<http://www.defenseofreason.com>

Vasilaki, 2017 <https://arxiv.org/abs/1710.04582>

EPICURUS & REINFORCEMENT LEARNING

“We say that pleasure is the beginning and the end of a happy life”

“For continual drinking and partying [...] do not produce a pleasant life, but sober reasoning which both examines the basis for every choice and avoidance [...]”

— Epicurus’ Letter to Menoeceus

Diogenes Laertius, Lives of Eminent Philosophers

REWARD FUNCTION

$$R = r_t + r_{t+1} + r_{t+2} + \dots + r_{t+N}$$

- Reward can be positive and negative.
- An action can bring a small immediate reward and a large future punishment.

Vasilaki, 2018 <http://bit.ly/RL-happiness>

Sutton & Barto, 2018 <http://incompleteideas.net/book/the-book.html>

DISCOUNT FACTOR

$$R_t = r_t + r_{t+1} + r_{t+2} + \dots + r_{t+N}$$

$$R_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

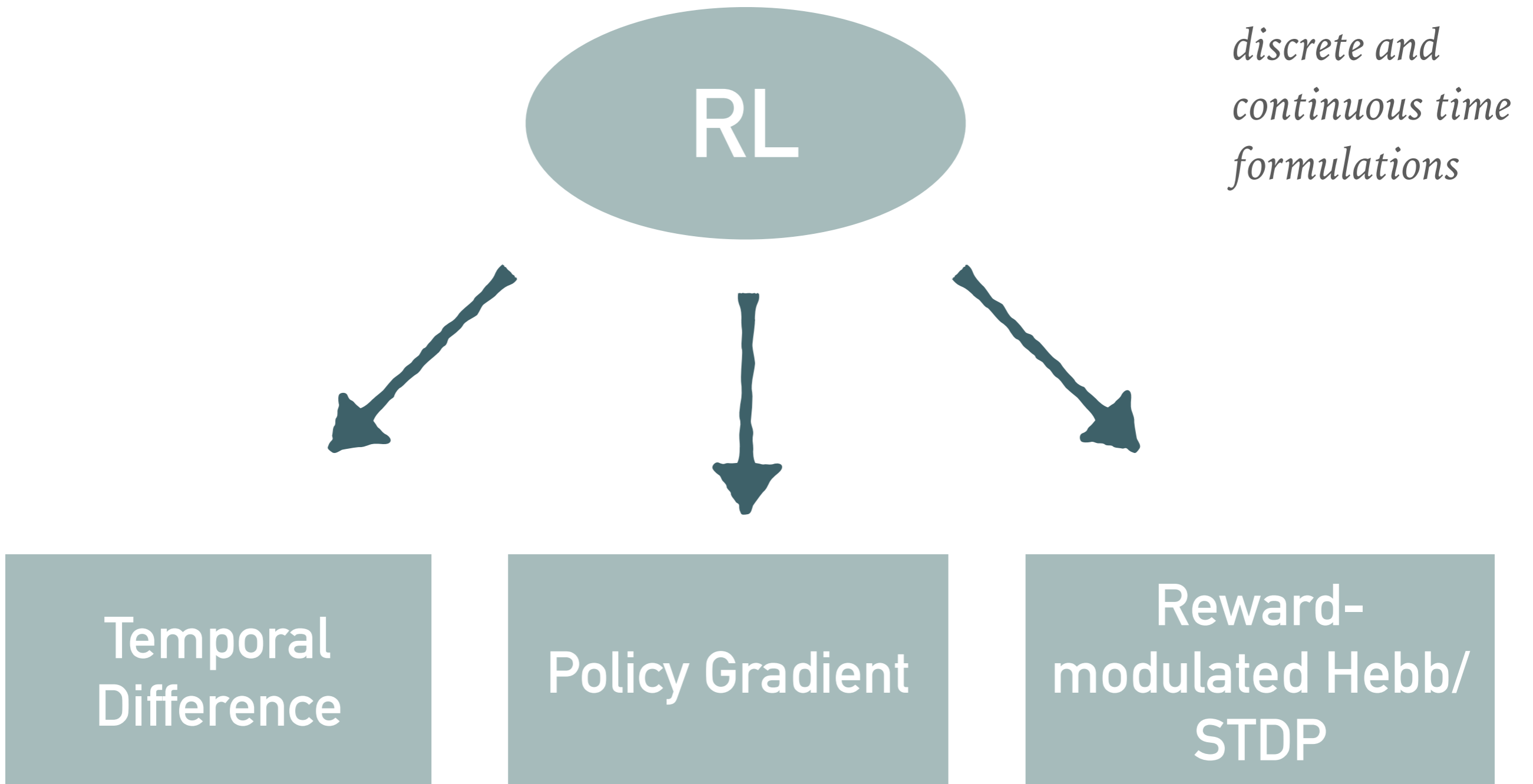
$$0 \leq \gamma < 1$$

This is why I am impatient.

Vasilaki, 2018 <http://bit.ly/RL-happiness>

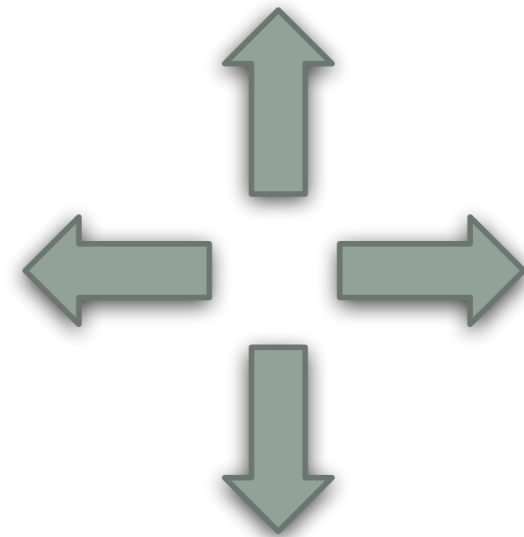
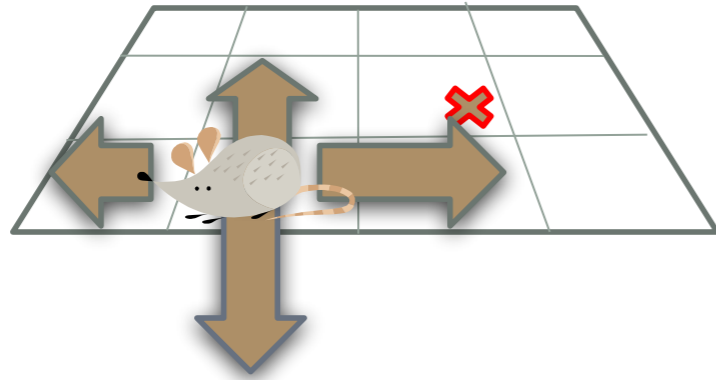
Sutton & Barto, 2018 <http://incompleteideas.net/book/the-book.html>

REINFORCEMENT LEARNING CATEGORIES



TEMPORAL DIFFERENCE LEARNING

States x Actions: $3 \times 4 \times 4 = 48$



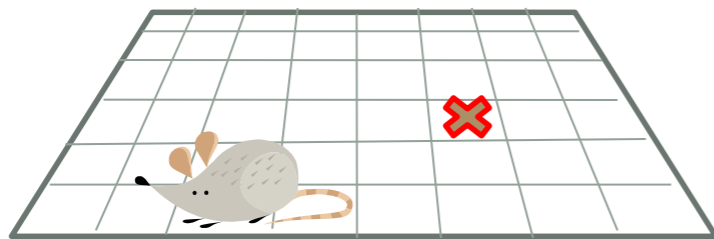
Maximize expected return

Q (state, action)

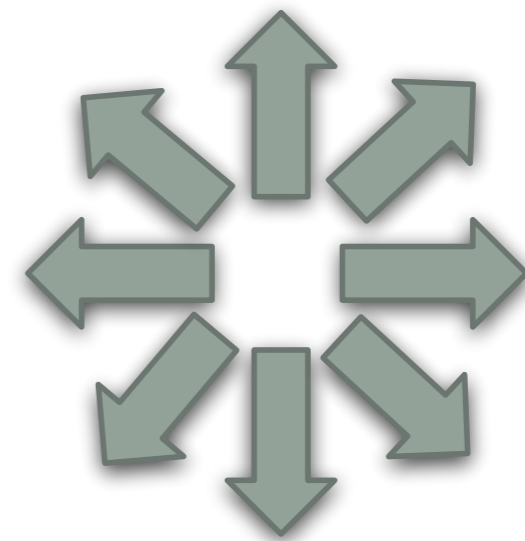
We do not know the Q values

TEMPORAL DIFFERENCE LEARNING

States x Actions: $6 \times 8 \times 8 = 384$



Q (state, action)



Learning process slows down

POLICY (π)

- How do I chose an action?
 - I know all the Q values
 - Exploit (max Q value, a.k.a Greedy policy)
 - I only have an estimate, and this may be wrong!
 - Explore (e.g. ϵ -Greedy, soft-max)

TEMPORAL DIFFERENCE LEARNING

$$R = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots$$

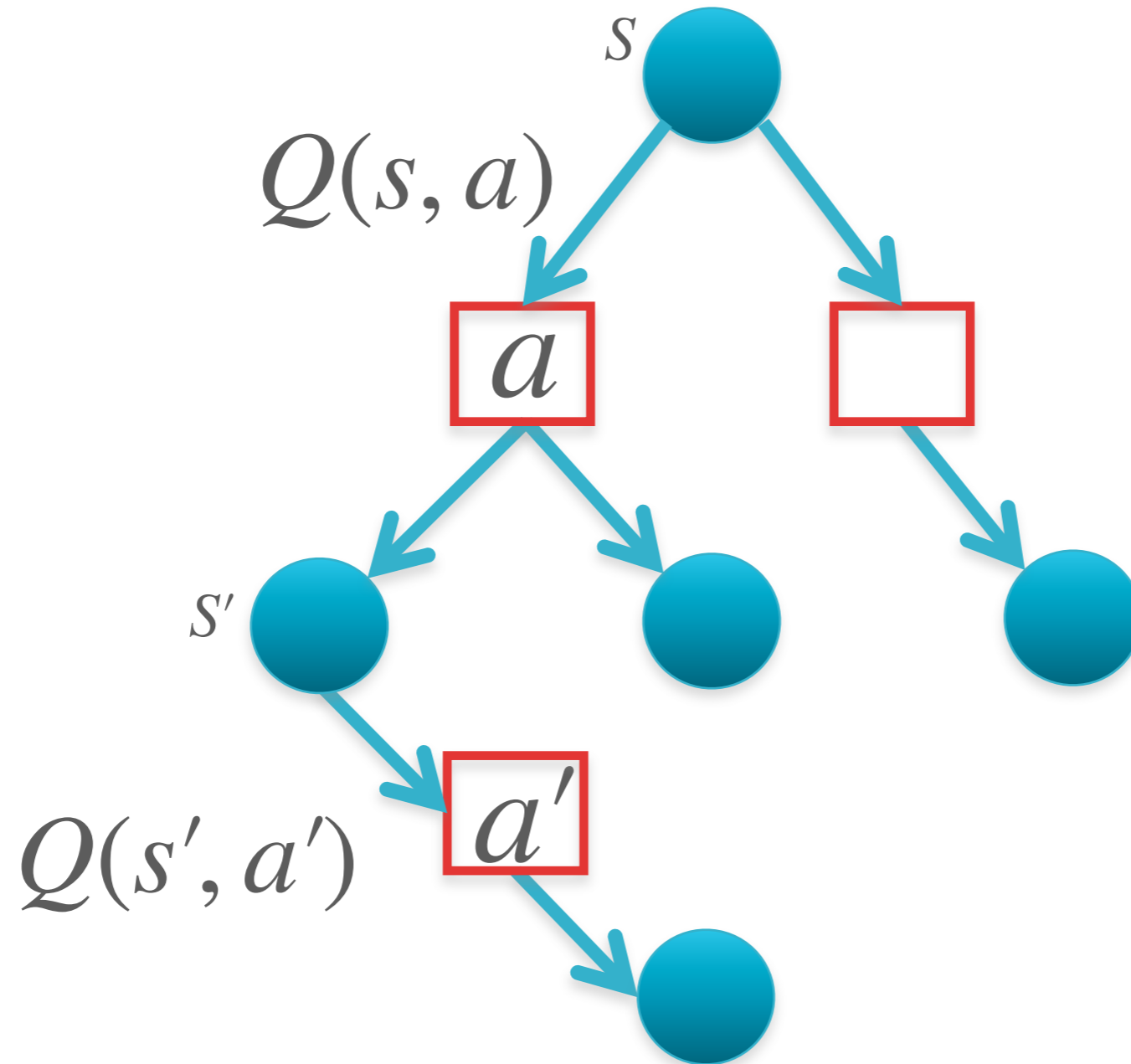
$$R = r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots)$$

$$Q^\pi(s, a) = E_\pi\{R_t \mid s_t = s, a_t = a\}$$

Bellman Equations,

Markovian property

SARSA $\Delta Q(s, a) = \eta \left((r + \gamma Q(s', a')) - Q(s, a) \right)$



SARSA AND HAPPINESS

$$\Delta Q(s, a) = \eta \left(\underbrace{(r + \gamma Q(s', a'))}_{\text{What I "actually" get}} - \underbrace{Q(s, a)}_{\text{Anticipated reward}} \right)$$

What I "actually" get

Anticipated reward

A positive reward may feel like punishment

A negative reward may feel like reward

DEEP REINFORCEMENT LEARNING



$$L(\theta_i) = \left(\underbrace{(r + \gamma Q(s', a', \theta_{i-1}))}_{\text{Target}} - Q(s, a, \theta_i) \right)^2$$

Target

DEEP REINFORCEMENT LEARNING



<https://deepmind.com/research/dqn/>

Mnih et al, Nature 2015
Mnih et al, NIPS 2013

REINFORCEMENT LEARNING CATEGORIES

RL

```
graph TD; RL([RL]) --> TD[Temporal Difference]; RL --> PG[Policy Gradient]; RL --> RMH[Reward-modulated Hebb/STDP];
```

Temporal
Difference

Policy Gradient

Reward-
modulated Hebb/
STDP

POLICY GRADIENT METHODS

$$\langle R \rangle_{x,y} = \sum_{x,y} R(x,y)P(y|x)P(x) \qquad \langle \Delta w \rangle_{x,y} = \alpha \frac{\partial \langle R \rangle_{x,y}}{\partial w}$$

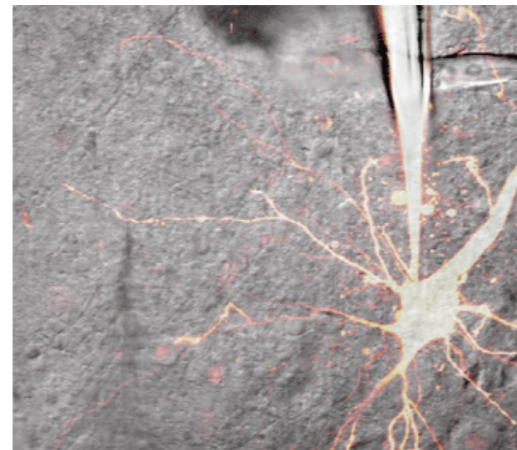
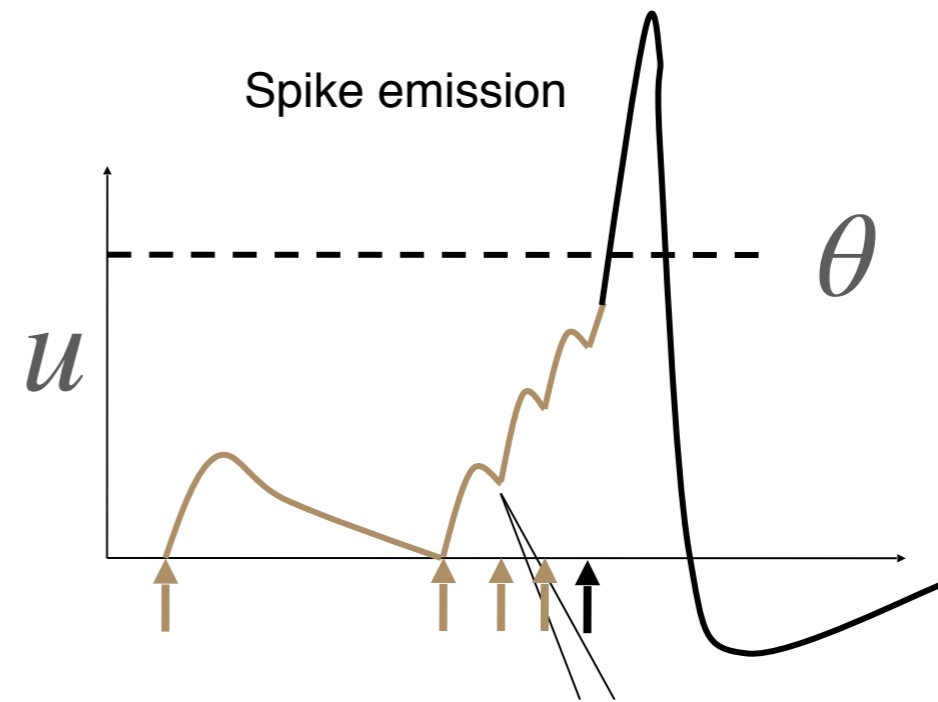
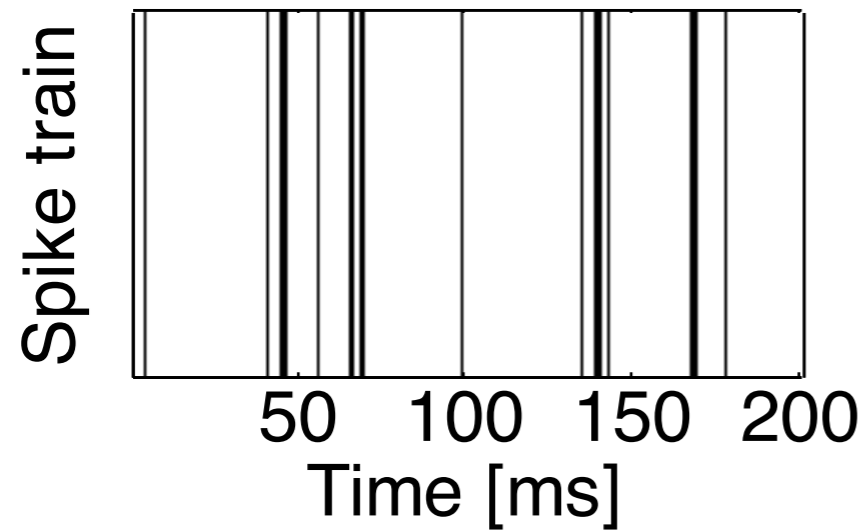
$$\frac{dw_{ij}}{dt} = \underbrace{\alpha}_{\text{Learning rate}} \underbrace{R(t)}_{\text{Reward}} \underbrace{e_{ij}(t)}_{\text{Eligibility trace}}$$

$$\frac{de_{ij}}{dt} = -\frac{e_{ij}}{\underbrace{\tau_e}_{\text{Time constant}}} + \left(\underbrace{Y_i(t)}_{\text{Spike}} - \underbrace{\rho(t)}_{\text{Instantaneous probability of firing}} \right) \sum_{t_j^f} \underbrace{\epsilon(t - t_j^f)}_{\text{Input at synapse ij}}$$

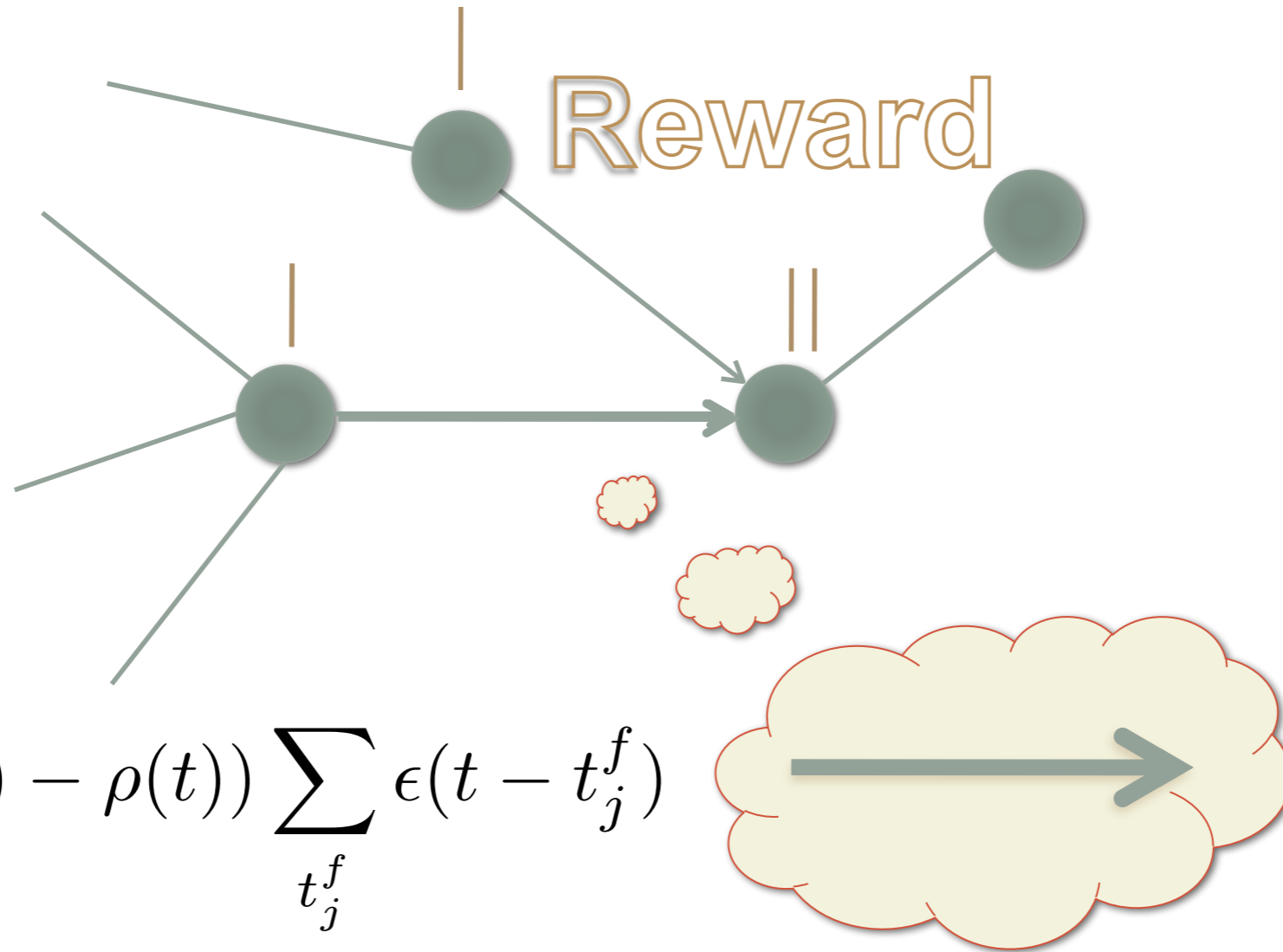
Williams 1992
 Xie and Seung, 2004
 Pfister et al. 2006
 Florian 2007
 Vasilaki et al, 2009

SPIKING NEURONS

- Spikes
- Threshold

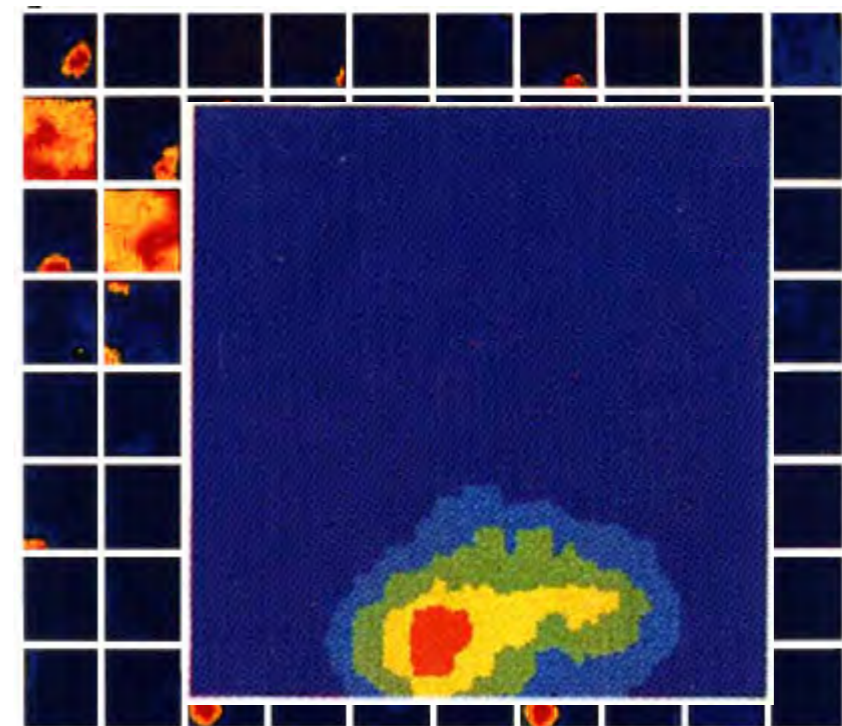
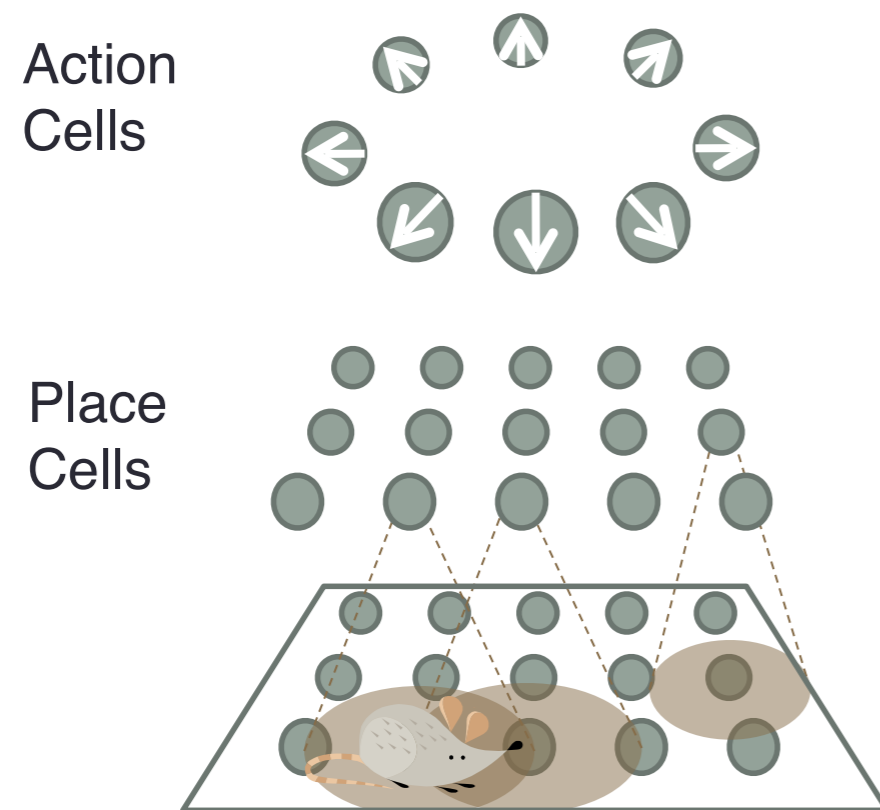


ELIGIBILITY TRACE



SPIKE BASED REINFORCEMENT LEARNING

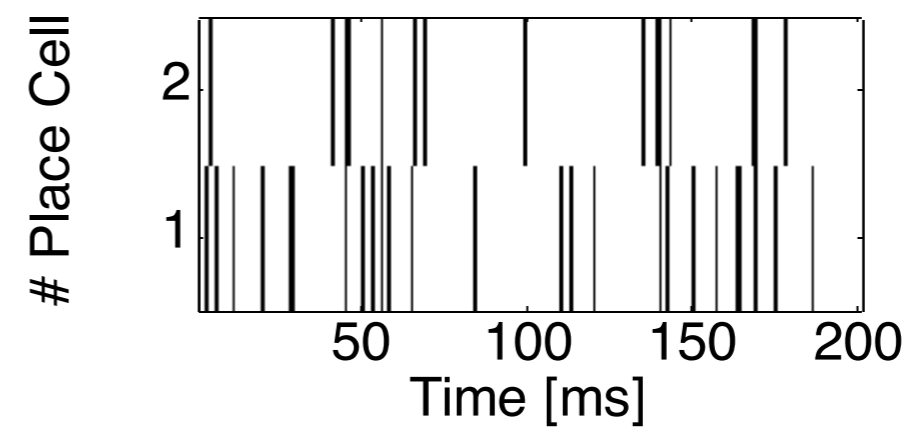
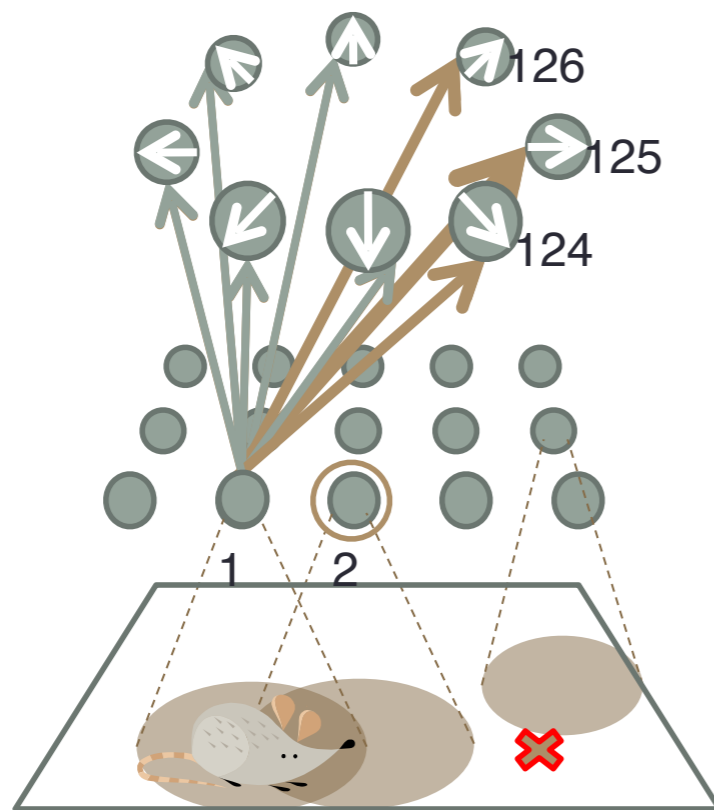
Action encoded by the centre of mass of the activity.



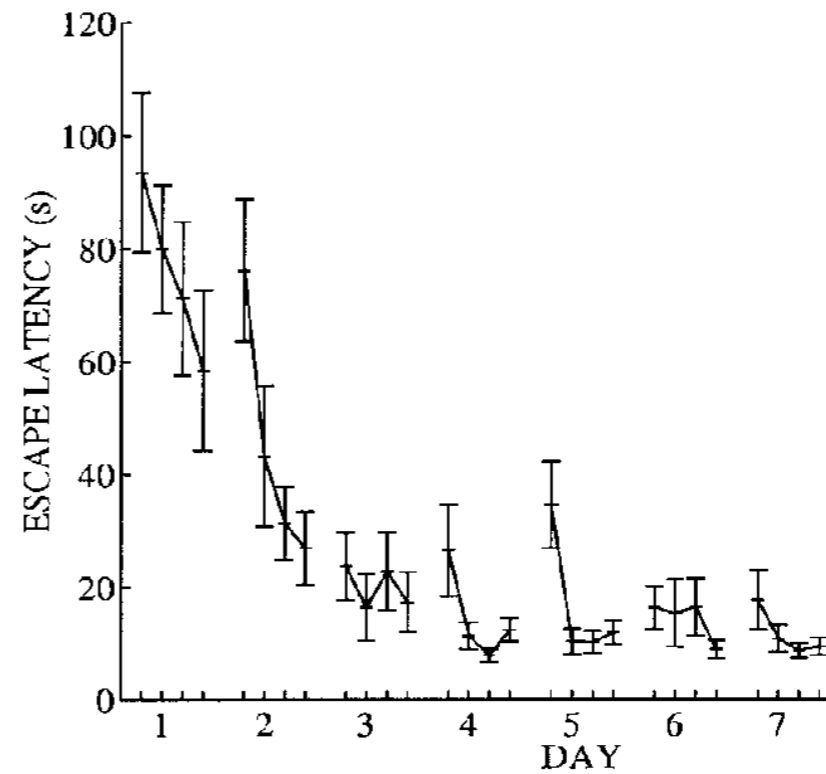
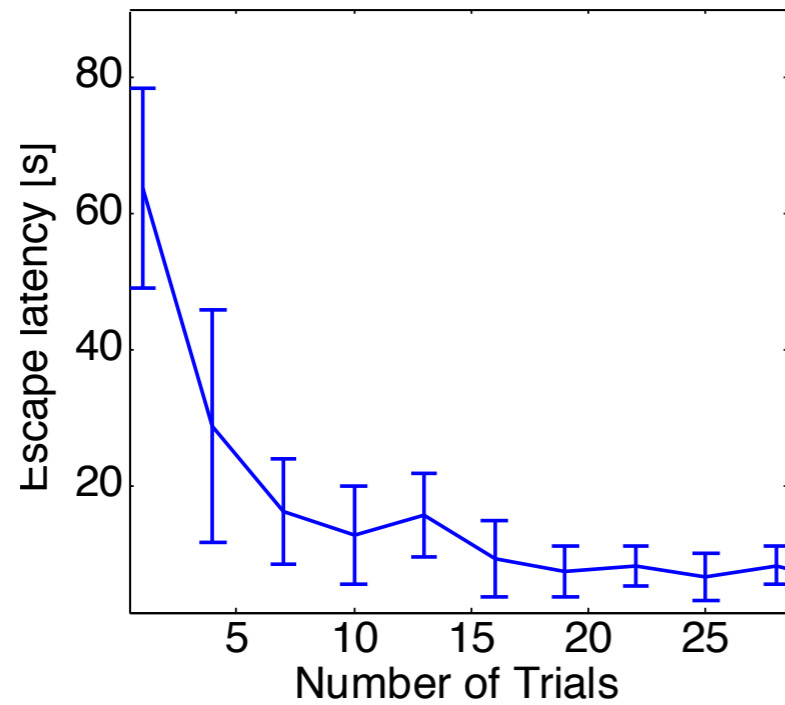
Nakazawa et al 2004 Nature Reviews | Neuroscience

O'Keefe & Dostrovsky , 1971

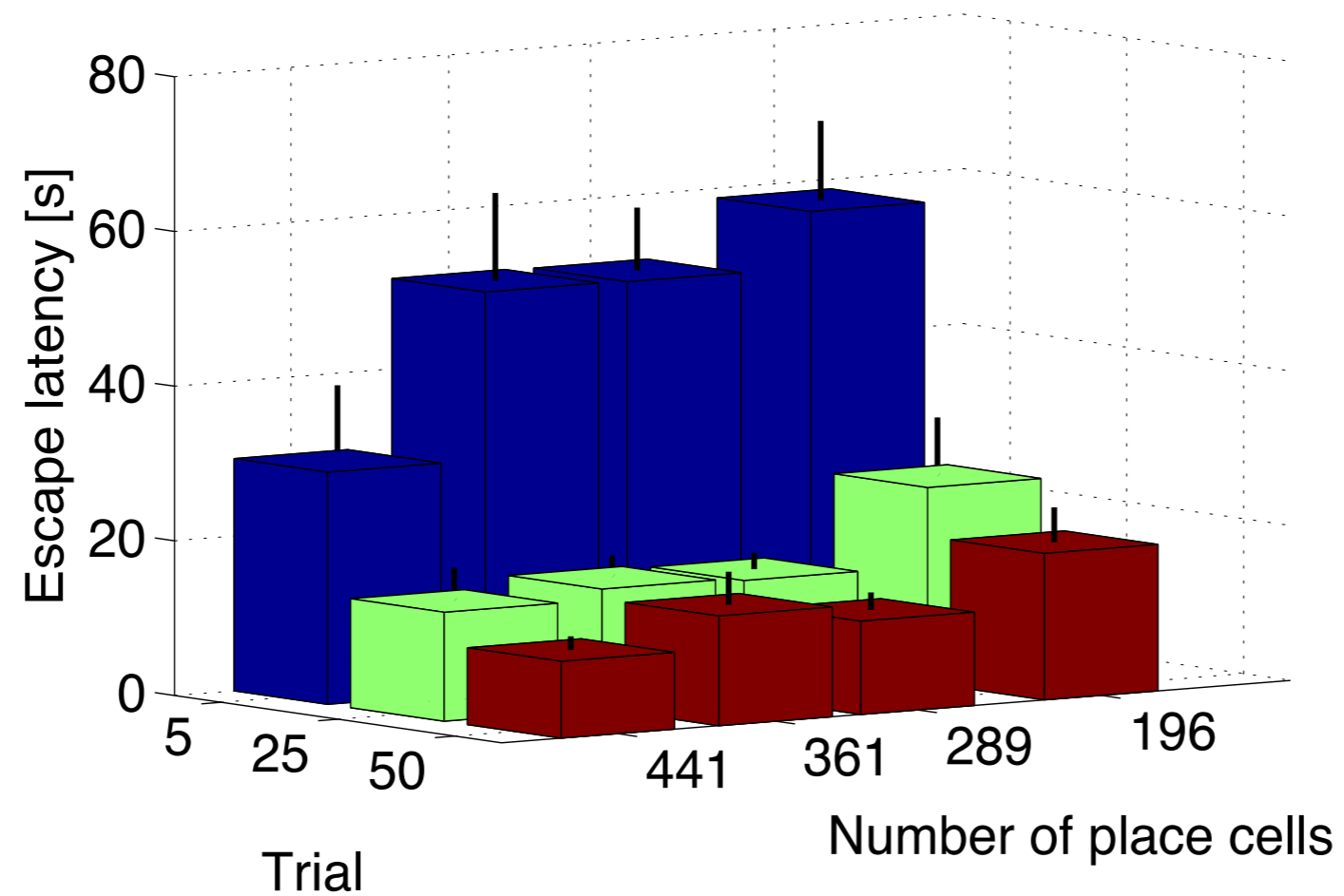
SPIKE BASED REINFORCEMENT LEARNING



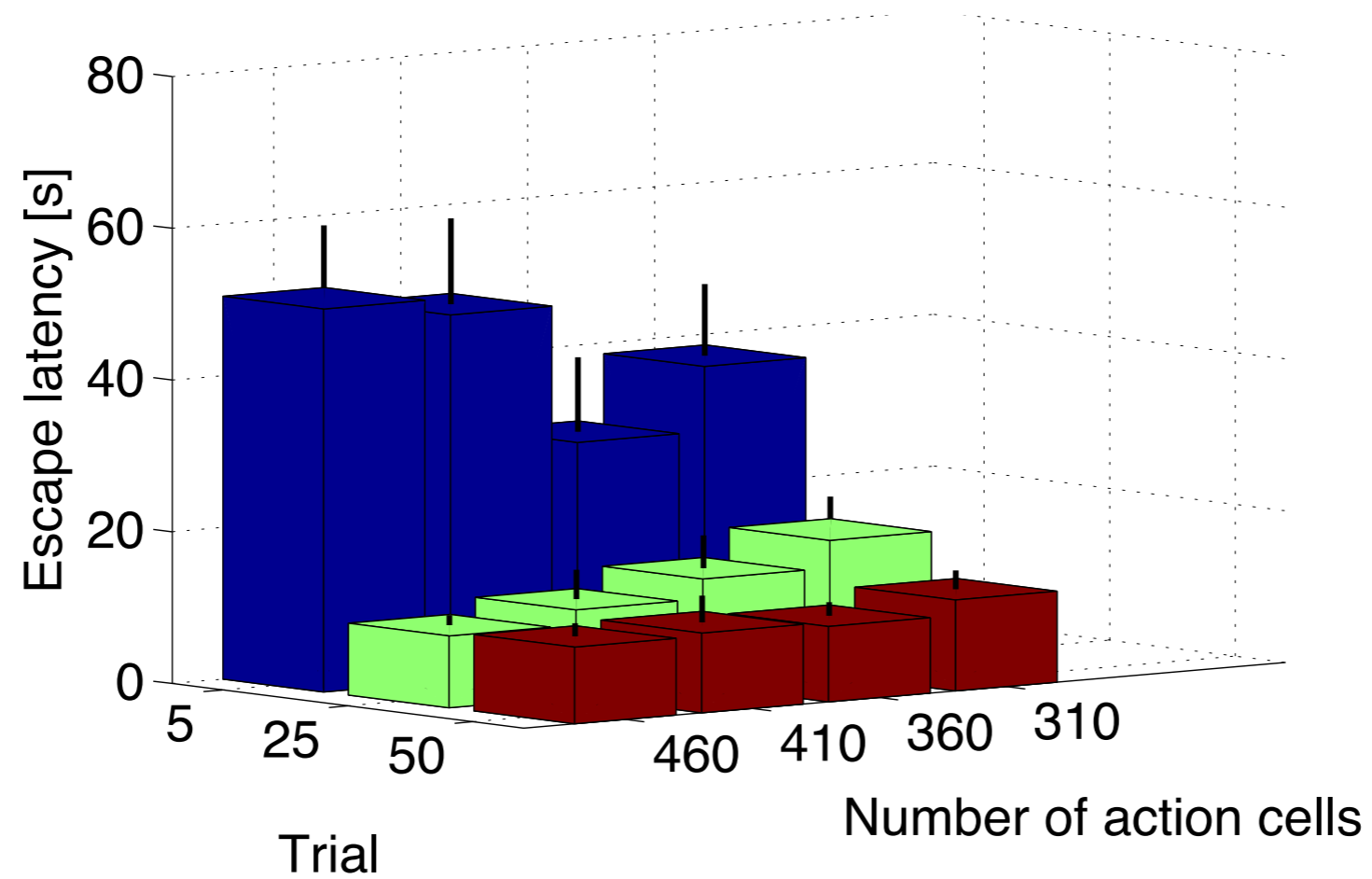
SIMULATION VS EXPERIMENT



SCALING PROPERTIES



SCALING PROPERTIES



REINFORCEMENT LEARNING CATEGORIES

RL

```
graph TD; RL([RL]) --> TD[Temporal Difference]; RL --> PG[Policy Gradient]; RL --> RMH[Reward-modulated Hebb/STDP];
```

Temporal
Difference

Policy Gradient

Reward-
modulated Hebb/
STDP

REWARD MODULATED HEBB

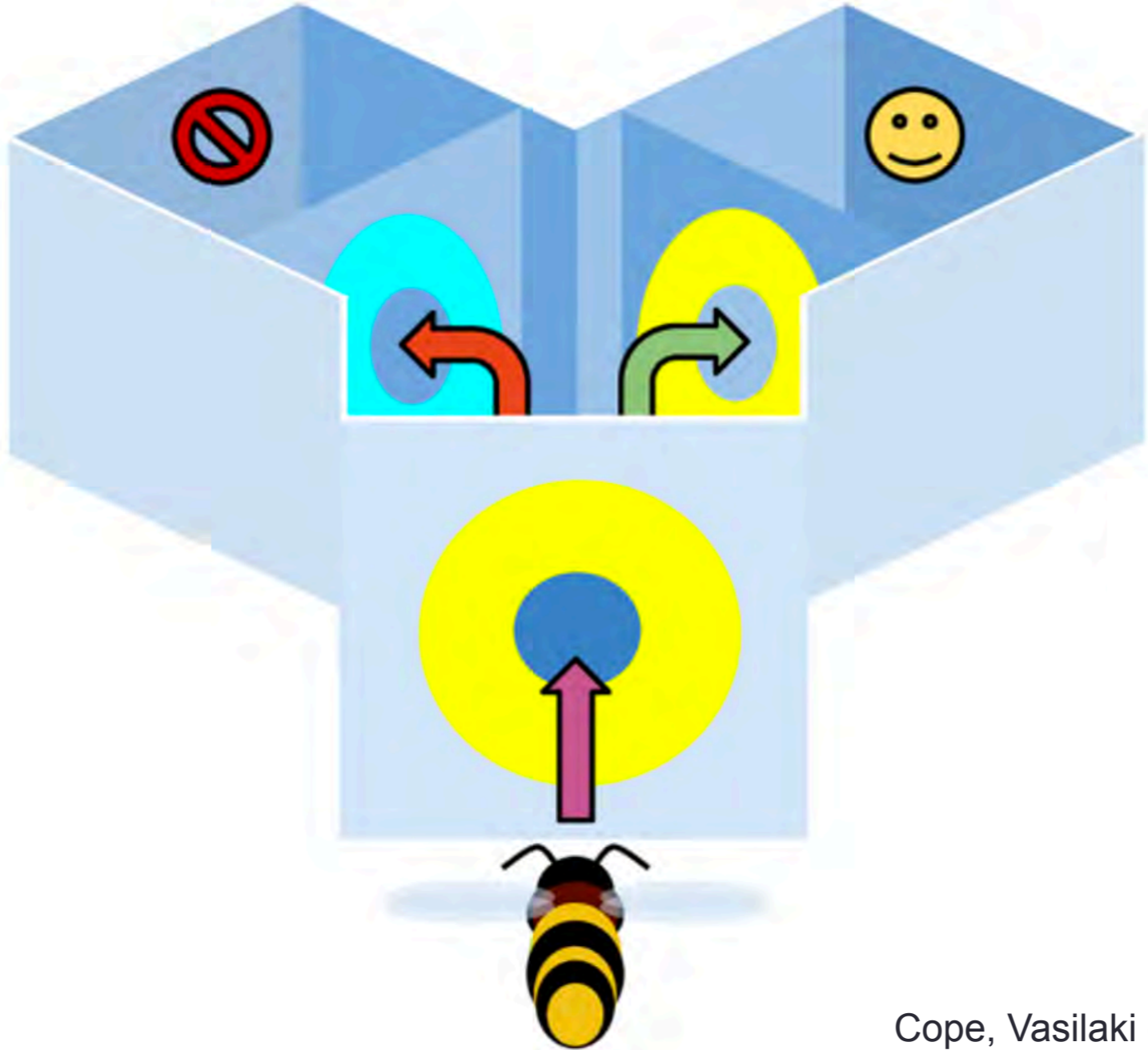


“Neurons that fire together wire together”

Correlations

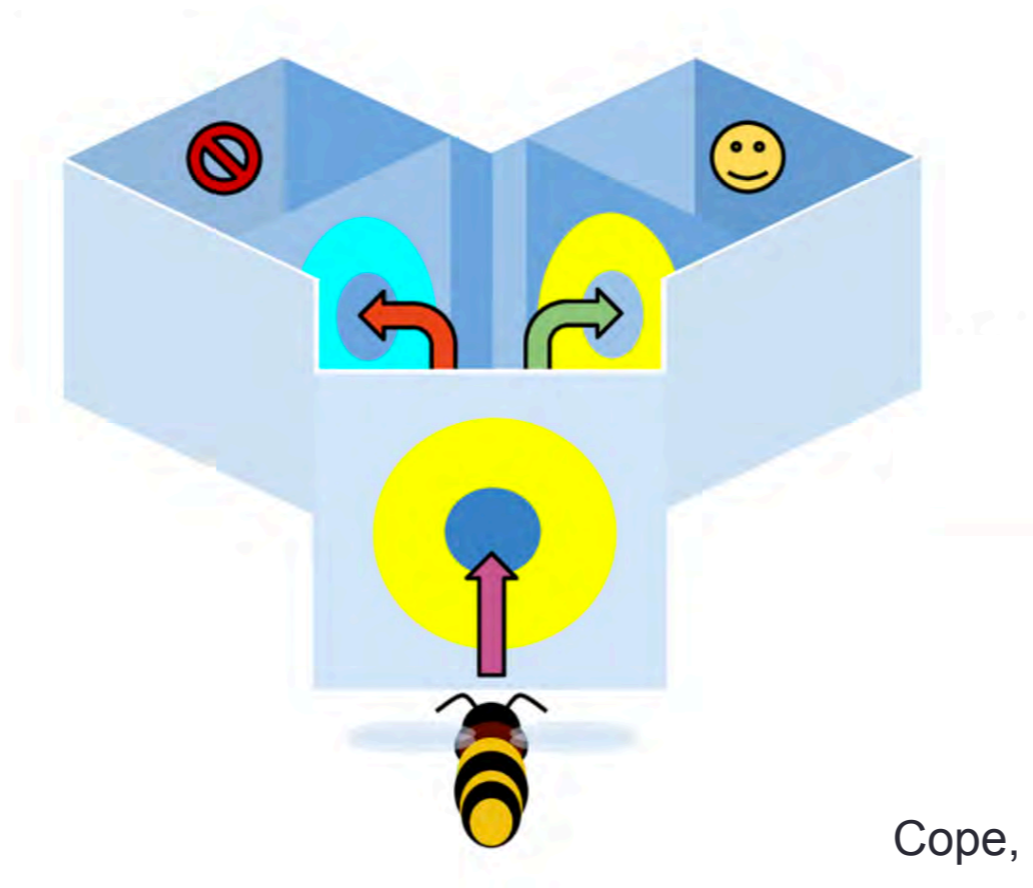
$$\Delta w = g(\text{reward}) \times f(\text{presynaptic activity}, \text{postsynaptic activity})$$

TO BEE OR NOT TO BEE



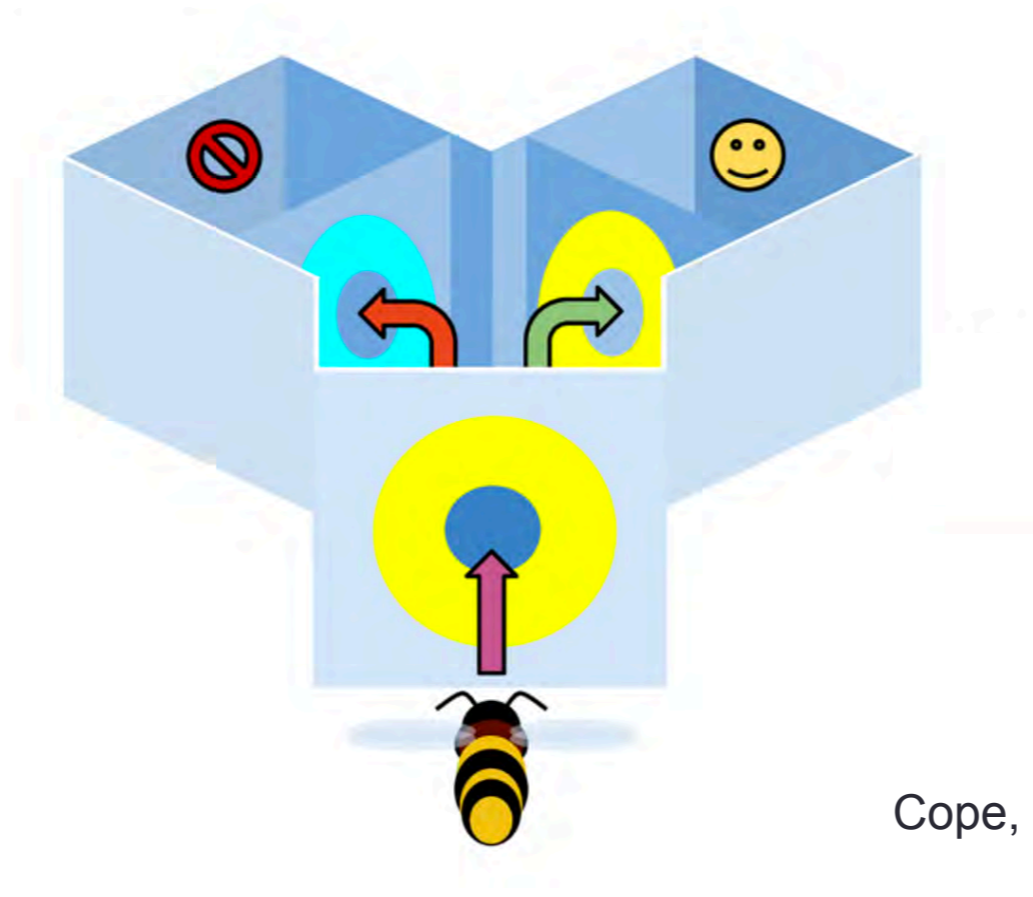
TO BEE OR NOT TO BEE

- Learns with visual stimuli
- Generalises to odours
 - Does leaning “sameness” require higher order cognition?



TO BEE OR NOT TO BEE

- Key idea: it learns to go to the pattern of the highest/lowest activity.
- “Stimulus-specific adaptation” - like observations. Repetitive stimulus results in lower activity.

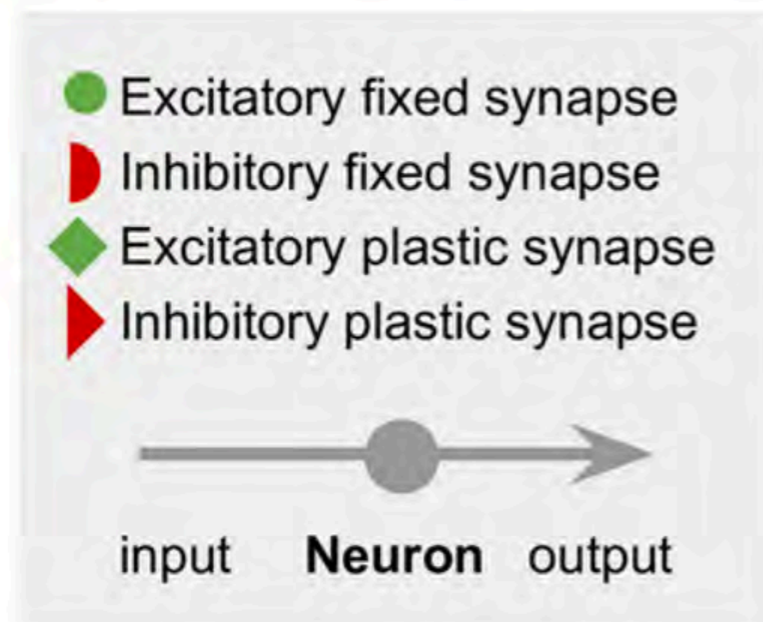


SAMENESS, DIFFERENCE AND TRANSFER OF LEARNING



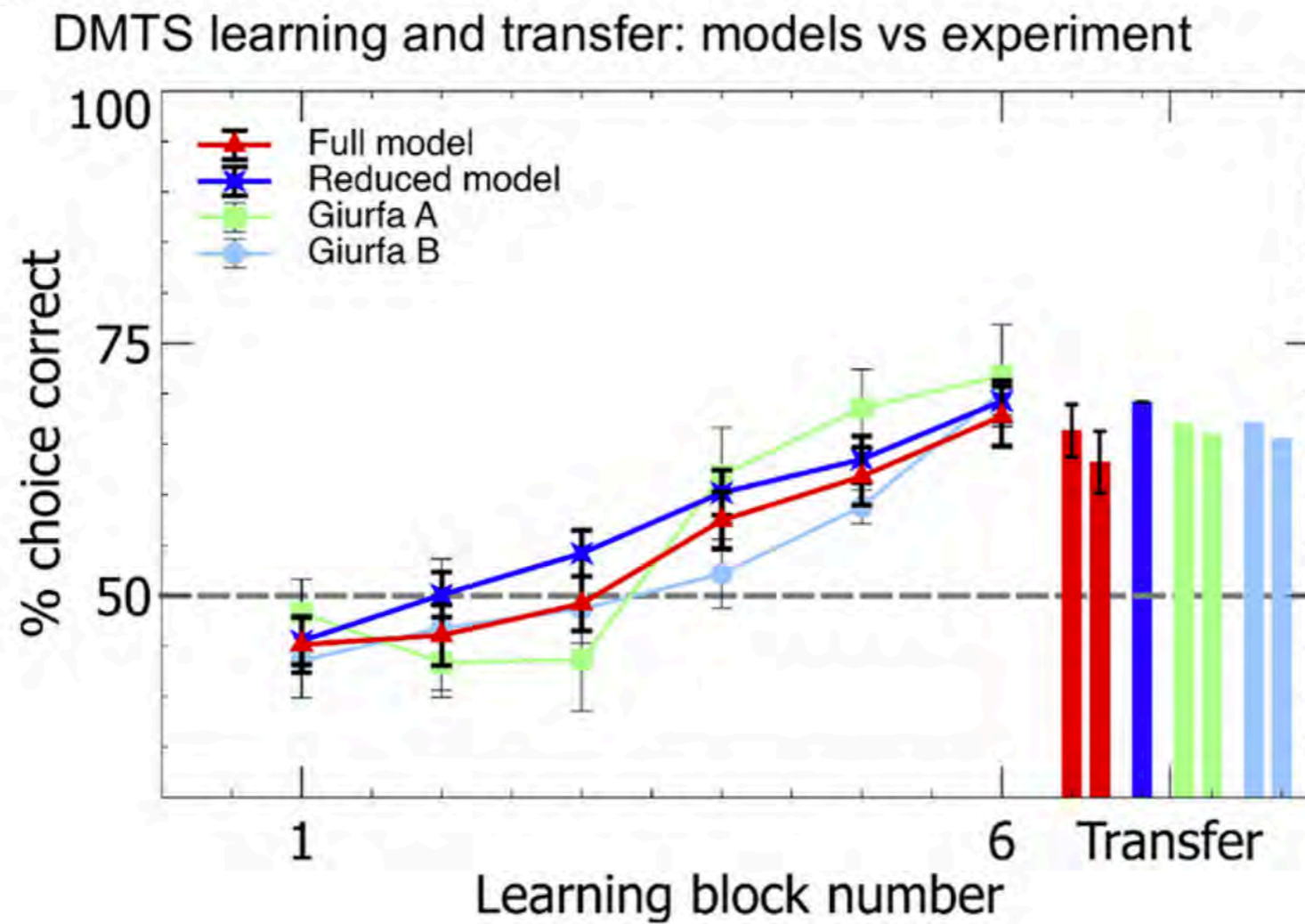
$$P(GO) = \frac{1}{1 + e^{-(c-d)(GO-NOGO)}}$$

$$P(NOGO) = 1 - P(GO)$$

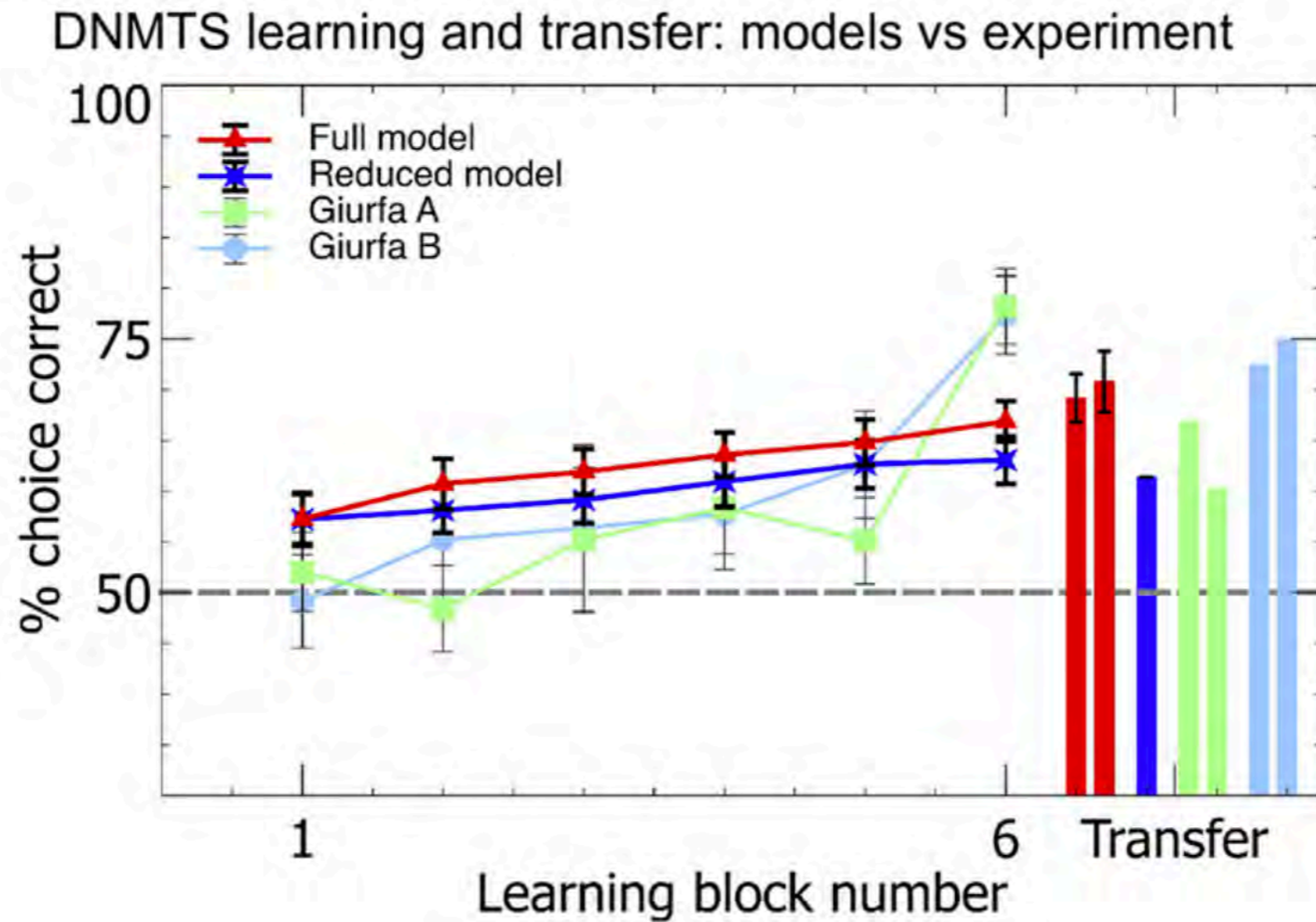


$$\Delta w^i = -\eta (R - R_b) \times \text{presynaptic activity} \times \text{postsynaptic activity}$$

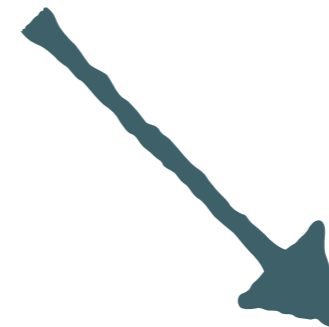
SAMENESS, DIFFERENCE AND TRANSFER OF LEARNING



SAMENESS, DIFFERENCE AND TRANSFER OF LEARNING



SUMMARY



Temporal
Difference

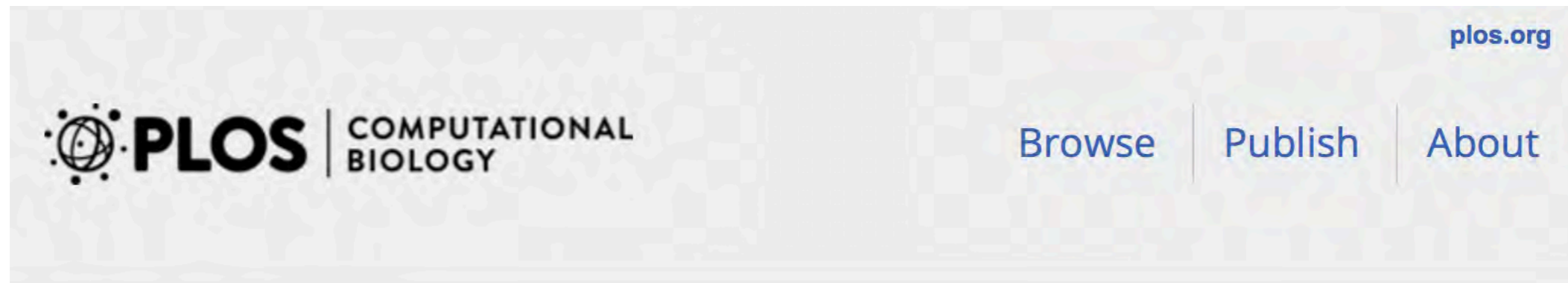
Policy Gradient

Reward-
modulated Hebb/
STDP

SUMMARY

- Eligibility & Future Rewards
 - Not necessary in TD learning but can speed up (depending on the task)
 - May be necessary in Policy Gradient learning (depending on the formalism)
 - May be necessary in reward-modulated Hebb/STDP (depending on the formalism)

ACKNOWLEDGEMENTS



 OPEN ACCESS  PEER-REVIEWED

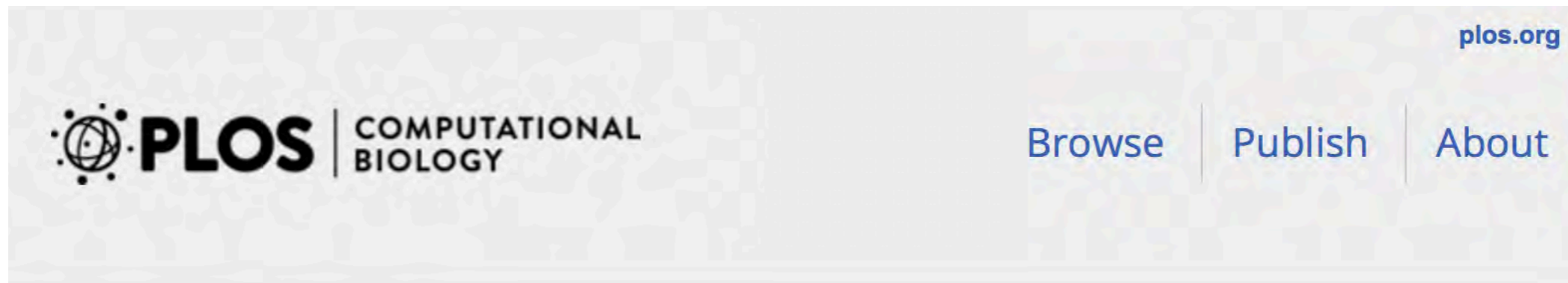
RESEARCH ARTICLE

Spike-Based Reinforcement Learning in Continuous State and Action Space: When Policy Gradient Methods Fail

Eleni Vasilaki , Nicolas Frémaux, Robert Urbanczik, Walter Senn, Wulfram Gerstner

Published: December 4, 2009 • <https://doi.org/10.1371/journal.pcbi.1000586>

ACKNOWLEDGEMENTS



 OPEN ACCESS  PEER-REVIEWED

RESEARCH ARTICLE

Abstract concept learning in a simple neural network inspired by the insect brain

Alex J. Cope , Eleni Vasilaki, Dorian Minors, Chelsea Sabo, James A. R. Marshall, Andrew B. Barron

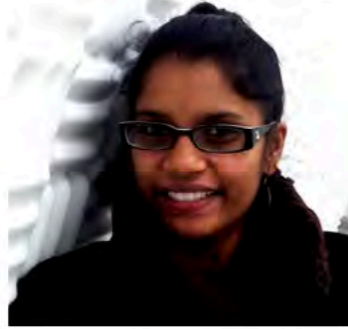
Version 2



Published: September 17, 2018 • <https://doi.org/10.1371/journal.pcbi.1006435>

• >> See the preprint

Postdoctoral Fellows

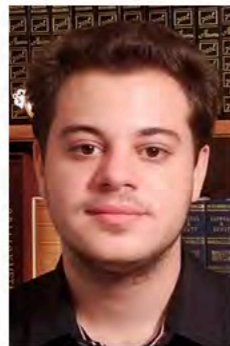


Dr Natacha Vanattou-Saïfoudine — with INI, Zurich



Dr Alex Cope — with Brains on Board project

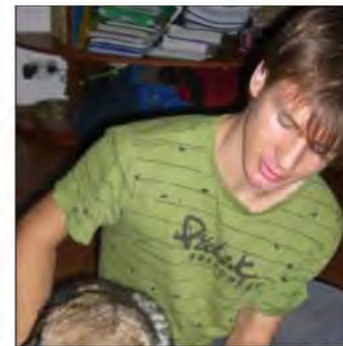
PhD students



Avgoustinos Vouros



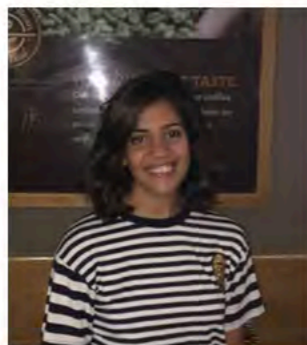
Chao Han



Luca Manneschi



Matthew Whelan — with Sheffield Robotics



Nada Abdelrahman — with Biomedical Sciences

