

# Learning with Opponent-Learning Awareness

*GTC 2018*

Full paper at AAMAS 18

Jakob N. Foerster<sup>1,2,†</sup>, Richard Y. Chen<sup>1,†</sup>, Maruan Al-Shedivat<sup>4</sup>,  
Shimon Whiteson<sup>2</sup>, Pieter Abbeel<sup>3</sup>, Igor Mordatch<sup>1</sup>

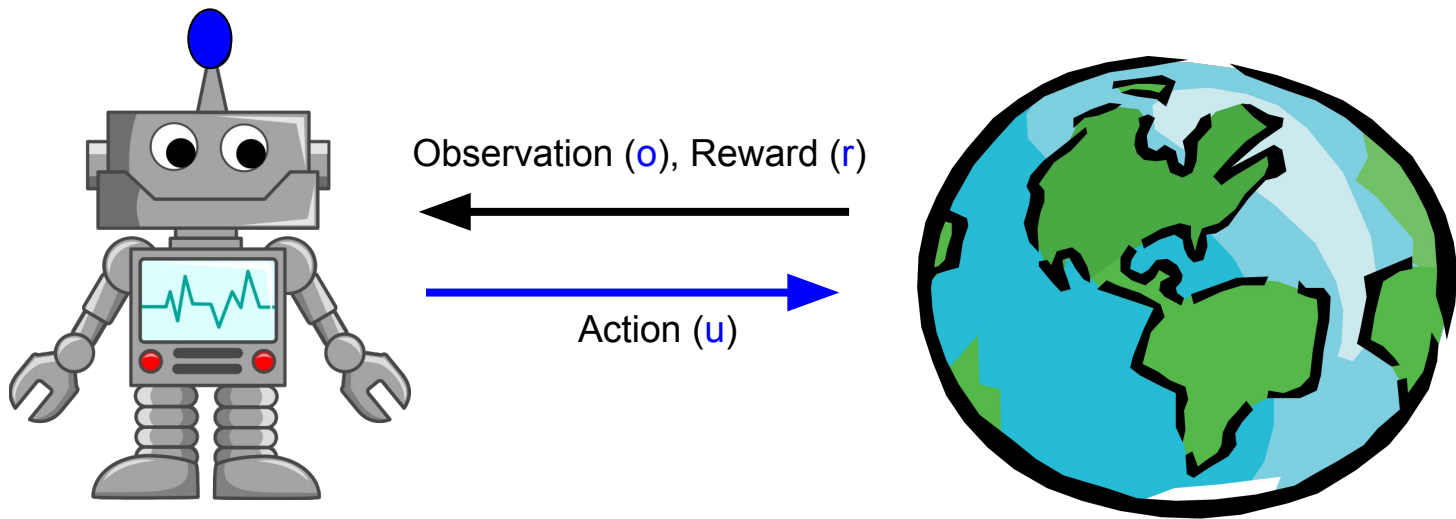


UNIVERSITY OF  
**OXFORD**



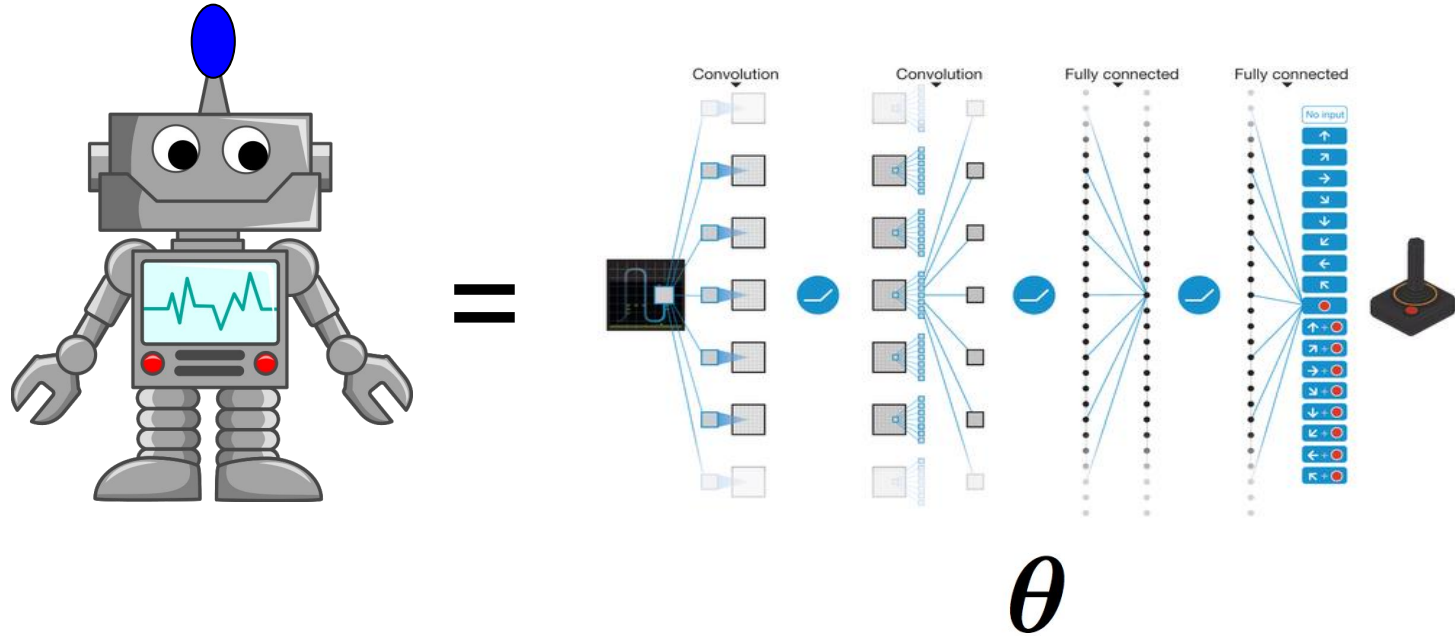


# Reinforcement Learning

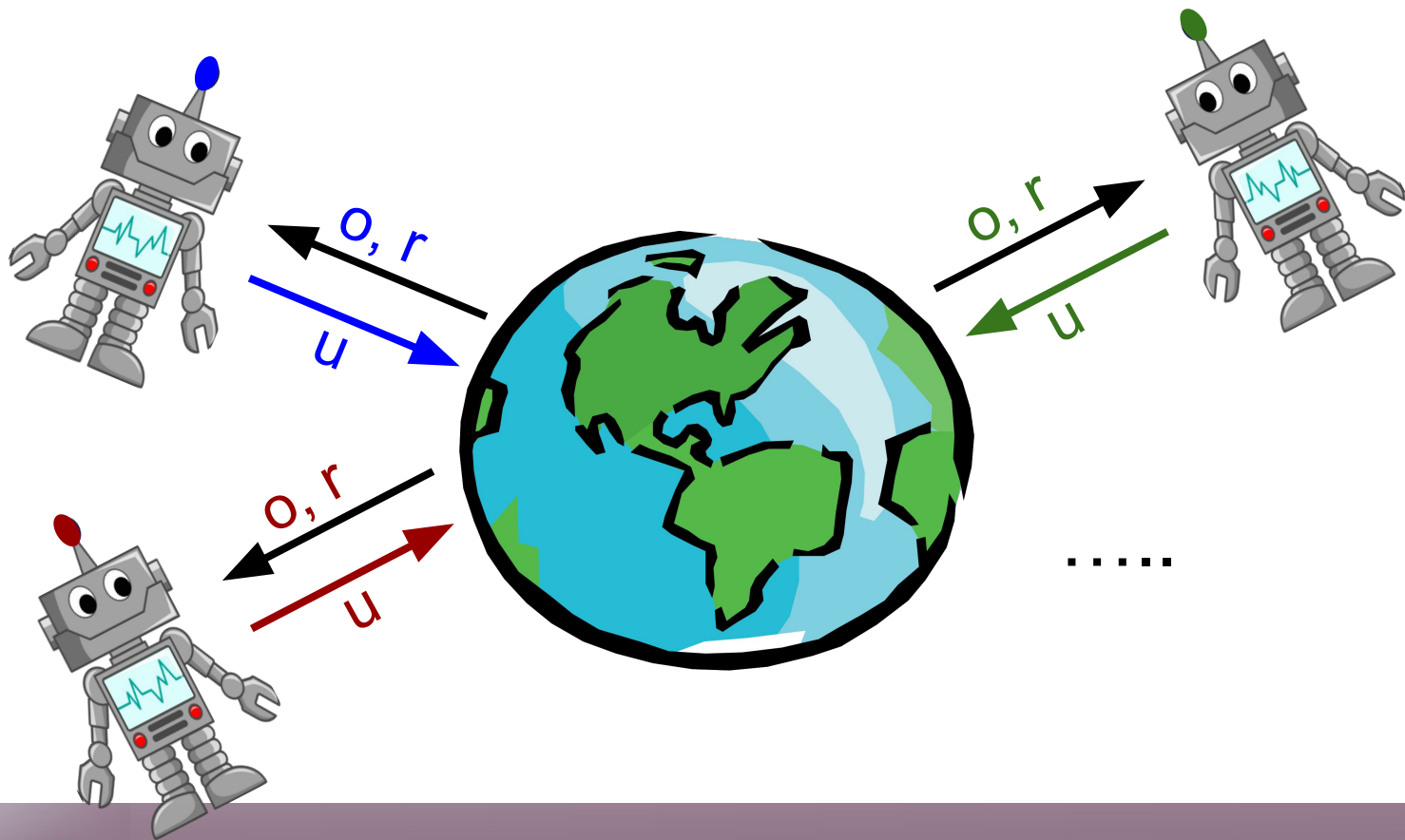


Goal is to maximise total return per episode:  $V = \sum \gamma^t r_t$

# Deep Reinforcement Learning



# Multi-Agent Reinforcement Learning [MARL]



# Some “great challenges” of MARL



- Communication
- Non-stationarity
- Credit Assignment
- Reciprocity

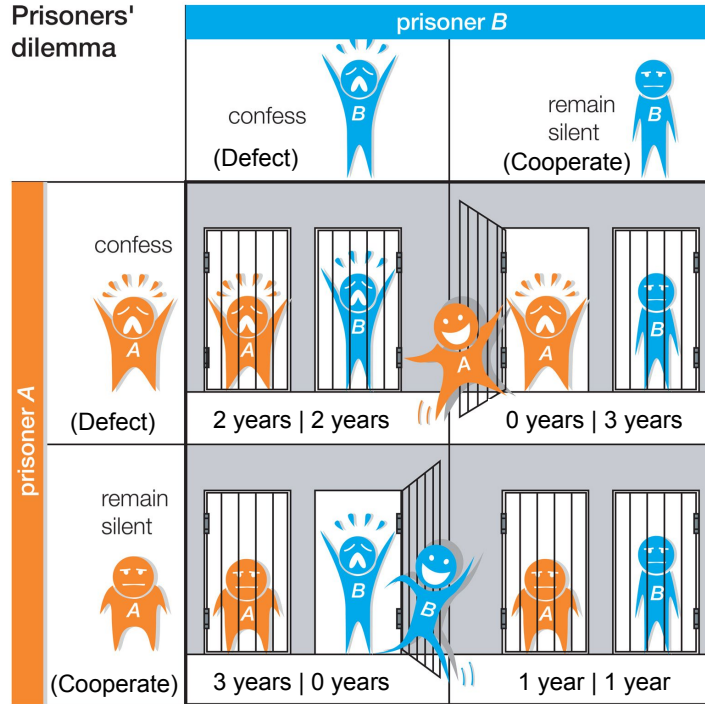


# LOLA Motivation



# Prisoners Dilemma

Prisoners' dilemma



	D	C
D	$(-2, -2)$	$(0, -3)$
C	$(-3, 0)$	$(-1, -1)$

© 2010 Encyclopædia Britannica, Inc.



# Prisoner's Dilemma

Payout matrix:

	D	C
D	$(-2, -2)$	$(0, -3)$
C	$(-3, 0)$	$(-1, -1)$

Background

- Single shot game:
  - Defection is only Nash equilibrium
- Repeated game (with high gamma):
  - Folk theorem says many equilibria

It's everywhere..:

## 4 Real-life examples

- 4.1 In environmental studies
- 4.2 In animals
- 4.3 In psychology
- 4.4 In economics
- 4.5 In sport
- 4.6 Multiplayer dilemmas
- 4.7 In international politics

# Related Work

## Non-cooperative Deep RL:

- Generalization of tit-for-tat with deep RL [Lerer & Peysakhovich, 2017]
- Investigation of pro-social Learners in generalised stag hunt [Peysakhovich & Lerer, 2017]
- Emergence of cooperation and competition [Leibo et al, 2017]
- Centralized actor-critic for training [Lowe et al, 2017]

## Opponent modeling:

- fictitious play [Brown, 1951],
- action prediction [Mealing & Shapiro, 2013]

## Opponent learning:

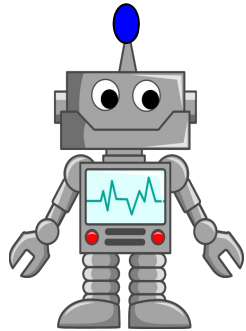
- Policy prediction under one-step learning dynamics [Zhang & Lester, 2017]
- Unrolled GAN [Metz et al, 2016] differentiates through opponent's update steps

## Human-Machine Interaction:

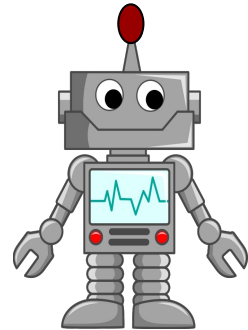
- “Planning for Autonomous Cars that Leverage Effects on Human Action” [Sadigh et al, 2016]

# Naive Learning

$$\theta_{i+1}^1 = \operatorname{argmax}_{\theta_1} V^1(\theta^1, \theta_i^2)$$



Naive

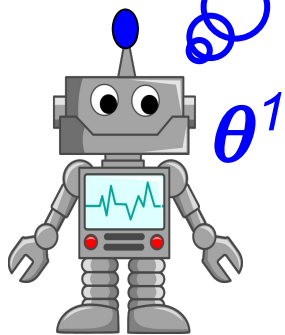
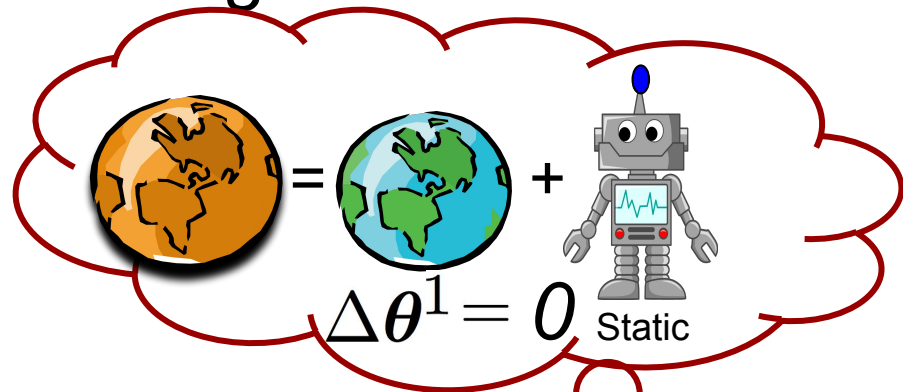
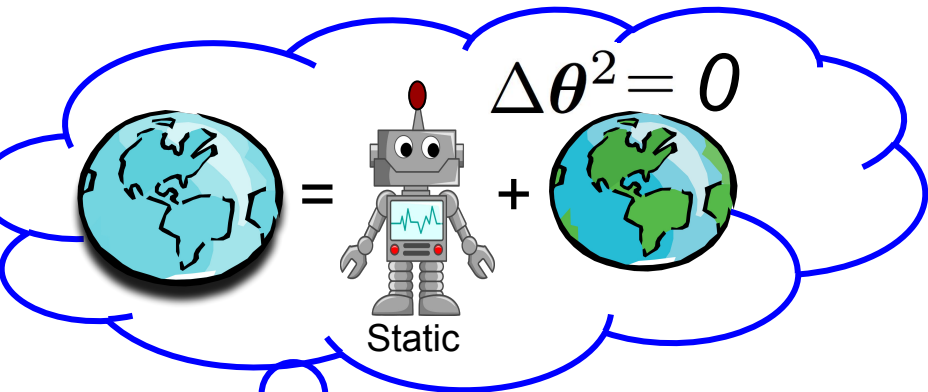


Naive



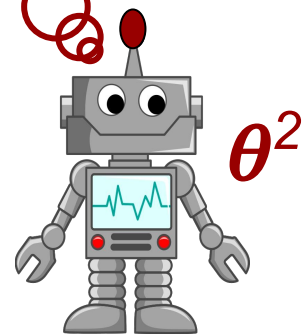
$$\theta_{i+1}^2 = \operatorname{argmax}_{\theta_2} V^2(\theta_i^1, \theta^2)$$

# Naive Learning



$\theta^1$

Naive

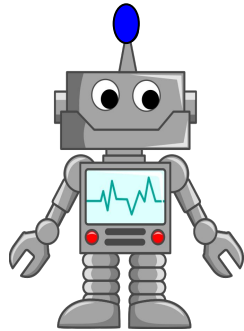


$\theta^2$

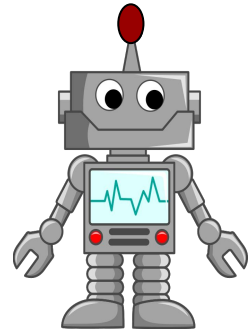
Naive

# Naive Learning with Gradients

$$\theta_{i+1}^1 = \theta_i^1 + \nabla_{\theta_i^1} V^1(\theta_i^1, \theta_i^2) \cdot \delta$$



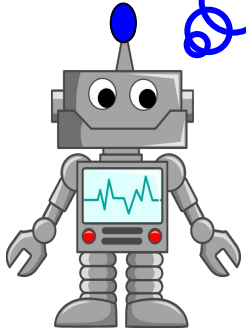
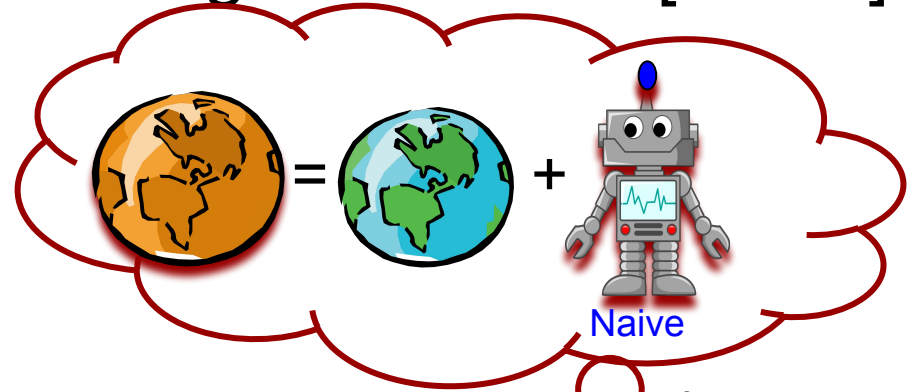
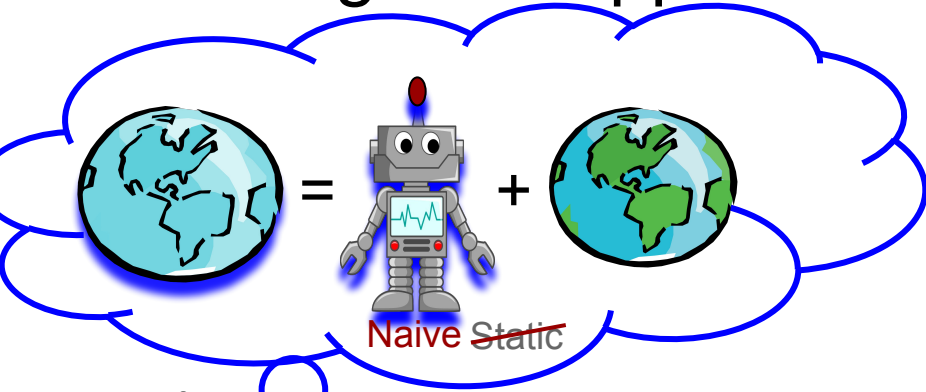
Naive



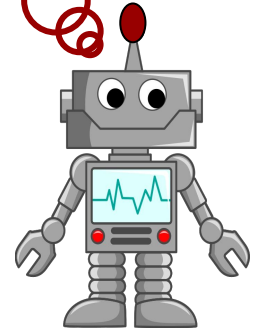
Naive

$$\theta_{i+1}^2 = \theta_i^2 + \nabla_{\theta_i^2} V^2(\theta_i^1, \theta_i^2) \cdot \delta$$

# Learning with Opponent Learning Awareness [LOLA]



LOLA



LOLA

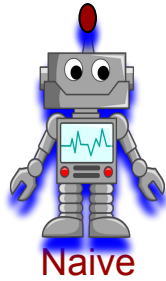
# LOLA with Gradients



=

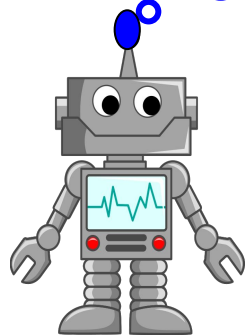


+



Naive

$$\Delta\theta^2 = \nabla_{\theta_i^2} V^2(\theta_i^1, \theta_i^2)$$



LOLA

$$\theta_{i+1}^1 = \theta_i^1 + \nabla_{\theta_i^1} V^1(\theta^1, \theta^2 + \Delta\theta^2) \cdot \delta$$

$$V^1(\theta^1, \theta^2 + \Delta\theta^2) \approx$$

$$V^1(\theta^1, \theta^2) + (\Delta\theta^2)^T \nabla_{\theta^2} V^1(\theta^1, \theta^2)$$



# LOLA Maths

Optimize Return after one step of opponent learning:

$$V^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2 + \Delta\boldsymbol{\theta}^2) \approx V^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) + (\Delta\boldsymbol{\theta}^2)^T \nabla_{\boldsymbol{\theta}^2} V^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2)$$

$$\Delta\boldsymbol{\theta}^2 = \nabla_{\boldsymbol{\theta}^2} V^2(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \cdot \eta$$

LOLA learning rule:

$$\begin{aligned} f_{\text{lola}}^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) &= \nabla_{\boldsymbol{\theta}^1} V^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \\ &+ \left( \nabla_{\boldsymbol{\theta}^2} V^1(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \right)^T \nabla_{\boldsymbol{\theta}^1} \nabla_{\boldsymbol{\theta}^2} V^2(\boldsymbol{\theta}^1, \boldsymbol{\theta}^2) \cdot \delta\eta \end{aligned}$$

Health warning:

This requires access to true value function and derivatives

# LOLA Policy Gradient

Can use Policy Gradients to estimate all gradients

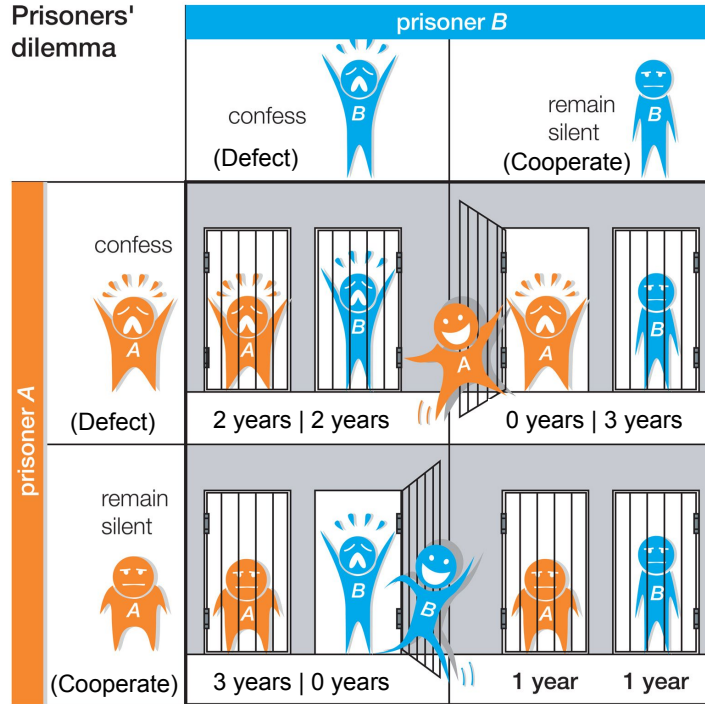
$$\begin{aligned} \mathbf{f}_{\text{lola, pg}}^1 &= \nabla_{\boldsymbol{\theta}^1} \mathbb{E} R_0^1(\tau) \cdot \delta + \\ &(\nabla_{\boldsymbol{\theta}^2} \mathbb{E} R_0^1(\tau))^T \nabla_{\boldsymbol{\theta}^1} \nabla_{\boldsymbol{\theta}^2} \mathbb{E} R_0^2(\tau) \cdot \delta \eta. \end{aligned}$$

LOLA term is still tractable and exact (in expectation):

$$\begin{aligned} &\nabla_{\boldsymbol{\theta}^1} \nabla_{\boldsymbol{\theta}^2} \mathbb{E} R_0^2(\tau) \\ &= \mathbb{E} \left[ R_0^2(\tau) \nabla_{\boldsymbol{\theta}^1} \log \pi^1(\tau) (\nabla_{\boldsymbol{\theta}^2} \log \pi^2(\tau))^T \right] \\ &= \mathbb{E} \left[ \sum_{t=0}^T \gamma^t r_t^2 \cdot \left( \sum_{l=0}^t \nabla_{\boldsymbol{\theta}^1} \log \pi^1(u_l^1 | s_l) \right) \right. \\ &\quad \left. \left( \sum_{l=0}^t \nabla_{\boldsymbol{\theta}^2} \log \pi^2(u_l^2 | s_l) \right)^T \right]. \end{aligned}$$

# Prisoners Dilemma

Prisoners' dilemma

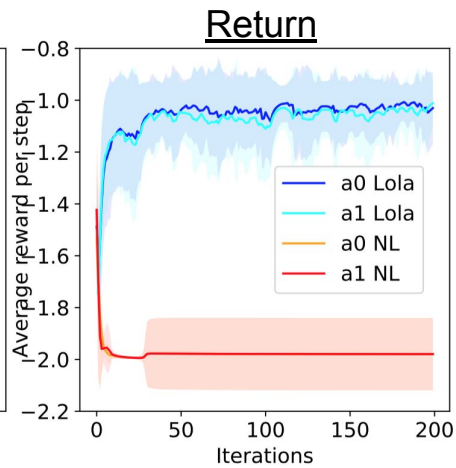
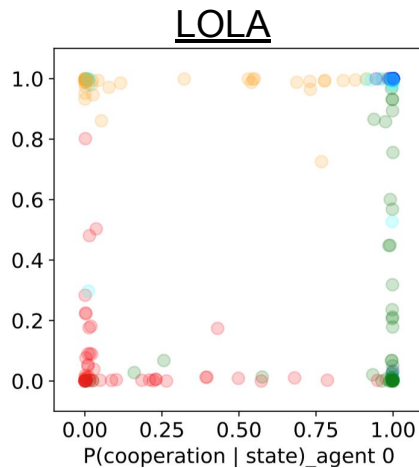
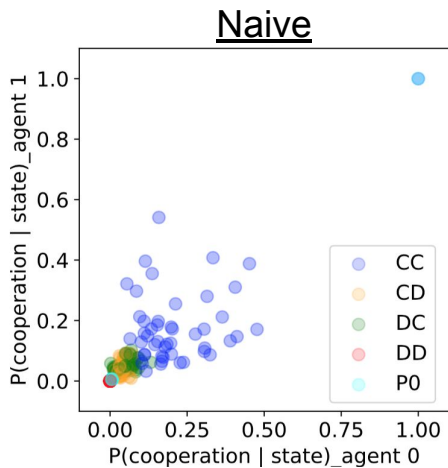


	D	C
D	$(-2, -2)$	$(0, -3)$
C	$(-3, 0)$	$(-1, -1)$

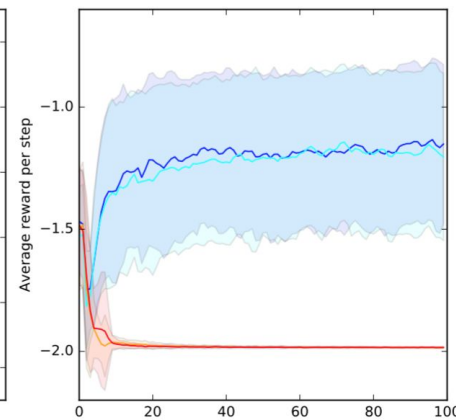
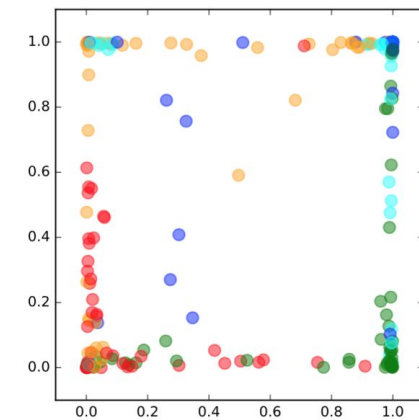
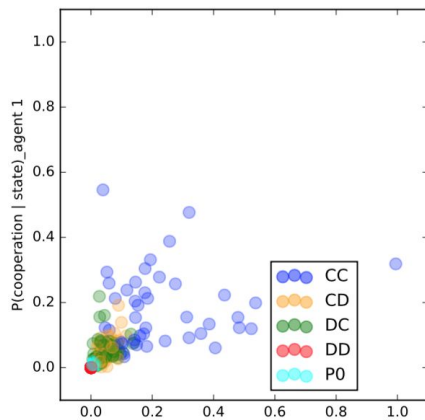
© 2010 Encyclopædia Britannica, Inc.

# Iterated Prisoner's Dilemma

Exact



Policy Gradient



# Matching Pennies

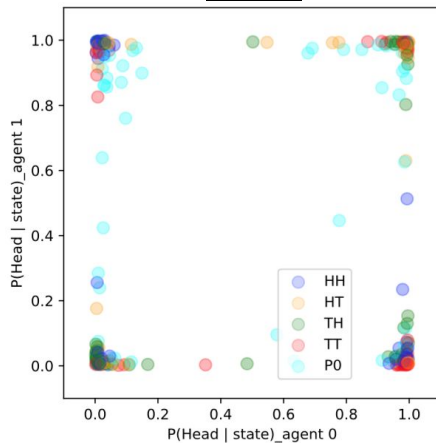


	Head	Tail
Head	$(+1, -1)$	$(-1, +1)$
Tail	$(-1, +1)$	$(+1, -1)$

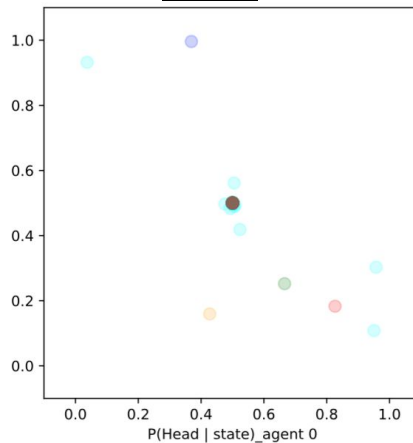
# Iterated Matching Pennies

Exact

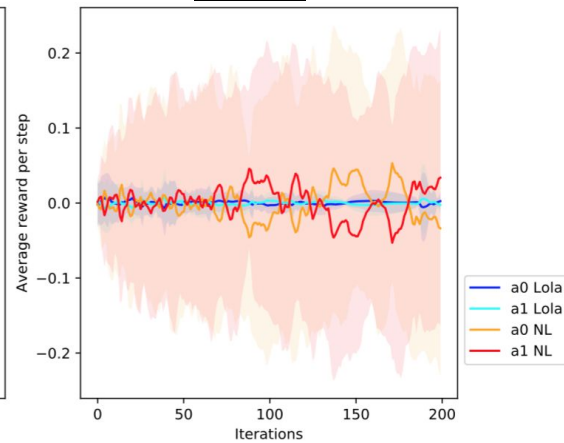
Naive



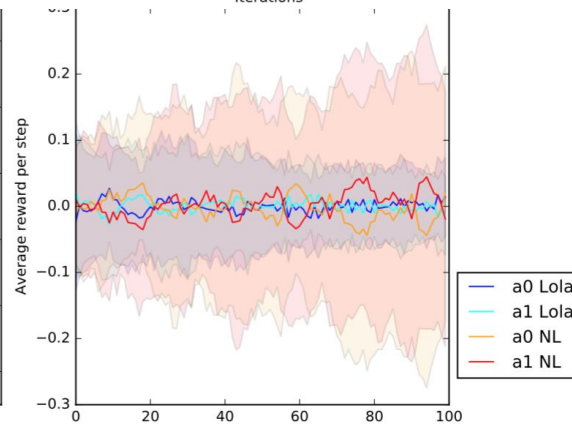
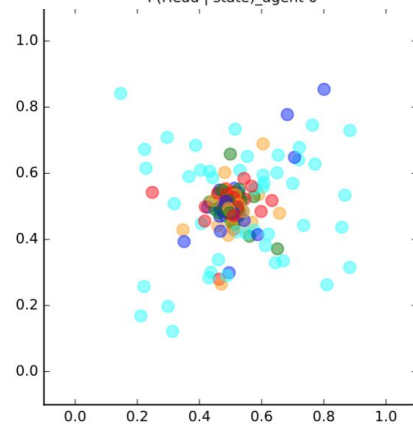
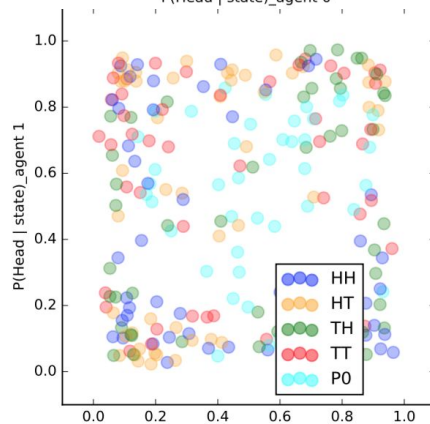
LOLA



Return



Policy Gradient



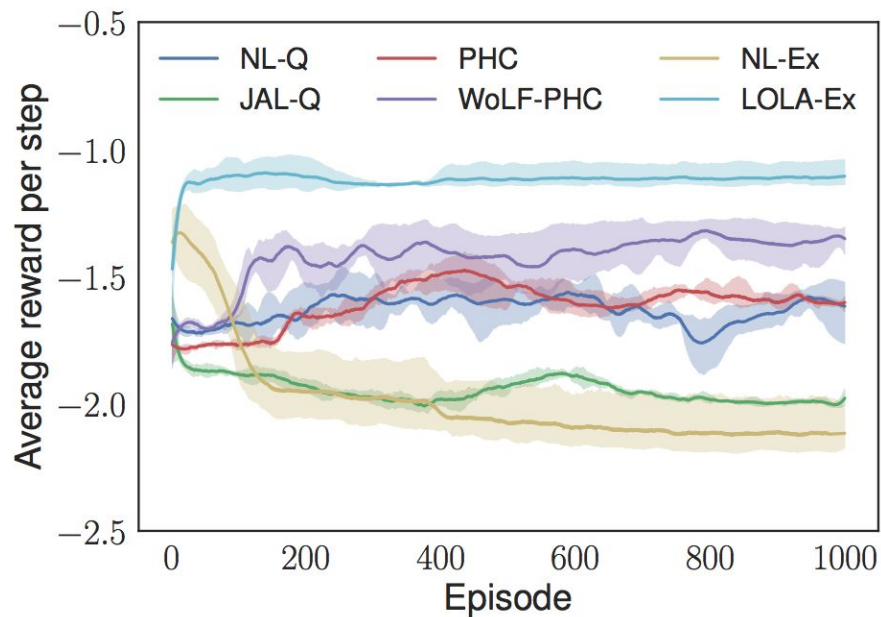
# Results

	IPD		IMP	
	%TFT	R(std)	%Nash	R(std)
NL-Ex.	20.8	-1.98(0.14)	0.0	0(0.37)
LOLA-Ex.	81.0	-1.06(0.19)	98.8	0(0.02)
NL-PG	20.0	-1.98(0.00)	13.2	0(0.19)
LOLA-PG	66.4	-1.17(0.34)	93.2	0(0.06)

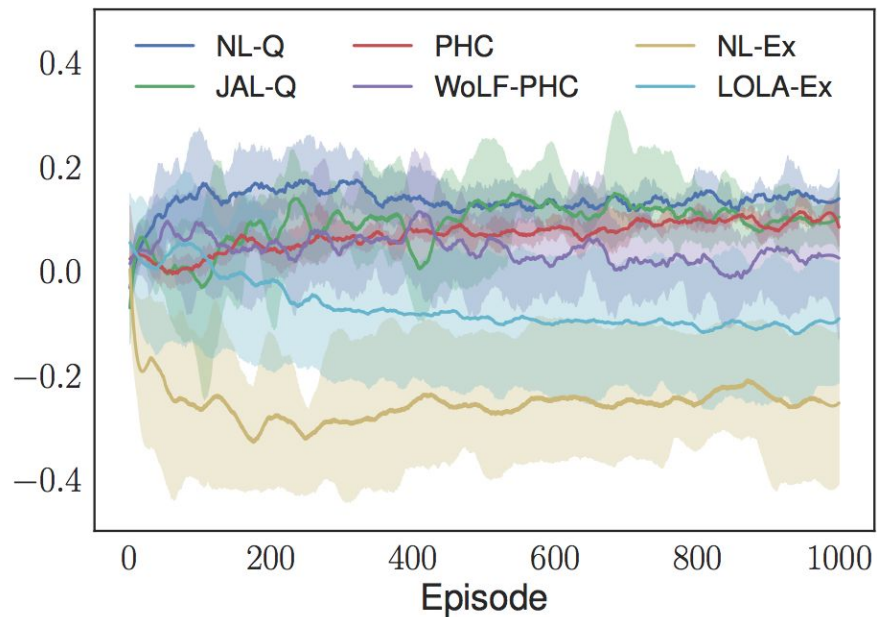


# Round Robin Tournament

Iterated Prisoner's Dilemma



Iterated Matching Pennies



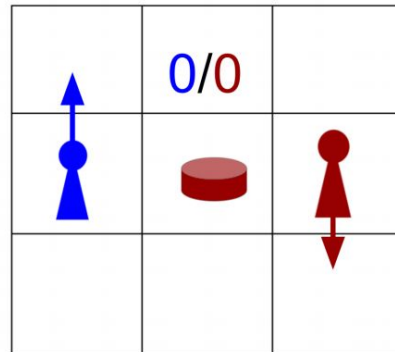
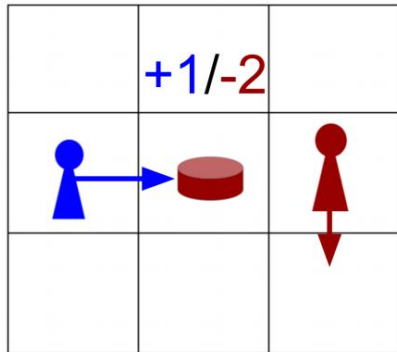
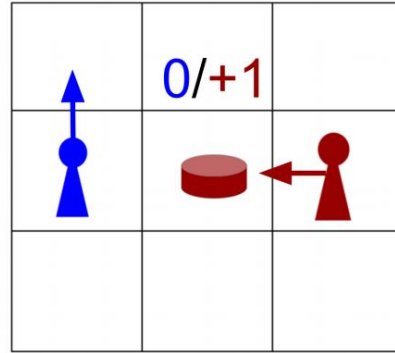
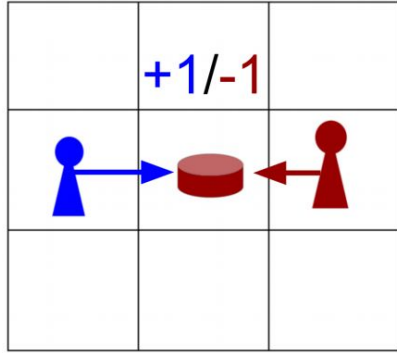
# LOLA with Opponent Modelling (LOLA-OM)



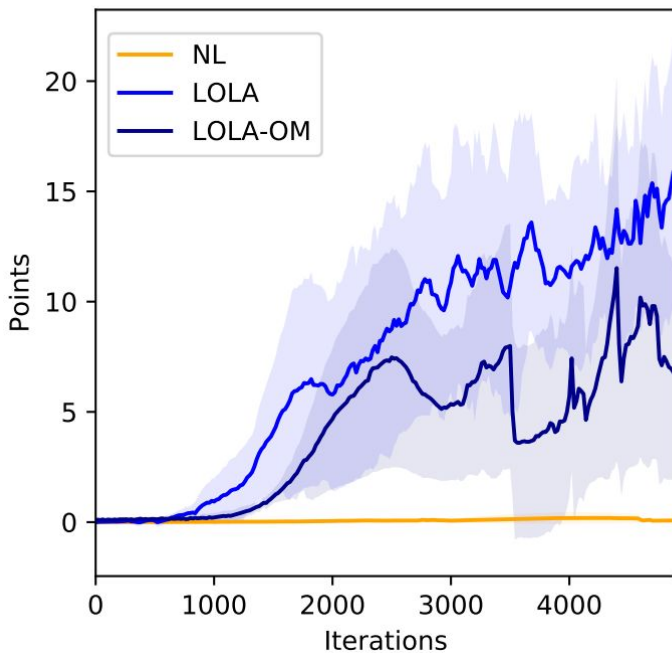
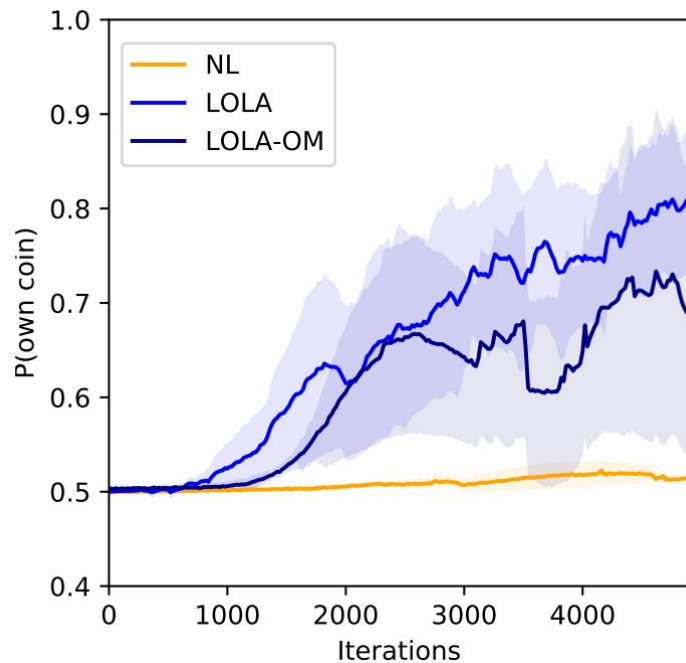
$$\theta^2 \rightarrow \hat{\theta}^2$$

$$\hat{\theta}^2 = \operatorname{argmax}_{\theta^2} \sum_t \log \pi_{\theta^2}^2(u_t^2 | s_t)$$

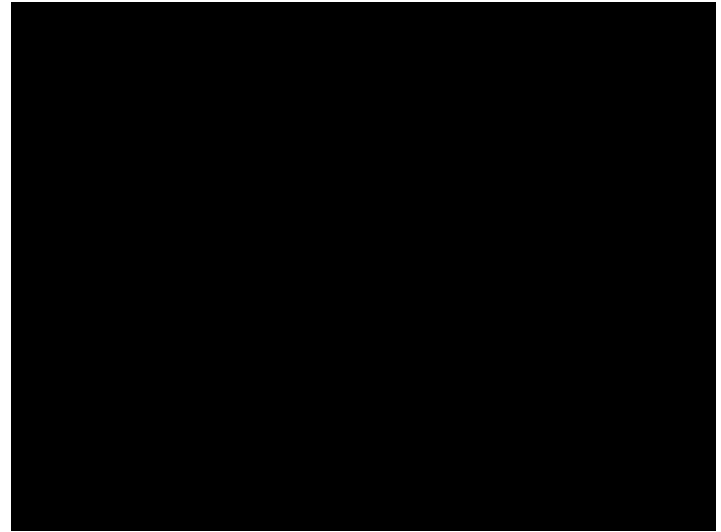
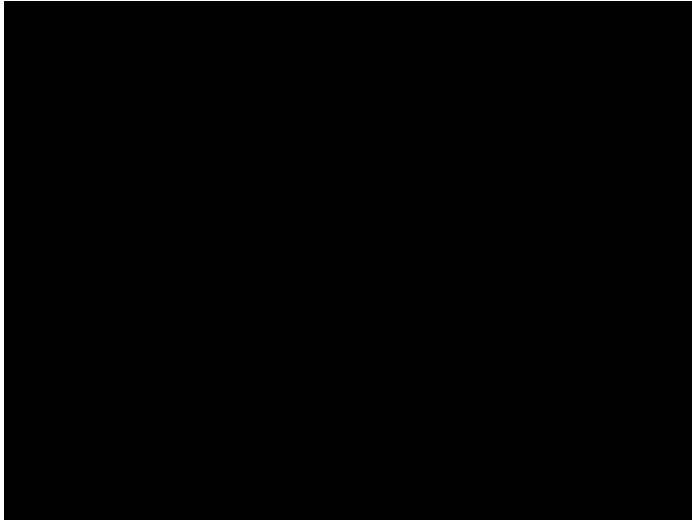
# LOLA with Recurrent Deep RL



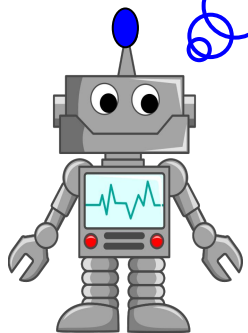
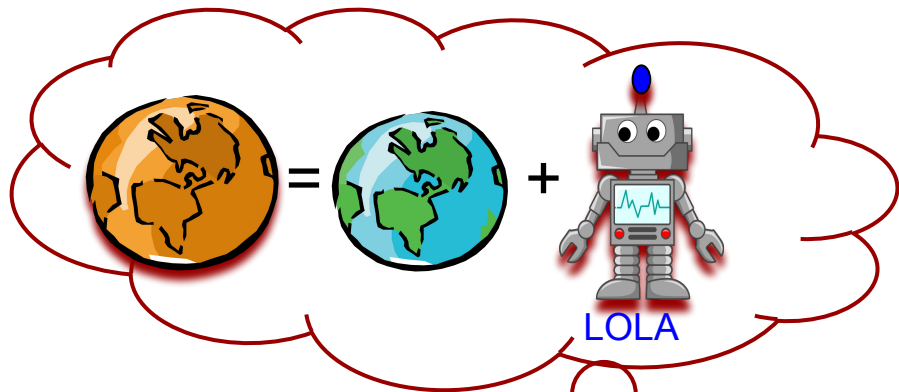
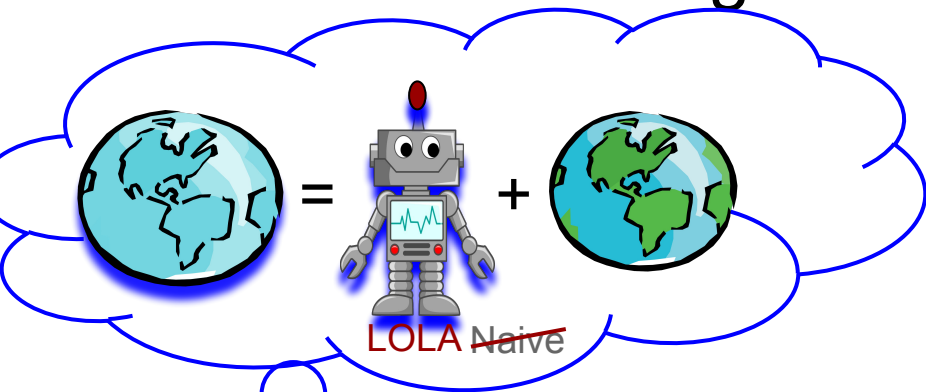
# LOLA with Recurrent Deep RL



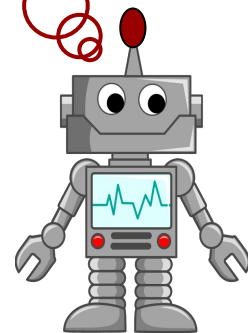
# LOLA with Recurrent Deep RL



# Higher Order LOLA



2nd Order  
LOLA



2nd Order  
LOLA

# Higher Order LOLA results

	NL	1st order	2nd Order
NL	$(-1.99, -1.99)$	$(-1.54, -1.28)$	-
1st	$(-1.28, -1.54)$	$(-1.04, -1.04)$	$(-1.14, -1.17)$



# LOLA Open Challenges



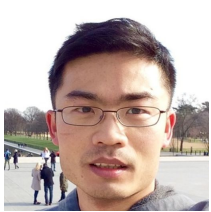
- Unknown update rules?
- Adversarial update rules?
- Proofs

# LOLA Conclusion

- State of the art Deep-MARL methods lead to defection
- LOLA leads to emergent reciprocity
- Cooperation arises out of selfish interest, considering learning of the opponent
- Works both in an exact setting and in Deep RL using policy gradients



# Acknowledgements



UNIVERSITY OF  
OXFORD



**nVIDIA**



***Thank you for listening!!***



***Any questions?***



**THE END**