# COLA: Consistent Learning with Opponent-Learning Awareness

Timon Willi*, Alistair Letcher*, Johannes Treutlein*, Jakob Foerster

UNIVERSITY OF OXFORD
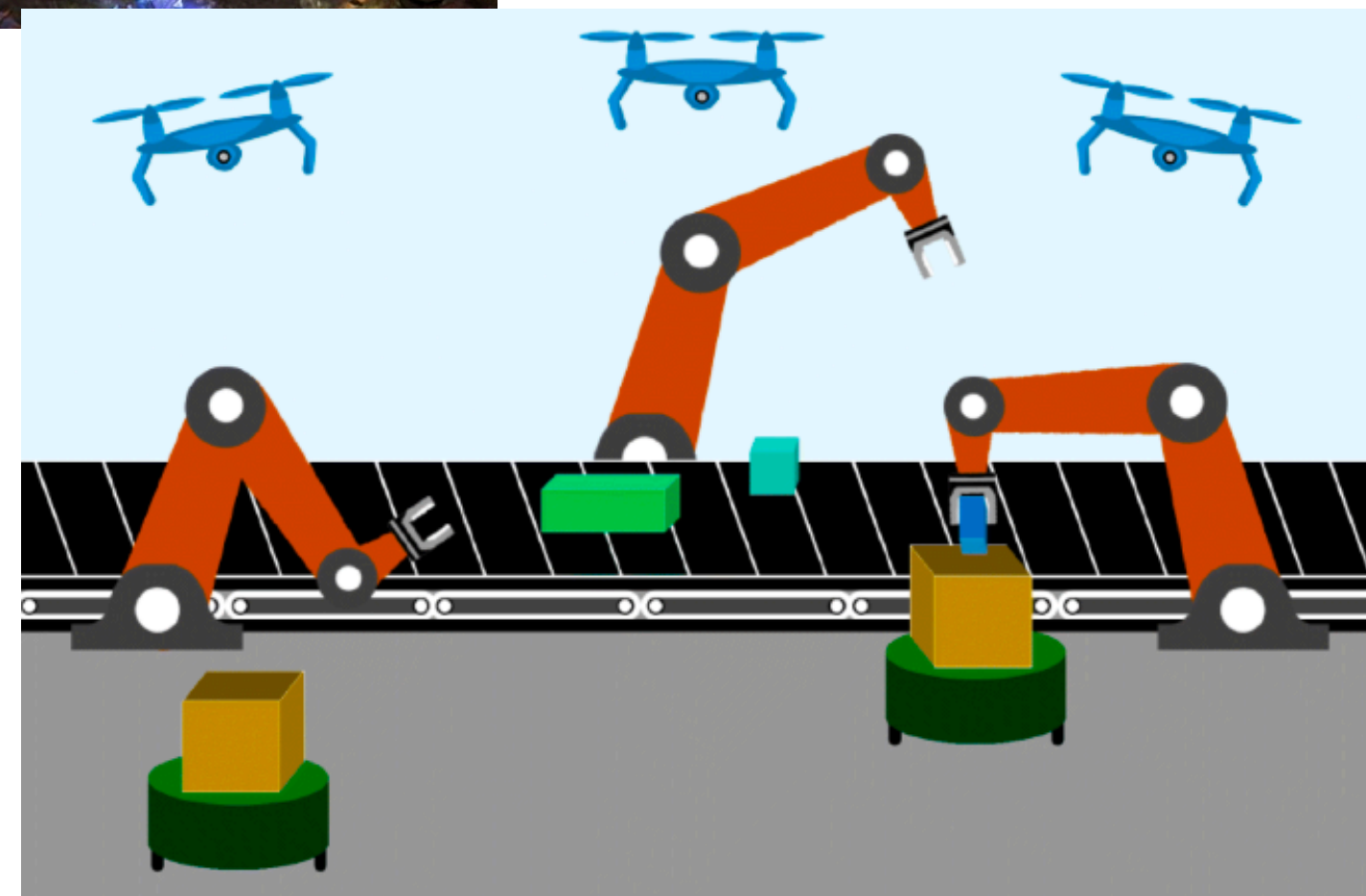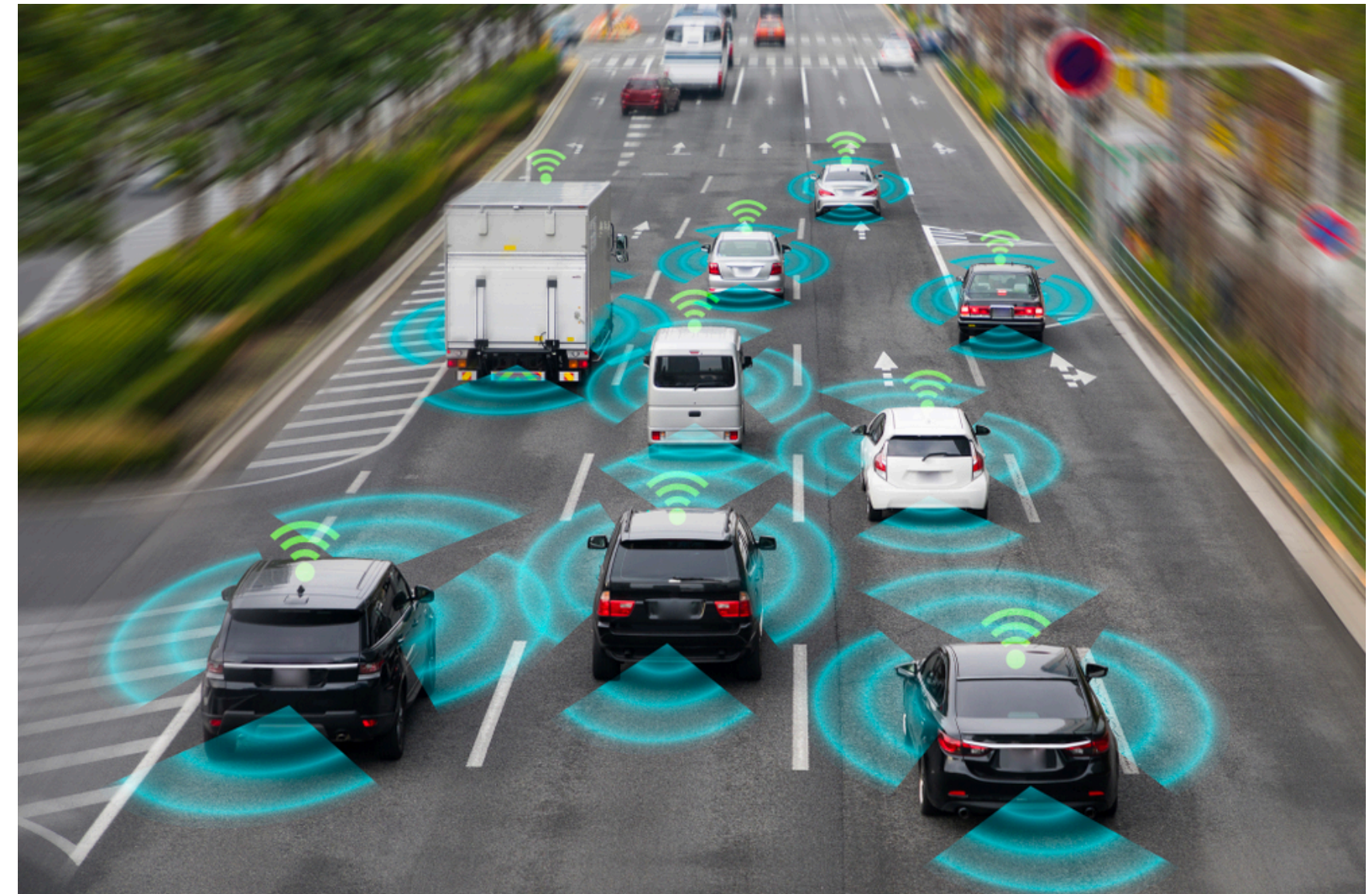
FLAIR Foerster Lab for AI Research

VECTOR INSTITUTE

UNIVERSITY OF TORONTO

* Equal Contribution

# Multi-Agent Reinforcement Learning

# General-Sum Games

# Opponent Shaping & LOLA

- **LOLA:**[1] $\Delta\theta_1 = -\alpha\nabla_1\left(L^1\left(\theta_1, \theta_2 + \widehat{\Delta\theta_2}\right)\right)$, where $\widehat{\Delta\theta_2} = -\alpha\nabla_2 L^2$

- LOLA finds **Tit-for-Tat** in the Iterated Prisoner's Dilemma.

- LOLA is **inconsistent:** it assumes that the opponent is a **Naive Learner.**

- LOLA does not find Stable Fixed Points (SFPs).

- Stable Fixed Points are a popular solution concept.

1 Foerster et al., 2017

# Contributions

# Higher-Order LOLA (HOLA) & iLOLA

- **HOLAn:**

$$h_1^{n+1} := -\alpha \nabla_1 \left( L^1 \left( \theta_1, \theta_2 + h_2^n \right) \right)$$

$$h_2^{n+1} := -\alpha \nabla_2 \left( L^2 \left( \theta_1 + h_1^n, \theta_2 \right) \right)$$

**where** $h_1^{-1} = h_2^{-1} = 0$

- **iLOLA** := $\displaystyle\lim_{n \to \infty} \begin{pmatrix} h_1^n \\ h_2^n \end{pmatrix}$ **if** $\mathrm{HOLA}n = \left( h_1^n, h_2^n \right)$ **converges pointwise as** $n \to \infty$

# Consistency

- **Consistency:** Any update functions $f_1 : \mathbb{R}^d \to \mathbb{R}^{d_1}$ and $f_2 : \mathbb{R}^d \to \mathbb{R}^{d_2}$ are consistent if they satisfy:

$$f_1\left(\theta_1, \theta_2\right) = -\alpha \nabla_1 \left( L^1 \left( \theta_1, \theta_2 + f_2\left(\theta_1, \theta_2\right) \right) \right)$$

$$f_2\left(\theta_1, \theta_2\right) = -\alpha \nabla_2 \left( L^2 \left( \theta_1 + f_1\left(\theta_1, \theta_2\right), \theta_2 \right) \right)$$

- **Proposition:** *iLOLA is consistent under mutual opponent shaping*

# Competitive Gradient Descent ≠ iLOLA

The authors of CGD[1] (2019) claim that:

1. the higher-order ``series-expansion [of CGD] would recover higher-order LOLA'' (page 4)

2. ``LCGD [Linearized CGD] coincides with first-order LOLA'' (page 6)

- **Proposition:** CGD does not coincide with iLOLA and does not solve the inconsistency problem

**What I think that they think that I think ... that they do**: Another game-theoretic interpretation of CGD follows from the observation that its update rule can be written as

$$\begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix} = - \begin{pmatrix} \mathrm{Id} & \eta D_{xy}^2 f \\ \eta D_{yx}^2 g & \mathrm{Id} \end{pmatrix}^{-1} \begin{pmatrix} \nabla_x f \\ \nabla_y g \end{pmatrix}. \qquad (4)$$

Applying the expansion $\lambda_{\max}(A) < 1 \Rightarrow (\mathrm{Id} - A)^{-1} = \lim_{N \to \infty} \sum_{k=0}^{N} A^k$ to the above equation, we observe that the first partial sum ($N = 0$) corresponds to the optimal strategy if the other player's strategy stays constant (GDA). The second partial sum ($N = 1$) corresponds to the optimal strategy if the other player thinks that the other player's strategy stays constant (LCGD, see Figure 1). The third partial sum ($N = 2$) corresponds to the optimal strategy if the other player thinks that the other player thinks that the other player's strategy stays constant, and so forth, until the Nash equilibrium is recovered in the limit. For small enough $\eta$, we could use the above series expansion to solve for $(\Delta x, \Delta y)$, which is known as Richardson iteration and would recover high order LOLA (Foerster et al., 2018). However, expressing it as a matrix inverse will allow us to use optimal Krylov subspace methods to obtain far more accurate solutions with fewer gradient evaluations.

the case of linear quadratic GANs.
Daskalakis et al. (2017) proposed to modify GDA as

$$\Delta x = - (\nabla_x f(x_k, y_k) + (\nabla_x f(x_k, y_k) - \nabla_x f(x_{k-1}, y_{k-1})))$$
$$\Delta y = - (\nabla_y g(x_k, y_k) + (\nabla_y g(x_k, y_k) - \nabla_y g(x_{k-1}, y_{k-1}))),$$

which we will refer to as optimistic gradient descent ascent (OGDA). By interpreting the differences appearing in the update rule as finite difference approximations to Hessian vector products, we see that (to leading order) OGDA corresponds to yet another second order correction of GDA (see Figure 1). It will also be instructive to compare the algorithms to linearized competitive gradient descent (LCGD), which is obtained by skipping the matrix inverse in CGD (which corresponds to taking only the leading order term in the limit $\eta D_{xy}^2 f \to 0$) and also coincides with first order LOLA (Foerster et al., 2018). As illustrated in Figure 1, these six algorithms amount to different subsets of the following four terms.

# COLA: Consistent Learning with Opponent-Learning Awareness

- Naive Solution to consistency: Iteratively compute higher orders of LOLA until convergence.

  - May diverge and requires arbitrarily high derivatives: expensive!

- We propose **COLA:** Learn $f_1$ and $f_2$ using neural networks

$$C_1\left(\phi_1, \phi_2, \theta_1, \theta_2\right) = \left\| f_1 + \alpha \nabla_1\left(L^1\left(\theta_1, \theta_2 + f_2\right)\right) \right\|$$

$$C_2\left(\phi_1, \phi_2, \theta_1, \theta_2\right) = \left\| f_2 + \alpha \nabla_2\left(L^2\left(\theta_1 + f_1, \theta_2\right)\right) \right\|$$

# COLA: Theoretical Results

**COLA's solutions are not necessarily unique**

**Proposition 4.5.** *Solutions to the consistency equations are not unique, even when restricted to linear solutions; more precisely, there exist several linear consistent solutions to the Tandem game.*

**Consistency does not imply preservation of SFPs**

**Proposition 4.6.** *Consistency does not imply preservation of SFPs: there is a consistent solution to the Tandem game with $\alpha = 1$ that fails to preserve **any** SFP. Moreover, for any $\alpha > 0$, there are **no** linear consistent solutions to the Tandem game that preserve more than one SFP.*

**COLA is more robust than LOLA**

**Proposition 4.7.** *For any non-zero initial parameters and any $\alpha > 1$, LOLA and SOS have divergent iterates in the Hamiltonian game. By contrast, any linear solution to the consistency equations converges to the origin for any initial parameters and **any** look-ahead rate $\alpha > 0$; moreover, the speed of convergence strictly increases with $\alpha$.*

# Experiments

# COLA's Consistency

- COLA finds consistent update functions *even when HOLA diverges.*

- COLA's updates are *similar to iLOLA when HOLA converges*.

- COLA tends to find similar solutions over different runs (despite the theoretical result)

| $\alpha$ | LOLA | HOLA3 | HOLA6 | COLA |
|---|---|---|---|---|
| 1.0 | 128.0 | 512 | 131072 | 3e-14±2e-15 |
| 0.5 | 12.81 | 14.05 | 12.35 | 2e-14±5e-15 |
| 0.3 | 2.61 | 2.05 | 0.66 | 4e-14±3e-15 |
| 0.1 | 0.08 | 9e-3 | 2e-6 | 6e-14±9e-15 |
| 0.01 | 1e-5 | 2e-8 | 4e-14 | 1e-14±4e-14 |

**Table 1: Consistency Loss on Tandem**

| $\alpha$ | LOLA | HOLA3 | HOLA6 |
|---|---|---|---|
| 1.0 | 0.94±0.04 | 0.94±0.04 | 0.94±0.04 |
| 0.5 | 0.88±0.12 | 0.88±0.12 | 0.08±0.13 |
| 0.3 | 0.92±0.01 | 0.91±0.01 | 0.80±0.01 |
| 0.1 | 0.95±0.01 | 0.99±0.01 | 0.99±0.01 |
| 0.01 | 0.99±0.01 | 1.00±0.00 | 1.00±0.00 |

**Table 2: Similarity Scores on Tandem**

| Game@LR | Cosine Sim |
|---|---|
| MP@10 | 0.97 ± 0.01 |
| MP@0.5 | 0.99 ± 0.01 |
| Tandem@0.1 | 1.00 ± 0.00 |
| Tandem@1.0 | 0.98 ± 0.01 |

**Table 3: Self-Similarity Scores**

# Tandem Game

$$L^1(x, y) = (x + y)^2 - 2x$$

$$L^2(x, y) = (x + y)^2 - 2y$$

- COLA and HOLA converge to similar solutions

- CGD converges to a different solution than COLA and HOLA8

- COLA does not recover SFP

# Iterated Prisoner's Dilemma

**CGD does not coincide with iLOLA and does not solve the inconsistency problem**

Table 1. Payoff Matrix for the Prisoner's Dilemma

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



- COLA finds prosocial solution

- COLA's policy is similar to Tit-for-Tat

- CGD does not find prosocial solution

# Matching Pennies

Table 6. Payoff Matrix for the Matching Pennies game.

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

- COLA converges robustly and fast to a wide range of hyperparameter values



MatchingPennies



COLA on MP



HOLA4 on MP

# Ultimatum Game



$$L_1 = - \left( 5p_{\textbf{fair}} + 8 \left( 1 - p_{\textbf{fair}} \right) p_{\textbf{accept}} \right)$$

$$L_2 = - \left( 5p_{\textbf{fair}} + 2 \left( 1 - p_{\textbf{fair}} \right) p_{\textbf{accept}} \right)$$

- Player 1 has access to 10$

  - Option 1: Split money 50/50

  - Option 2: Split money 80/20

- Player 2 can accept or decline

- *COLA is the only algorithm consistently converging to the fair solution*

# Conclusion

- Corrected a claim made in prior work

- iLOLA solves part of the consistency problem of LOLA

- Even with consistency, opponent shaping does not preserve SFPs

- We introduced **COLA**

- COLA tends to find prosocial solutions

# Thank you!