# Multi-Agent Equilibria: From Verification to Modification and Beyond

Muhammad Najib

School of Mathematical and Computer Sciences
Heriot-Watt University, UK

# Five Trends in Computing

1. Ubiquity
2. Interconnection
3. Delegation
4. Human Orientation
5. Intelligence

## Five Trends in Computing

1. **Ubiquity**
2. Interconnection
3. Delegation
4. Human Orientation
5. Intelligence

- Computing systems are everywhere (Moore's law: small, low-power, inexpensive CPUs).

- Computing systems embedded in devices around us: Roomba, smart fridge, Alexa,...
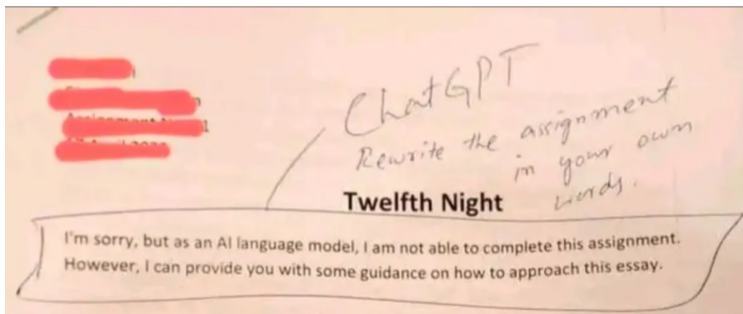
## Five Trends in Computing

1. Ubiquity
2. **Interconnection**
3. Delegation
4. Human Orientation
5. Intelligence

- Computer systems connected with one and another.
- e.g., internet

## Five Trends in Computing

- Computers do things for us (we let them take control).
- Fly-by-wire planes, autonomous cars, ...

1. Ubiquity
2. Interconnection
3. **Delegation**
4. Human Orientation
5. Intelligence

- Computers do things for us (we let them take control).
- Fly-by-wire planes, autonomous cars, …

1. Ubiquity
2. Interconnection
3. **Delegation**
4. Human Orientation
5. Intelligence

## Five Trends in Computing

- Many computer systems are designed to interact with humans.
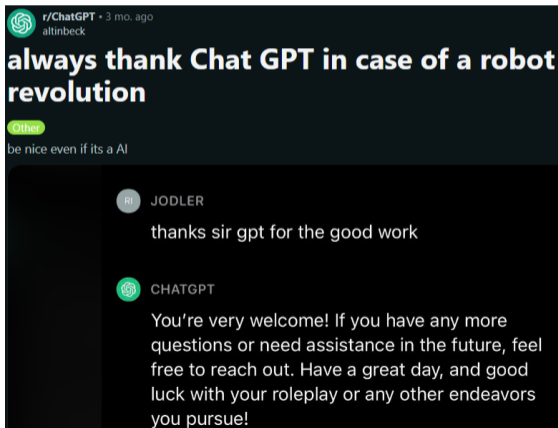- We interact with them like with humans (Alexa, Siri,...).

1. Ubiquity
2. Interconnection
3. Delegation
4. **Human Orientation**
5. Intelligence

- Many computer systems are designed to interact with humans.
- We interact with them like with humans (Alexa, Siri,…).

1. Ubiquity
2. Interconnection
3. Delegation
4. **Human Orientation**
5. Intelligence



r/ChatGPT • 3 mo. ago
altinbeck

**always thank Chat GPT in case of a robot revolution**

Other

be nice even if its a AI

RI · JODLER

thanks sir gpt for the good work

CHATGPT

You're very welcome! If you have any more questions or need assistance in the future, feel free to reach out. Have a great day, and good luck with your roleplay or any other endeavors you pursue!

# Five Trends in Computing

1. Ubiquity
2. Interconnection
3. Delegation
4. Human Orientation
5. **Intelligence**

- Data + Compute Power + Algorithm & Engineering
- AI systems become smarter, more capable.

## Five Trends in Computing

1. Ubiquity
2. Interconnection
3. Delegation
4. Human Orientation
5. Intelligence

Manifestations:

- Cloud computing
- Internet of Things
- Ubiquitous computing
- Semantic Web
- ...
- **Multi-agent systems**

## What is an Agent?

> "... a computer system that is capable of independent (**autonomous**) **action on behalf of its user**."[a]
>
> ---
> [a] Michael Wooldridge. *An Introduction to Multiagent Systems*. 2nd ed. Chichester, UK: Wiley, 2009.

> "... an **autonomous** entity which observes and **acts** upon an environment and directs its activity **towards achieving goals**."[a]
>
> ---
> [a] Stuart J. Russell and Peter Norvig. *Artificial Intelligence: A Modern Approach (4th Edition)*. Pearson, 2020. URL: http://aima.cs.berkeley.edu/.

9

# Example of an Agent

## Now ask Siri to ...

### Make a call

"Hey Siri, call Mom."

"Hey Siri, call Vivek's mobile on speakerphone."

Siri can also make and answer calls on HomePod ›

### Get directions

"Hey Siri, find coffee near me."

"Hey Siri, get directions home."

Use Siri with CarPlay ›

### Send a message

"Hey Siri, send a message to Ming Lu."

"Hey Siri, text Adrian and Sofia, 'Where are you?'"

Siri can read new messages on your AirPods ›

### Play music

"Hey Siri, play the hottest Taylor Swift tracks."

"Hey Siri, play the new Tame Impala album."

Learn more ways to play music ›

### Find information

"Hey Siri, what's the weather for today?"

"Hey Siri, how high is Mount Everest?"

Learn more things you can ask Siri ›

### Find your Apple device

"Hey Siri, where's my iPhone?"

"Hey Siri, find my AirPods."

Learn how to use Find My ›

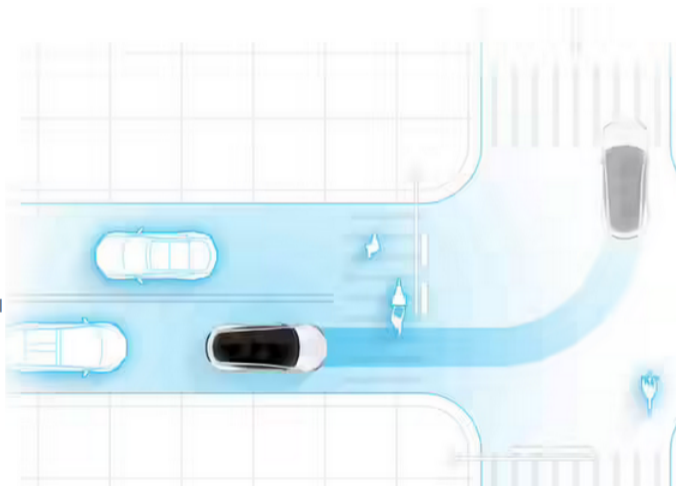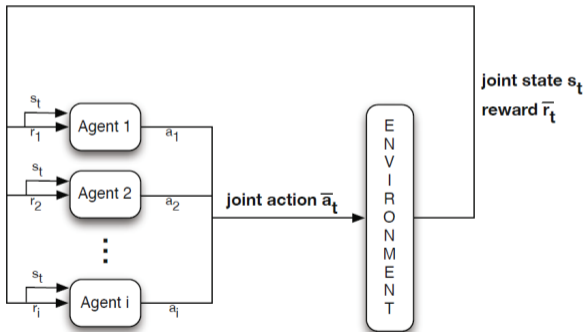**TESLA**

## From Home

All you will need to do is get in and tell your car where to go. If you don't say anything, your car will look at your calendar and take you there as the assumed destination. Your Tesla will figure out the optimal route, navigating urban streets, complex intersections and freeways.

# What is a Multi-Agent System?

- A system consists of **multiple agents** that **interact** with one another.
- Agents **act** on behalf of users/stakeholders with **different goals and preferences**.
- They interact and act upon the **environment**.



Source: Nowe, Ann & Vrancx, Peter & De Hauwere, Yann-Michaël. (2012). Game Theory and Multi-agent Reinforcement Learning.

## Example of a Multi-Agent System

- Algorithmic/high-frequency trading.
- Trading softwares **buy & sell** stocks to **generate as much money as possible.**

# Problem with Multi-Agent Systems

- MASs are prone to **instability** and might have **unpredictable dynamics**.
- Or, some stable behaviour gives rise to **bad outcomes**.
- 2010 Flash Crash[a]: over a 30 minutes period, Dow Jones lost (momentarily) over a trillion dollars of valuation.
  - "...the interaction between automated execution programs and algorithmic trading strategies can quickly erode liquidity and result in disorderly markets."[b]

[a]https://www.theguardian.com/business/2015/apr/22/2010-flash-crash-new-york-stock-exchange-unfolded
[b]U.S. Securities and Exchange Commission; Commodity Futures Trading Commission. "Findings Regarding the Market Events of May 6, 2010"



PREVIOUS CLOSE: 10,868.10

Dow industrials

Close
10,520.32
−3.2%

**Momentary Lapse**
Stock markets plunged suddenly yesterday afternoon and gained speed as computer programs prevented losses. But almost as quickly, the market recovered much of the decline.

2:46
9,869.62
−9.2%

10 A.M.   11 A.M.   12 P.M.   1 P.M.   2 P.M.   3 P.M.
Source: Bloomberg                          THE NEW YORK TIMES

15

# Problem with Multi-Agent Systems



- With *safety critical* systems (e.g., autonomous cars), not only we risk losing money but human lives.

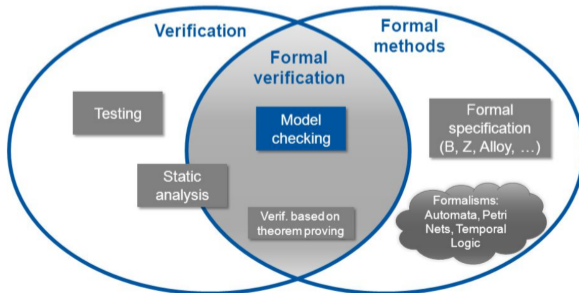# Problem with Multi-Agent Systems



- With *safety critical* systems (e.g., autonomous cars), not only we risk losing money but human lives.

We want our AI (multi-agent) systems to be **'CORRECT'**

# Part I: Verification

## Correctness in Computer Science

- The **correctness problem** has been one of the most widely studied problems in computer science over the past fifty years, and remains a topic of fundamental concern to the present day
- the correctness problem: checking that computer systems behave as their designer intends
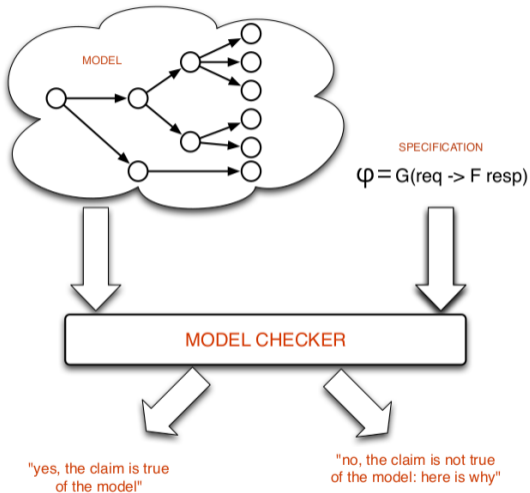- **Formal verification** is the problem of checking that a system $P$ is *correct* with respect to a formal specification $\varphi$ (e.g., LTL)

## Linear Temporal Logic (LTL)

- Standard formal language for talking about (infinite) state sequences
- Has been around for more than four decades[1]
- Propositional logic ($\land, \lor, \neg, \dots$) + temporal modalities ($\mathbf{G}, \mathbf{F}, \mathbf{X}, \dots$)
    - $\mathbf{G}p$: is always the case that $p$
    - $\mathbf{F}q$: will eventually the case that $q$
- We can express something like:
    - "*it is always not hot in Aberdeen*": $\mathbf{G}\neg hot$
    - "*eventually will be sunny in Aberdeen*": $\mathbf{F}sunny$

---

[1] Amir Pnueli. "The temporal logic of programs". In: *18th Annual Symposium on Foundations of Computer Science (sfcs 1977)*. ieee. 1977, pp. 46–57.
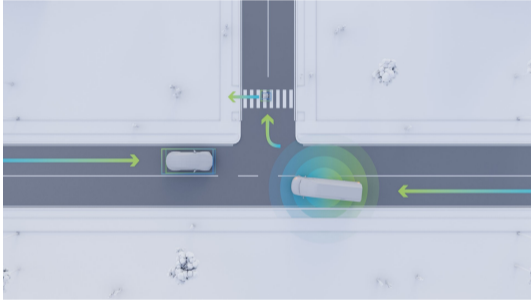
# (LTL) Model Checking



MODEL

SPECIFICATION

$\varphi = G(req \rightarrow F\ resp)$

MODEL CHECKER

"yes, the claim is true of the model!"

"no, the claim is not true of the model: here is why"

Very influential: 4 Turing Award Winners

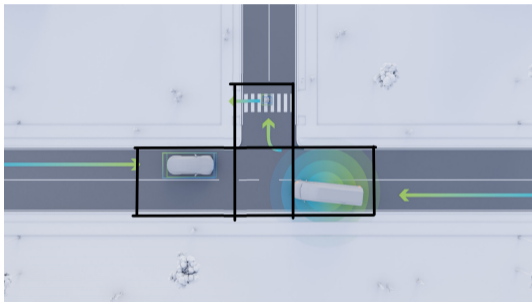| 1996 | Amir Pnueli | | For seminal work introducing temporal logic into computing science and for outstanding contributions to program and systems verification.[35] |
|------|-------------|---|---|
| 2007 | Edmund M. Clarke | | For their roles in developing model checking into a highly effective verification technology, widely adopted in the hardware and software industries.[38] |
| | E. Allen Emerson | | |
| | Joseph Sifakis | | |

## From Scenario to Model Checking



Source: https://www.digitrans.expert/en

- Two autonomous vehicles are approaching a junction.
- One is turning, the other one is going straight.
- We want: *"avoid collisions"*
- Once a collision occurs, the vehicles cannot continue their journey

Source: https://www.digitrans.expert/en
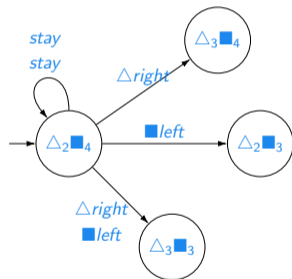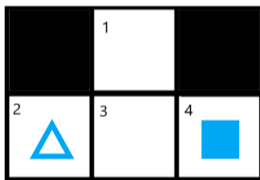
- Abstracting $\rightarrow$ discretising
- *"avoid collisions"*: **G**¬*collide*

# From Scenario to Model Checking

*"avoid collisions"*: $\mathbf{G}\neg collide$, where *collide* means $\triangle$ and $\blacksquare$ are **in the same location**
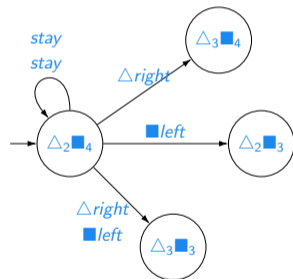
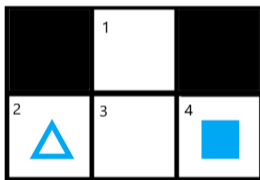$$\varphi := \mathbf{G}\neg \bigvee_{i \in \{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i)$$

# From Scenario to Model Checking

*"avoid collisions"*: **G**¬*collide*, where *collide* means △ and ■ are **in the same location**

$$\varphi := \mathbf{G}\neg \bigvee_{i \in \{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i)$$



$\varphi$ is **violated** since it is *possible* to reach the state $\triangle_3\blacksquare_3$

23

*"avoid collisions"*: **G**¬*collide*, where *collide* means △ and ■ are **in the same location**

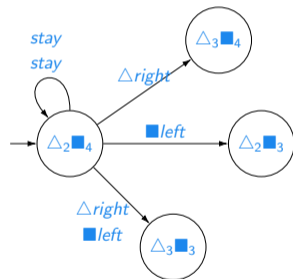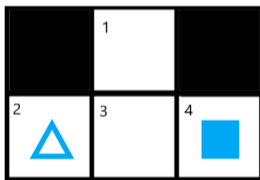$$\varphi := \mathbf{G}\neg \bigvee_{i \in \{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i)$$



$\varphi$ is **violated** since it is *possible* to reach the state $\triangle_3\blacksquare_3$

Is this **reasonable?**

**Not All Behaviours Are Equal, but Some Are More Unequal Than Others**



Source: https://www.digitrans.expert/en

- A collision is a **possible** behaviour.
- However, not a **rational** behaviour.
- The vehicles would **prefer** to **avoid** a collision: wait for the other vehicle to pass, then continue to its destination
- Classical verification is not a good/reasonable approach to check the correctness of such a scenario.

**Problem with Classical Notion of Correctness Problem**

How should we define correctness in MASs?



Classical notion of correctness ignores agents **goals/preferences**

# A New Notion of Correctness Problem

How should we define correctness in MASs?



Correctness with respect to **rational choices** of agents

## Rational Verification[2]

**Classical Verification**

Is the system correct?

$$\Downarrow$$

**Rational Verification**

Is the system correct wrt behaviours that can be **sustained by rational choices** of agents?

- Use **game theory** to model/analyse rational behaviours.
- Turn MASs into **multi-player games**.

---

[2]Alessandro Abate et al. "Rational verification: game-theoretic verification of multi-agent systems". In: *Applied Intelligence* 51.9 (2021), pp. 6569–6584.

# Why games?

- Games serves as **abstractions** for *strategic interactions* between self-interested players/agents
- Various settings: *turn-based vs concurrent, zero-sum vs general-sum, cooperative vs non-cooperative, ...*
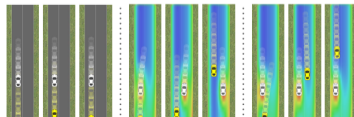- Relevant for many scenarios in autonomous/AI systems
  - e.g., zero-sum: DeepMind AlphaZero (go, chess, shogi playing), concurrent: resource sharing/allocation (server, GPU power),...
  - even autonomous vehicles

## Hierarchical Game-Theoretic Planning for Autonomous Vehicles

Jaime F. Fisac[*1]   Eli Bronstein[*1]   Elis Stefansson[2]   Dorsa Sadigh[3]   S. Shankar Sastry[1]   Anca D. Dragan[1]

*Abstract*— The actions of an autonomous vehicle on the road affect and are affected by those of other drivers, whether overtaking, negotiating a merge, or avoiding an accident. This mutual dependence, best captured by dynamic game theory, creates a strong coupling between the vehicle's planning and its predictions of other drivers' behavior, and constitutes an open problem with direct implications on the safety and viability of

28

**What is a Game?**

Ingredients:

1. Several decision makers (**the players/agents**)
2. Players have different goals (**the goals**)
3. Each player can affect the outcome for all (**the actions**)

**What is a Game?**

Ingredients:

1. Several decision makers (**the players/agents**)
2. Players have different goals (**the goals**)
3. Each player can affect the outcome for all (**the actions**)

### Game theory

the methodology of using mathematical tools to model and analyse situations of interactive decision making.

## How to model rational behaviours?

- What kind of behaviour is **rational**?
- Game theory proposes many "solution concepts", i.e., a formal rule for 'predicting' how a game will be played
- The most influential is **Nash equilibrium**: Nobel prize in Economics 1994

# How to model rational behaviours?

- What kind of behaviour is **rational**?
- Game theory proposes many "solution concepts", i.e., a formal rule for 'predicting' how a game will be played
- The most influential is **Nash equilibrium**: Nobel prize in Economics 1994

## How to model rational behaviours?

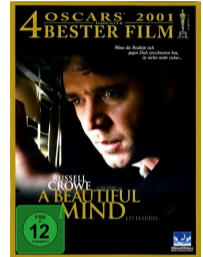- What kind of behaviour is **rational**?
- Game theory proposes many "solution concepts", i.e., a formal rule for 'predicting' how a game will be played
- The most influential is **Nash equilibrium**: Nobel prize in Economics 1994



### Nash equilibrium

A situation where no player in a game would want to change their strategy, while keeping the other players' strategies constant

# From Scenario to Game: T-Junction Game



- the players: $\triangle, \blacksquare$
- the goals:
  - Player $\triangle$ wants to go straight: $\gamma_\triangle := \mathbf{F}\triangle_4$
  - Player $\blacksquare$ wants to turn: $\gamma_\blacksquare := \mathbf{F}\blacksquare_1$
- the actions: players can move to adjacent locations

$$\varphi := \mathbf{G}\neg \bigvee_{i \in \{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i) \quad \gamma_\triangle := \mathbf{F}\triangle_4 \quad \gamma_\blacksquare := \mathbf{F}\blacksquare_1$$



$\triangle$ moves: right, right and $\blacksquare$ moves: left, up

Not a NE, since (for example) $\triangle$ can stay put and wait for $\blacksquare$ to go up, then proceed to move right, right

# From Verification to Rational Verification

# From Verification to Rational Verification



MODEL

SPECIFICATION

$\varphi = G(req \rightarrow F\ resp)$

MODEL CHECKER

"yes, the claim is true of the model"

"no, the claim is not true of the model: here is why"

MODEL

QUERY

$\varphi = G(req \rightarrow F\ resp)$

PLAYER PREFERENCES

RATIONAL MODEL CHECKER

"the claim is true in some equilibrium"

"no, the claim does not hold in any equilibrium

**Rational Verification: Decision Problems**

- **Safety:** **all** stable outcomes (e.g., NE) do **not violate** a desirable property $\varphi$ (A-NASH)
- **Liveness:** there **exists** a stable outcome that satisfies a desirable property $\varphi$ (E-NASH)
- **Stability:** Is there any stable outcome? (NON-EMPTINESS)

# Rational Verification Tool: EVE

- Equilibrium Verification Environment (EVE)[3]
- Automata-theoretic techniques
- Support memoryful strategies; players can fully implement LTL goals
- EVE online: http://eve.cs.ox.ac.uk/



| Home | EVE | Download | People | Publications | Others |

**Welcome to EVE Website**

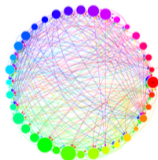EVE (Equilibrium Verification Environment) is a formal verification tool for the automated analysis of temporal equilibrium properties of concurrent and multi-agent systems represented as multi-player games. Systems are modelled using the Simple Reactive Module Language (SRML) as a collection of independent system components (players/agents in a game), which are assumed to have goals expressed using Linear Temporal Logic (LTL) formulae. In particular, EVE checks for the existence of Nash equilibria in such systems and can be used to do rational synthesis and verification automatically.

---

[3] Julian Gutierrez et al. "Automated temporal equilibrium analysis: Verification and synthesis of multi-player games". In: *Artificial Intelligence* (2020).

## Tractability

- Decision Problems (A/E-Nash, Non-Emptiness) with LTL are expensive: 2EXPTIME
- What can we do to improve?
- Use different goals and properties: GR(1) and mean-payoff value

# GR(1)[6]

The language of *General Reactivity of rank 1*, denoted GR(1), is the fragment of LTL of formulae written in the following form:

$$(\mathbf{GF}\psi_1 \wedge \ldots \wedge \mathbf{GF}\psi_m) \rightarrow (\mathbf{GF}\varphi_1 \wedge \ldots \wedge \mathbf{GF}\varphi_n),$$

each $\psi_i$ and $\varphi_i$ is a Boolean combination of atomic propositions.

$$(\mathbf{GF}req_1 \wedge \mathbf{GF}req_2) \rightarrow \mathbf{GF}ack$$

GR(1) synthesis has been used for controllers of ground robots[4], UAVs[5].

---

[4]Hadas Kress-Gazit, Georgios E. Fainekos, and George J. Pappas. "Where's Waldo? Sensor-Based Temporal Logic Motion Planning". In: *ICRA*. 2007.
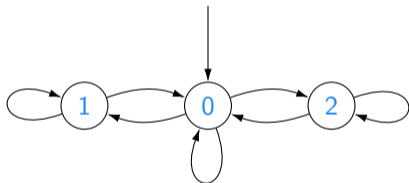
[5]Thomas B. Apker, Benjamin Johnson, and Laura Humphrey. "LTL Templates for Play-Calling Supervisory Control". In: *AIAA Infotech @ Aerospace*. 2016.

[6]Roderick Bloem et al. "Synthesis of Reactive(1) designs". In: *J. Comput. Syst. Sci.* 78.3 (2012), pp. 911–938.

# Mean-payoff value

For an infinite sequence $\beta \in \mathbb{R}^\omega$ of real numbers, let $\mathrm{mp}(\beta)$ be the *mean-payoff* value of $\beta$, defined as follows:

$$\mathrm{mp}(\beta) = \lim_{n \to \infty} \inf \frac{1}{n} \sum_{i=0}^{n-1} \beta[i]$$



$$\beta^1 = 00000000000\ldots \qquad \mathrm{mp}(\beta^1) = 0$$
$$\beta^2 = 01010101010\ldots \qquad \mathrm{mp}(\beta^2) = 0.5$$
$$\beta^3 = 01020102010\ldots \qquad \mathrm{mp}(\beta^3) = 3/4$$

# Cases

**E-Nash**

Given: Game $\mathcal{G}$, temporal property $\varphi$.

Quest: Is there any Nash Equilibrium $\vec{\sigma}$ in $\mathcal{G}$ such that $\pi(\vec{\sigma}) \models \varphi$?

|  | $\gamma_i$ | $\varphi$ | E-NASH |
|---|---|---|---|
|  | LTL | LTL | 2EXPTIME-complete |
| GR(1) games | GR(1) | LTL | ? |
|  | GR(1) | GR(1) | ? |
| mp games | mp | LTL | ? |
|  | mp | GR(1) | ? |

## Complexity Results

| $\gamma_i$ | $\varphi$ | E-Nash |
|-----|-----|-----|
| LTL | LTL | 2EXPTIME-complete |
| GR(1) | LTL | PSPACE-complete |
| GR(1) | GR(1) | FPT |
| mp | LTL | PSPACE-complete |
| mp | GR(1) | NP-complete |

- Non-Emptiness (E-Nash when $\varphi = \top$):
  - LTL games: 2EXPTIME-complete
  - GR(1) games: FPT
  - mp games: NP-complete
- A-Nash: 2EXPTIME, PSPACE, FPT, PSPACE, coNP.

# Part II: Modification

$$\varphi := \mathbf{G}\neg \bigvee_{i\in\{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i) \quad \gamma_\triangle := \mathbf{F}\triangle_4 \quad \gamma_\blacksquare := \mathbf{F}\blacksquare_1$$

$\triangle$ moves: right, right and $\blacksquare$ stays in 4 forever

this is a NE, but a bad one: nobody achieves their goal

# Dealing with missing or bad equilibria

### Problem
Individually rational choices can cause outcomes that are highly undesirable, *e.g.*, there is no equilibrium or the temporal specification is not satisfied.

### Question
The problem with this is intrinsic in the system. Can we modify it in order to gain (desirable) equilibria?
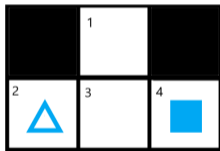
### Solution
Equilibrium Design: redesign the game such that individually rational behaviour leads to desired outcomes.

## Modifying Games

- Norms
- Modify goals
- Provide incentives

# Modifying Games

$$\varphi := \mathbf{G}\neg \bigvee_{i\in\{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i) \quad \gamma_\triangle := \mathbf{F}\triangle_4 \quad \gamma_\blacksquare := \mathbf{F}\blacksquare_1$$



- Introduce a norm: $\blacksquare$ cannot stay in the same place for 2 consecutive time steps
- Modify the goal: $\gamma_\triangle := \mathbf{F}\triangle_4 \wedge \mathbf{X}\neg\triangle_3$
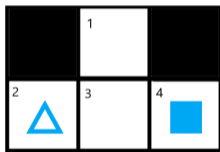
$$\varphi := \mathbf{G}\neg \bigvee_{i \in \{1,2,3,4\}} (\triangle_i \wedge \blacksquare_i) \quad \gamma_\triangle := \mathbf{F}\triangle_4 \quad \gamma_\blacksquare := \mathbf{F}\blacksquare_1$$
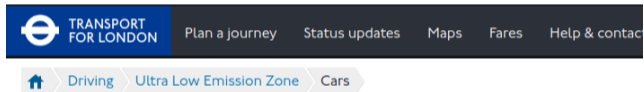


- Introduce a norm: $\blacksquare$ cannot stay in the same place for 2 consecutive time steps
- Modify the goal: $\gamma_\triangle := \mathbf{F}\triangle_4 \wedge \mathbf{X}\neg\triangle_3$

Every NE satisfies $\varphi$

# Modification via Incentives

- Sometimes, designer cannot prohibit actions (e.g., according to some laws)
- Designer can only incentivise players to take/avoid some actions



**Cars**

Cars need to meet minimum emissions standards when travelling within the Ultra Low Emission Zone (ULEZ) or the daily £12.50 charge must be paid.

# Equilibrium Design via Incentives[7]

Given a mean-payoff game $\mathcal{G}$, a temporal specification $\varphi$ and a budget $\beta \in \mathbb{N}$

**Definition (Weak Implementation)**

find an incentive scheme $\kappa$ with $\mathrm{cost}(\kappa) \leq \beta$ such that $(\mathcal{G}, \kappa, \varphi)$ solves E-NASH positively.

**Definition (Strong Implementation)**

find an incentive scheme $\kappa$ with $\mathrm{cost}(\kappa) \leq \beta$ such that $(\mathcal{G}, \kappa, \varphi)$ solves A-NASH positively.

---

[7] Julian Gutierrez et al. "Equilibrium Design for Concurrent Games". In: *CONCUR*. 2019, 22:1–22:16.

# Complexity results

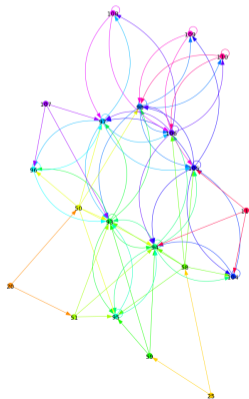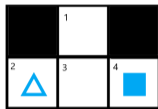|  | LTL Spec. | GR(1) Spec. |
|---|---|---|
| Weak Implementation | PSPACE-complete | NP-complete |
| Strong Implementation | PSPACE-complete | $\Sigma_2^P$-complete |

# Part III: Beyond

# From Verification to Explanation



- Equilibria can be complex and opaque
- There are 22 states and 82 edges
- Not easy to understand
- ☺: "Why do we have to wait? It's wasting my time.[a]"
- △: "We have to wait to avoid crashing to another car."
- ☺: "But why us, it's *unfair*![b]"
- △: "This is the most reasonable choice, because..."
- ☺: "OK!"

---

[a]Value alignment problem.
[b]Fairness problem.

**From Verification to Explanation**

Explainability and Transparency

- Is the agent's goal aligned with user's?
- If so, how can we extract (synthesise) strategies and present them in a human-friendly way?
- Are the strategies fair?

## Other avenues

- Decision Problems with LTL are expensive: 2EXPTIME
  **Statistical methods:** can these make it more practical? E.g., model checking with the Monte Carlo method[8]

- **Learning agents**: What if the players use some learning element, e.g., reinforcement learning?[9]

- **Privacy & security**: So far the setting has been *perfect information*. What if this is not a viable setting? For instance, we might not want other vehicles to know our home address.

---

[8] Radu Grosu and Scott A Smolka. "Monte carlo model checking". In: *TACAS*. 2005.

[9] Lewis Hammond et al. "Multi-Agent Reinforcement Learning with Temporal Logic Specifications". In: *AAMAS*. 2021.

## Conclusion

- The future looks increasingly more and more multi-agent
- Want and need these multi-agent systems to be safe and correct
- Verification of Multi-Agent Systems
  - A new and more appropriate notion of correctness: rational verification
  - Modelling systems as games
  - Tool: EVE
- Challenges
  - Practicality and scalability
  - Incorporating agents who learn
  - How to ensure privacy and security?
  - How to make decisions transparent to human?

# Thank you!