

Introduction to Game Theory and Multi-Agent Systems

Agnieszka Mensfelt

Kostas Stathis

Vince Tencsenyi

<https://dicelab-rhul.github.io/Strategic-AI-Autoformalization>

ESSAI, 30/06/25





Outline

Introduction

Game Theory

Multi-agent
Systems

1 Introduction

2 Foundations of Strategic Interaction

3 Multi-agent Systems

Introduction



Motivation

Why strategic interactions matter in AI

Introduction

Game Theory

Multi-agent
Systems

Many applications are naturally multi-agent:



Human teams and
companies



Markets and economies



Transportation networks



Distributed software
systems



Communication networks



Robotic teams

How should agents act in the presence of other agents?



Strategic interaction

Why reasoning is essential

Introduction

Game Theory

Multi-agent
Systems

One way to answer this question is by expecting agents to think strategically.

- ▶ **Interdependent Decisions:** Agent choices affect and depend on others — strategic reasoning enables intelligent interaction.
- ▶ **Lack of information:** Strategic reasoning allows agents to model beliefs about others and act under information that is incomplete or imperfect.
- ▶ **Dynamic environment:** Strategic reasoning allows agents respond effectively to these ongoing changes.
 - **Collaboration vs Competition**

Key Message

To build intelligent and interactive agents, strategic reasoning is indispensable.

Games as Microcosms of Intelligence

Introduction

Game Theory

Multi-agent
Systems

Games provide a natural setting for studying and evaluating strategy.

- ▶ Success in games captures key elements of **strategic reasoning**:
 - **Problem-solving** in dynamic, rule-based environments
 - **Plan** sequences of actions to achieve long-term objectives
 - **Adaptation** to opponents' tactics and evolving situations
 - **Learning** from feedback to improve future decisions
- ▶ Quantifiable performance in games enables a **comparative assessment of strategies** and, as a result, **intelligence**



Why Game Playing Matters to AI Research

Game Playing as a Testbed for strategic AI

Introduction

Game Theory

Multi-agent
Systems

- ▶ Games have long been used to test AI:



Deep Blue (Chess)



AlphaGo (Go)

- ▶ These systems were successful but **not general** - they only played one game.
- ▶ **Can we build AI that plays *any* game?**

General Game Playing for Strategic AI

Game Playing as a Testbed for Strategic AI

Introduction

Game Theory

Multi-agent
Systems

- ▶ **Goal:** Build agents capable of **General Game Playing (GGP)** - excel in multiple games from their formal rules alone.
- ▶ **Key Idea:** Use a shared, declarative formalism for representing a game - the **Game Description Language (GDL)**.
- ▶ **Approach:**
 - Provide agents with new game rules expressed in GDL, without additional code or training.
 - When encountering a new game:
 - The agent parses and understands its rules and objectives.
 - It uses this understanding to plan, adapt, and compete strategically.
 - Shift away from game-specific logic toward **domain-general intelligence**.

Foundations of Strategic Interaction

Game Theory Foundations

Introduction

Game Theory

Multi-agent
Systems

A *game* models strategic interactions between decision-makers (players) where outcomes are interdependent on players' choices.¹

¹R. B. Myerson, "An introduction to game theory," Northwestern University, Center for Mathematical Studies in Economics and Management Science, Discussion Papers 623, 1984.

Games: Definition and Components

Introduction

Game Theory

Multi-agent
Systems

$$G = (N, \{A_i\}_{i \in N}, \{U_i\}_{i \in N})$$

- **Players** ($N = \{1, \dots, n\}$): Strategic decision-makers involved in the game.
- **Actions** (A_i): The set of immediate choices available to player i .
- **Strategy Profile** ($o = (a_1, \dots, a_n)$): A tuple representing the selected strategy of each player – the game's *outcome*.
- **Utility Functions** ($U_i : A_1 \times \dots \times A_n \rightarrow \mathbb{R}$): An ordinal/cardinal representation of player i 's preferences over outcomes.



Strategies

Introduction

Game Theory

Multi-agent
Systems

A strategy is a complete specification of play; a plan for a course of actions.

- ▶ A **pure strategy** for a player i is a single deterministic action $a_i \in A_i$.
- ▶ A **mixed strategy** is a probability distribution over pure strategies $\sigma_i \in \Delta(A_i)$, $\sum_{a_i \in A_i} \sigma_i(a_i) = 1$, where $\Delta(A_i)$ is the set of all probability distributions over A_i .



Representative Forms

Introduction

Game Theory

Multi-agent
Systems

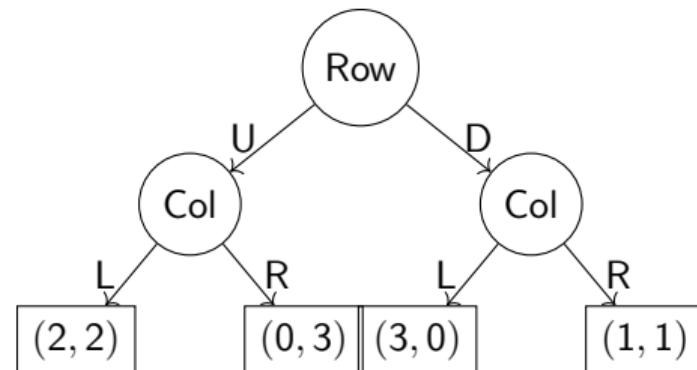
Normal Form

- Captures simultaneous actions
- Captures strategic structure
- $(\pi_{\text{row}}, \pi_{\text{col}})$ are the payoffs of an o

	L	R
U	(2,2)	(0,3)
D	(3,0)	(1,1)

Extensive Form

- Captures sequential moves
- Captures temporal/hierarchical structure



Non-cooperative Games^{2,3}

Introduction

Game Theory

Multi-agent
Systems

Non-Cooperative Games

- ▶ Players act independently
- ▶ Focus on individual strategic behaviour
- ▶ Communication and binding agreements disallowed

Cooperative Games

- ▶ Players can form coalitions
- ▶ Focus on joint gains and payoff division
- ▶ Assumes binding agreements

²J. F. Nash, "Non-cooperative games," in *The Foundations of Game Theory* Vol 4, Routledge, 2024, pp. 329–340.

³A. Rubinstein, H. W. Kuhn, O. Morgenstern, et al., *Theory of Games and Economic Behavior: 60th Anniversary Commemorative Edition*. Princeton university press, 2007.



Non-cooperative Participatory Game

Introduction

Game Theory

Multi-agent
Systems

A group of hunters awaits a stag known to follow a certain path. If they all cooperate, they can kill it and feast. If they act alone or are discovered, the stag flees, and all go hungry.

Hours pass without a sign. A day goes by. The stag is not guaranteed to come today, but the hunters believe it will. Suddenly, a hare appears.

If one hunter breaks cover to catch the hare, he eats – but ruins the chance at the stag. The others starve.

Do you shoot the Hare or wait for the Stag?

menti.com : 7649 2610



Payoff Structures

Introduction

Game Theory

Multi-agent
Systems

Matching Pennies

	Head	Tails
Head	(1, -1)	(-1, 1)
Tails	(-1, 1)	(1, -1)

Stag Hunt

	Stag	Hare
Stag	(5, 5)	(0, 3)
Hare	(3, 0)	(1, 1)

Prisoner's Dilemma

	C	D
C	(3, 3)	(0, 5)
D	(5, 0)	(1, 1)

Battle of the Sexes

	Ballet	Football
Ballet	(2, 1)	(0, 0)
Football	(0, 0)	(1, 2)



Payoff Structures

Introduction

Game Theory

Multi-agent
Systems

Zero-Sum Game

	L	R
U	(A, -A)	(-C, C)
D	(-B, B)	(D, -D)

Non-Zero-Sum Game

	L	R
U	(A, A)	(C, C)
D	(B, B)	(D, D)

Symmetric Payoffs

	L	R
U	(A, A)	(C, C)
D	(B, B)	(D, D)

Asymmetric Payoffs

	L	R
U	(A _i , A _j)	(C _i , C _j)
D	(B _i , B _j)	(D _i , D _j)

Solution Concepts: Dominant Strategies

Introduction

Game Theory

Multi-agent
Systems

A **dominant strategy** is one that always yields a higher payoff, regardless of what the opponent does.

Strict dominance:

$a_i \in A_i$ is strictly dominant for player i if

$$\forall a_{-i} \in A_{-i}, \forall a'_i \neq a_i : \\ U_i(a_i, a_{-i}) > U_i(a'_i, a_{-i})$$

Prisoner's Dilemma

	C	D
C	(3, 3)	(0, 5)
D	(5, 0)	(1, 1)

Stag Hunt

	S	H
S	(4, 4)	(0, 3)
H	(3, 0)	(3, 3)

Solution Concepts: Nash Equilibrium

Introduction

Game Theory

Multi-agent
Systems

A **Nash equilibrium** is a strategy profile where no player can benefit by unilaterally deviating.

(a_1^*, \dots, a_n^*) is a Nash equilibrium if $\forall i \in N$:

$$U_i(a_i^*, a_{-i}^*) \geq U_i(a_i, a_{-i}^*) \quad \forall a_i \in A_i$$

Prisoner's Dilemma

	C	D
C	(3, 3)	(0, 5)
D	(5, 0)	(1, 1)

Stag Hunt

	S	H
S	(4, 4)	(0, 3)
H	(3, 0)	(3, 3)

Solution Concepts: Pareto Efficiency

Introduction

Game Theory

Multi-agent
Systems

An outcome is **Pareto efficient** if no player can be made better off without making another worse off.

$o \in A_1 \times \cdots \times A_n$ is Pareto efficient if

$\nexists o'$ such that $U_i(o') \geq U_i(o) \ \forall i$
 and $\exists j, \ U_j(o') > U_j(o)$

Prisoner's Dilemma

	C	D
C	(3, 3)	(0, 5)
D	(5, 0)	(1, 1)

Stag Hunt

	S	H
S	(4, 4)	(0, 3)
H	(3, 0)	(3, 3)

Solution Concepts: Social Welfare

Introduction

Game Theory

Multi-agent
Systems

Social welfare is the total utility across all players.

An outcome is *socially optimal* if it maximizes:

$$o^* = \arg \max_{o \in A_1 \times \dots \times A_n} \sum_{i \in N} U_i(o)$$

Prisoner's Dilemma

	C	D
C	(3, 3)	(0, 5)
D	(5, 0)	(1, 1)

Stag Hunt

	S	H
S	(4, 4)	(0, 3)
H	(3, 0)	(3, 3)



Participatory Guessing Game

Introduction

Game Theory

Multi-agent
Systems

2 person game instruction:

"Your task is to guess how many balls are hidden in this black box. The box can be empty and can hold up to 100 balls. However, there is a twist! The winning guess is not the one closest to the actual number of balls hidden in the box. The closest guess to the two players' choices' mean multiplied by 0.66 will be the winner."

menti.com : 3613 7700

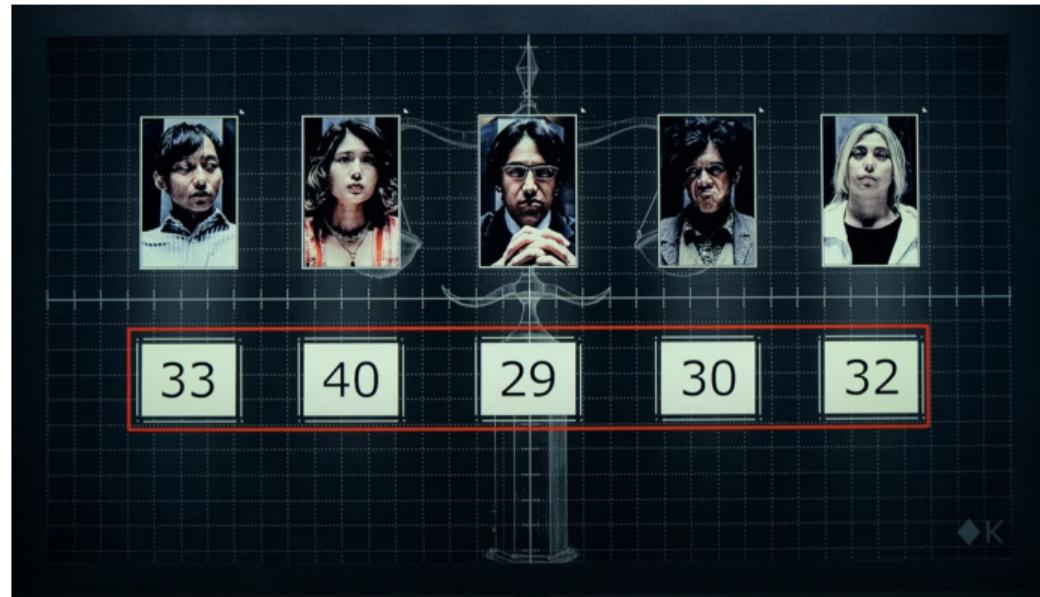


Strategic Reasoning in Guessing Games

Introduction

Game Theory

Multi-agent
Systems



Strategic Reasoning in the p -Beauty Contest

Introduction

Game Theory

Multi-agent
Systems

The p -Beauty Contest^{4,5}

- ▶ Each player selects a number in a known interval (e.g., $[0, 100]$).
- ▶ The winner is the player whose guess is closest to p times the average of all guesses.

Why is it interesting?

- ▶ Models *strategic anticipation* through **recursive reasoning**:
 - “What others think I will think...”
 - Formally captured by k-level theory
- ▶ No strictly dominant strategy
- ▶ Nash equilibrium: all players choose 0
- ▶ Special case⁶ with $n = 2$

⁴J. M. Keynes, “The general theory of employment,” *The quarterly journal of economics*, vol. 51, no. 2, pp. 209–223, 1937.

⁵R. Nagel, “Unraveling in guessing games: An experimental study,” *The American Economic Rev.*, vol. 85, no. 5, pp. 1313–1326, 1995.

⁶B. Grosskopf and R. Nagel, “The two-person beauty contest,” *Games and Economic Behavior*, vol. 62, no. 1, pp. 93–99, 2008.



Information Structures

Introduction

Game Theory

Multi-agent
Systems

Complete vs Incomplete Information

- ▶ Can players fully observe others' payoffs?

Perfect vs Imperfect Information

- ▶ Can players fully observe others' actions?

Symmetric vs Asymmetric Information

- ▶ Can players access the same quality/quantity of information?

Multi-agent Systems

Intelligent Agents and Strategic Reasoning

Introduction

Game Theory

Multi-agent
Systems

Real-world conflicts entail practical scenarios, where societies of actors are involved in cooperative or competitive interactions. For such cases, agent-based approaches are deemed a natural metaphor⁷.

⁷M. Wooldridge, *An introduction to multiagent systems*. John Wiley & Sons, 2009.



Agents and the Environment

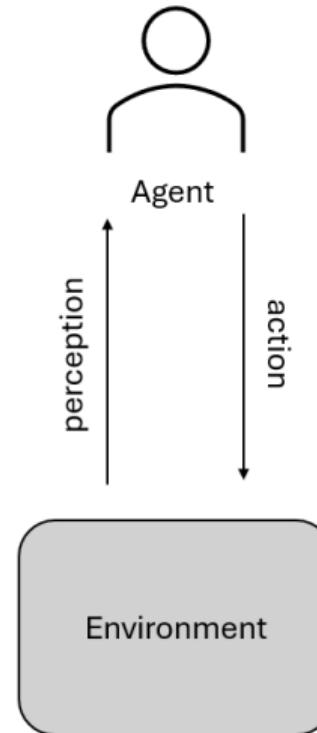
Introduction

Game Theory

Multi-agent
Systems

What is an agent?

An agent perceives the environment and produces actions that affect it.





Agents and the Environment

Introduction

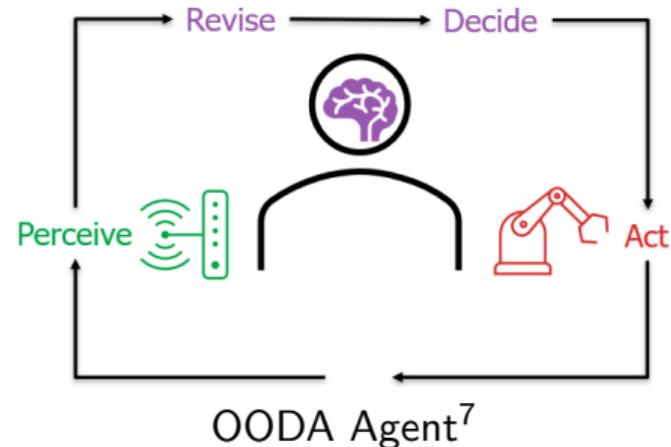
Game Theory

Multi-agent
Systems

What is an agent?

Agents are characterized by⁶:

- ▶ Autonomy (operating without direct intervention)
- ▶ Social ability (communicating with other agents)
- ▶ Reactivity (responding to environmental changes)
- ▶ Proactivity (exhibiting goal-directed behaviour)



⁶M. Wooldridge and N. R. Jennings, "Intelligent agents: Theory and practice," *The Knowledge Engineering Review*, vol. 10, no. 2, pp. 115–152, 1995

⁷V. Trencsenyi, A. Mensfelt, and K. Stathis, "Approximating human strategic reasoning with IIM-enhanced recursive reasoners leveraging multi-agent hypergames," *arXiv preprint arXiv:2502.07443*, 2025



Introduction

Game Theory

Multi-agent
Systems

How do agents work internally?

- ▶ **Reactive Architectures** Rely on simple stimulus-response mechanisms without internal models or planning.
- ▶ **Deliberative Architectures** Maintain internal representations, and support reasoning, memory, and planning, often inspired by cognitive processes.
- ▶ **Hybrid Architectures** Combine reactive and deliberative layers to balance reactive responsiveness with symbolic representations and goal-driven behaviour.

BDI Agents: Human-like Practical Reasoning

Introduction

Game Theory

Multi-agent
Systems

Belief–Desire–Intention (BDI) agents model human-like deliberative processes by decoupling reasoning into:

- **Beliefs (B)**: informational state (agent's internal model of the world)
- **Desires (D)**: motivational state (goals the agent might want to achieve)
- **Intentions (I)**: deliberative commitments to plans of action

BDI architecture functions:

- \mathcal{B} : Belief revision from new percepts
- \mathcal{O} : Option generation based on B, I
- \mathcal{F} : Intention filtering from B, D, I
- \mathcal{P} : Planning over B, I to action sequence Π

BDI agents cycle through these functions to continuously update beliefs, deliberate, and act — paralleling human practical reasoning.



Are LLMs Agents?

Introduction

Game Theory

Multi-agent
Systems

Are LLMs agents?

What kind of agent are they?





Are LLMs agents?

Introduction

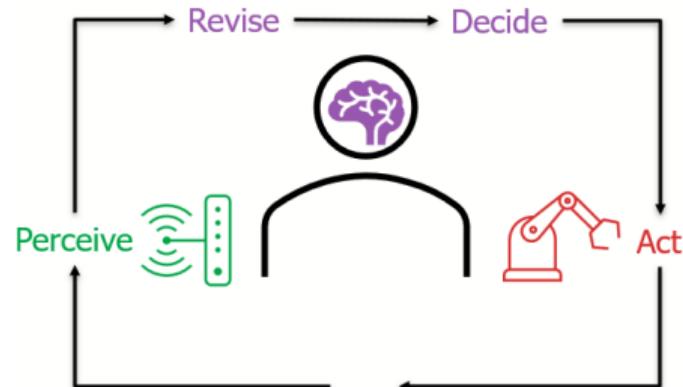
Game Theory

Multi-agent
Systems

An agent perceives the environment and produces actions that affect it.

LLM-as-agent⁸ builds on this weaker notion, where the prompt is a perception and the response is the action.

LLM-as-mind⁹



⁸X. Liu, H. Yu, H. Zhang, et al., "Agentbench: Evaluating llms as agents," *arXiv preprint arXiv:2308.03688*, 2023

⁹V. Trencsenyi, A. Mensfelt, and K. Stathis, "Approximating human strategic reasoning with llm-enhanced recursive reasoners leveraging multi-agent hypergames," *arXiv preprint arXiv:2502.07443*, 2025

Multi-Agent Systems

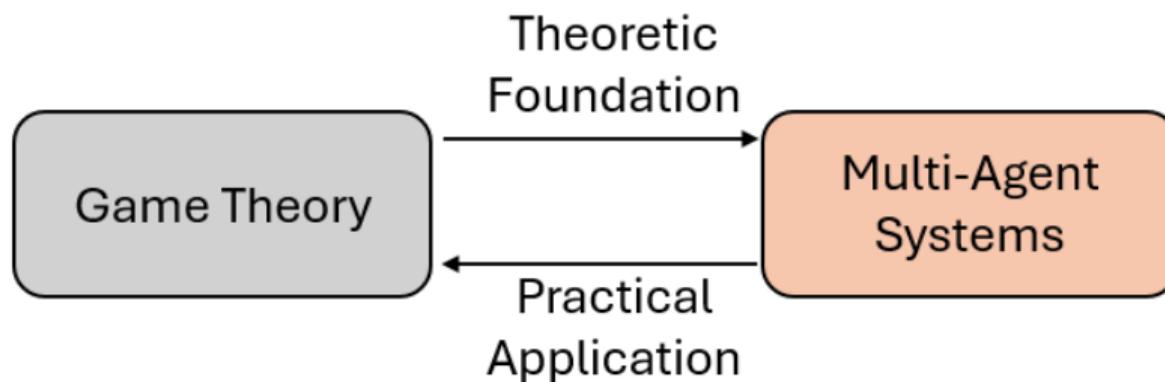
Introduction

Game Theory

Multi-agent
Systems

Multi-Agent Systems (MAS) are societies of autonomous agents interacting in a shared environment.

- ▶ Each agent pursues individual or shared goals
- ▶ Agents coordinate, compete, and communicate
- ▶ Strategic interactions naturally emerge



Uncertainty in MAS

Introduction

Game Theory

Multi-agent
Systems

Realistic MAS must handle uncertainty beyond ideal game-theoretic assumptions.

- **Environmental uncertainty:** Partial observability, dynamic environments, non-stationarity
- **Action and Strategic uncertainty:** Non-deterministic effects of actions, unknown opponent strategies
- **Bounded rationality and belief misalignment:** Cognitive limits, mismatched models of the world or other agents