# Master
# Reinforcement Learning 2022
# Lecture 7:
# Multi-Agent

Aske Plaat

# Different Approaches

- Model-free

  - Value-based [2,3]

  - Policy-based [4]

- Model-based

  - Learned [5]

  - Perfect; Two-Agent [6]

- Multi-agent [7]

- Hierarchical Reinforcement Learning (Sub-goals) [8]

- Meta Learning [9]

# Motivation

# Motivation

- Real-world decision making: model interaction

- Poker

- StarCraft

- Football

# Overview

- 1: Competitive

- 2: Cooperative

- 3: Mixed


- 1: CFR

- 2: Centr/Decentr, Opponent

- 3: Evo, Swarm, Population Based Teams


- 1: Poker

- 2: Hide and Seek
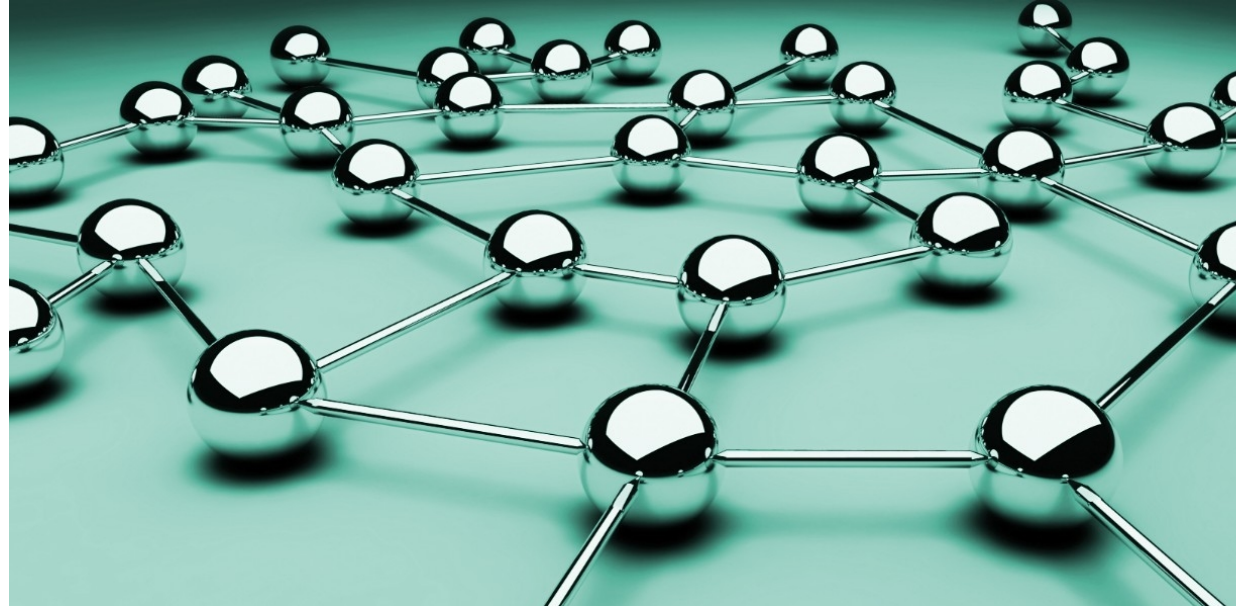
- 3: Capture The Flag

# Social Behavior

- Modeling competition; egoism

- Modeling cooperation; altruism

- Emergent social behavior
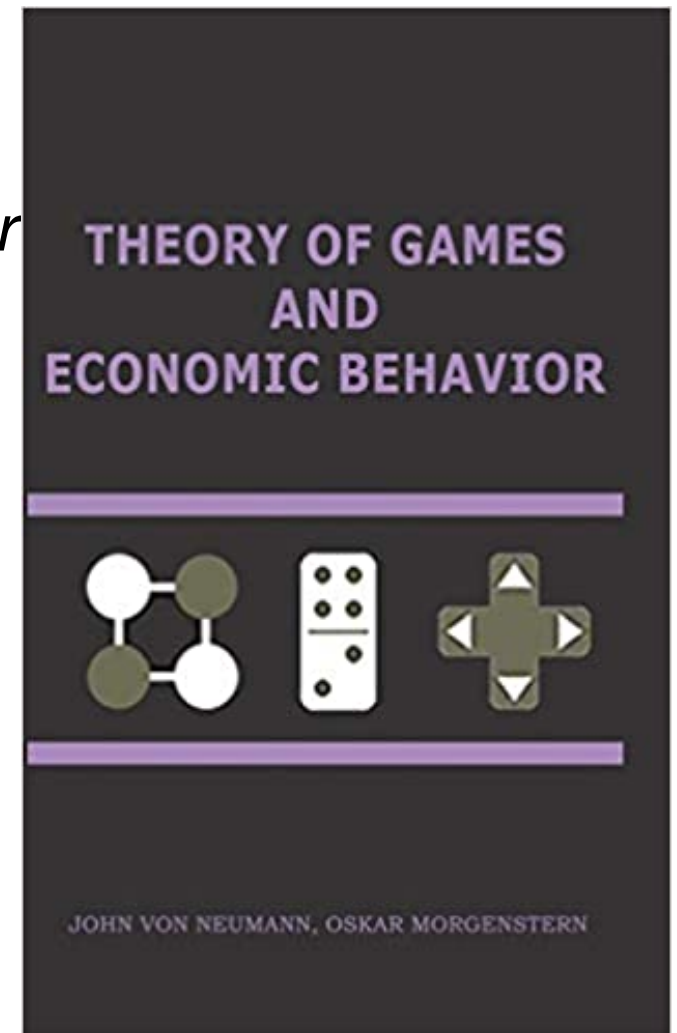
# Related Fields



- Multi-agent Systems

- Swarm computing, Evolutionary Algorithms

- Complex Networks, Real World Networks

# Multi-agent problems

# Game Theory

- Von Neumann & Morgenstern, 1944
  *Theory of Games and Economic Behavior*
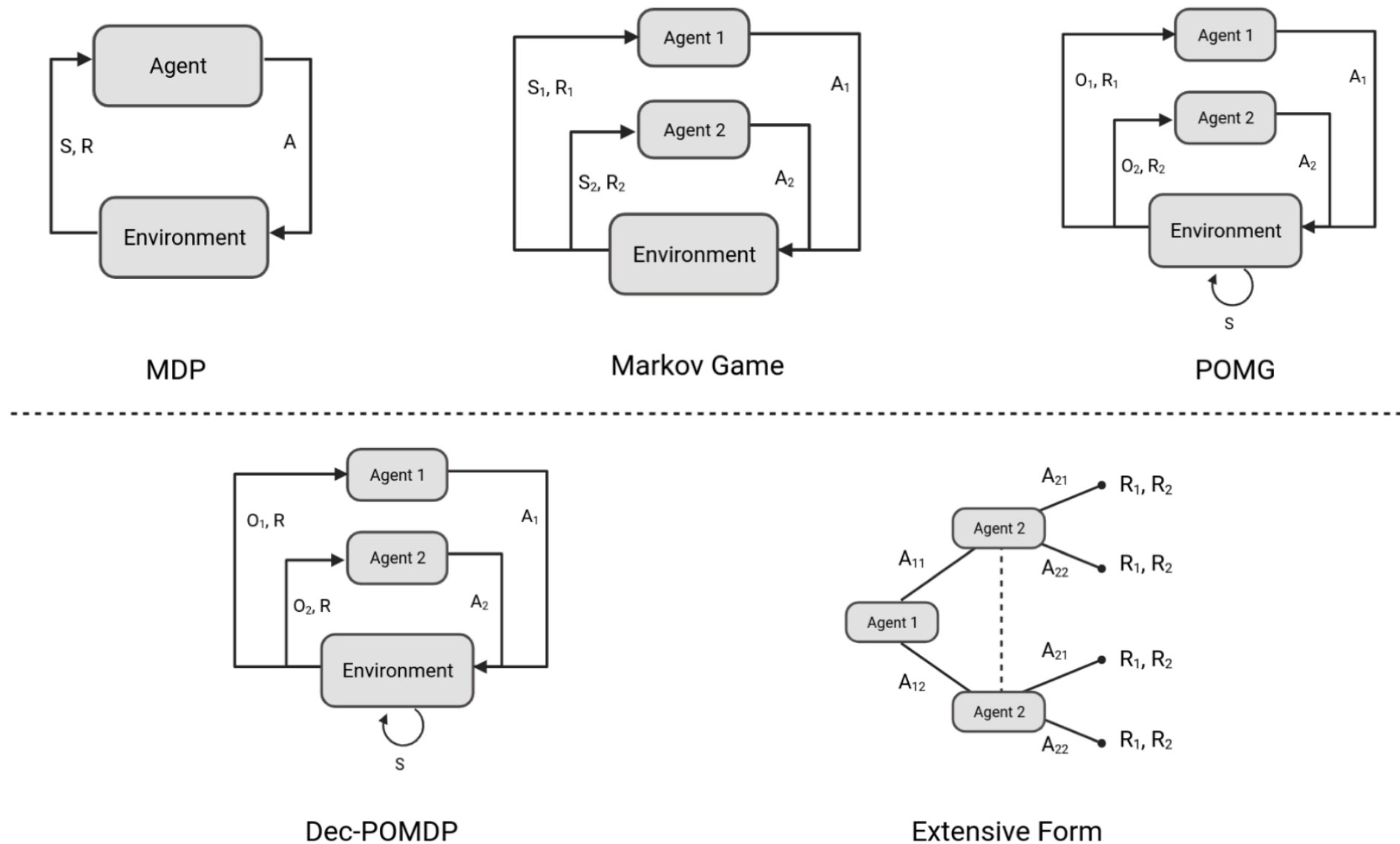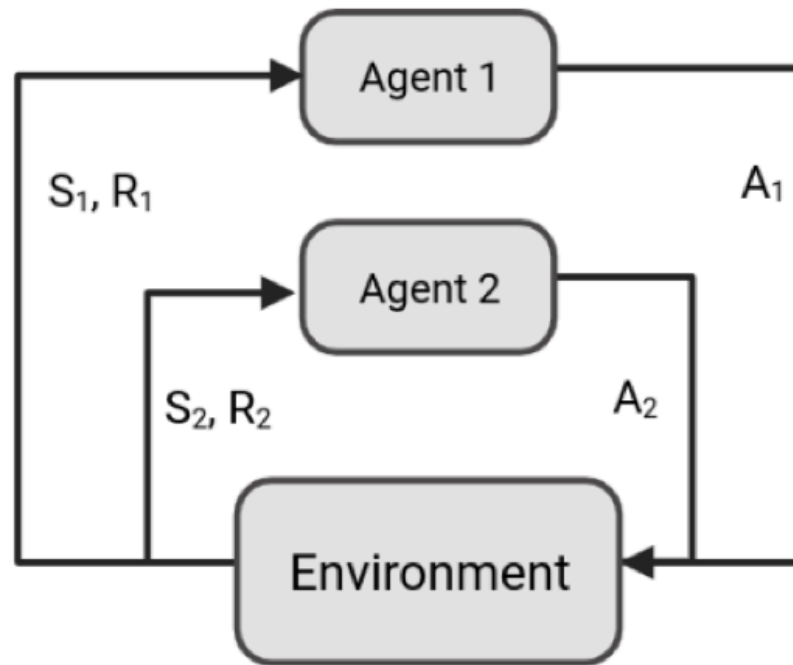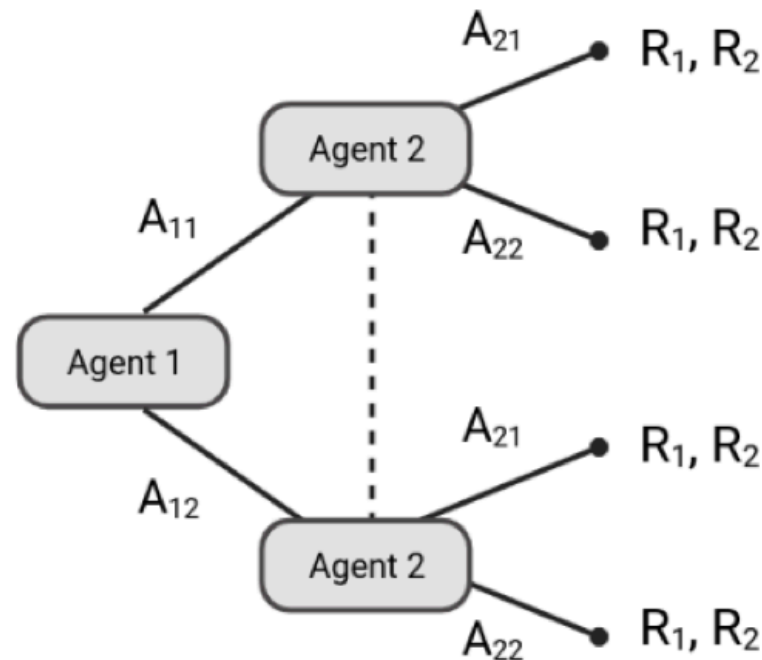
- MDP

- Partial information: POMDP

**Fig. 2 Visual depiction of the main problem representations in multiagent reinforcement learning** The MDP is the primary framework used in the single-agent setting. An agent is in some state $S$, performs action $A$, and receives a reward $R$ from the environment. In partially observable environments, the agent cannot view the true state $S$ and receives an observation $O$ instead. For simplicity, all figures display the interaction between two agents $i = 1, 2$ but can be extended to more agents.

# Stochastic Games



- Stochastic Games

- Markov Games

# Extensive-form Games



- Imperfect information games

- Possible outcomes: information set

# Competition

# Competition



- Zero sum; win/loss

- John Nash:

  The Nash equilibrium is point $\pi^\star$ from which in a non-collaborative setting none of the agents has any incentive to deviate.

- It is the optimal competitive strategy; each agent chooses best actions for themselves assuming others do the same

# Nash equilibrium

- "Multi-agent minimax"

- The Nash-policy for an agent is its best-response strategy

- It is guaranteed to do no worse than tie against any opponent strategy

- For games of imperfect information the Nash equilibrium is an expected outcome

# Nash equilibrium

Firm B

|  | Hold down output | Increase output |
|---|---|---|
| **Hold down output** | A gets $1,000<br>B gets $1,000 | A gets $200<br>B gets $1,500 |
| **Increase output** | A gets $1,500<br>B gets $200 | A gets $400<br>B gets $400 |

Firm A

← **Nash**

# Counterfactual Regret Minimization

- Multi-agent, partial information, competition

- Algorithm: Counterfactual regret minimization

- Minimize the regret of not having taken the right action, playing many "what-ifs" (counterfactuals)

- CFR is probabilistic multi-agent version of competitive minimax

- Works quite well in Poker
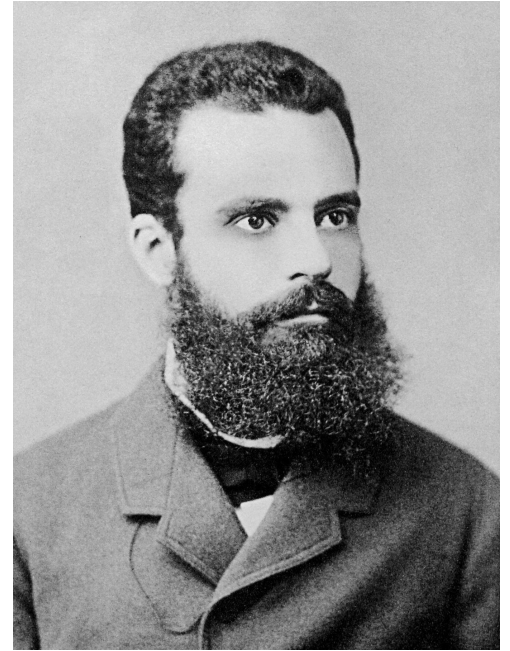
- Complicated code, see paper

# Poker

# Pluribus

Science

AAAS

PLYR 6

PLYR 2

76%

PLYR 1

PLYR 5

FOLD

PLYR 3

PLYR 4

CALLING OUR BLUFF

AI masters multiplayer poker
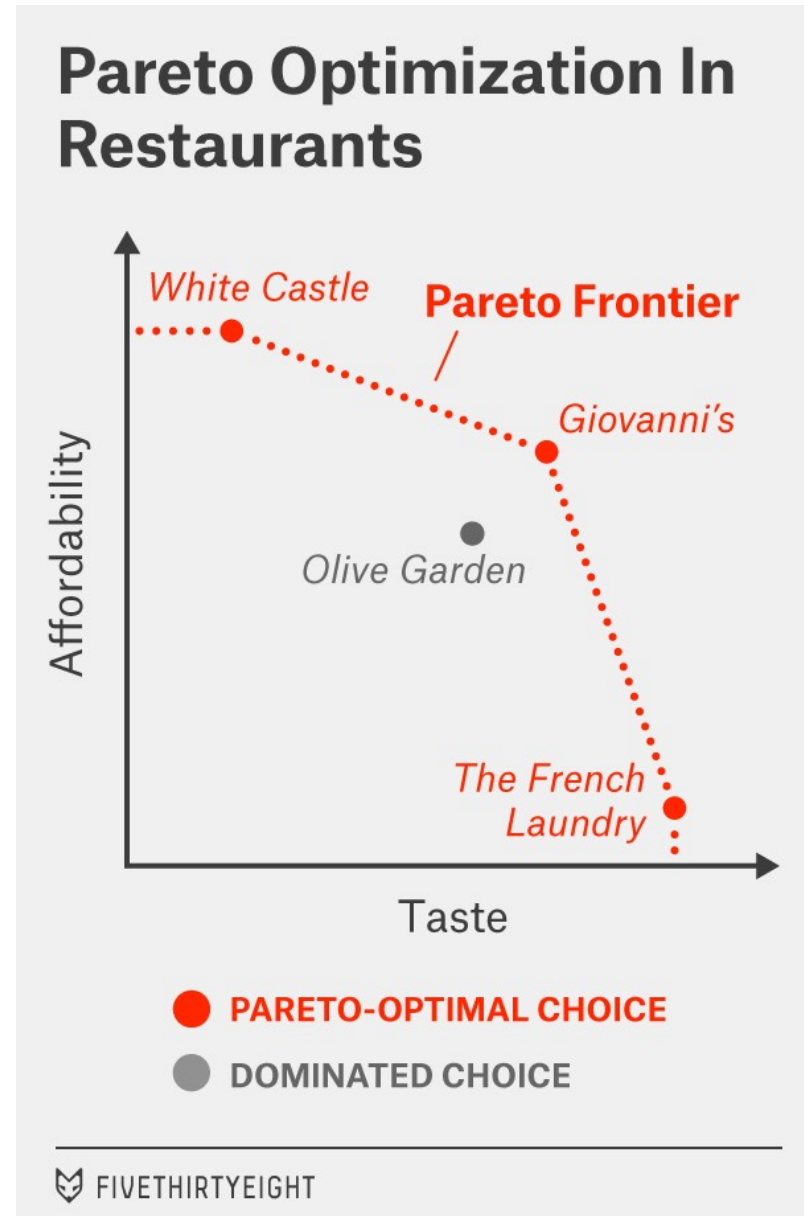*pp. 864 & 885*

# Cooperation

# Cooperation



- Non zero sum; win/win

- Vilfredo Pareto
  Pareto front is, in a cooperative setting, the combination of choices where no agent can be better off without at least making one other agent worse off

- It is the optimal cooperative strategy, the best outcome without hurting others.

# Pareto front



Pareto Optimization In Restaurants

White Castle

**Pareto Frontier**

Giovanni's

Olive Garden

The French Laundry

Affordability

Taste

● PARETO-OPTIMAL CHOICE
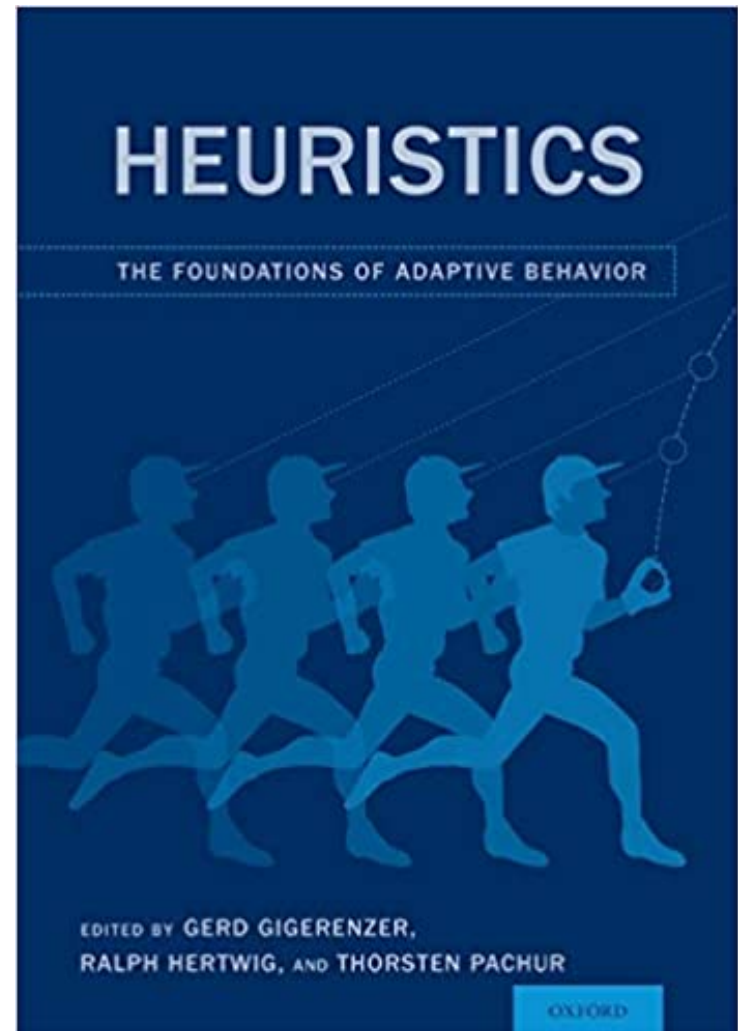
● DOMINATED CHOICE

FIVETHIRTYEIGHT

# Cooperative Behavior

- Dealing with nonstationarity and partial observability can be done (ignored) by separate training, no communication

- Realism can be improved with Centralized Training/Decentralized Execution -> Centralized controller, or interaction graphs

- Active field of research; overview

  - Value based: VDN, QMIX

  - Policy based: COMA, MADDPG

  - Opponent modeling: DRON, LOLA

  - Communication: Diplomacy game

  - Psychology: Heuristics
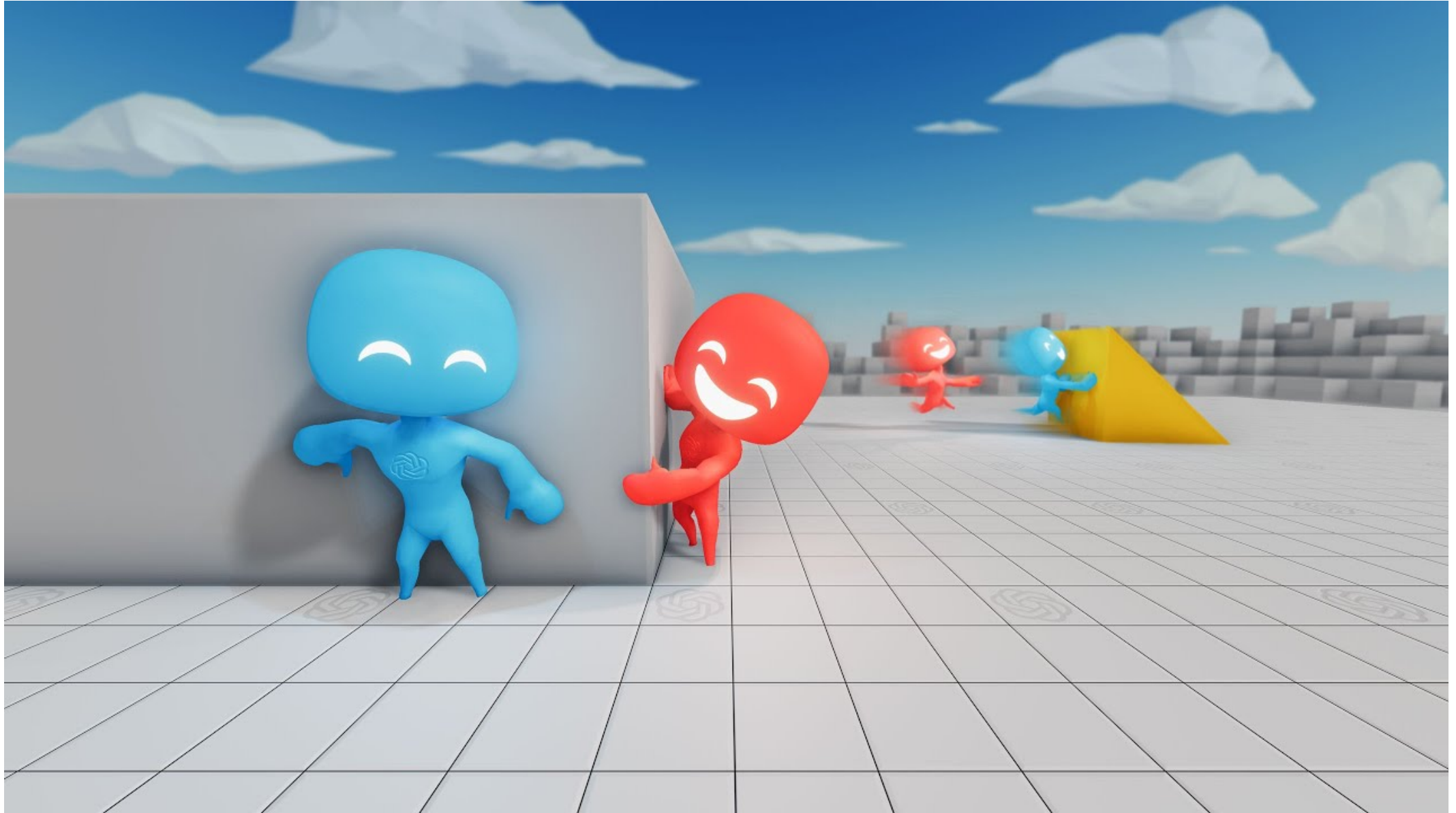
# Heuristics

SIMPLE
HEURISTICS
THAT MAKE US
SMART

GERD GIGERENZER, PETER M. TODD,
AND THE ABC RESEARCH GROUP

HEURISTICS

THE FOUNDATIONS OF ADAPTIVE BEHAVIOR

EDITED BY GERD GIGERENZER,
RALPH HERTWIG, AND THORSTEN PACHUR

OXFORD

# Emergent Cooperation

- [Baker, 2019]

- The agents can **move** by setting a force on themselves in the x and y directions as well as rotate along the z-axis.

- The agents can **see** objects in their line of sight and within a frontal cone.

- The agents can **sense** distance to objects, walls, and other agents around them using a lidar-like sensor.

- The agents can **grab and move** objects in front of them.

- The agents can **lock** objects in place. Only the team that locked an object can unlock it.

# Hide and Seek

# Mixed

- Prisoner's dilemma

- Iterated prisoner's dilemma

- Emerging social norms

# Prisoner's Dilemma

|  | **Confess** Defect | **Silent** Cooperate |
|---|---|---|
| **Confess** Defect | $(-5, -5)$ *Nash* | $(0, -10)$ |
| **Silent** Cooperate | $(-10, 0)$ | $(-2, -2)$ *Pareto* |

# Iterated Prisoner's Dilemma

- You remember "opponent's" behavior

- You will continue to meet your "opponents"

- Famous Experiment by Axelrod

- Rappoport introduced Tit for Tat

- You start being nice (Cooperating) and then do what the other did the previous round
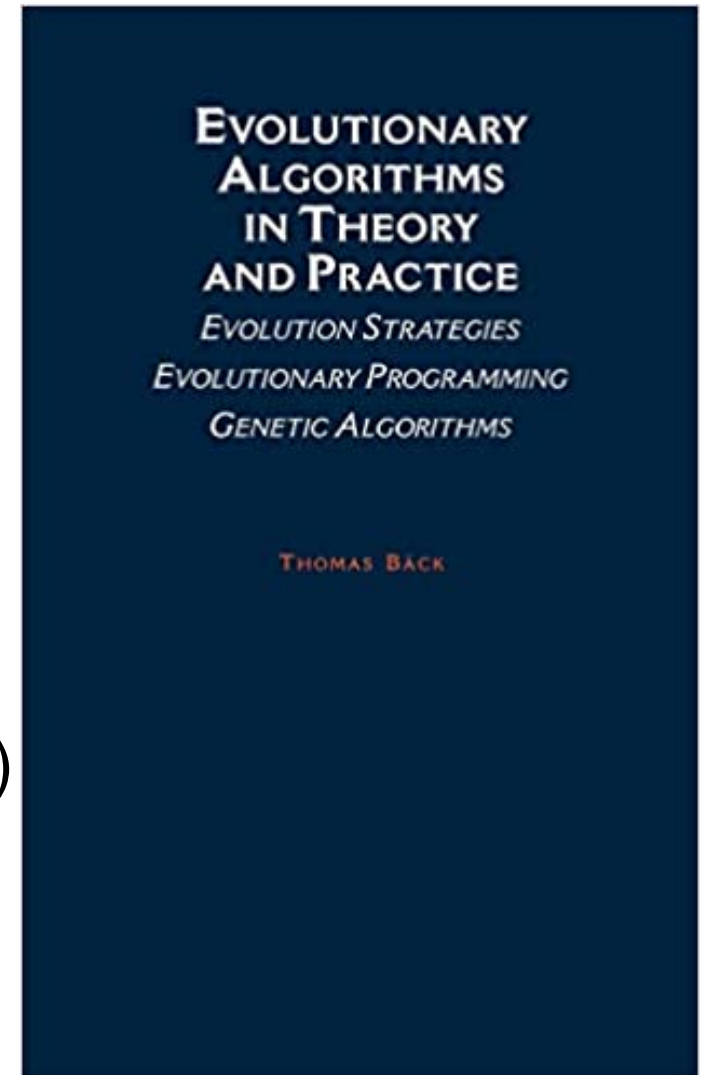
# Tit for Tat

# Algorithms

# Challenges

- Partial Observability -> Large State Space
  (Information sets)

- Nonstationary Environments -> Large State Space
  (Calculate all configurations)

- Multiple Agents -> Large State Space
  (Esp. with simultaneous actions)

# Evolutionary Approaches

- Evolutionary Algorithms

- Swarm Computing

- Population based training (teams, HRL)

**EVOLUTIONARY ALGORITHMS IN THEORY AND PRACTICE**
EVOLUTION STRATEGIES
EVOLUTIONARY PROGRAMMING
GENETIC ALGORITHMS

THOMAS BÄCK

# Evolutionary Framework

---

**Algorithm 7.1** Evolutionary Framework [36]

1: Generate the initial population randomly
2: **repeat**
3:      Evaluate the fitness of each individual of the population
4:      Select the fittest individuals for reproduction
5:      Through crossover and mutation generate new individuals
6:      Replace the least fit individuals by the new individuals
7: **until** terminated

---

# Evo

| Multi-Agent Reinforcement Learning | Evolutionary Computation |
|:---:|:---:|
| agent | individual |
| some | many |
| all agents | population |
| environment | problem |
| reward | fitness |
| policy | genes |
| adaptation | mutation and combination |
| time step | generation |
| feedback | selection |

- Highly parallel

- Multi-agent population based optimization

- Single-agent deep network policy optimization instead of backpropagation

- Single fitness function, determines cooperation or competition

# Swarm Intelligence Algorithms

## A Tutorial

Edited by

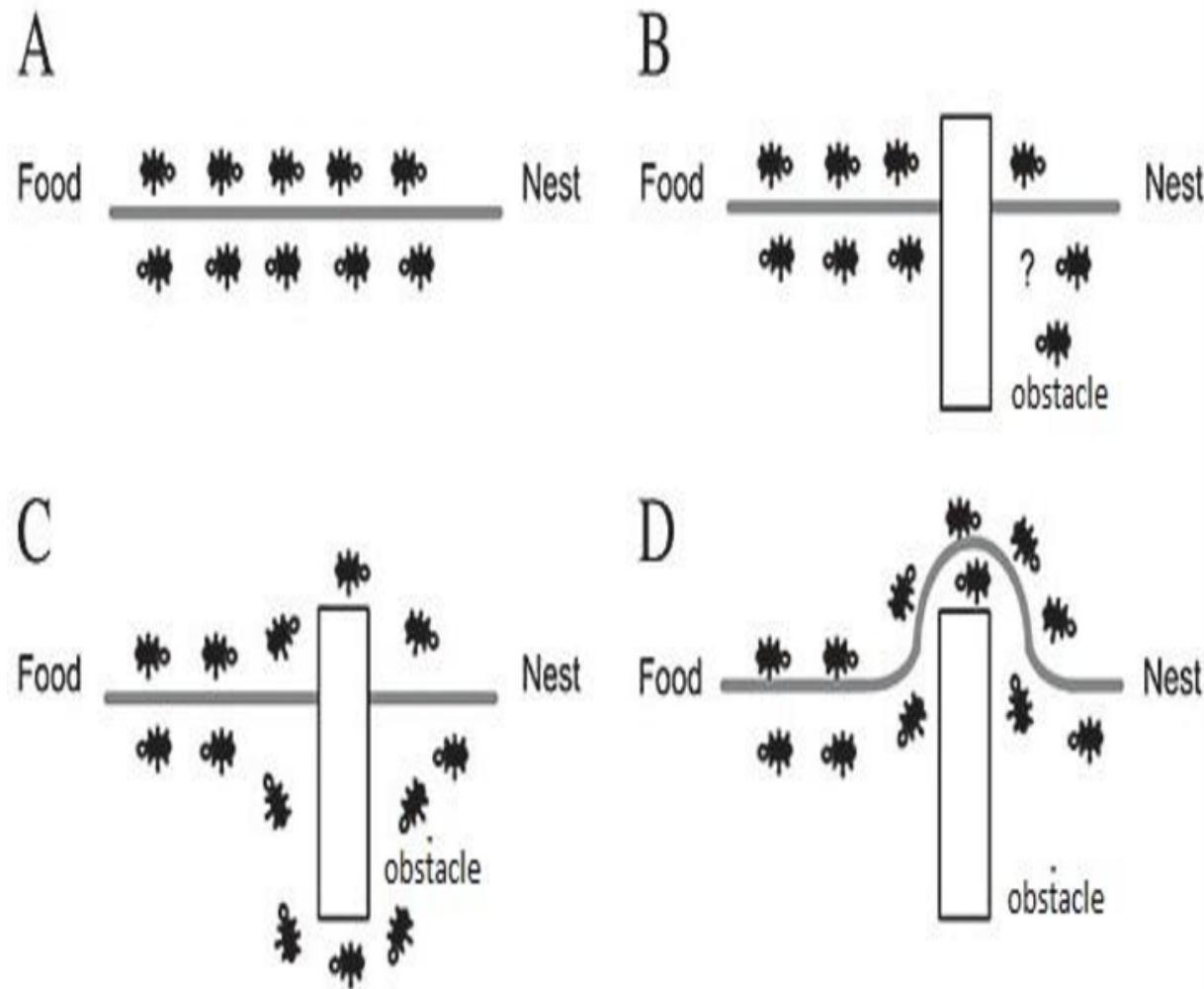Adam Slowik

# Ant Colony Optimization

- ACO



**Fig 1.**
A: Ants in a pheromone trail between nest and food.
B: an obstacle interrupts the trail.
C: Ants find two paths to go around the obstacle

# Population-based training

- Teams

- Hierarchical

- Cooperation, competition

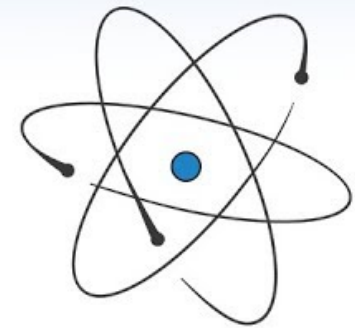- Within Teams, between teams

- Blends RL and Evo

# Population-based training

**Algorithm 7.2** Population Based Training [352]

**procedure** TRAIN($\mathcal{P}$)                                                                                    ▷ initial population $\mathcal{P}$
    **for** $(\theta, h, p, t) \in \mathcal{P}$ (asynchronously in parallel) **do**
        **while** not end of training **do**
            $\theta \leftarrow \texttt{step}(\theta | h)$                          ▷ one step of optimisation using hyperparameters $h$
            $p \leftarrow \texttt{eval}(\theta)$                                     ▷ current model evaluation
            **if** $\texttt{ready}(p, t, \mathcal{P})$ **then**
                $h', \theta' \leftarrow \texttt{exploit}(h, \theta, p, \mathcal{P})$       ▷ use the rest of population for improvement
                **if** $\theta \neq \theta'$ **then**
                    $h, \theta \leftarrow \texttt{explore}(h', \theta', \mathcal{P})$          ▷ produce new hyperparameters $h$
                    $p \leftarrow \texttt{eval}(\theta)$                            ▷ new model evaluation
                **end if**
            **end if**
            update $\mathcal{P}$ with new $(\theta, h, p, t + 1)$                          ▷ update population
        **end while**
    **end for**
    **return** $\theta$ with the highest $p$ in $\mathcal{P}$
**end procedure**

# CTF



SUPERHUMAN QUAKE 3 AI TEAM

TWO MINUTE PAPERS

# StarCraft

# StarCraft

- Real Time Strategy

- 10^1685

- <u>AlphaStar</u>

- Population based multi agent methods

# Questions?