

Designing Consensus: Gamified Modeling and Simulation of Collaborative Decision-Making

Instructor: Jin Gao

IAP 2025 (Non-Credit)

2025.01.30

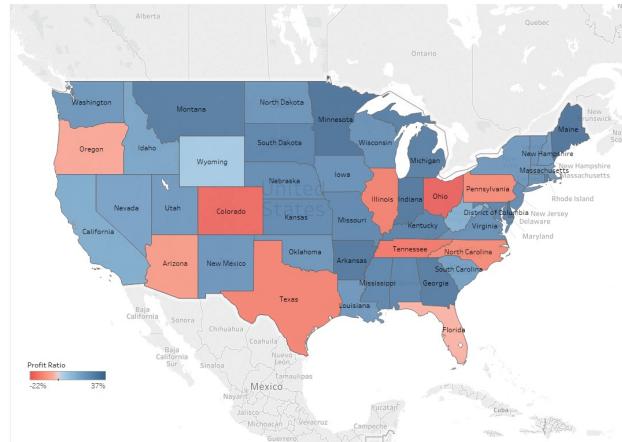
gaojin@mit.edu

Date	Topic	Contents
Jan 30	Multi-Agent Learning	<ul style="list-style-type: none">Step from the concept of game theories we built from the previous lectures, we explore methods to train agents to make decisions in the environment we created. Starting from reinforcement learning (RL), we will introduce multi-agent RL and explore the pros and cons of a variety of methods.Overview of advances in this field, including Deep Blue, AlphaGo, and Cicero.
	Workshop	<ul style="list-style-type: none">Present your final game design to the class (game rules, competition or collaboration mechanisms, balancing, strategies, etc.).
	Wrap up and conclude.	

Before making a decision, what info do we need?

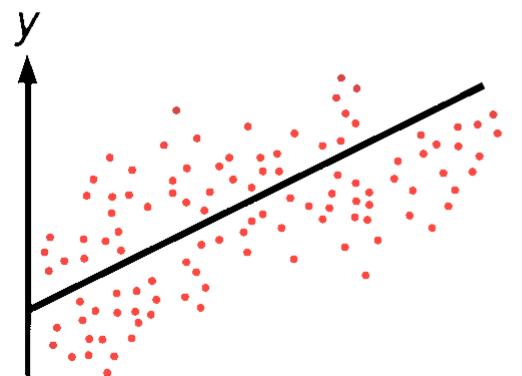
Descriptive

What happened?



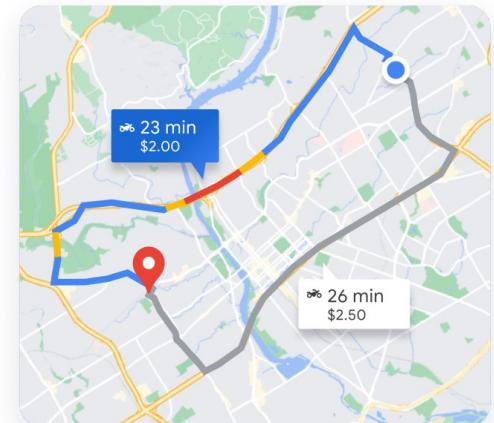
Predictive

What will happen?



Prescriptive

What should I do?



Data Management

Data Visualization

Forecasting

Machine Learning

Types of Analysis

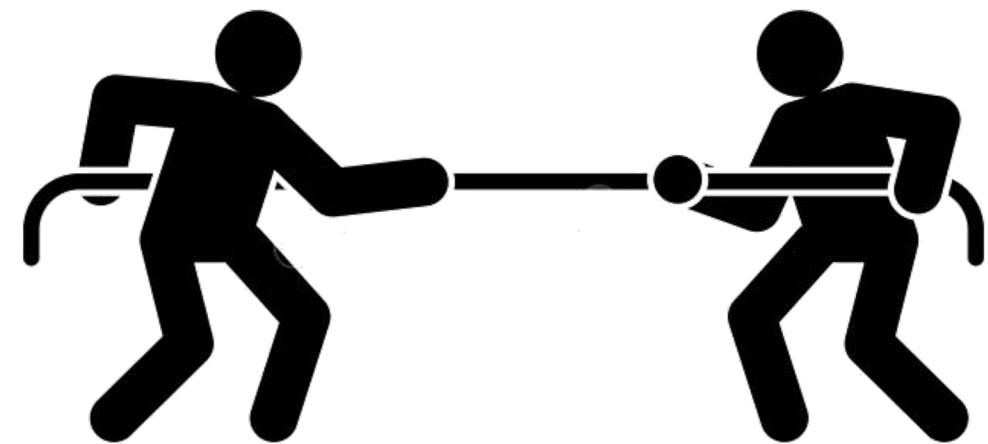
Optimization

Simulation

More than just yourself?

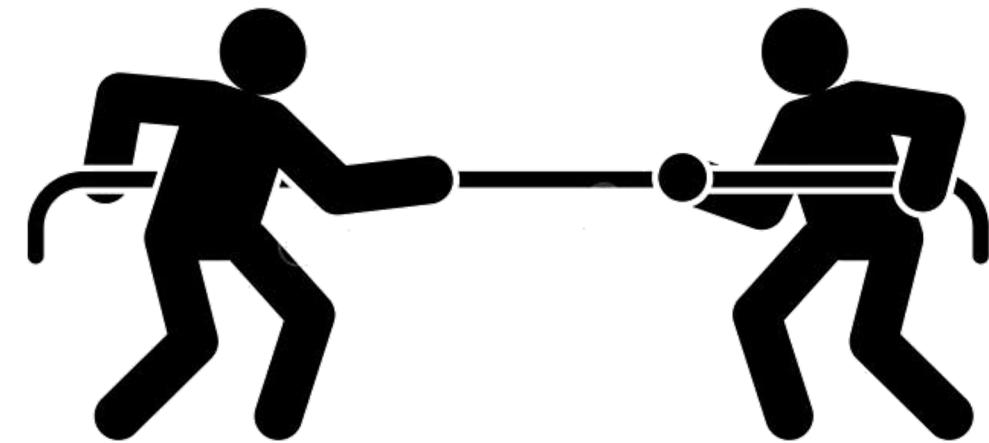
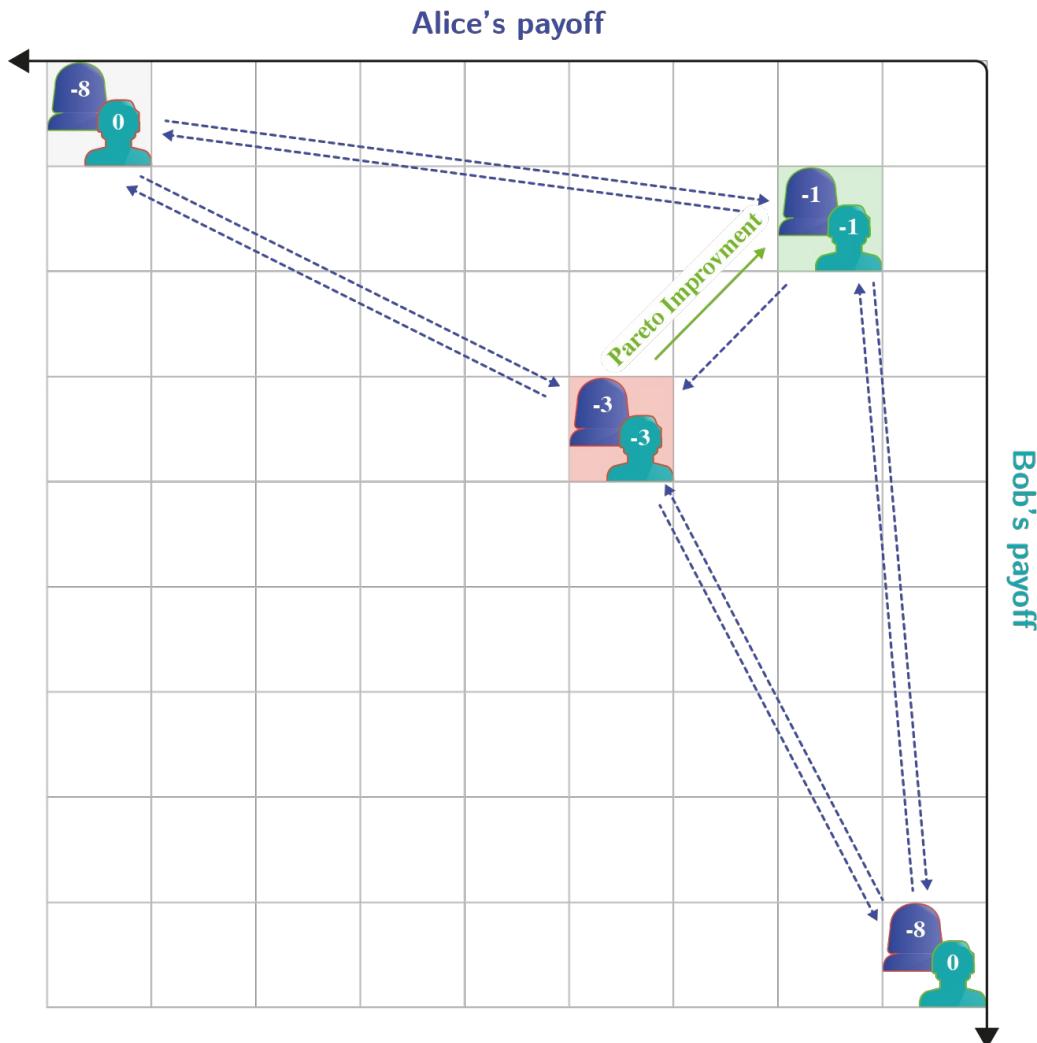


Corporative



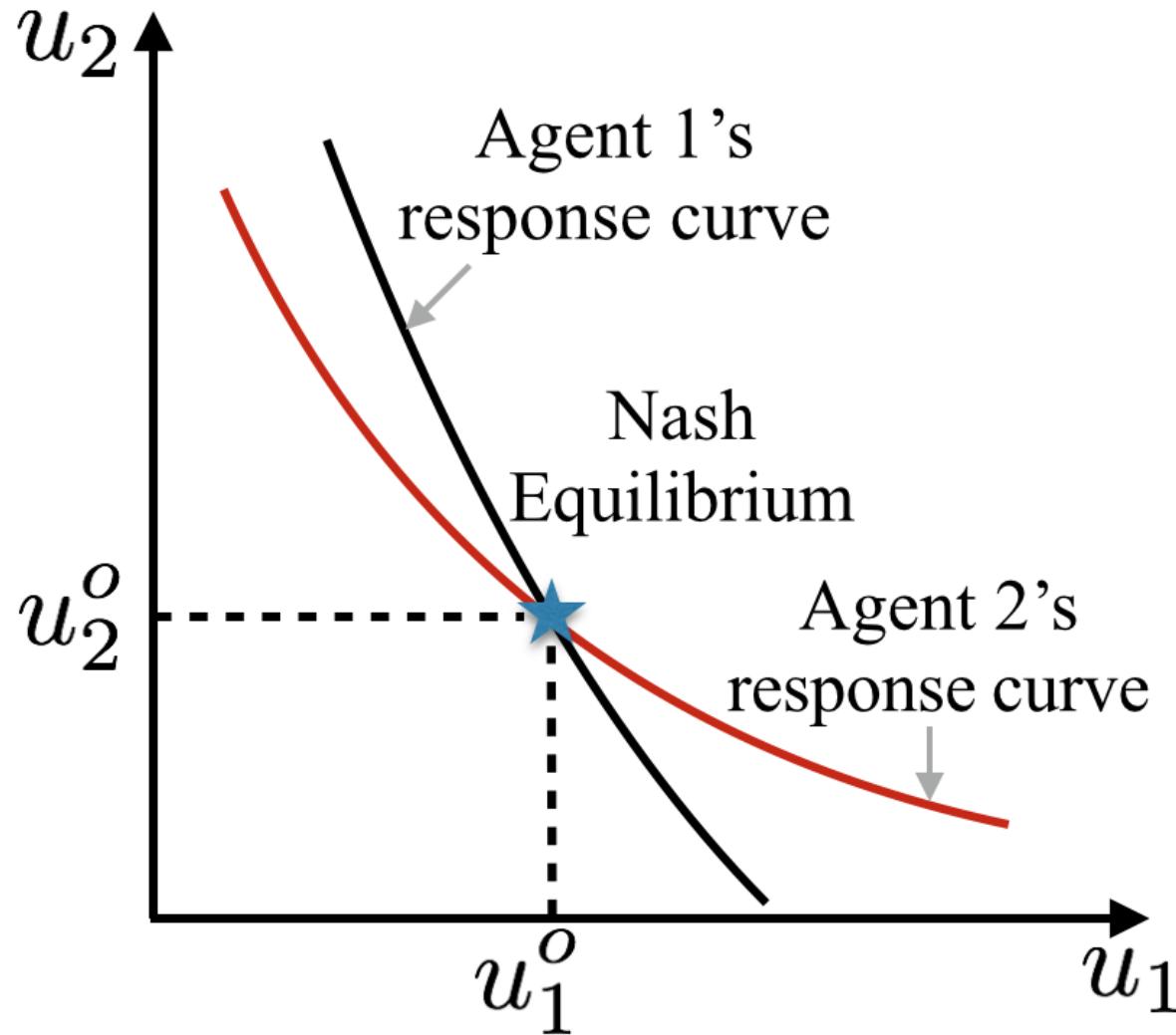
Competitive

Two-person situation – Zero-Sum Game



Competitive

Nash Equilibrium

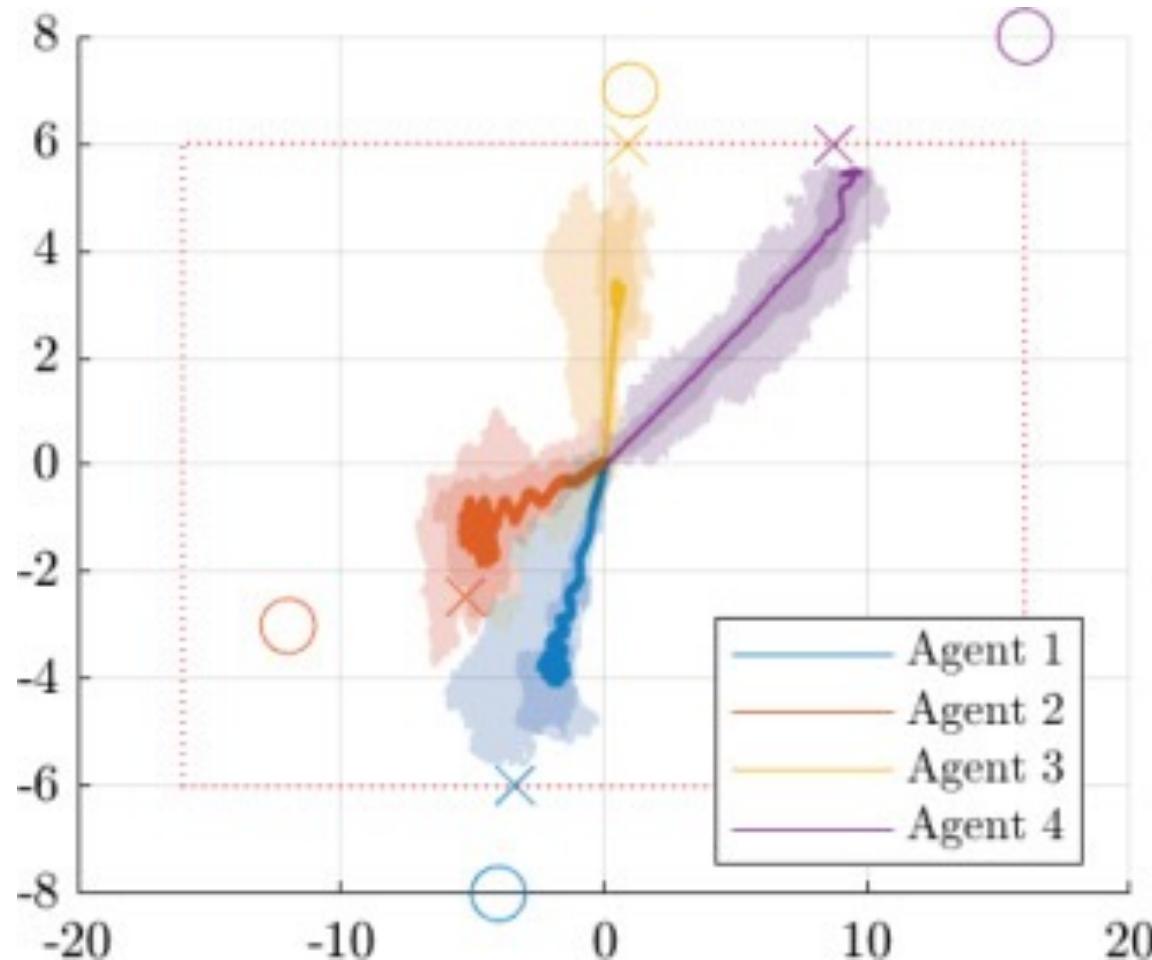


No player can improve their outcome by changing their strategy

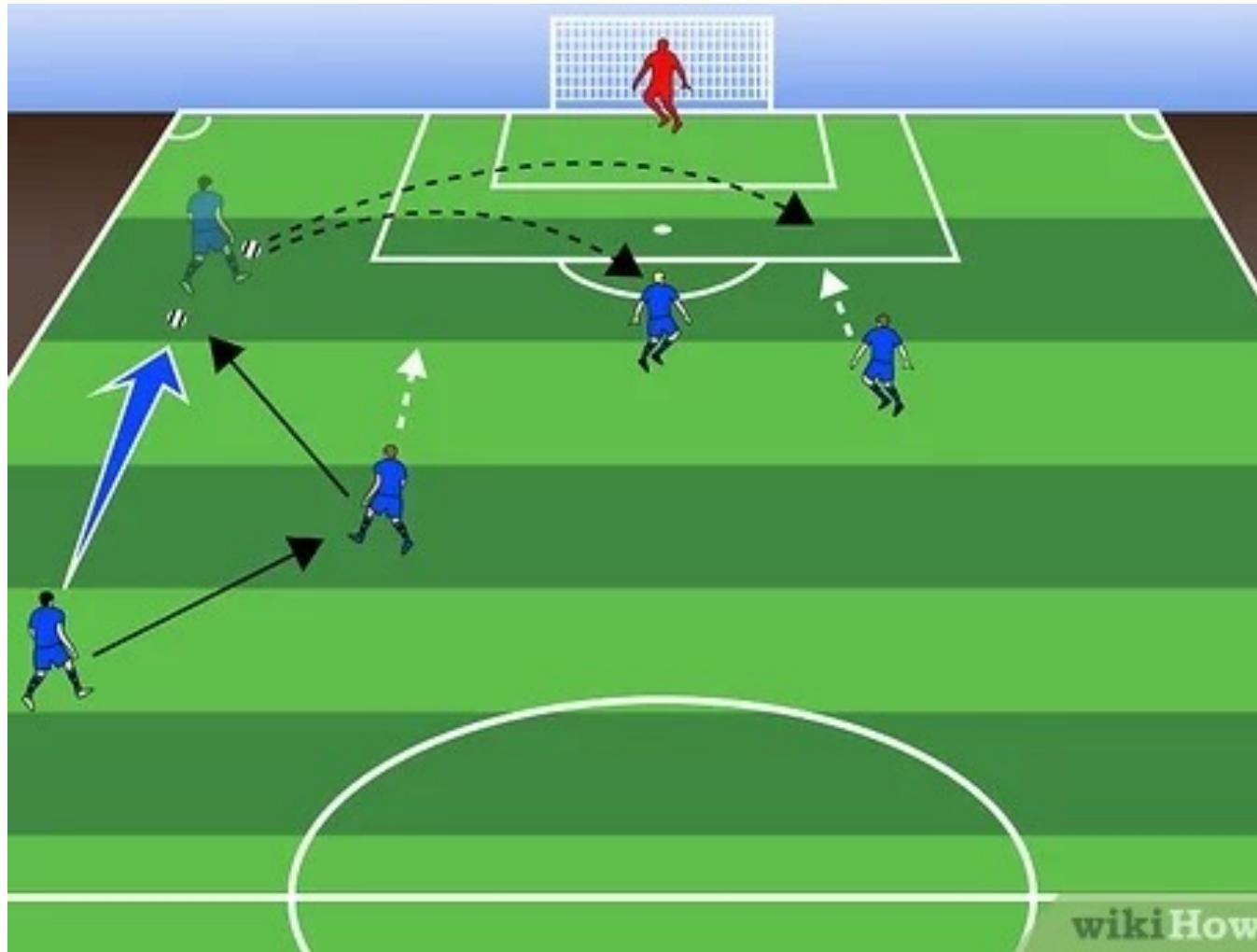
-- the solution of N-player noncooperative games.

-- hard to compute that *exact* point

What will happen if there are more than 2 people?

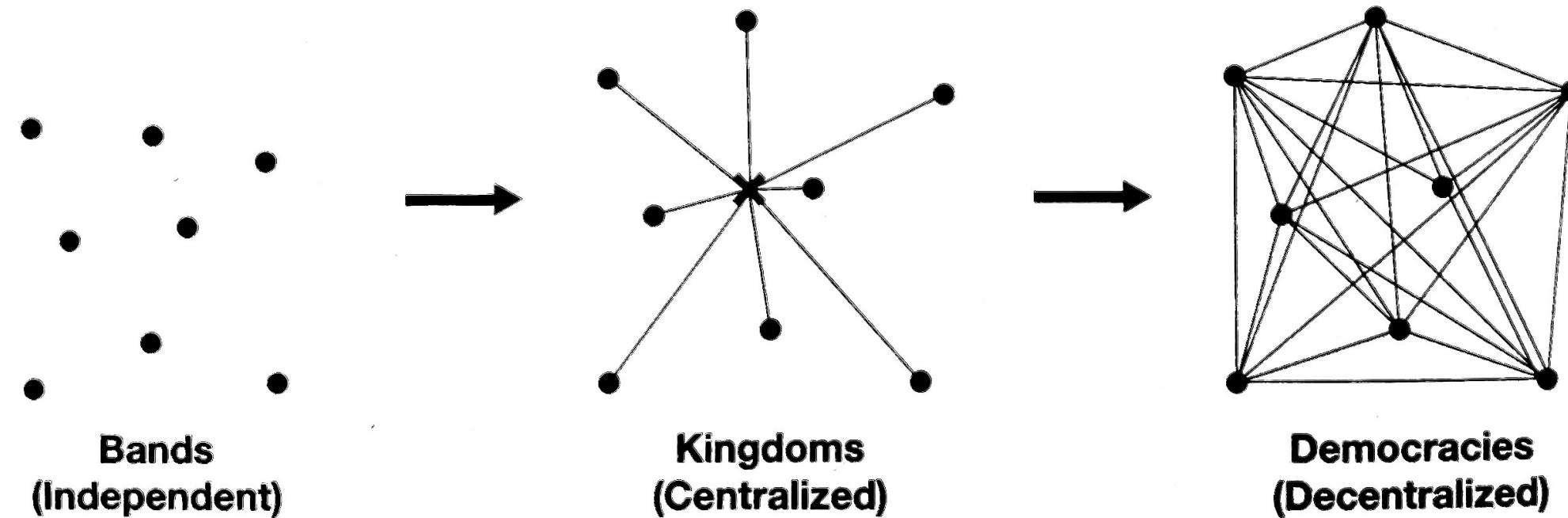


What will happen if the players needs collaborate?

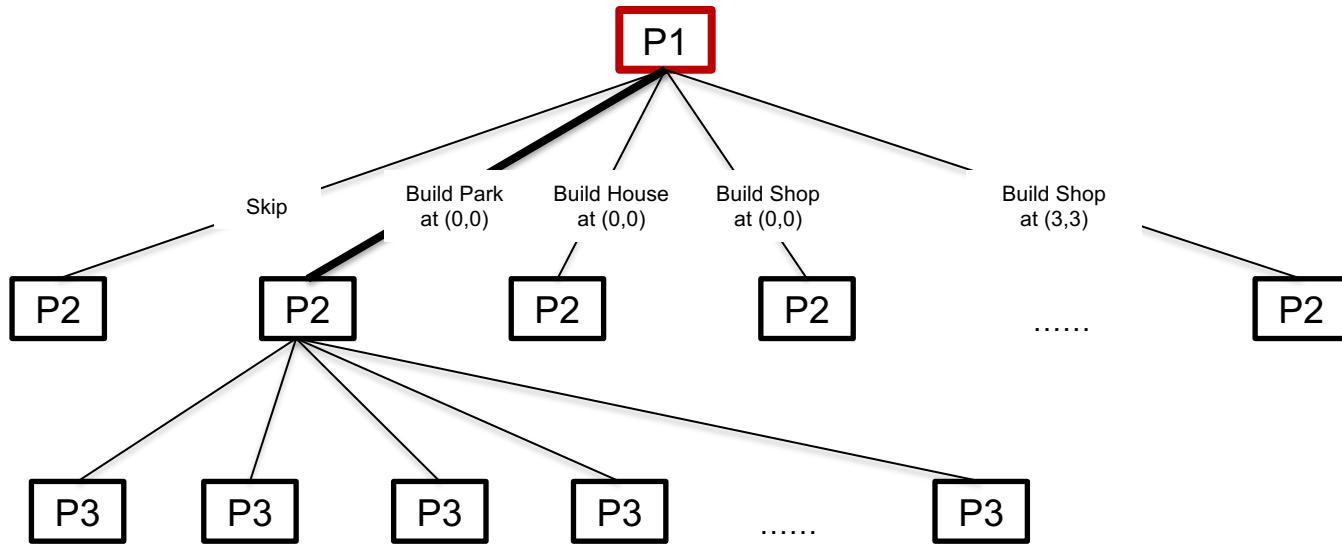


What will happen if the players have hierarchies?

The major ways human societies have been organized throughout history reveal a remarkably simple pattern that foreshadows how businesses are now changing.



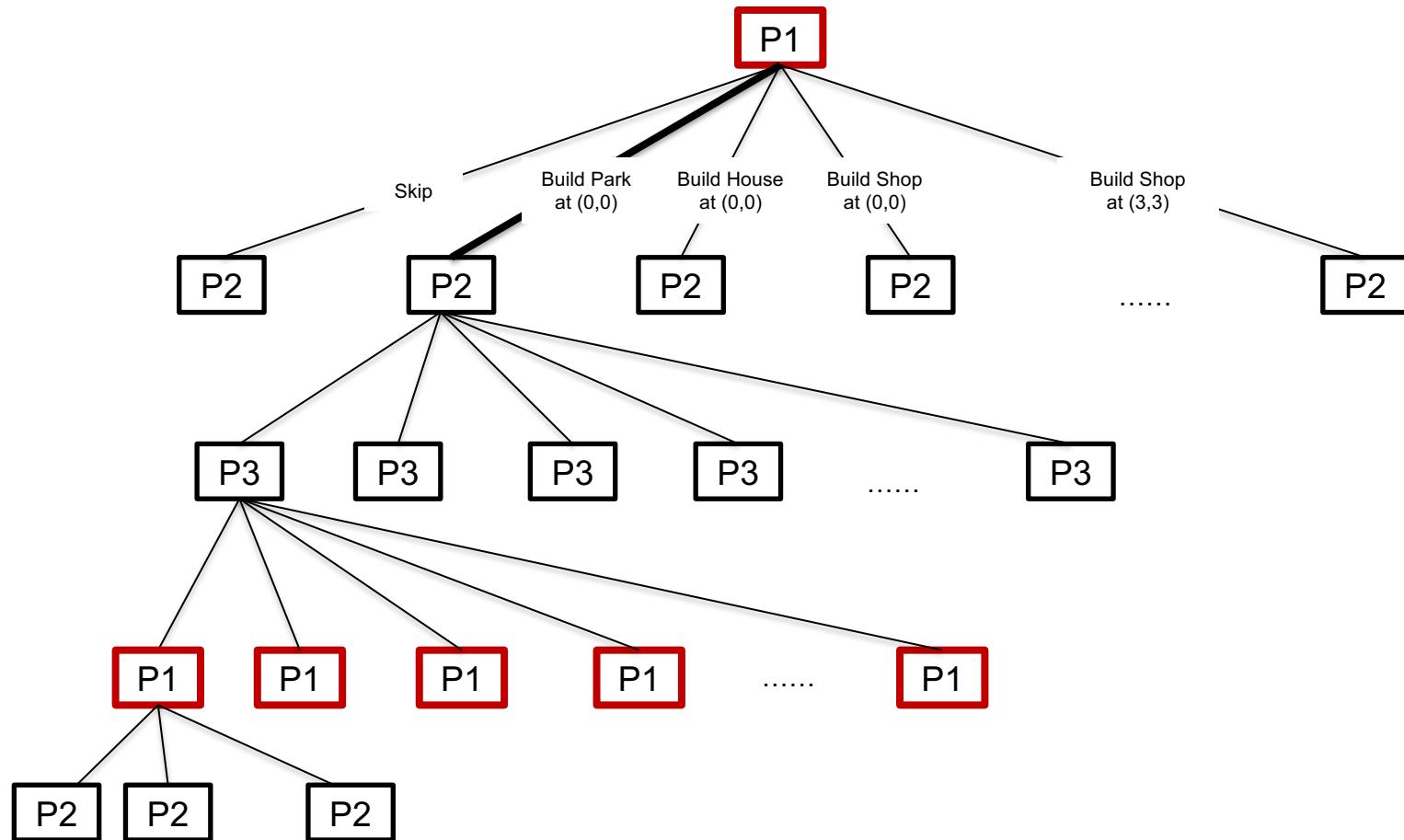
Intuition of making a series of decisions From Player 1's Perspective



Tree-form Decision Process (TFDP)

Intuition of making a series of decisions

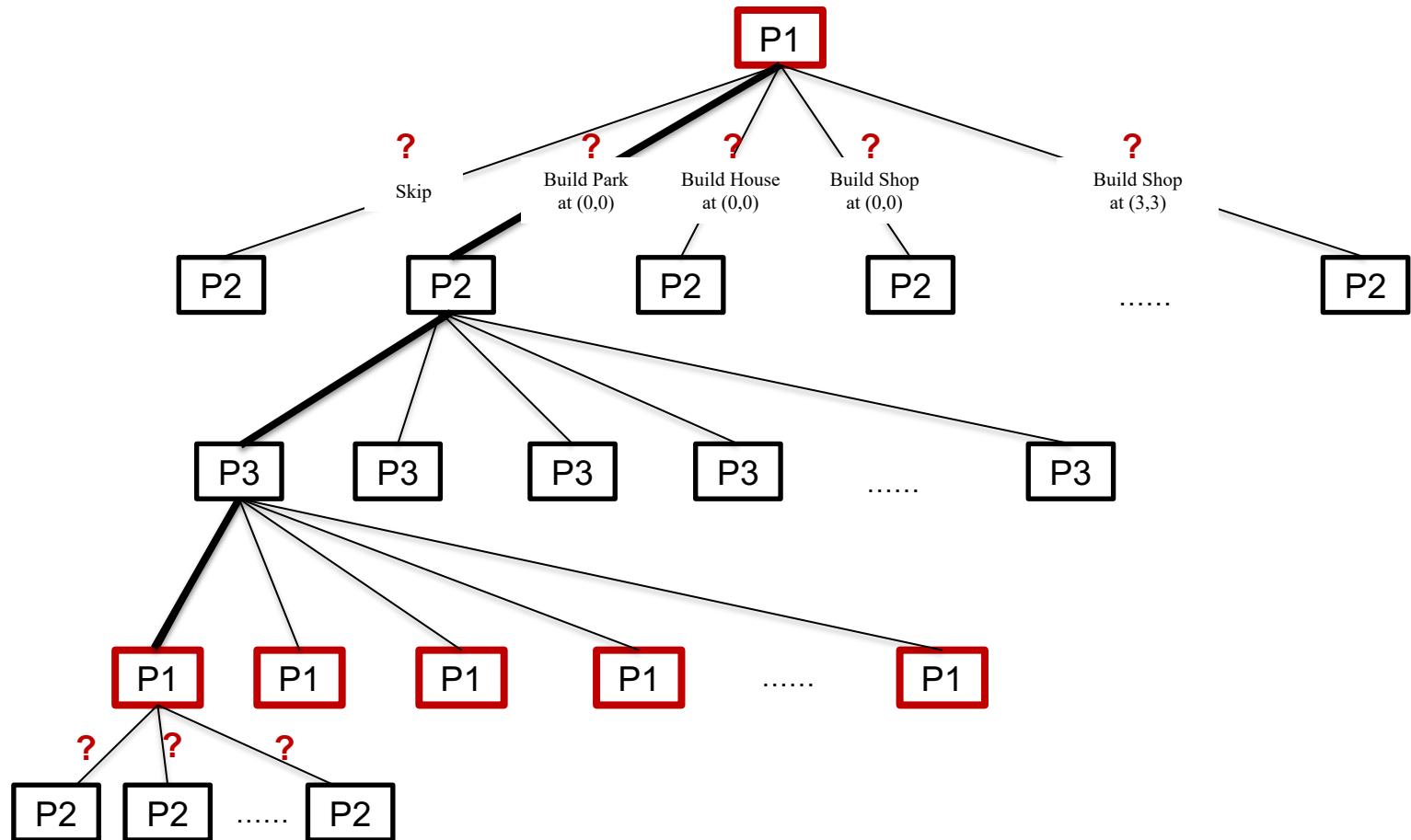
From Player 1's Perspective



Tree-form Decision Process (TFDP)

Intuition of making a series of decisions

From Player 1's Perspective



Tree-form Decision Process (TFDP)

Plotting Decisions with Hidden Information: Game Tree vs Tree Form Decision Process

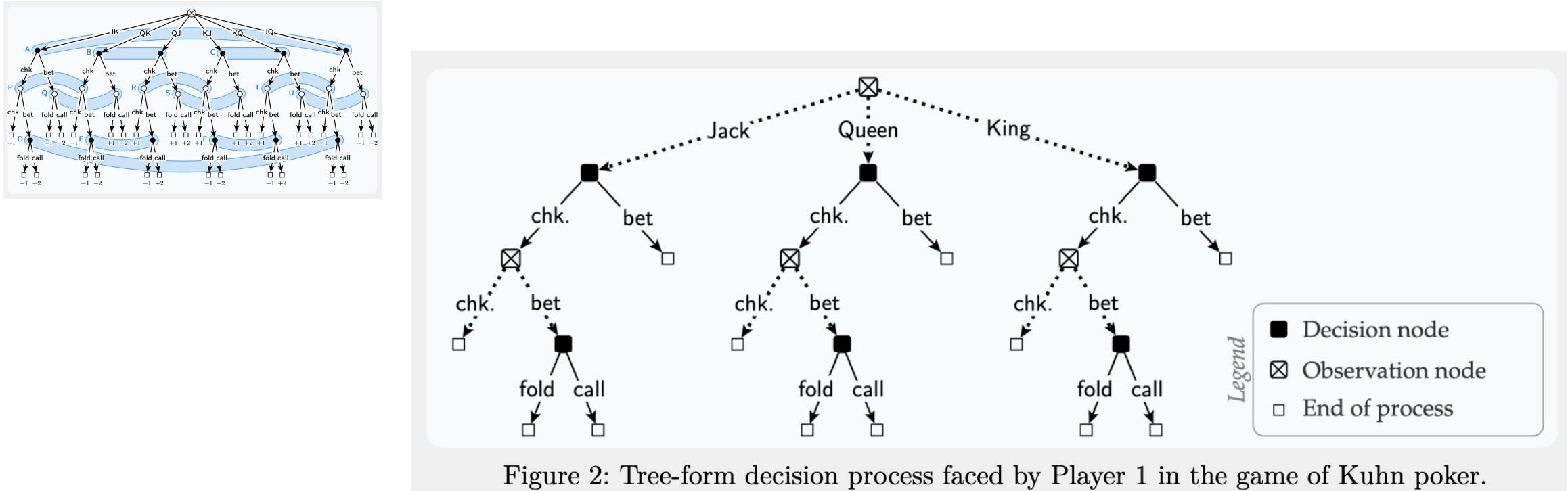
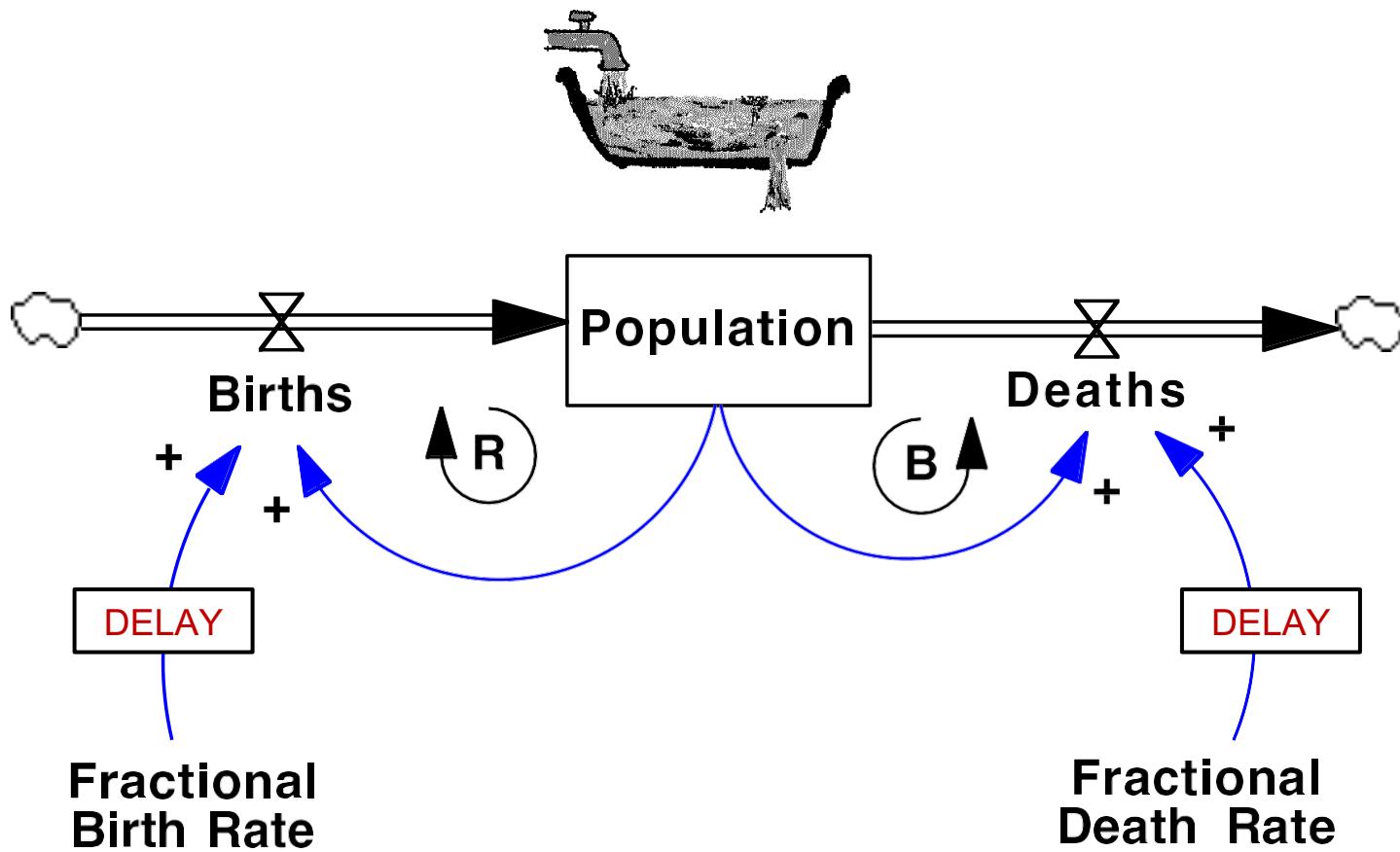


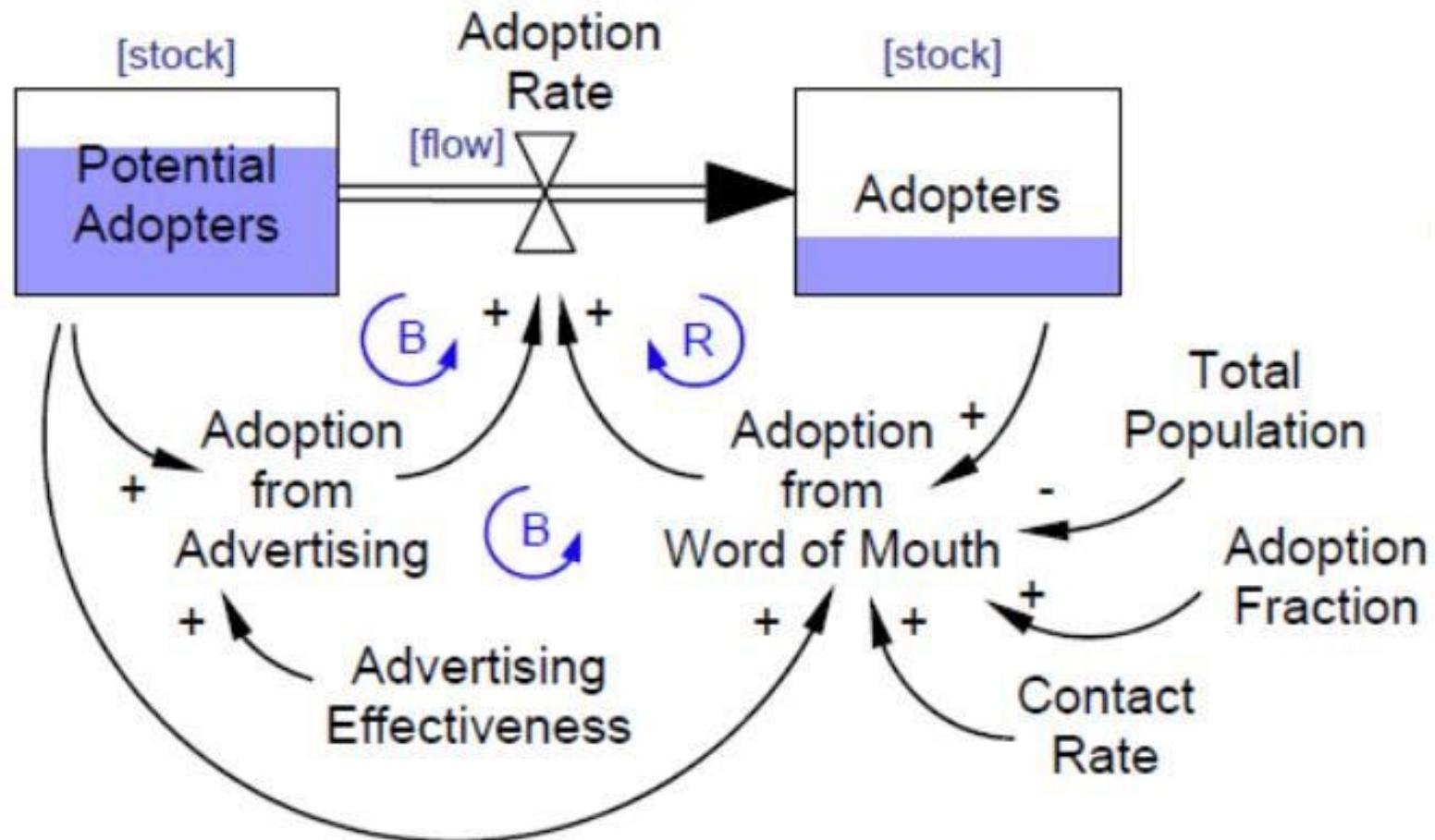
Figure 2: Tree-form decision process faced by Player 1 in the game of Kuhn poker.

Decision Tree

Basic Elements: Delay



Classic Models



Bass Diffusion Model

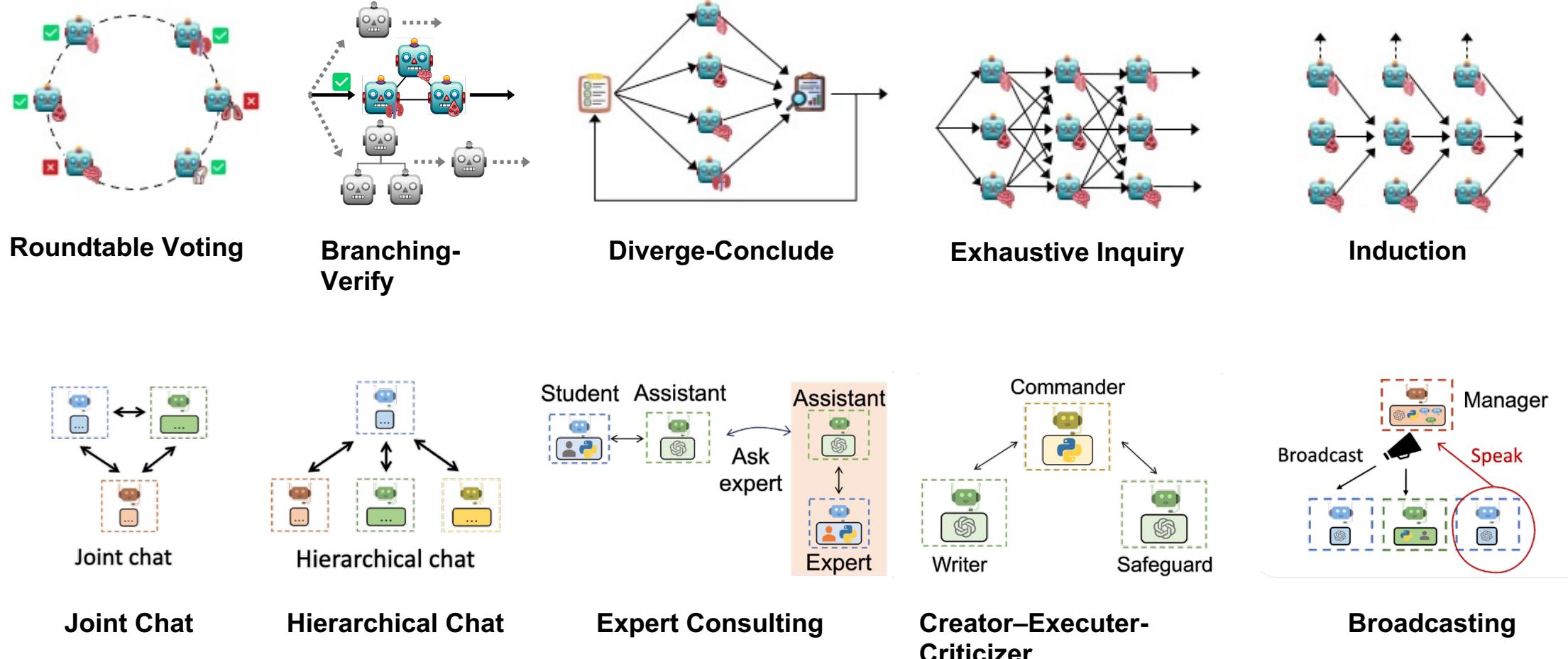


Multi-Generative Agent System

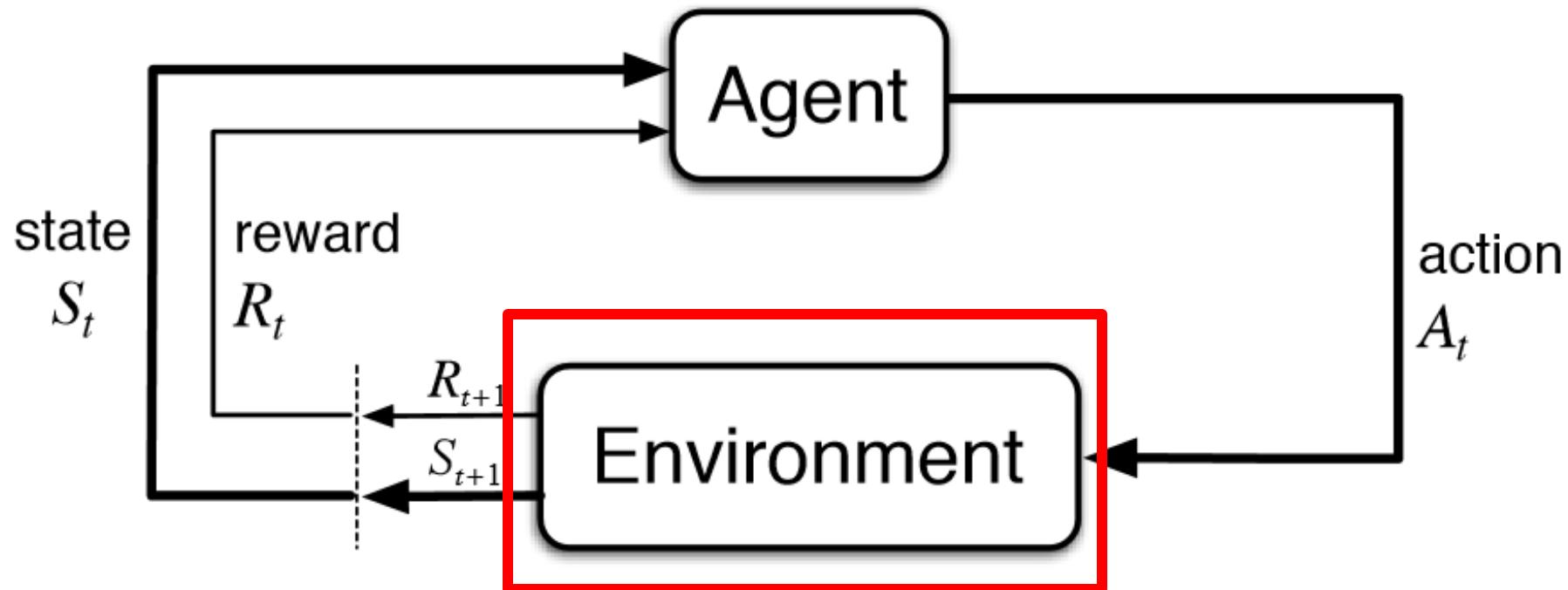


Multi-Generative Agent System proposed in recent researches [Park et al., 2023] (left) and [Chen et al., 2023] (right)

Agents-level Institution: Multi-Agent Debate Frameworks



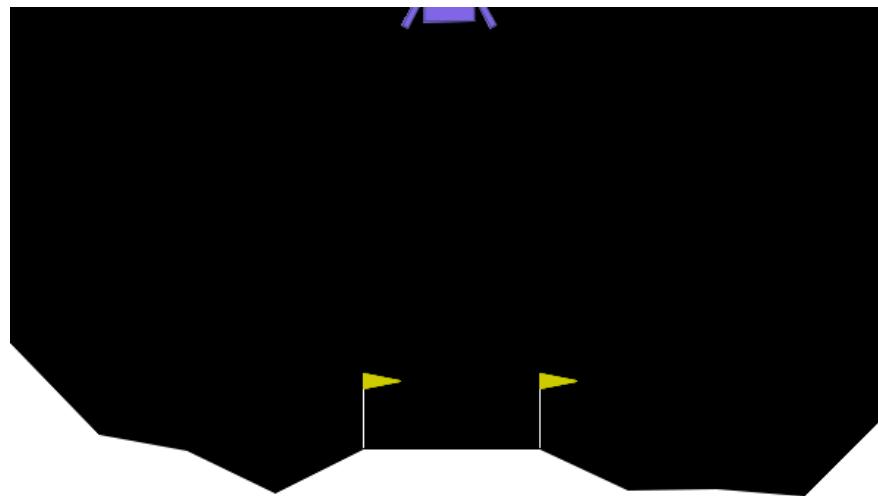
Where does agents learn?



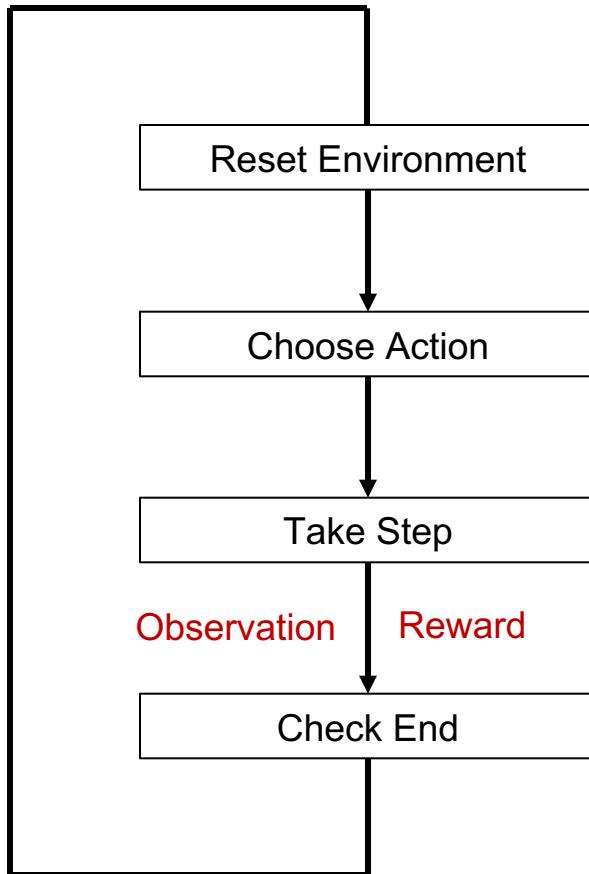
OpenAI Gymnasium



An API standard for reinforcement learning with a diverse collection of reference environments



Execution Loop



```
import gymnasium as gym

# Initialise the environment
env = gym.make("LunarLander-v3", render_mode="human")

# Reset the environment to generate the first observation
observation, info = env.reset(seed=42)
for _ in range(1000):
    # this is where you would insert your policy
    action = env.action_space.sample()

    # step (transition) through the environment with the action
    # receiving the next observation, reward and if the episode has terminated or truncated
    observation, reward, terminated, truncated, info = env.step(action)

    # If the episode has ended then we can reset to start a new episode
    if terminated or truncated:
        observation, info = env.reset()

env.close()
```

```
env_name: CartPole-v1  
steps: 69  
action: right  
reward: 1.00  
total_reward: 69.00  
terminated: False  
truncated: False
```



By Gregory

Episode
20,000

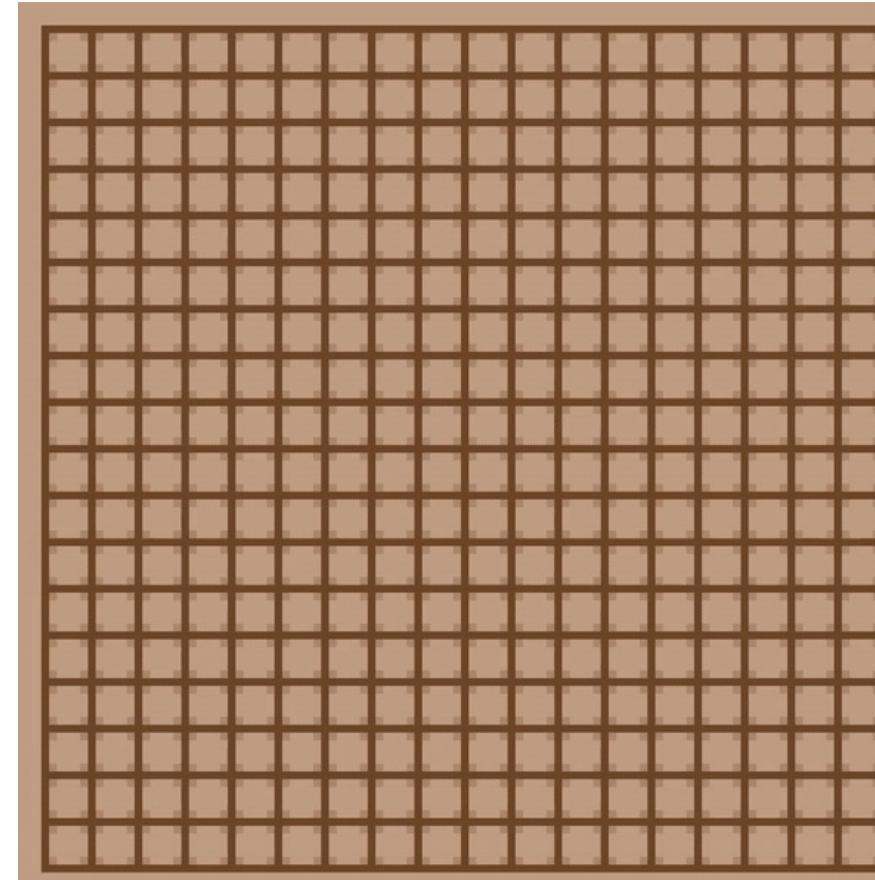


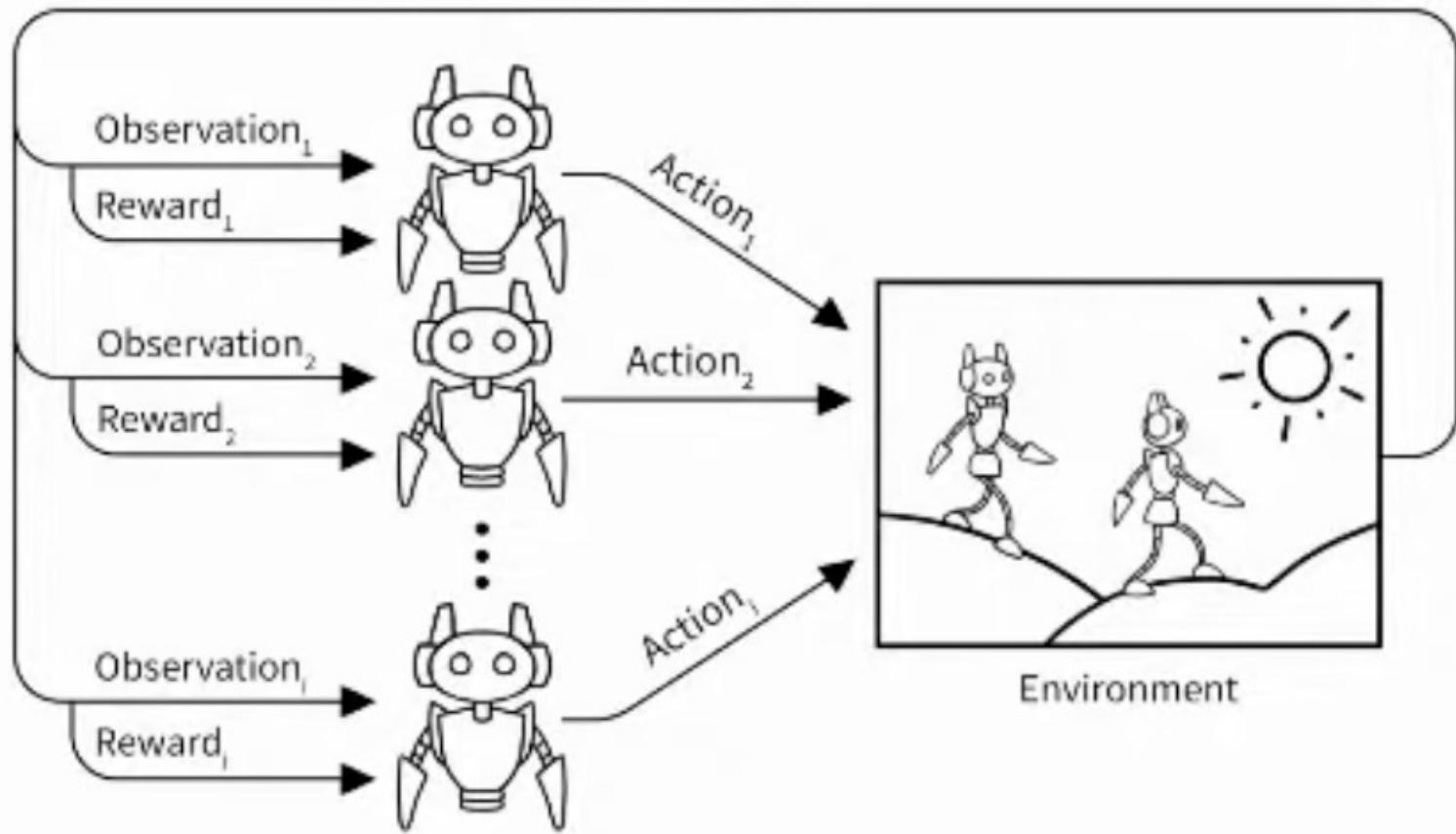
What about multi-agent?



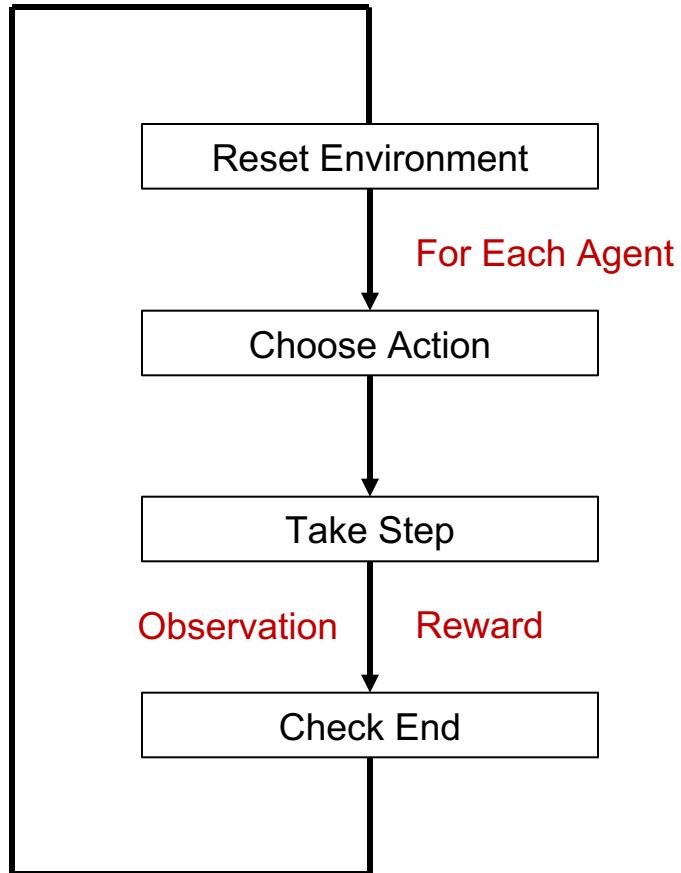


An API standard for multi-agent reinforcement learning.



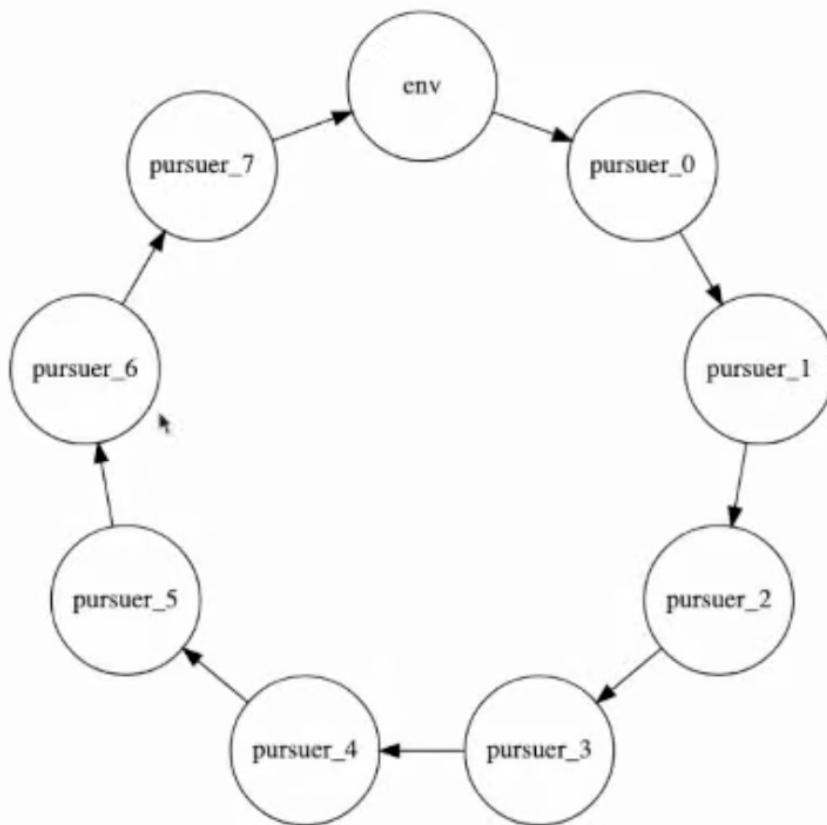


Execution Loop

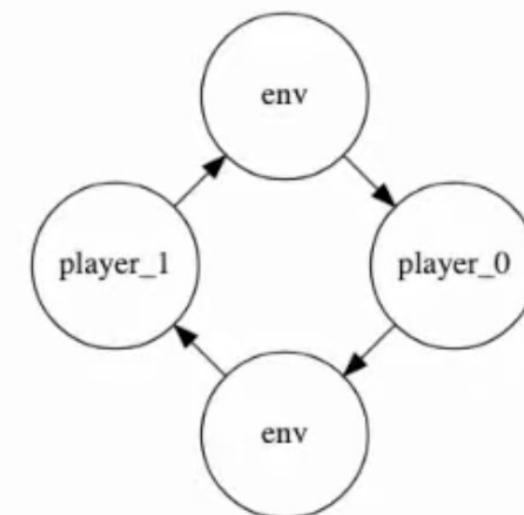


```
from pettingzoo.butterfly import knights_archers_zombies_v10
env = knights_archers_zombies_v10.env(render_mode="human")
env.reset(seed=42)
for agent in env.agent_iter():
    observation, reward, termination, truncation, info = env.last()
    action = policy(observation, agent)
    env.step(action)
```

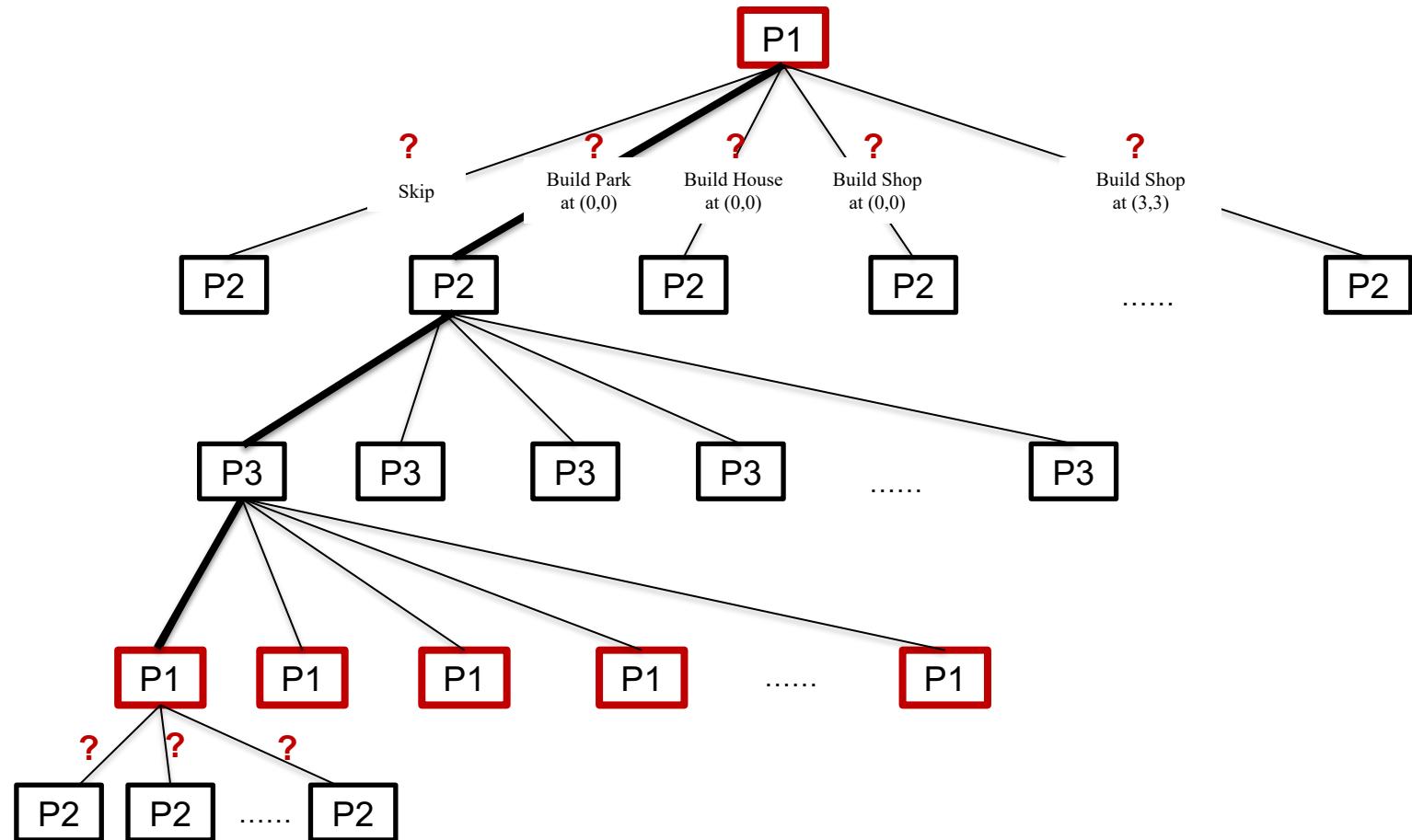
Agent-Environment Cycle



Pursuit



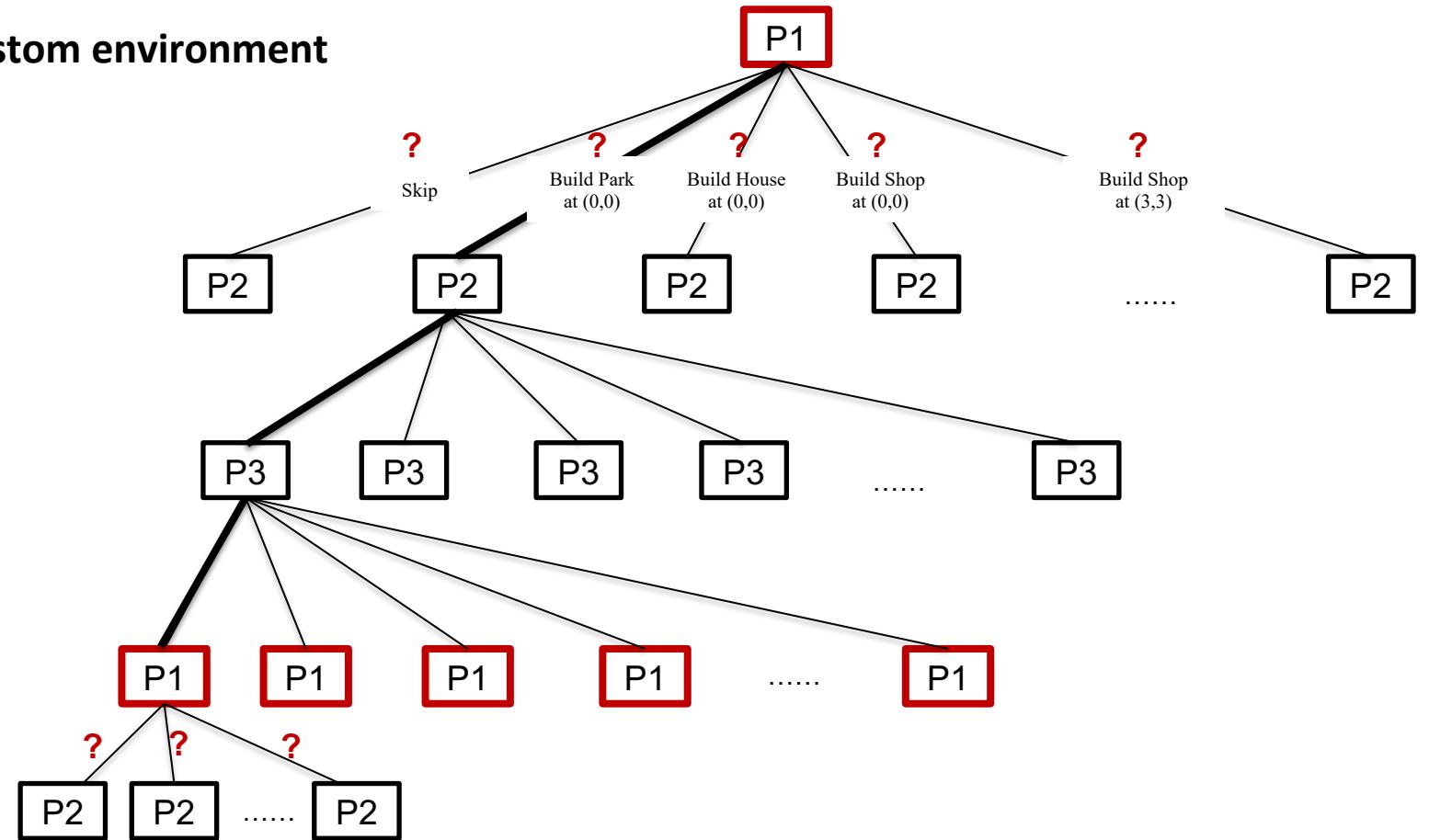
Chess



Tree-form Decision Process (TFDP)

Training in the environment

Using EPyMARL Library with custom environment



Training Outcome? The best game strategies

Cities Skylines



Training Outcome? The best game strategies

Cities Skylines



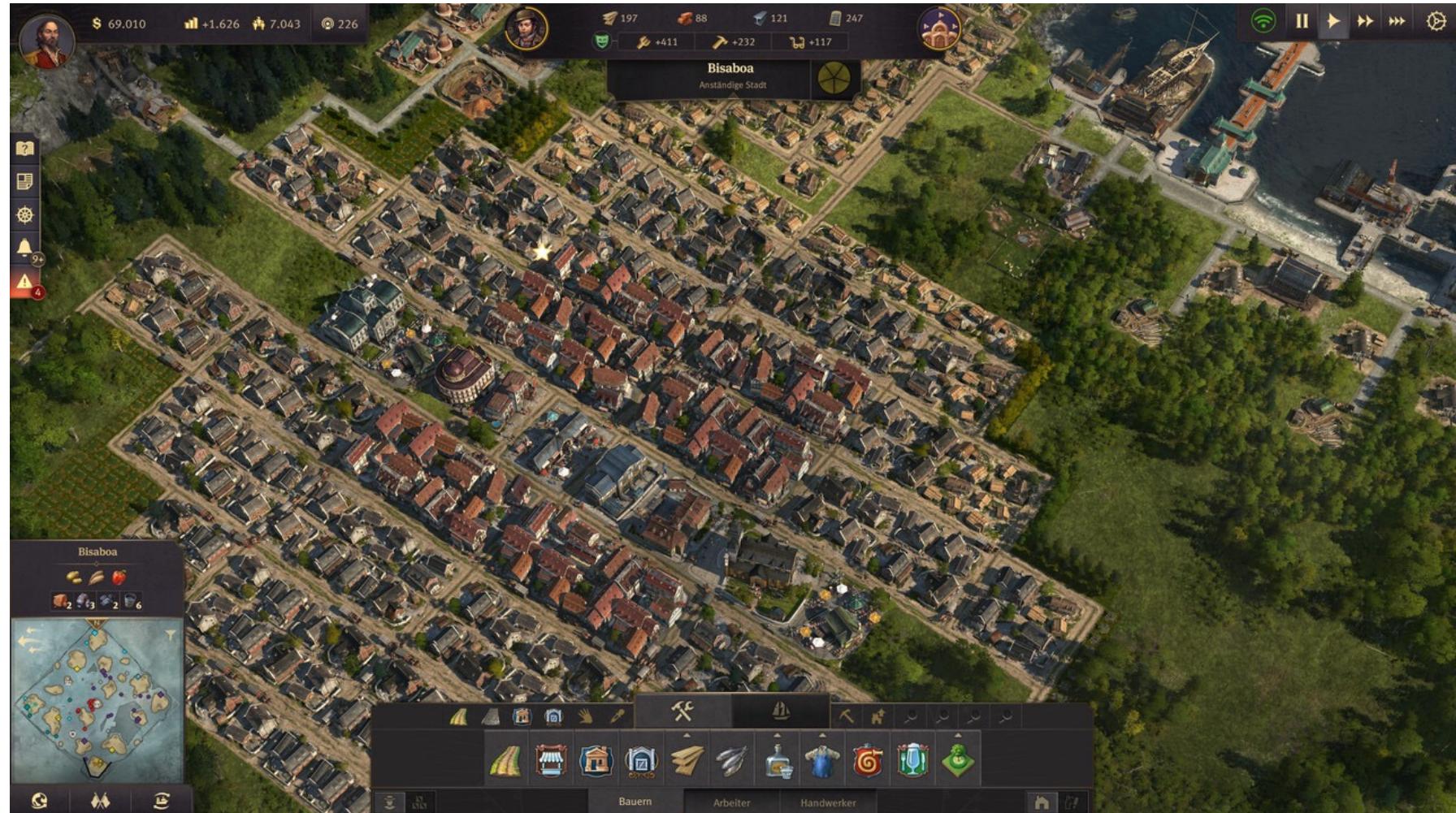
- Main Street
- Street
- Commercial and Office space
- Path
- Water
- Park



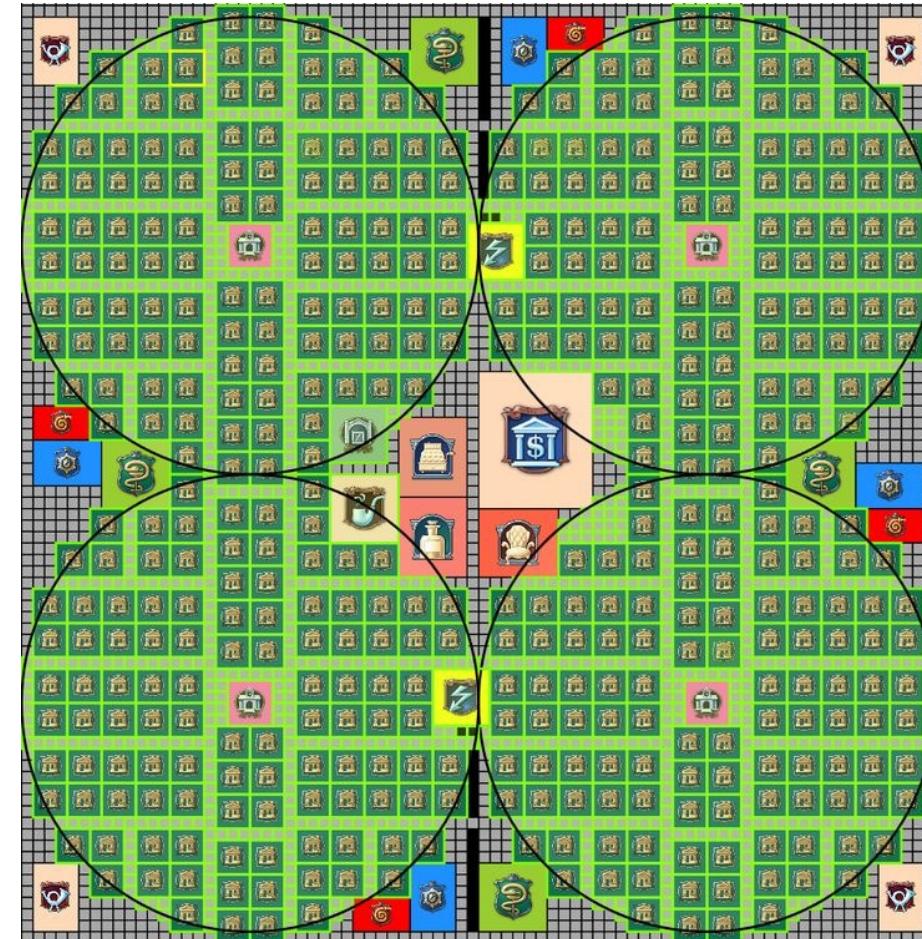
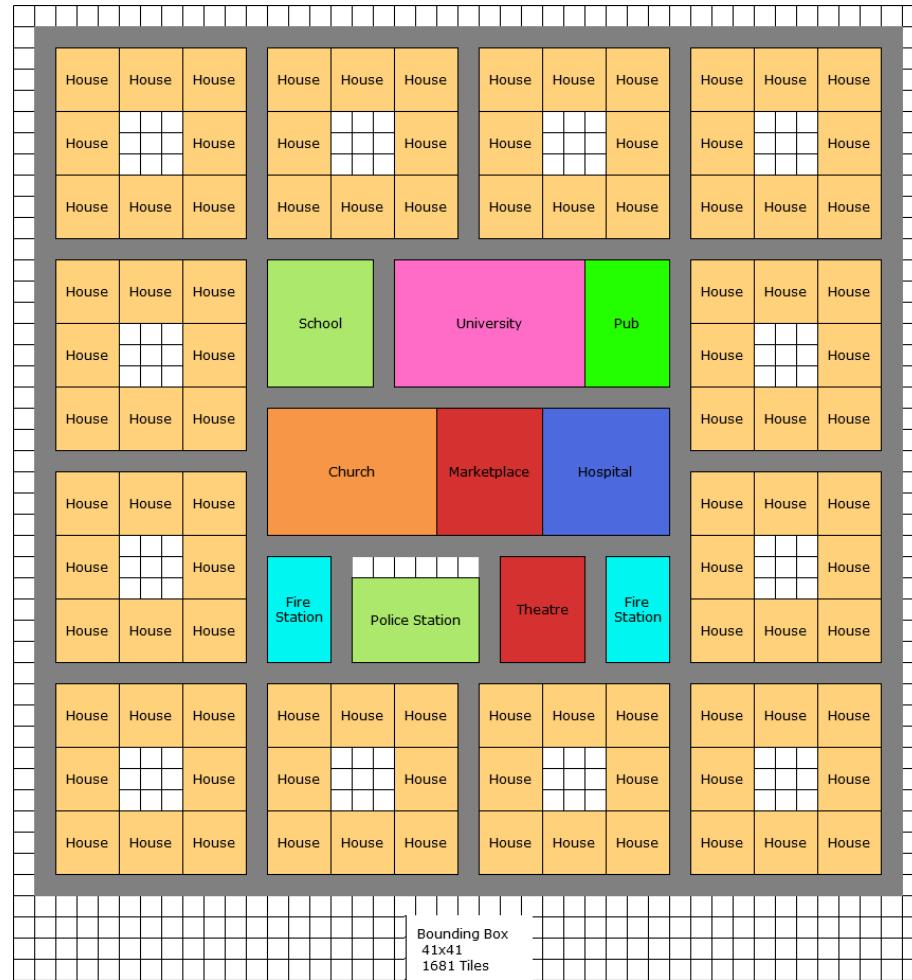
- Highway
- Main Street
- Street
- Commercial and Office space
- Path
- Water
- Park

Training Outcome? The best game strategies

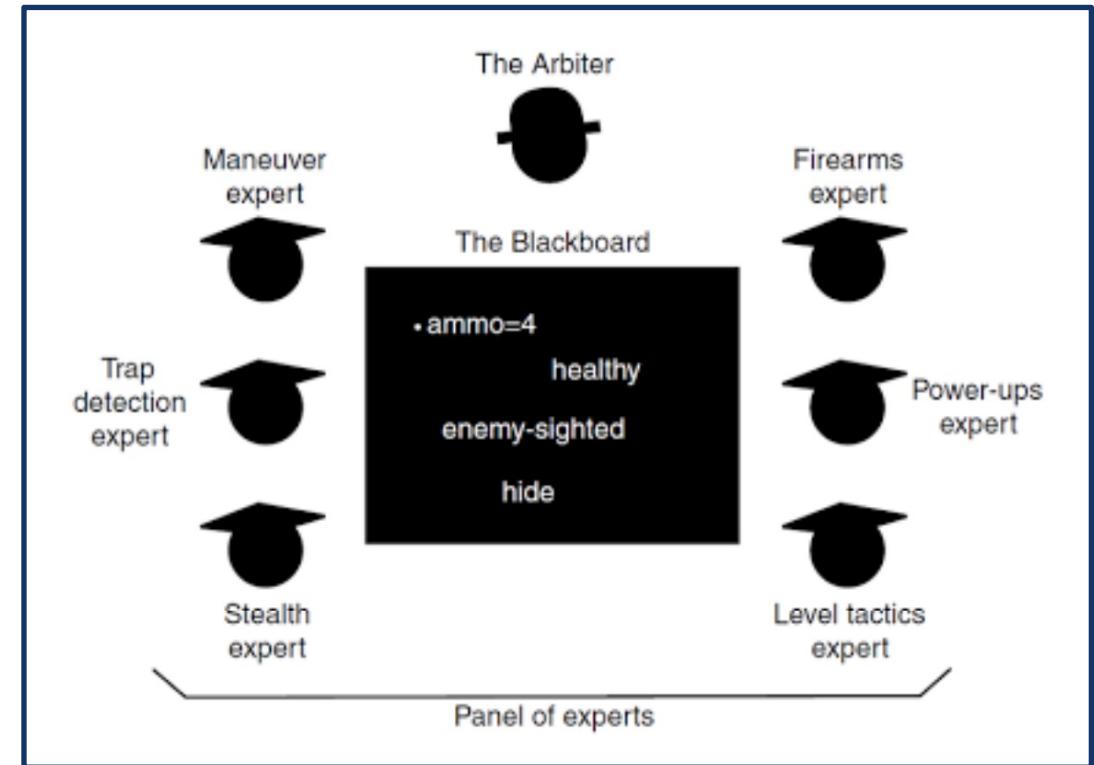
ANNO 1800



Training Outcome? The best game strategies



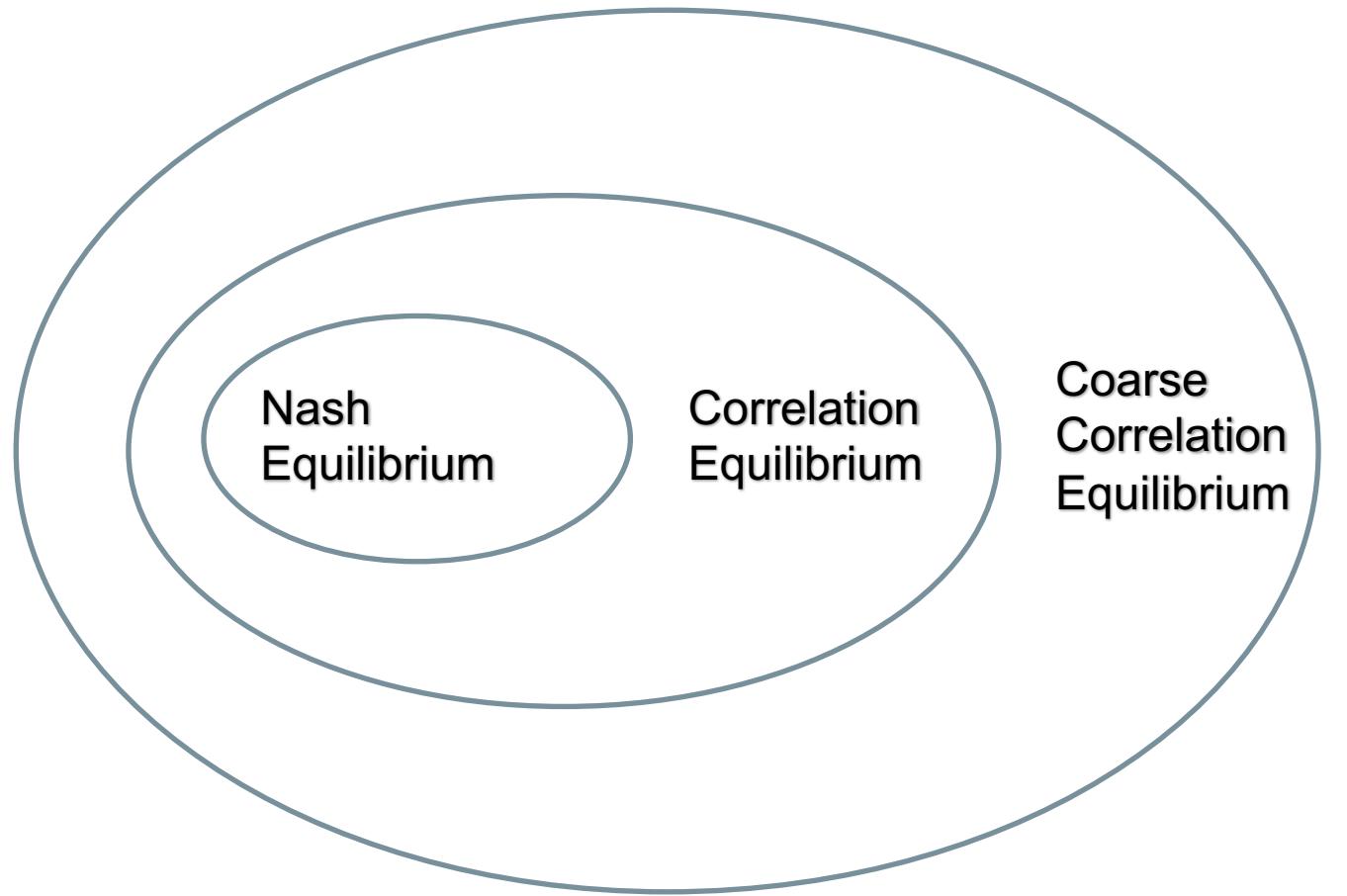
Collaboration between multiple players - Blackboard Model



Algorithms Inspired by Social Structures

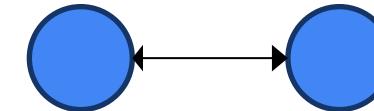
Category	Description	Illustrations	Algorithms
Independent Learning	Each agent learns independently, treating others as part of the environment.	Three separate circles, each labeled "Agent" below it, representing independent learning.	IQL, IPG, IA2C, ITRPO, IPPO
Centralized Critic	Uses a centralized critic to estimate value functions while policies remain decentralized for execution.	A central circle labeled "Central-critic" is connected to three separate circles labeled "Agent" by arrows, representing a centralized critic estimating values for decentralized agents.	COMA, MATRPO, MAPPO, MADDPG, MAA2C
Value Decomposition	Decomposes joint value functions into individual contributions to facilitate coordination among agents.	A group of four circles labeled "Joint Value" maps to three separate squares labeled "Individual contributions", representing decomposing joint value functions.	QMIX, VDN, QTRAN
Regret-Based Methods	Minimizing regret over iterations, commonly used in solving extensive-form games and multi-agent strategies.	Three circles inside an oval labeled "Min. Global Regret" with double-headed arrows between them, and a label "Min. Local Regret" below the oval, representing regret-based methods.	CFR, MC-CFR, MWU, Regret Matching, Regret Matching+

Existence of a good strategy : Game Theory 2.0

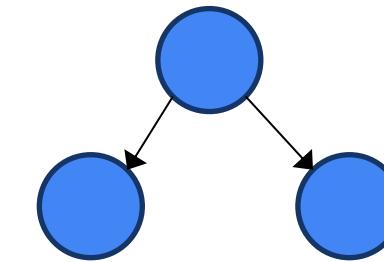


Easier to compute
(converge) !

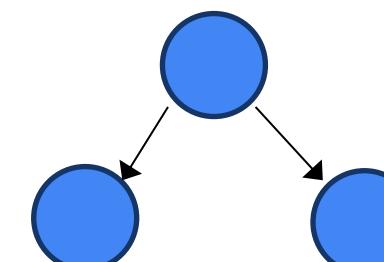
NE:



CE:



CCE:



G: 30	G: 30	G: 30	G: 30
V: 30	V: 30	V: 30	V: 30
D: 30	D: 30	D: 30	D: 30
G: 30	G: 30	G: 30	G: 30
V: 30	V: 30	V: 30	V: 30
D: 30	D: 30	D: 30	D: 30
G: 30	G: 30	G: 30	G: 30
V: 30	V: 30	V: 30	V: 30
D: 30	D: 30	D: 30	D: 30
G: 30	G: 30	G: 30	G: 30
V: 30	V: 30	V: 30	V: 30
D: 30	D: 30	D: 30	D: 30
G: 30	G: 30	G: 30	G: 30
V: 30	V: 30	V: 30	V: 30
D: 30	D: 30	D: 30	D: 30

Initial Game Board



Altruistic Player
Resources: Money: 20 Reputation: 20
Reward Function: $\alpha = 0.2, \beta = 0.8$



Balanced Player
Resources: Money: 20 Reputation: 20
Reward Weights: $\alpha = 0.5, \beta = 0.5$



Profit-Driven Player
Resources: Money: 20 Reputation: 20
Reward Function: $\alpha = 0.8, \beta = 0.2$

Players



Park
Cost: Money -1, Reputation -3
Revenue: Money -1/t, Reputation +3/t
Effects: Greenery +, Vitality -



House
Cost: Money -2, Reputation -2
Revenue: Money +2/t, Reputation 0/t
Effects: Greenery -, Density +

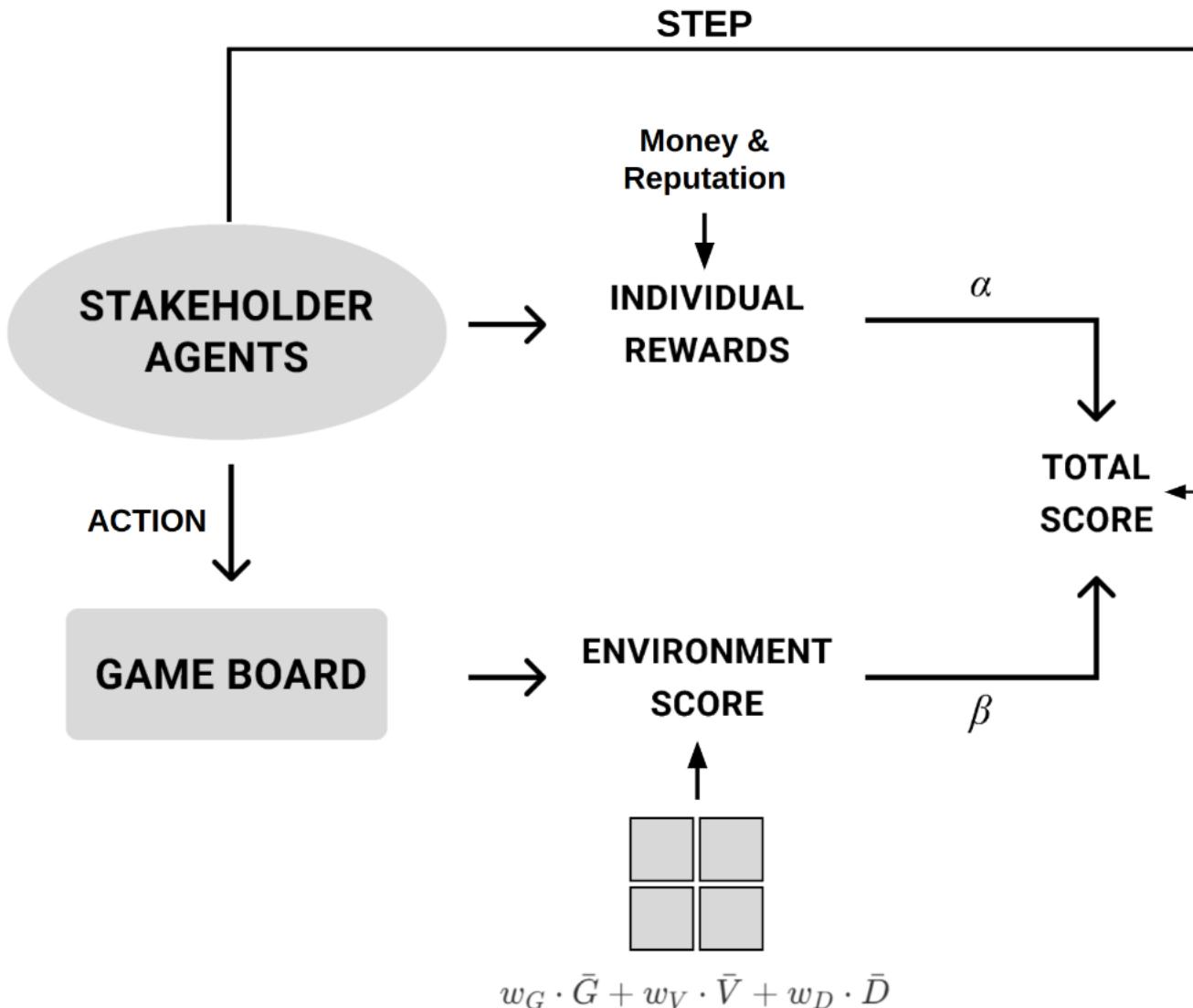


Shops
Cost: Money 3, Reputation 1
Revenue: Money 3/t, Reputation -1/t
Effects: Vitality +, Density -

Properties

Players and Properties

Case Study: Minimalist Simcity



Environment Parameters

Avg. Greenery : 30

Env Score:

Avg. Vitality: 30

30

Avg. Density : 30



Personal: 0

Combined: 24



Personal: 0

Combined: 15



Personal: 0

Combined: 6

G: 30
V: 30
D: 30

Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

30



Personal: 0

Combined: 24



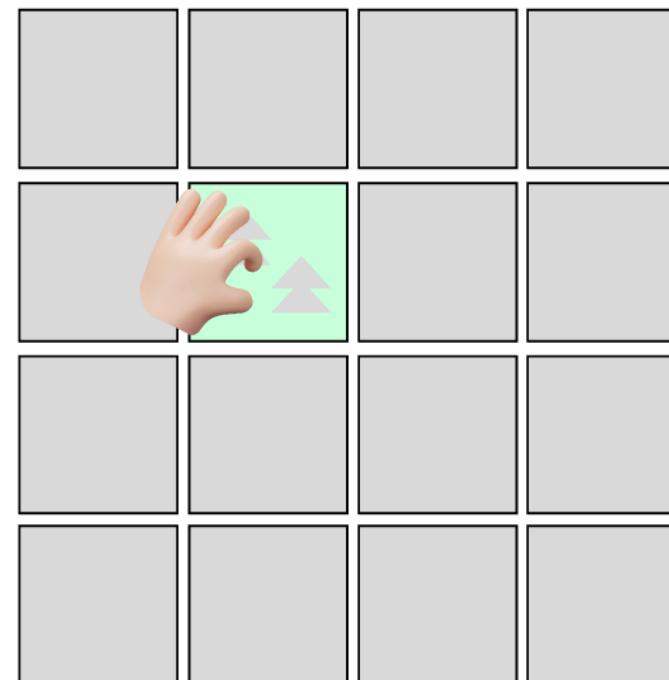
Personal: 0

Combined: 15



Personal: 0

Combined: 6



Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

35



Personal: 12

Combined: 30



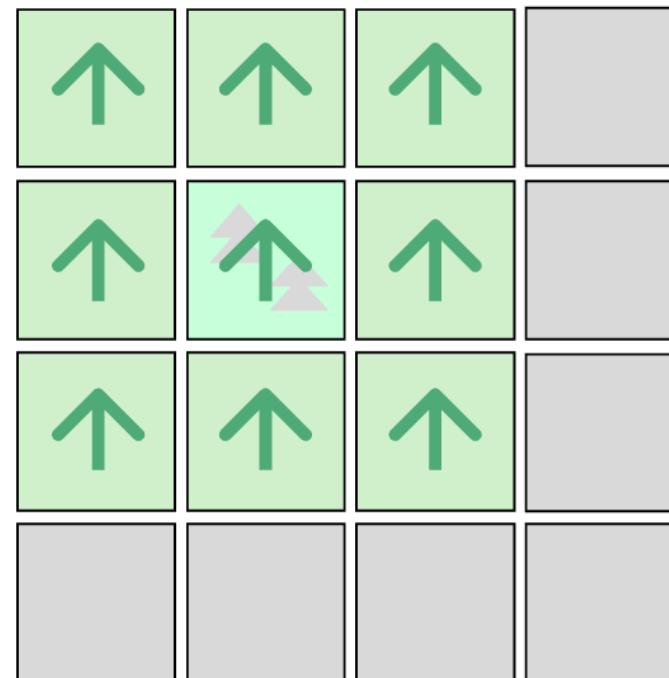
Personal: 10

Combined: 20



Personal: 8

Combined: 8



Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

35



Personal: 12

Combined: 30



Personal: 10

Combined: 20



Personal: 8

Combined: 8

G: 40
V: 30
D: 20

G: 50
V: 30
D: 10

G: 40
V: 30
D: 20

Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

35



Personal: 12

Combined: 30



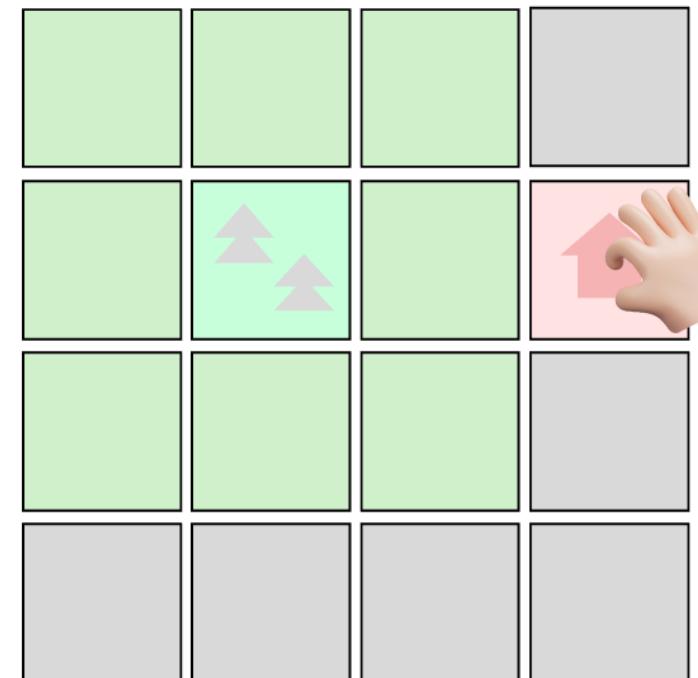
Personal: 10

Combined: 20



Personal: 8

Combined: 8



Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

38



Personal: 10

Combined: 26



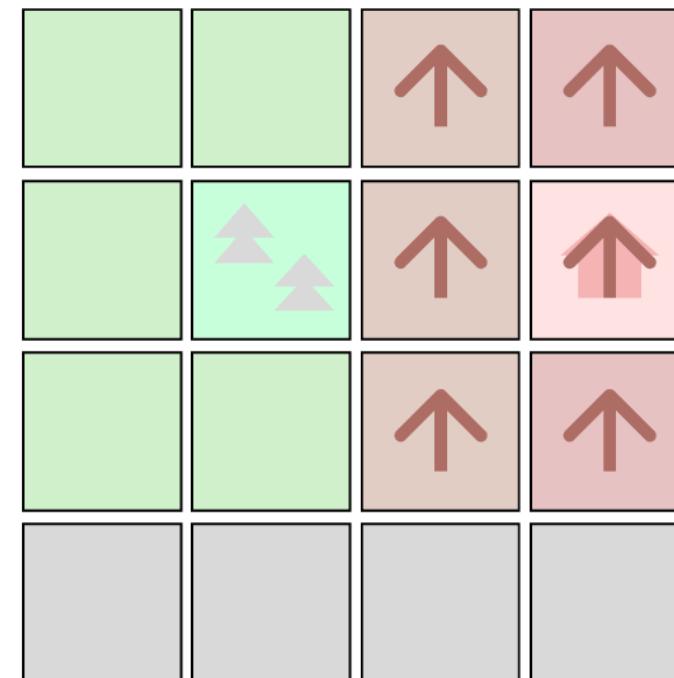
Personal: 12

Combined: 24



Personal: 20

Combined: 20



Environment Parameters

Avg. Greenery : 30

Avg. Vitality: 30

Avg. Density : 30

Env Score:

38



Personal: 10

Combined: 26



Personal: 12

Combined: 24



Personal: 20

Combined: 20



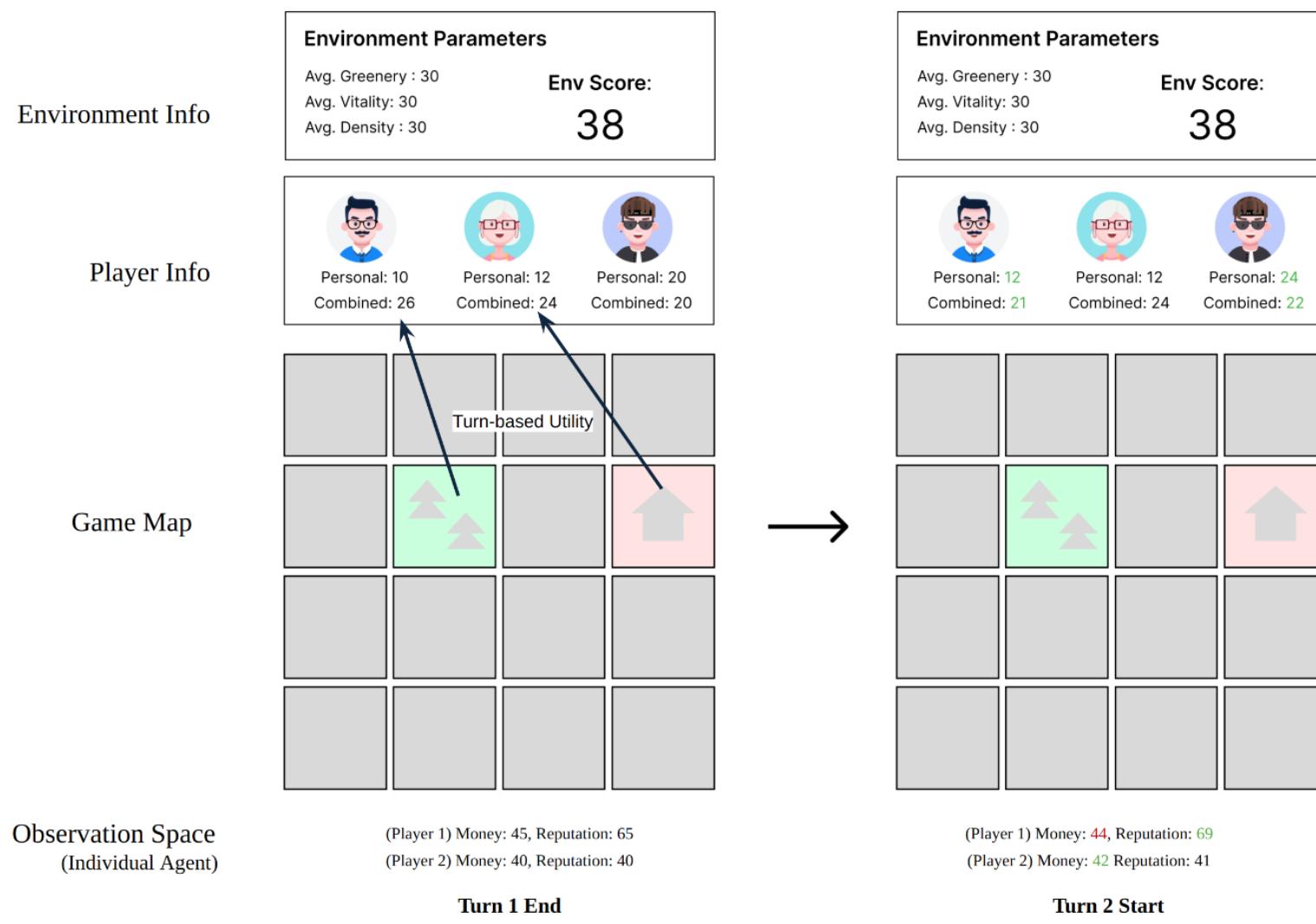


Figure 4: Turn Based Utilities

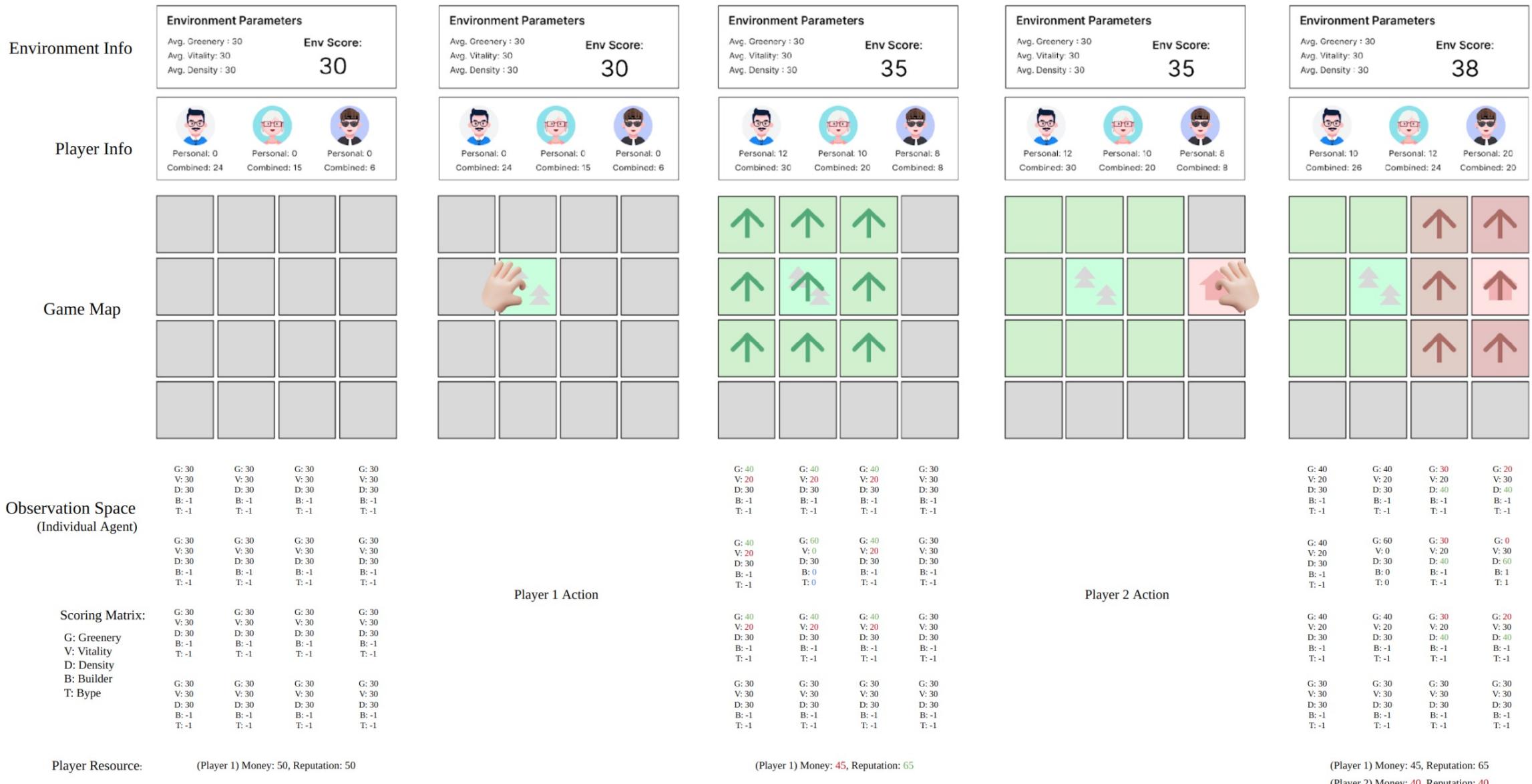


Figure 7: Observation Space for Agents

Grid Observations:

$$\begin{bmatrix} -35 & 70 & 0 \\ -15 & 50 & 20 \\ 5 & 30 & 30 \\ 35 & 0 & 20 \end{bmatrix} \quad \begin{bmatrix} -45 & 80 & -10 \\ -45 & 50 & 30 \\ 15 & 20 & 30 \\ 5 & 30 & 0 \end{bmatrix} \quad \begin{bmatrix} -55 & 90 & 0 \\ -35 & 70 & 20 \\ -15 & 50 & 30 \\ 25 & 10 & 30 \end{bmatrix} \quad \begin{bmatrix} -35 & 40 & 40 \\ -35 & 70 & 0 \\ 5 & 30 & 30 \\ -5 & 10 & 60 \end{bmatrix}$$

Builder Locations:

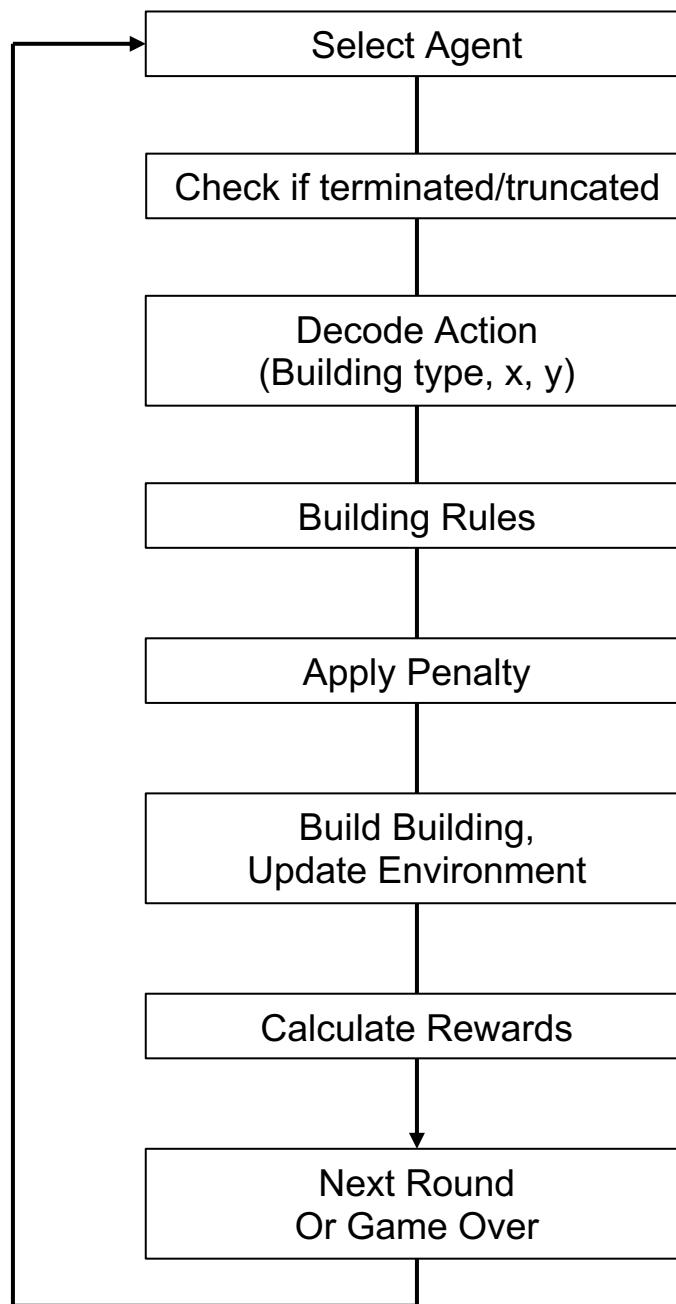
$$\begin{bmatrix} 0 & -1 & -1 & 2 \\ 0 & 2 & -1 & 1 \\ 0 & -1 & -1 & 2 \\ 1 & 1 & -1 & 1 \end{bmatrix}$$

Building Types:

$$\begin{bmatrix} 2 & -1 & -1 & 0 \\ 2 & 1 & -1 & 2 \\ 2 & -1 & -1 & 0 \\ 1 & 2 & -1 & 1 \end{bmatrix}$$

Player Resources:

$$\begin{bmatrix} \text{Money: 84} \\ \text{Reputation: 54} \end{bmatrix}$$

**Algorithm 1:** Step Function in SimCity Environment

```

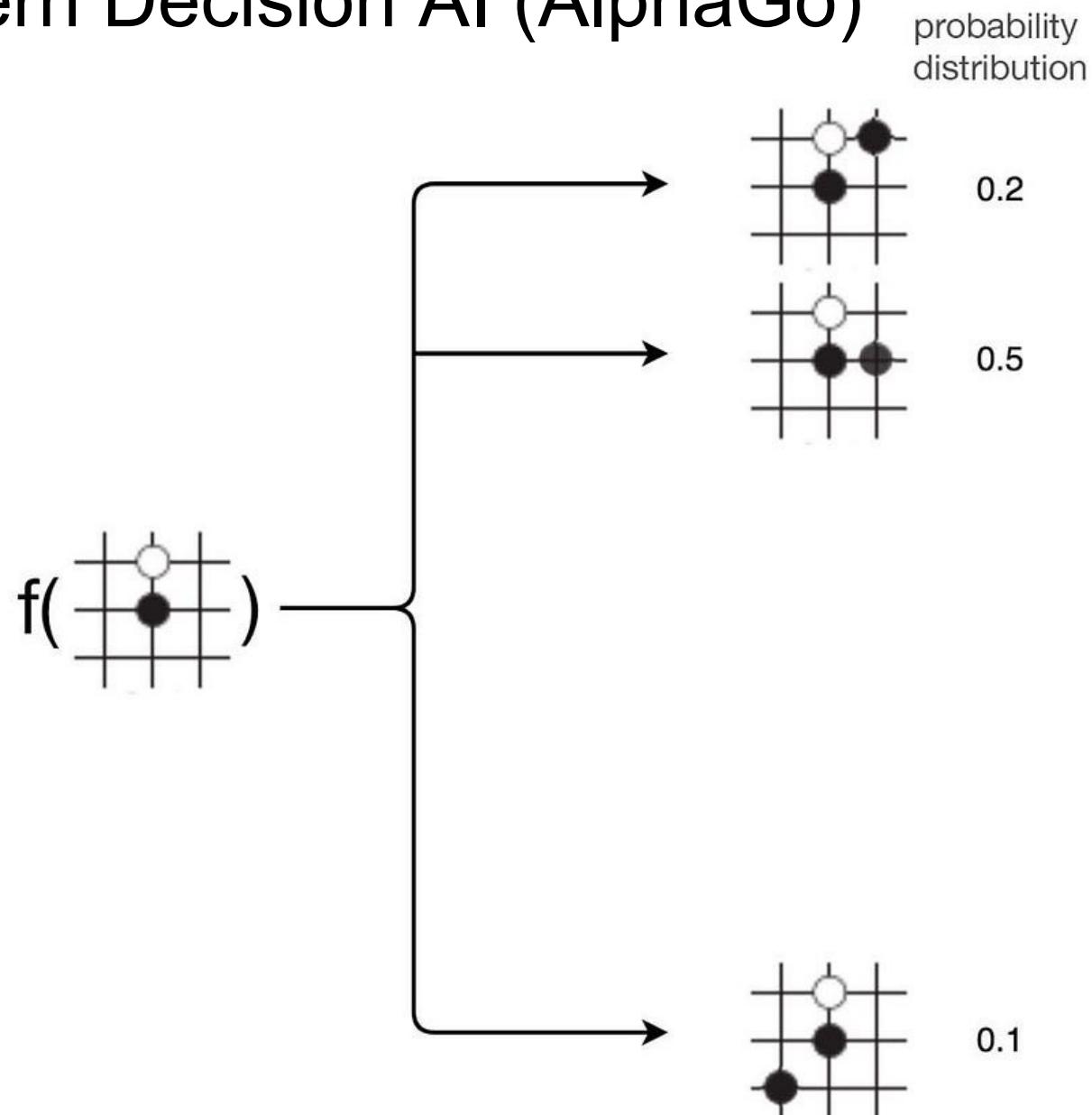
1 Input
2   ┌ Current Agent Action action, Training Mode (maximize individual reward or common reward)
3 Output
4   ┌ Updated Environment State(obs), Rewards, Terminated, Truncated, and Info

5 Initialize agent  $\leftarrow$  self.agent_selection;
6 if self.terminations[agent] or self.truncations[agent] then
7   ┌ return;
8 Decode Action: (building_type, x, y)  $\leftarrow$  decode_action(action);
9 if action = NO_OP then
10  ┌ No operation performed by agent;
11 else
12  ┌ if self.buildings[x][y] is not None then
13    ┌ Apply penalty for building on occupied cell;
14    ┌ reward  $\leftarrow$  reward + build_on_occupied_penalty;
15  else
16    ┌ if Agent lacks resources to build then
17      ┌ Apply penalty for insufficient resources;
18      ┌ reward  $\leftarrow$  reward + not_enough_resource_penalty;
19  else
20    ┌ Deduct resources from agent;
21    ┌ Update environment with new building;
22    ┌ Adjust grid and neighboring cells based on building effects;
23    ┌ reward  $\leftarrow$  reward + immediate_reward;

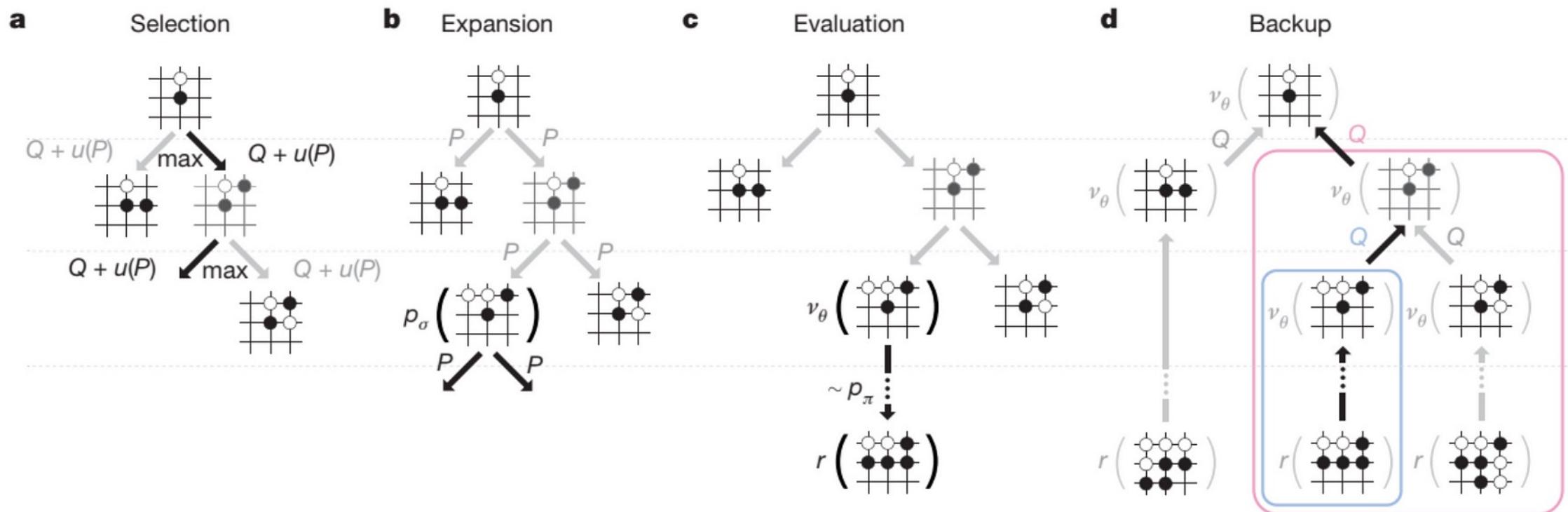
24 Update Utilities: Adjust player scores and resources based on building utilities;
25 Calculate Environment Score: env_score  $\leftarrow$  calculate_environment_score();
26 Determine Player Reward Weights Based on Type :  $(\alpha, \beta) \leftarrow$  weights based on player type;
27 Update Integrated Score:
28   self.players[agent].integrated_score  $\leftarrow$   $\alpha \times$  self.players[agent].self_score +  $\beta \times$  env_score;
29 Log Rewards: Compute individual reward for agent and common reward as the sum of individual reward.
30 Update Rewards: if mode = maximize individual reward then
31   ┌ reward  $\leftarrow$  each agent's reward (tuple);
32 else
33   ┌ reward  $\leftarrow$  common reward (float);
34 Increment Move Counter: self.num_moves  $\leftarrow$  self.num_moves + 1;
35 Check Termination Condition: self.terminations updated if game is over;
36 Select Next Agent: self.agent_selection  $\leftarrow$  self._agent_selector.next();
  
```

Category	Graph	Legend
Independent Learning		<ul style="list-style-type: none"> Actor-Critic (IA2C) – No Shared Information Actor-Critic (IA2C) – Shared Information IPPO – No Shared Information IPPO – Shared Information Q-Learning (IQL) – No Shared Information Q-Learning (IQL) – Shared Information
Centralized Critic		<ul style="list-style-type: none"> COMA – No Shared Information COMA – Shared Information MAA2C – No Shared Information MAA2C – Shared Information MAPPO – No Shared Information MAPPO – Shared Information
Value Decomposition		<ul style="list-style-type: none"> QMIX – No Shared Information QMIX – Shared Information VDN – No Shared Information VDN – Shared Information

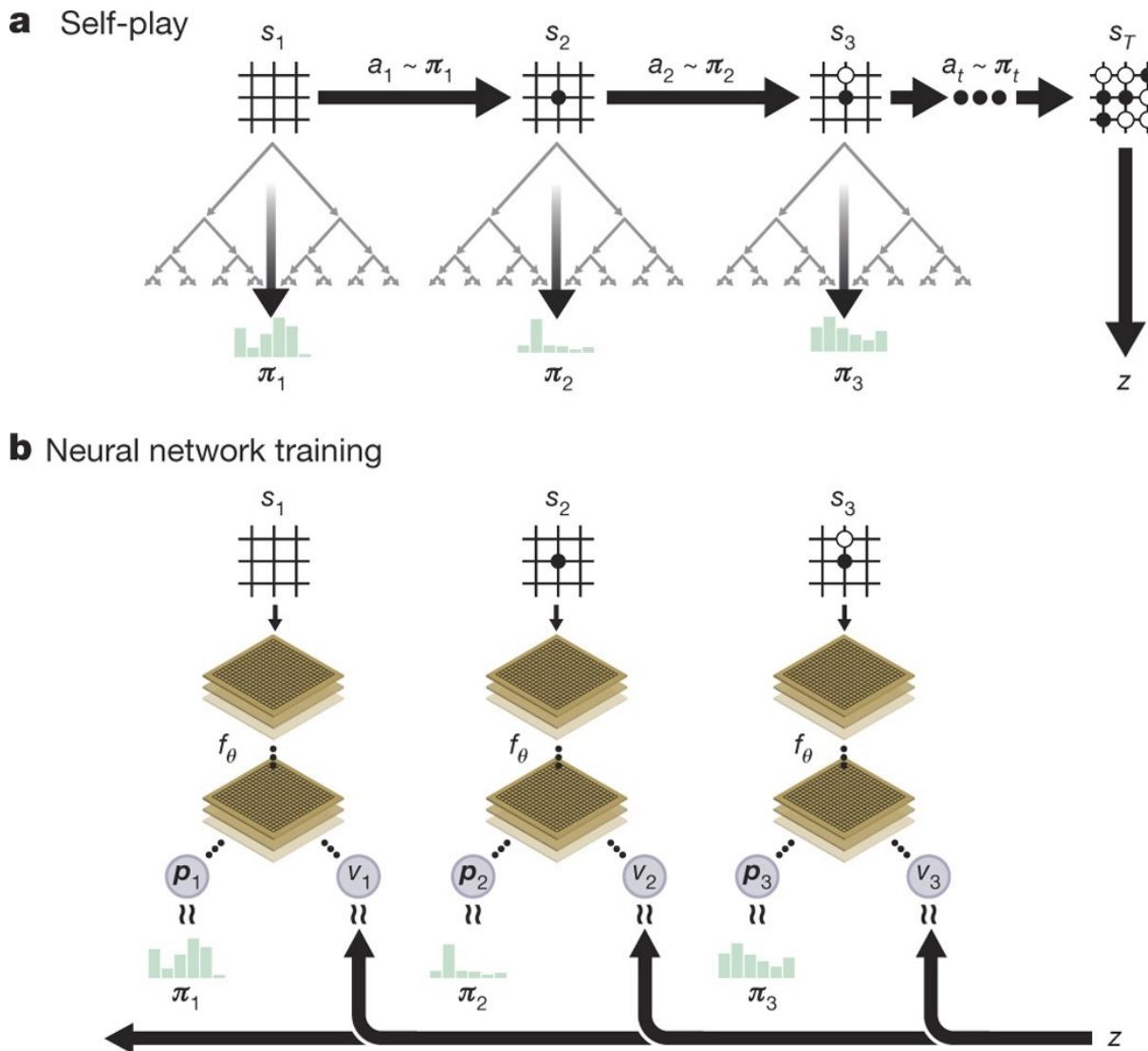
Reflect on Modern Decision AI (AlphaGo)



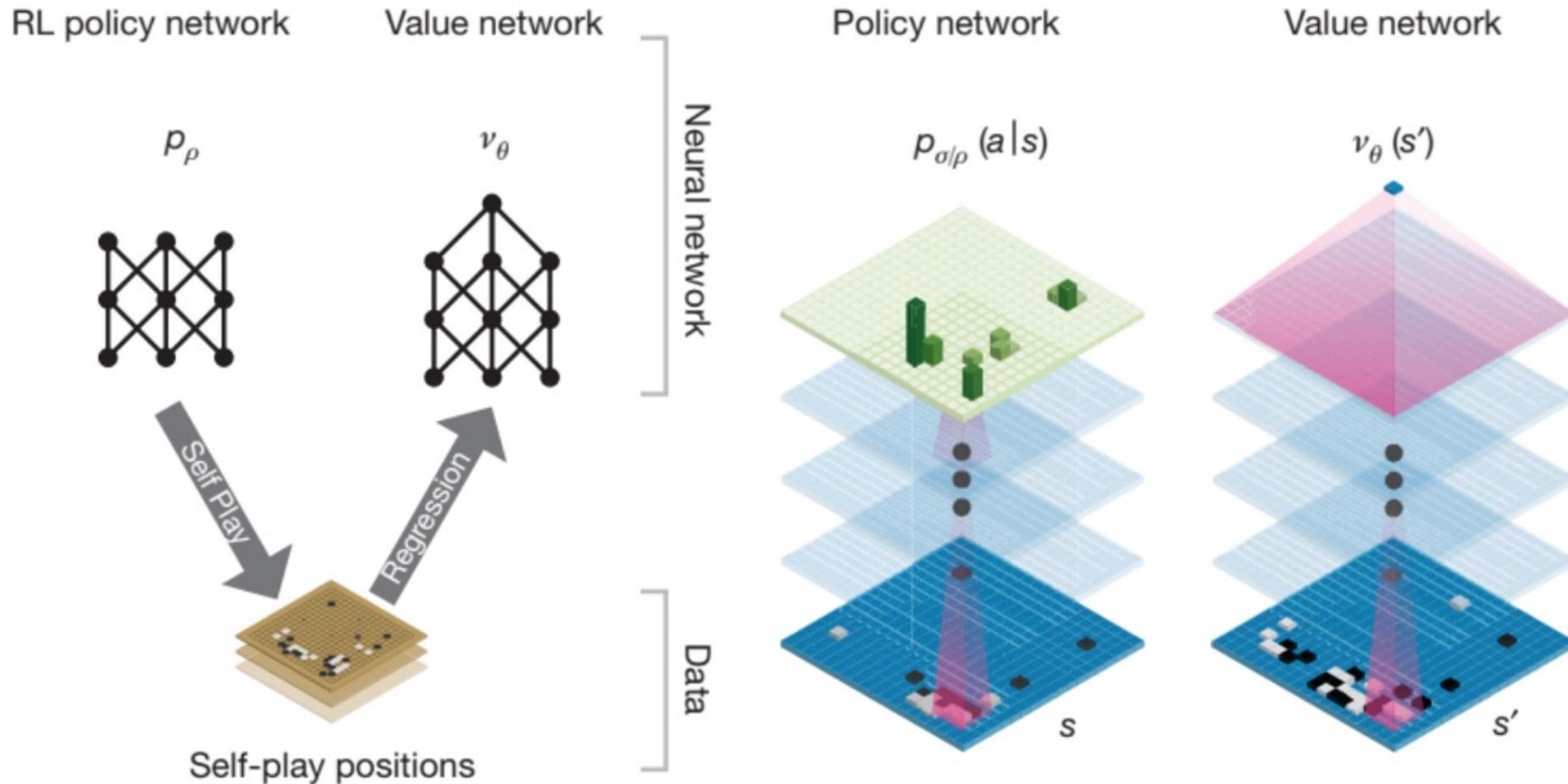
Reflect on Modern Decision AI (AlphaGo)



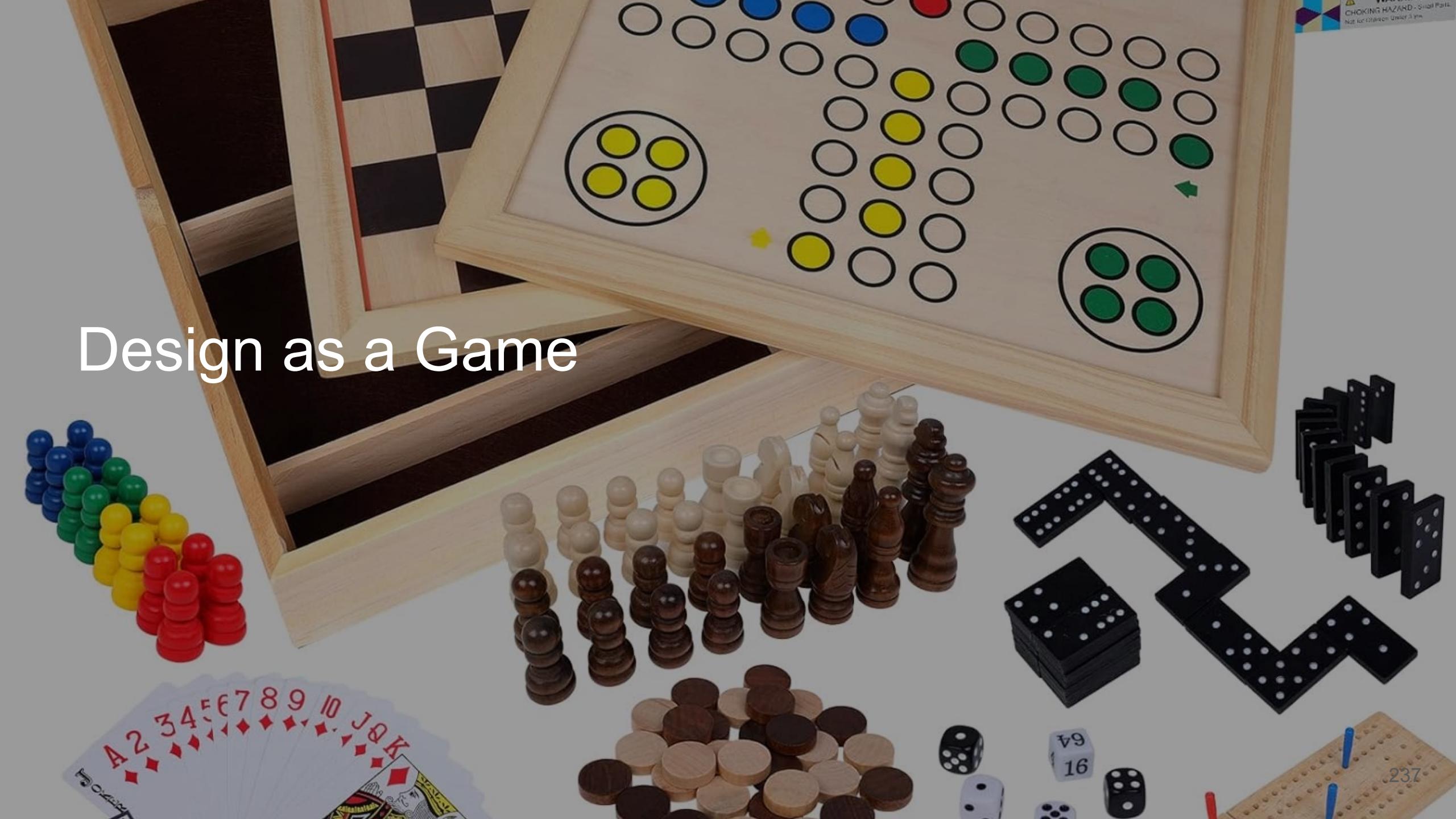
Reflect on Modern Decision AI (AlphaGo)



Reflect on Modern Decision AI (AlphaGo)



Design as a Game





Thank you for joining us!

Designing Consensus: Gamified Modeling and Simulation of Collaborative Decision-Making

Instructor: Jin Gao

IAP 2025 (Non-Credit)

gaojin@mit.edu

SCAN ME



Thank you for joining us!