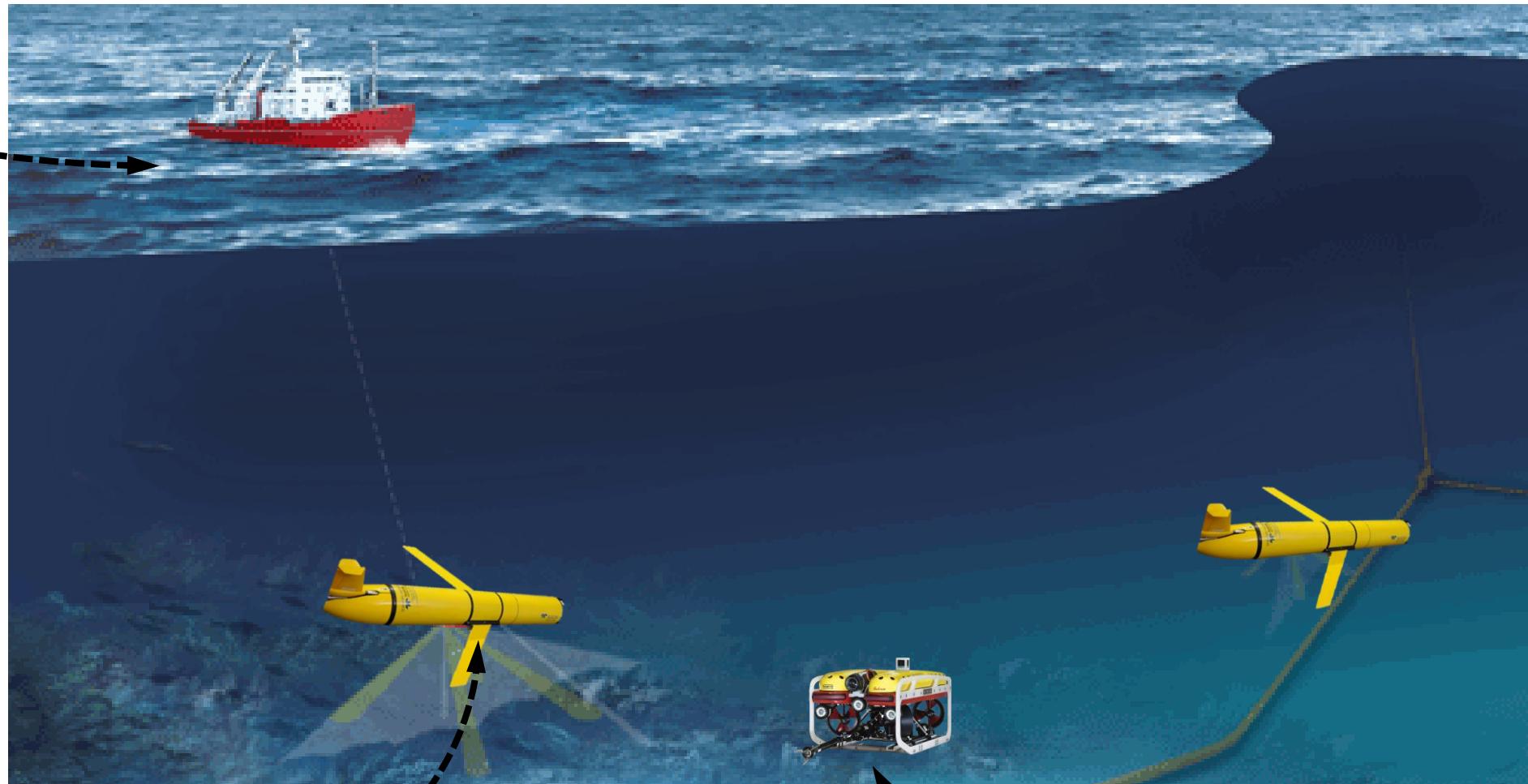


Monte Carlo Tree Search for Collaboration

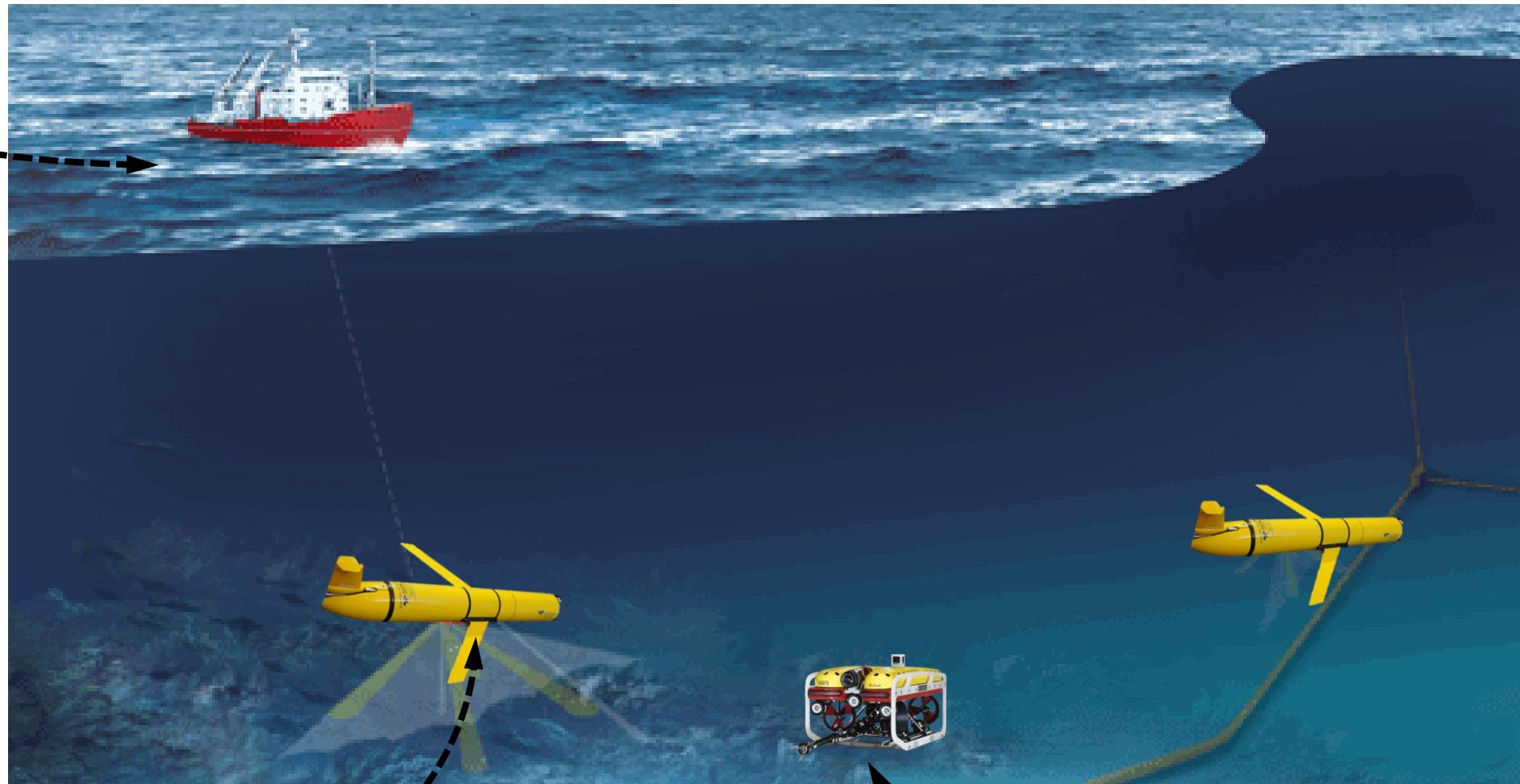
Can Pu, Filippos Sotiropoulos, Gabriel Margolis, Jacob Guggenheim, Zachary Duguid



**Autonomous Underwater
Vehicles**

**ROV exploration
vehicle**

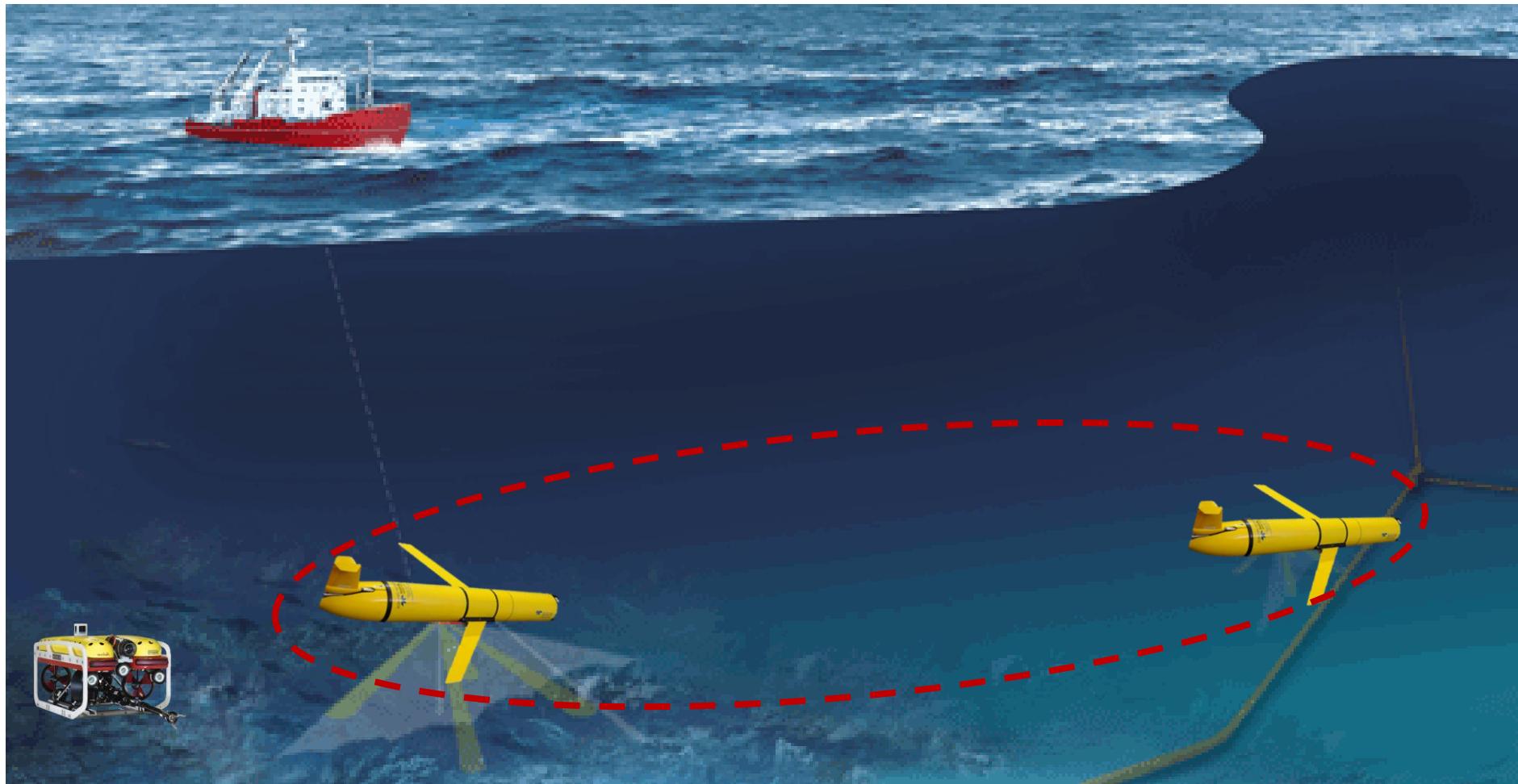
Rough
Mapping



Explore areas of
interest

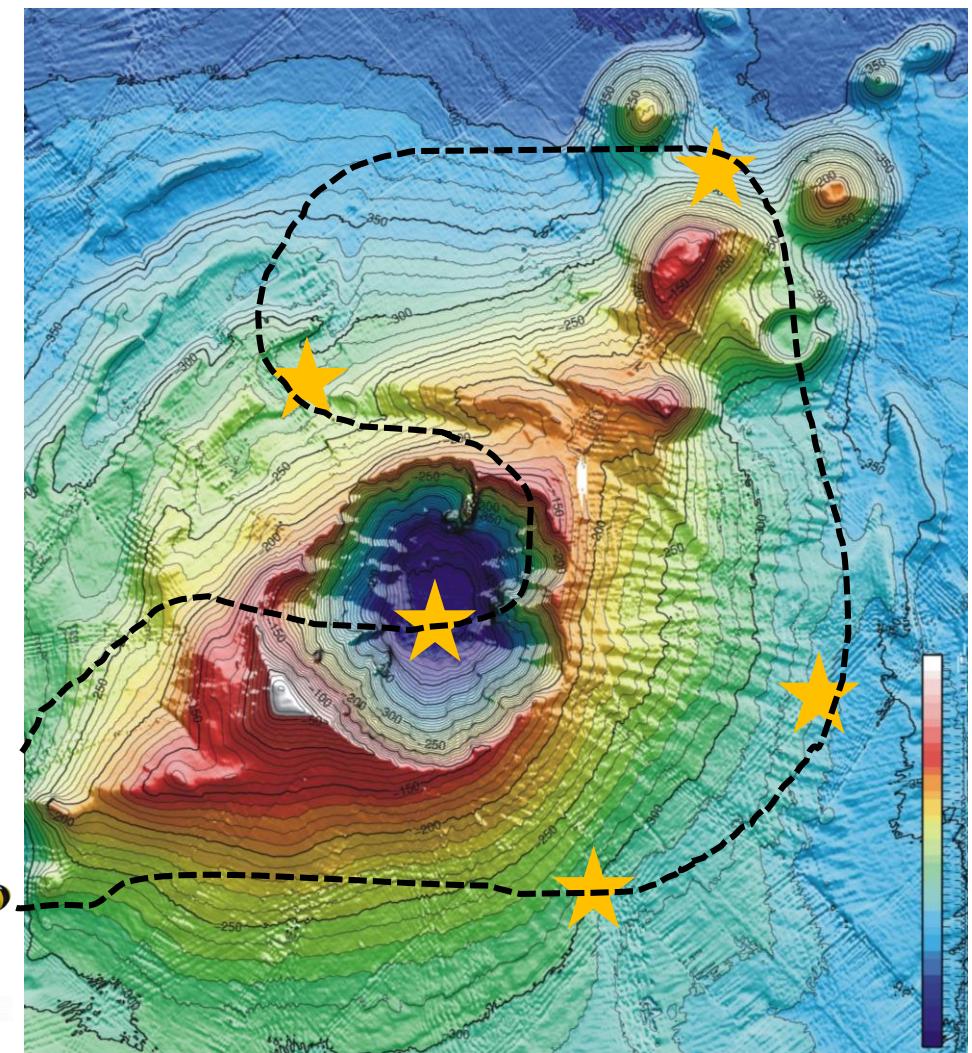
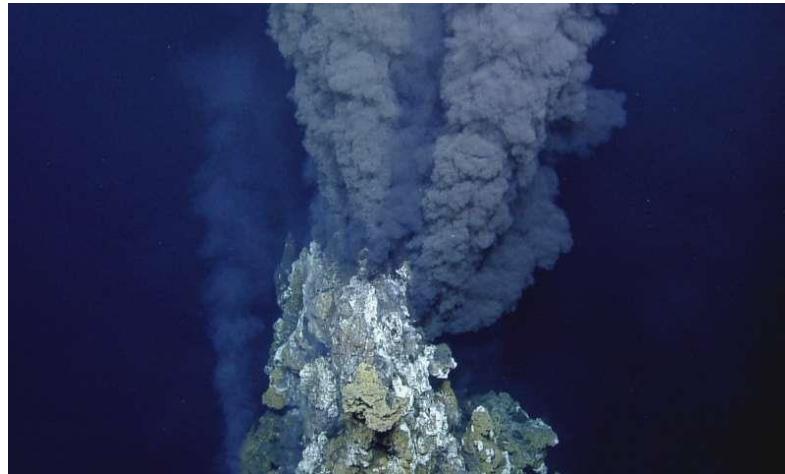
Gather samples and other
scientific tasks

Scenario

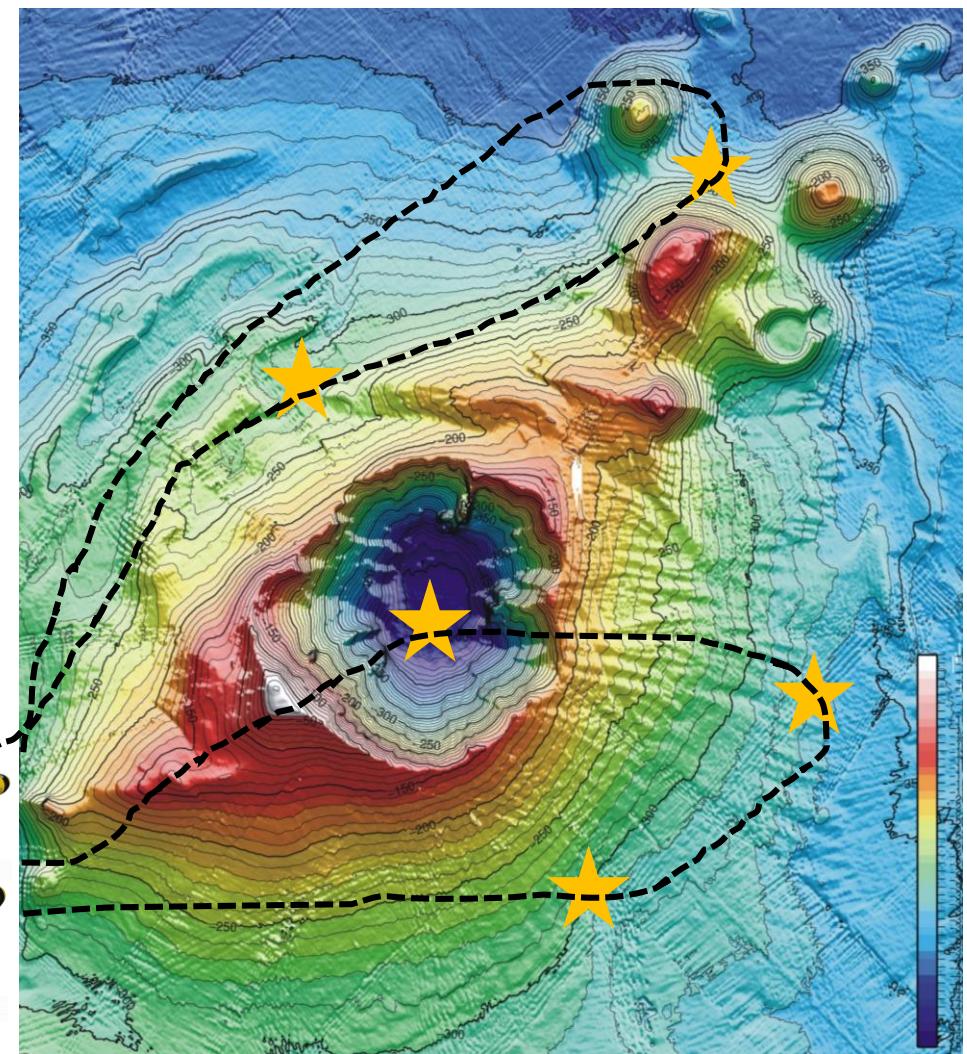
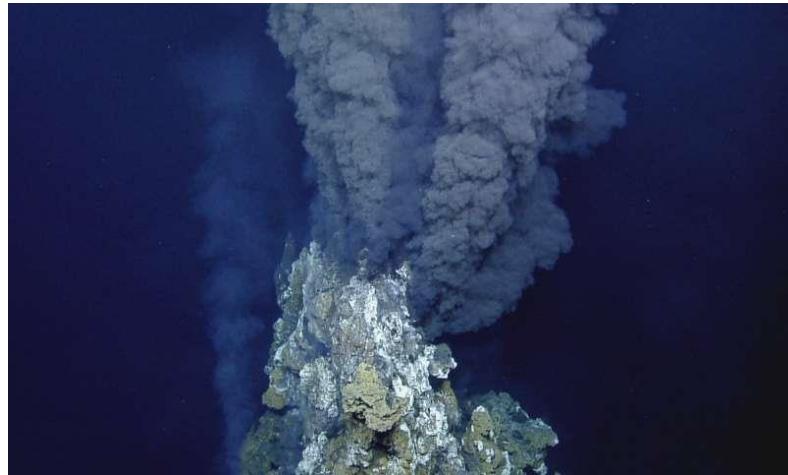


Focus on the AUV navigation

- AUV navigating through the Kolumbo volcanic crater.
- Trying to collect scientific information
- There may be many areas of value to visit



- AUV navigating through the Kolumbo volcanic crater.
- Trying to collect scientific information
- There may be many areas of value to visit
 - Using multiple AUV's can save time



Key features of the problem:

- Multiple agents coordinating their navigation to achieve a common objective
- Objective: explore areas with varying amounts of scientific interest.
 - Can be interpreted as a **reward**

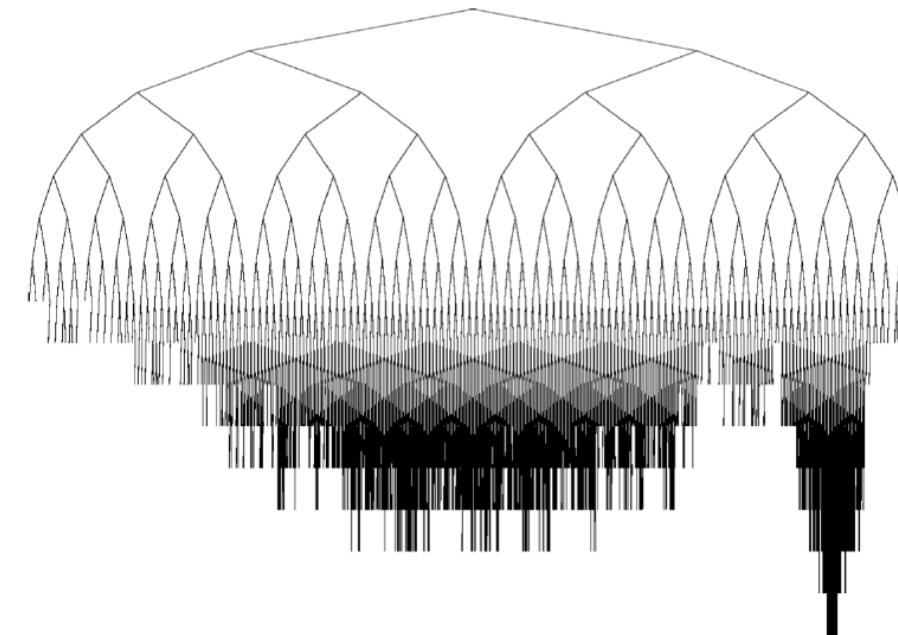
This problem can naturally be described in the *Markov Decision Process* framework but:

- State space to plan in very large because:
 - Many moves per agent (b) *and* many agents (n)
 - Distributed reward requires substantial planning depth (d)

$$O(b^{n \times d})$$



We will illustrate an algorithm, **Monte Carlo Tree Search**, (MCTS) which can help tackle these large dimensional problems.



Some examples of MCTS



Mario



Ms. Pacman

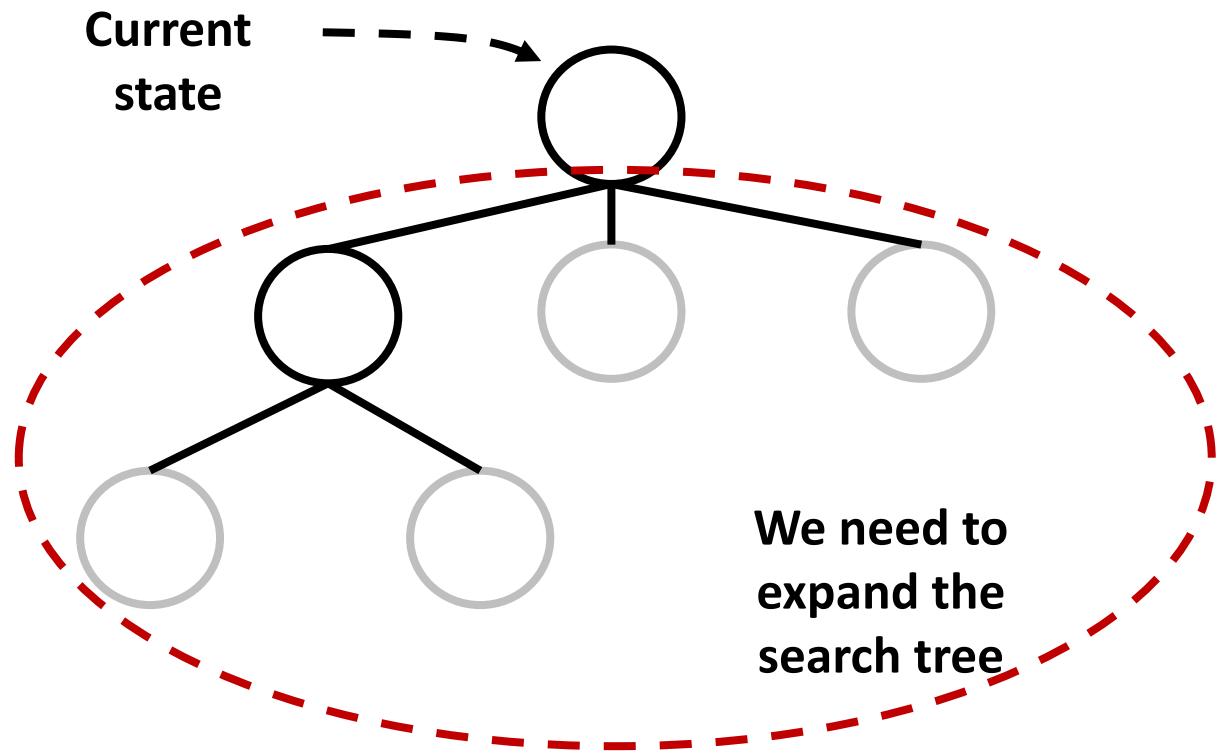
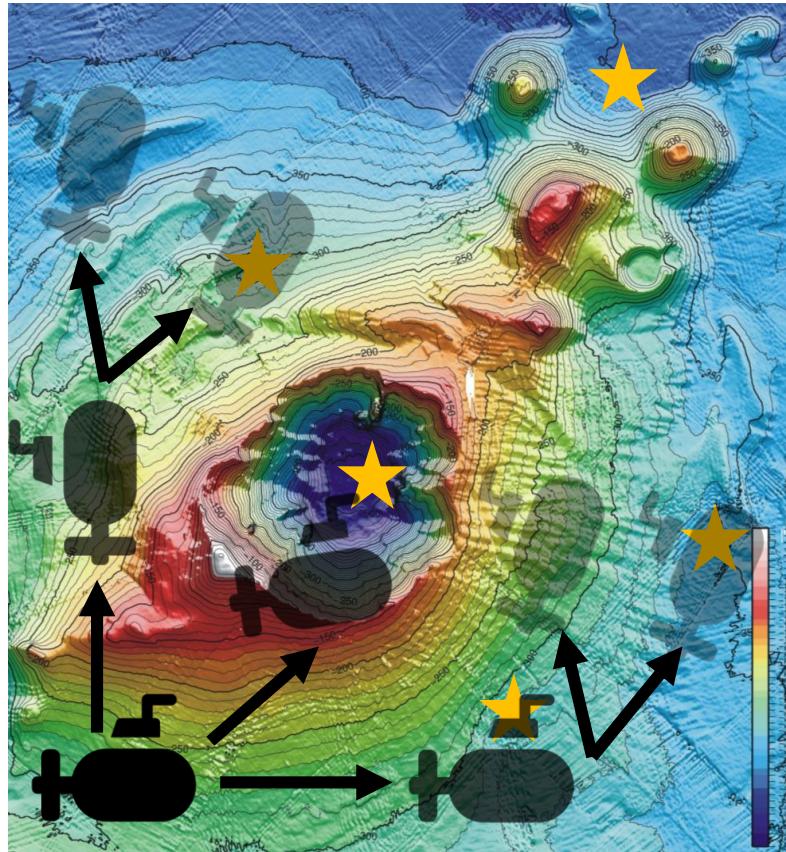


Rome: Total War

- Background
 - Monte Carlo Simulation
 - Upper Confidence Bounds
- Implementing MCTS with Upper Confidence Bounds for Trees
- Collaborative MCTS
 - Discuss Max-N, a variant of MCTS for multiple agents
 - Discuss how to decentralize MCTS
 - Discuss how to improve the default policy in decentralized MCTS

Path Planning Problem

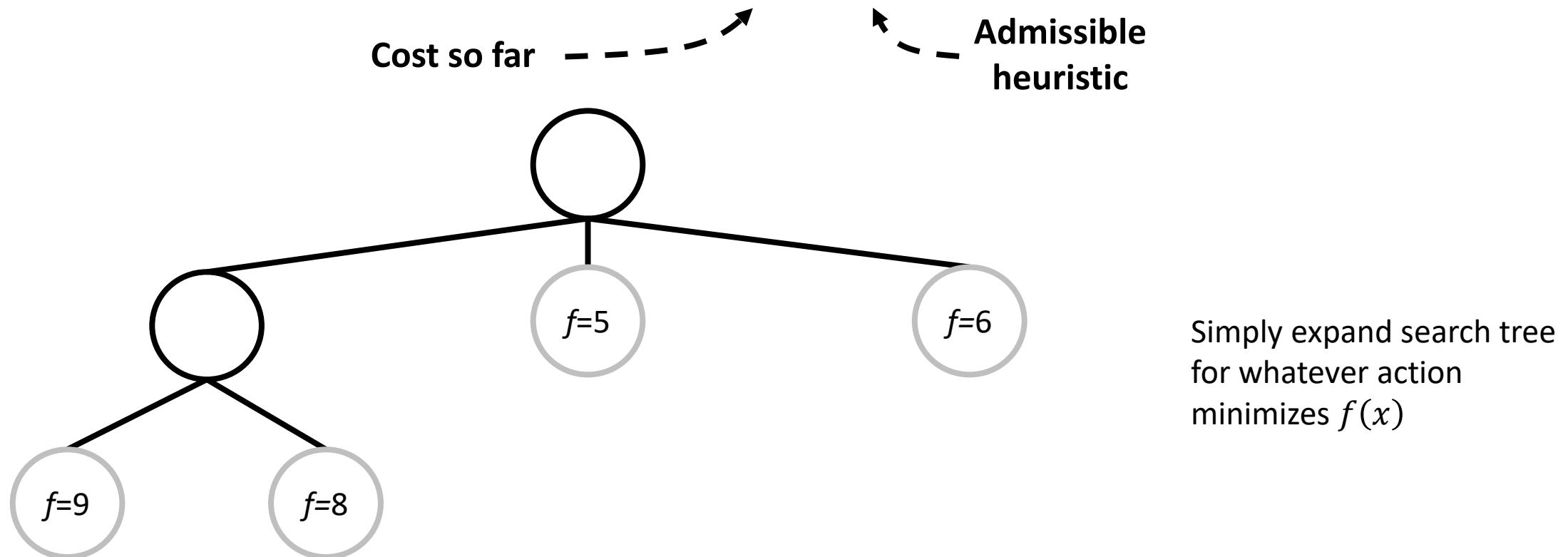
- The problem of planning can be described in terms of a decision tree



What happens if we try to solve a problem like this with a conventional heuristic search method: A*

Reminder: A* chooses a path to minimize the evaluation function:

$$f(x) = g(x) + h(x)$$

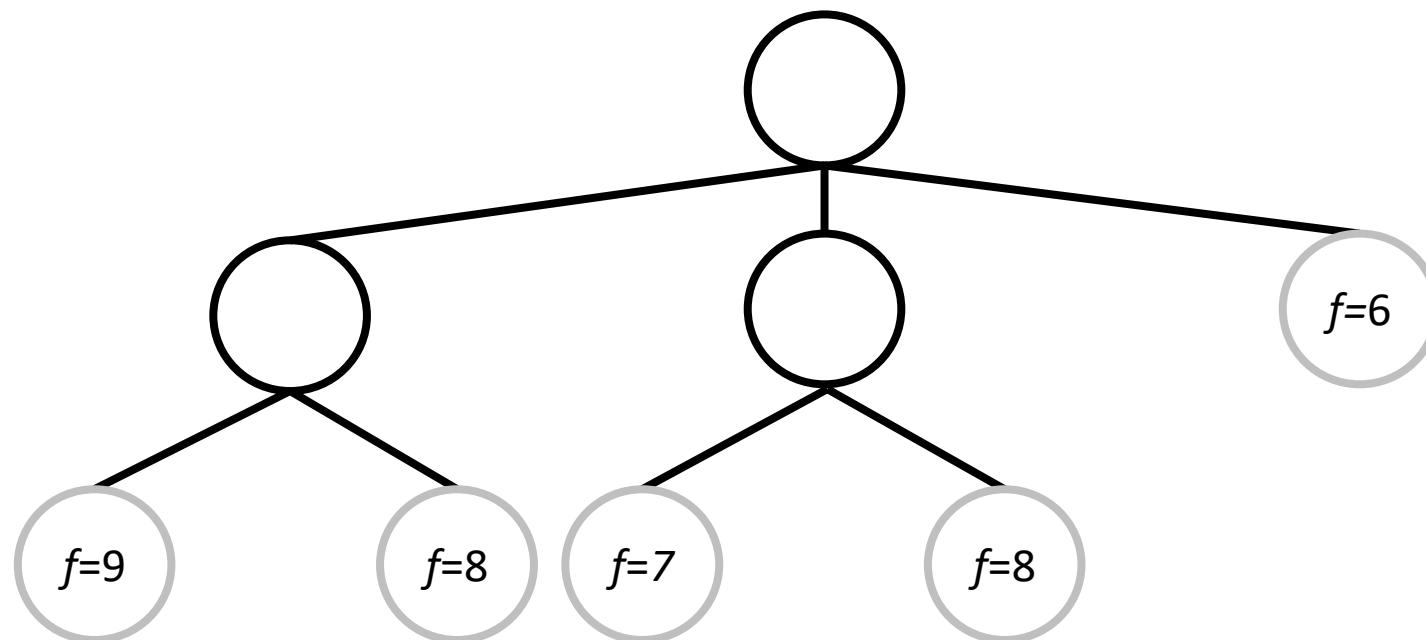


What happens if we try to solve a problem like this with a conventional heuristic search method: A*

Reminder: A* chooses a path to minimize the evaluation function:

$$f(x) = g(x) + h(x)$$

Cost so far Admissible heuristic



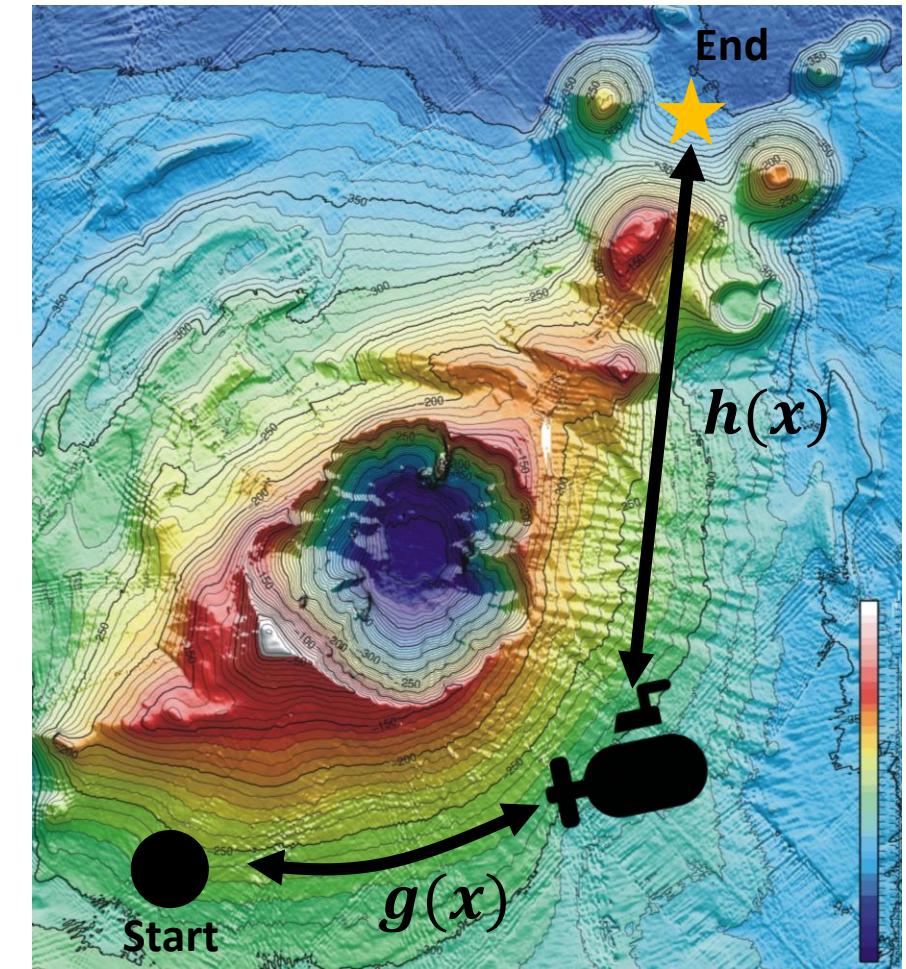
Simply expand search tree
for whatever action
minimizes $f(x)$

$$f(x) = g(x) + h(x)$$

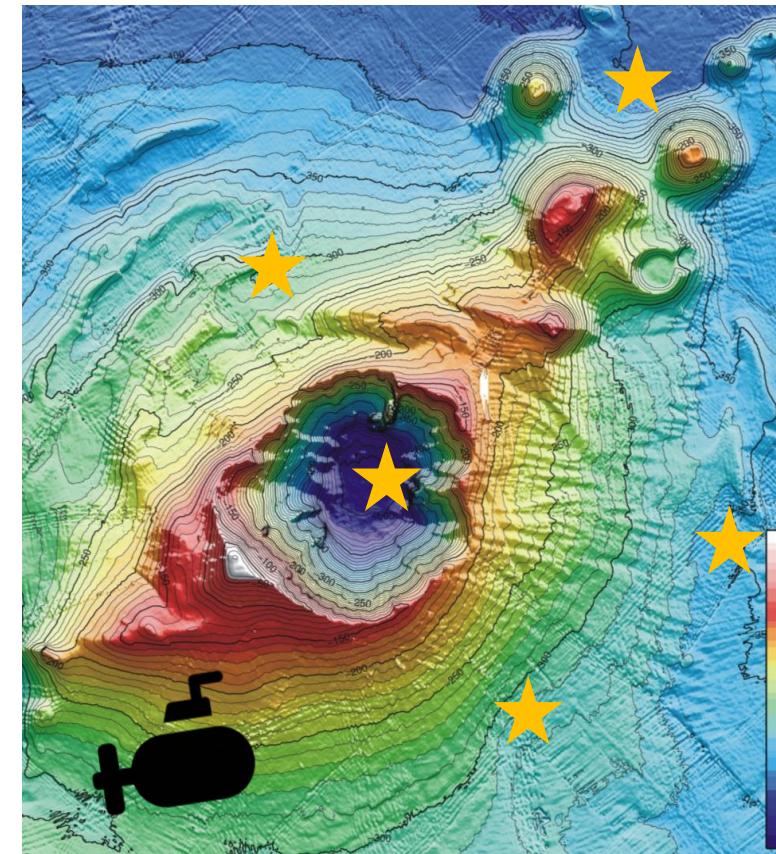
Cost so far Admissible heuristic

What heuristic ?

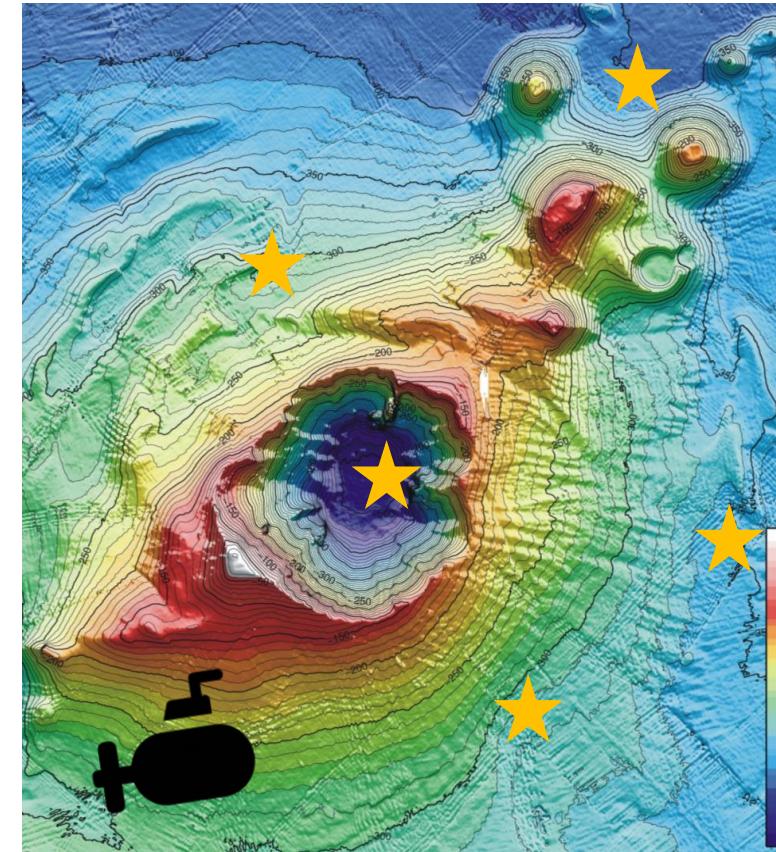
For a simple path planning problem one could use:
Euclidean distance to end goal



What heuristic should we use?



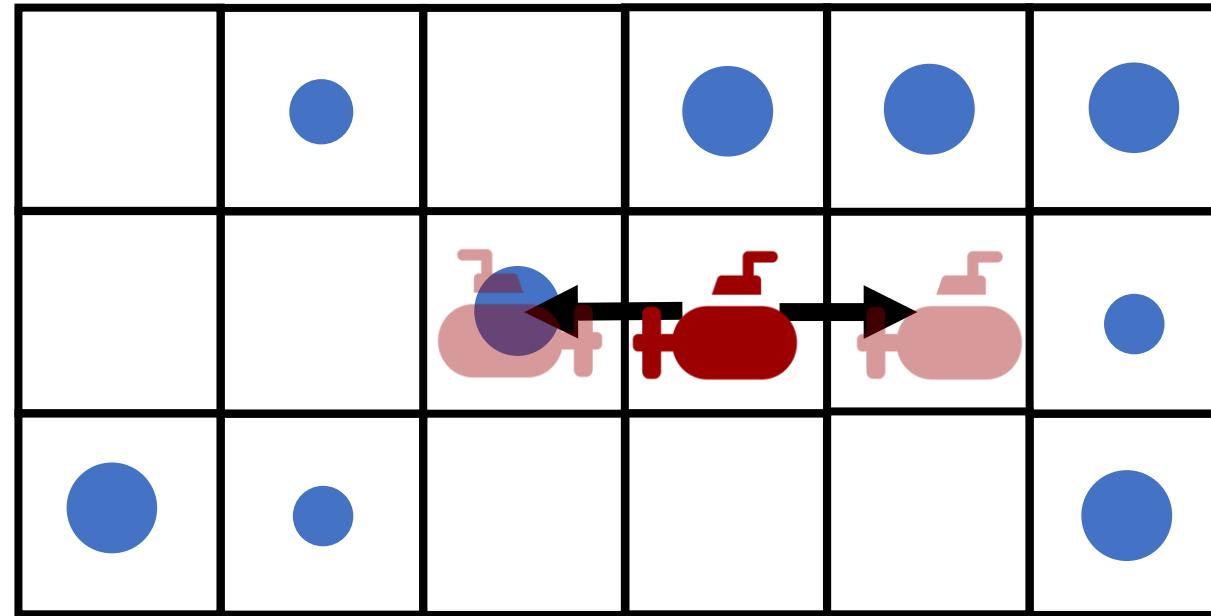
What heuristic should we use?



Finding a heuristic can be problematic!

Monte Carlo Simulation

Monte Carlo Methods: One can estimate the expected value of potential future reward from a state via random sampling of future moves.

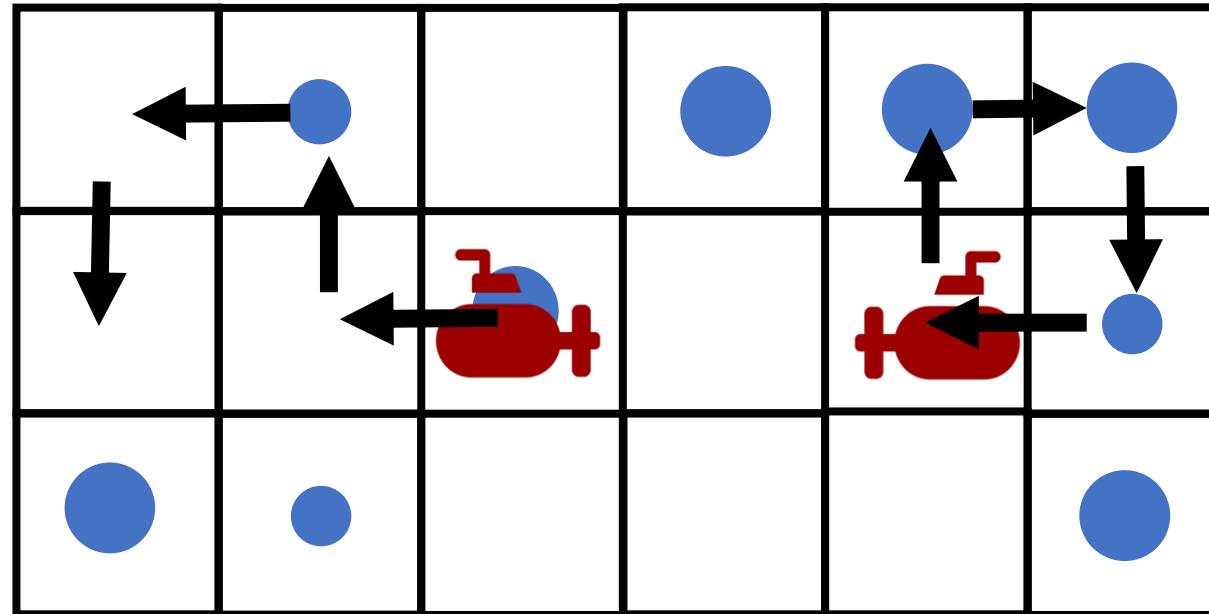


What move is a better strategy?

Monte Carlo Methods: One can estimate the expected value of potential future reward from a state via random sampling of future moves.



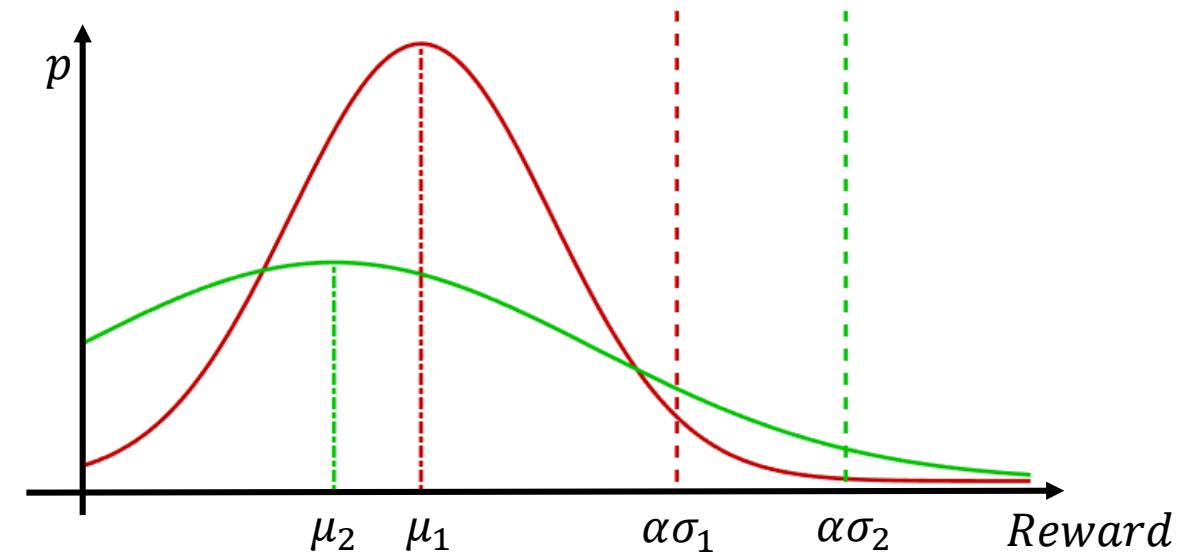
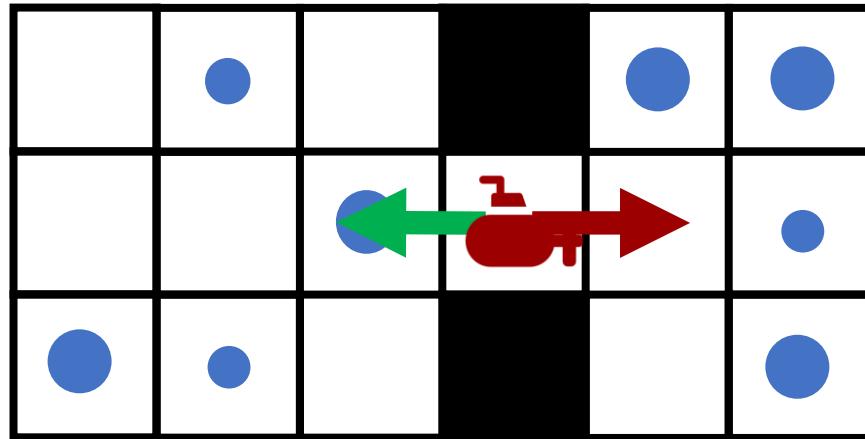
Random simulation can be used to ascertain the quality of a state: a heuristic!



What move is a better strategy?

But now random sampling introduces a new problem: Our sampled reward has uncertainty.
How do we pick the best way to expand the search tree?

Lets take the simplest case: two possible actions sampled with assumed Gaussian distribution.



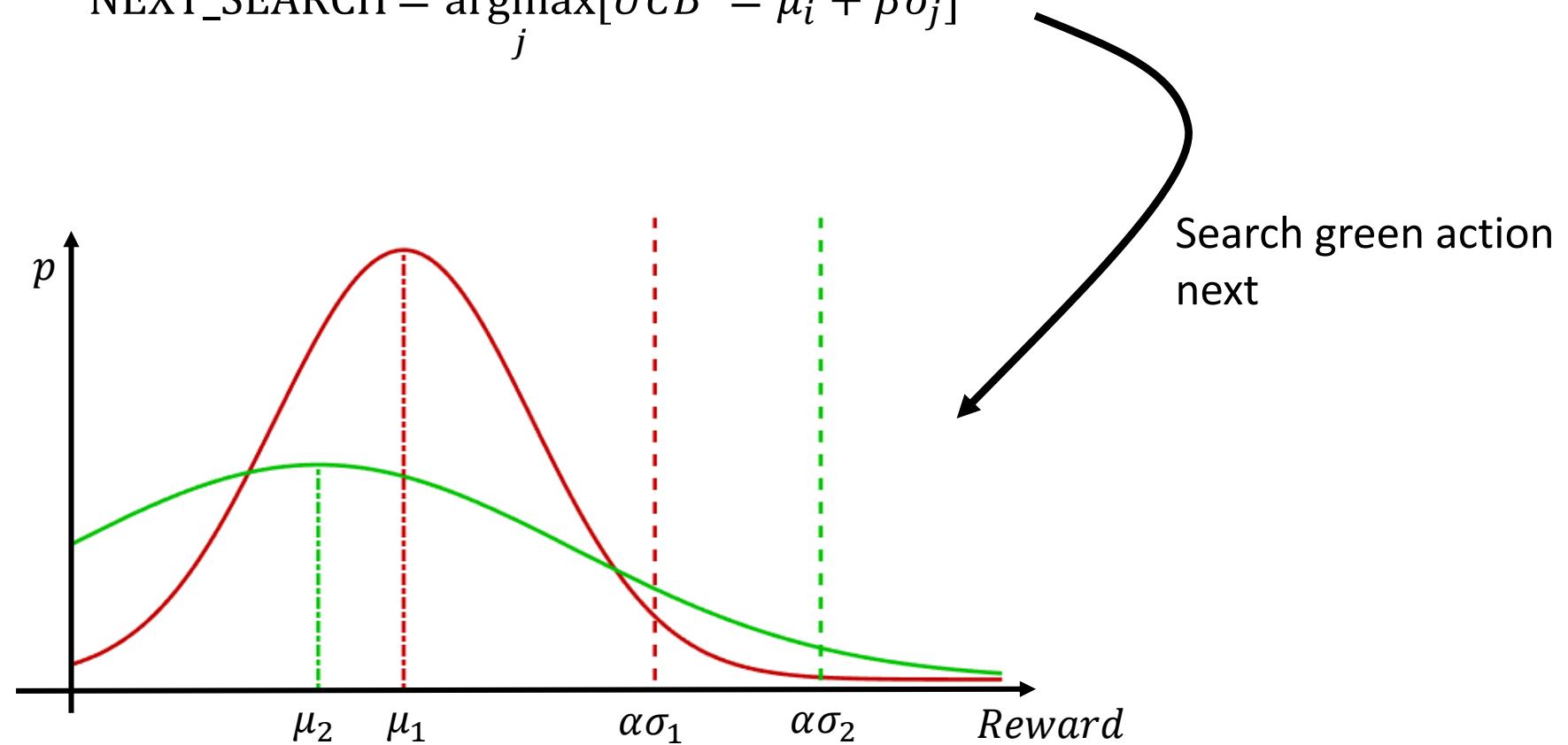
What should we sample next?

Balance searching in most promising direction (**Exploitation**) but also search where we are unsure (**Exploration**).

Upper Confidence Bounds

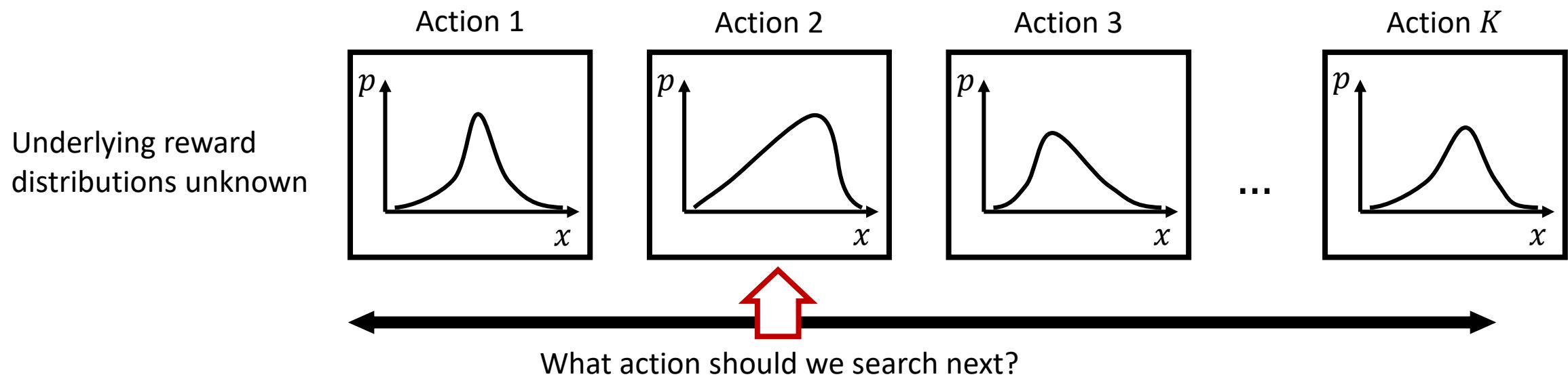
Idea: let's continue searching the action with the highest upper confidence bound. For gaussian:

$$\text{NEXT_SEARCH} = \underset{j}{\operatorname{argmax}} [UCB = \mu_j + \beta \sigma_j]$$



Upper Confidence Bounds

- One can sequentially sample one of K actions which yield a reward



- The policy should aim to minimize the *regret* after n searches:

$$R_N = \mu^* n - \mu_j \sum_{j=1}^K \mathbb{E}[T_j(n)]$$

- It represents the loss from not choosing the best action.

Upper Confidence Bounds

Since the distributions aren't necessarily Gaussian we can use the Chernoff-Hoeffding inequality:

$$P[\mathbb{E}[X_j] > \bar{X}_{j,t} + u_j] \leq e^{-2n_j u_j^2}$$



Set to a small threshold which we decrease with each sample



$$u_j = \sqrt{\frac{2 \ln n}{n_j}}$$

$$\text{UCB1} = \bar{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

subscript j : each possible action
 X_j : Random variable for reward
 u_j : Upper confidence bound
 n_j : number of samples

Sample action which maximizes:

$$\text{UCB1} = \bar{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

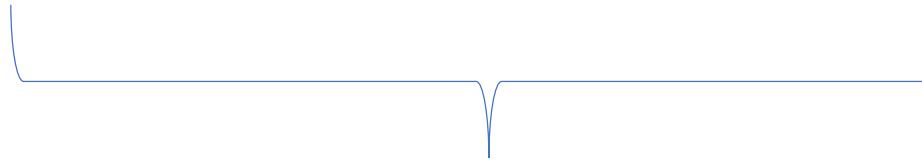
Exploitation

Mean reward from that action

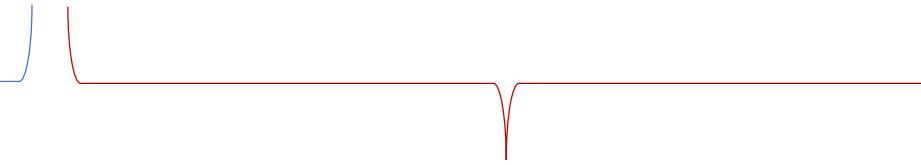
Exploration

Higher for poorly sampled actions

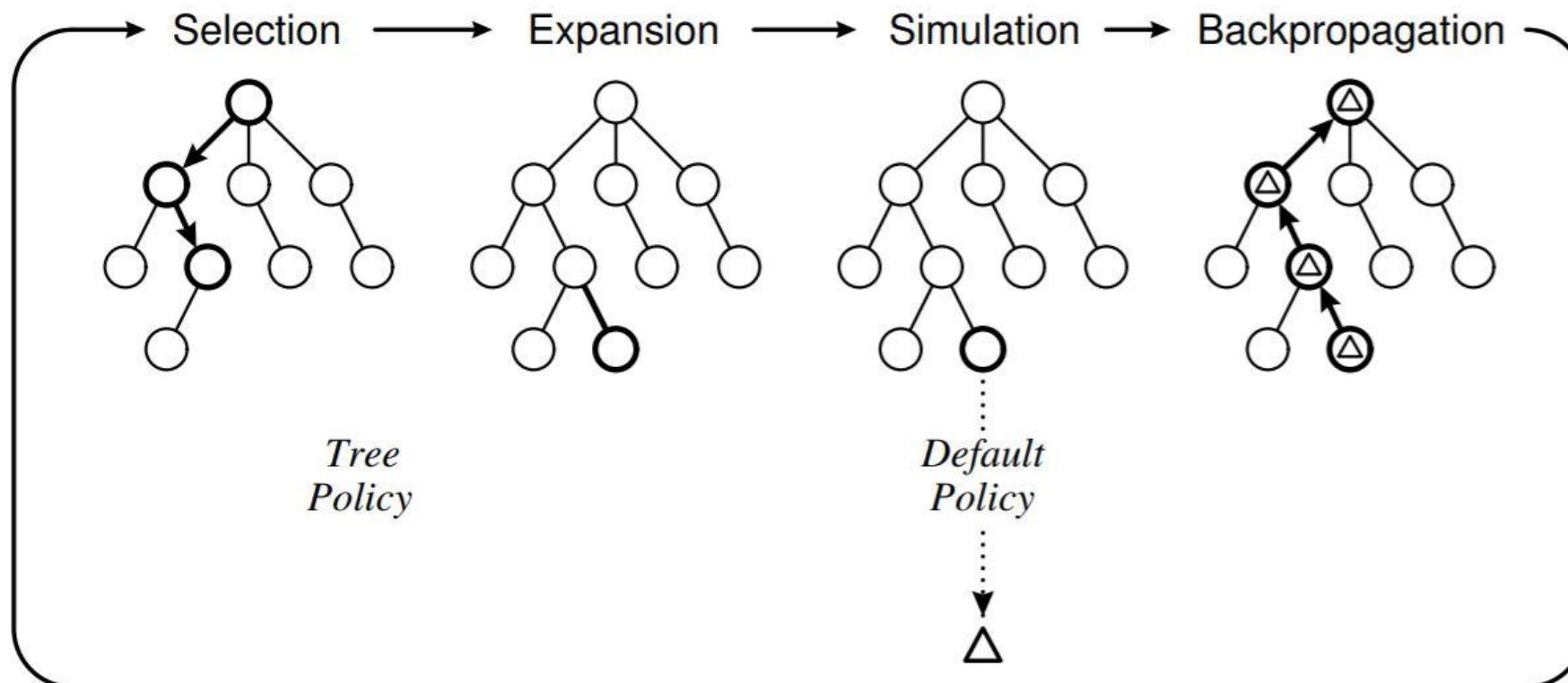
Monte Carlo Tree Search



Random Sampling



Decision Trees

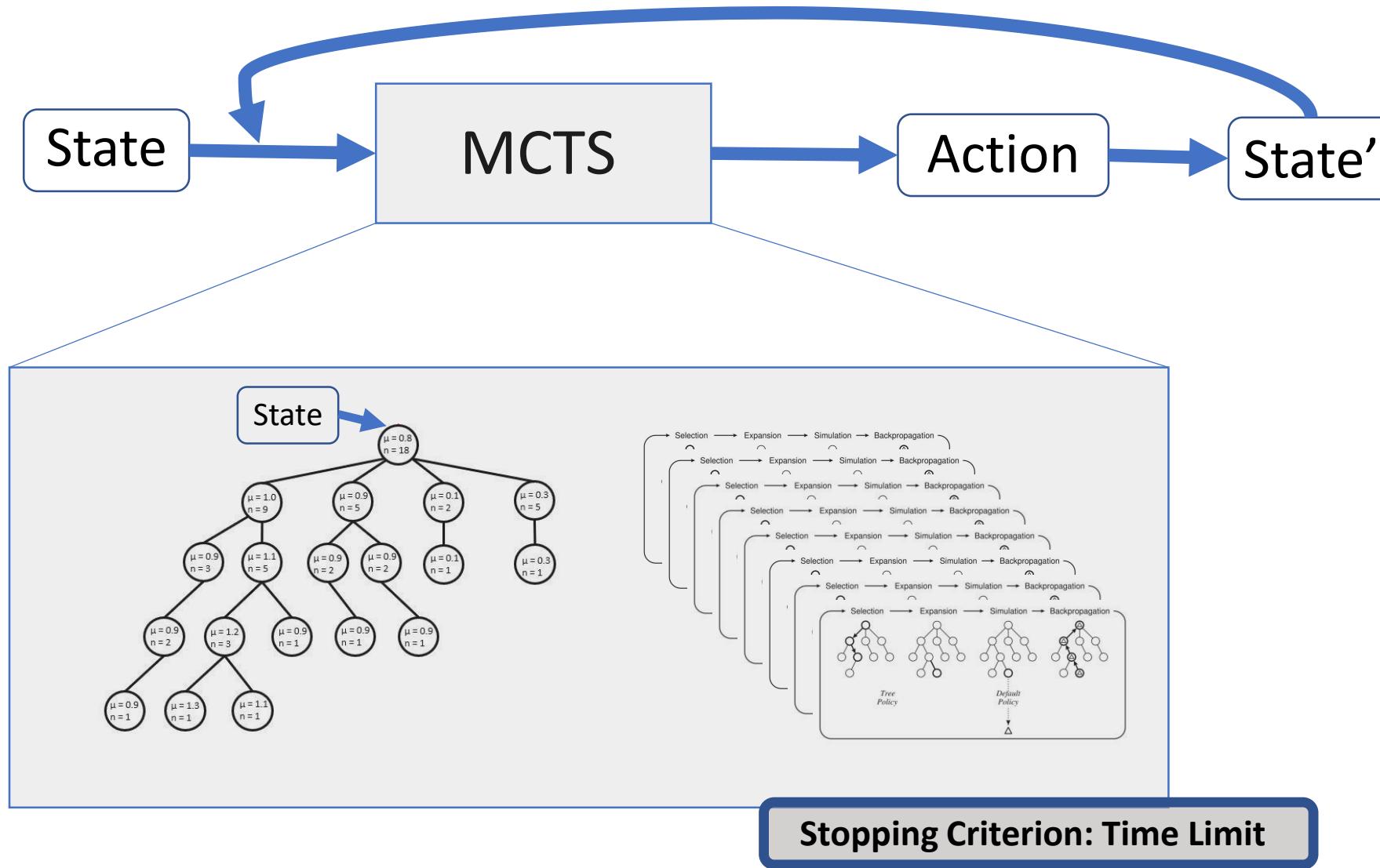
**Tree Policy**

Selection Phase
Example: UCB

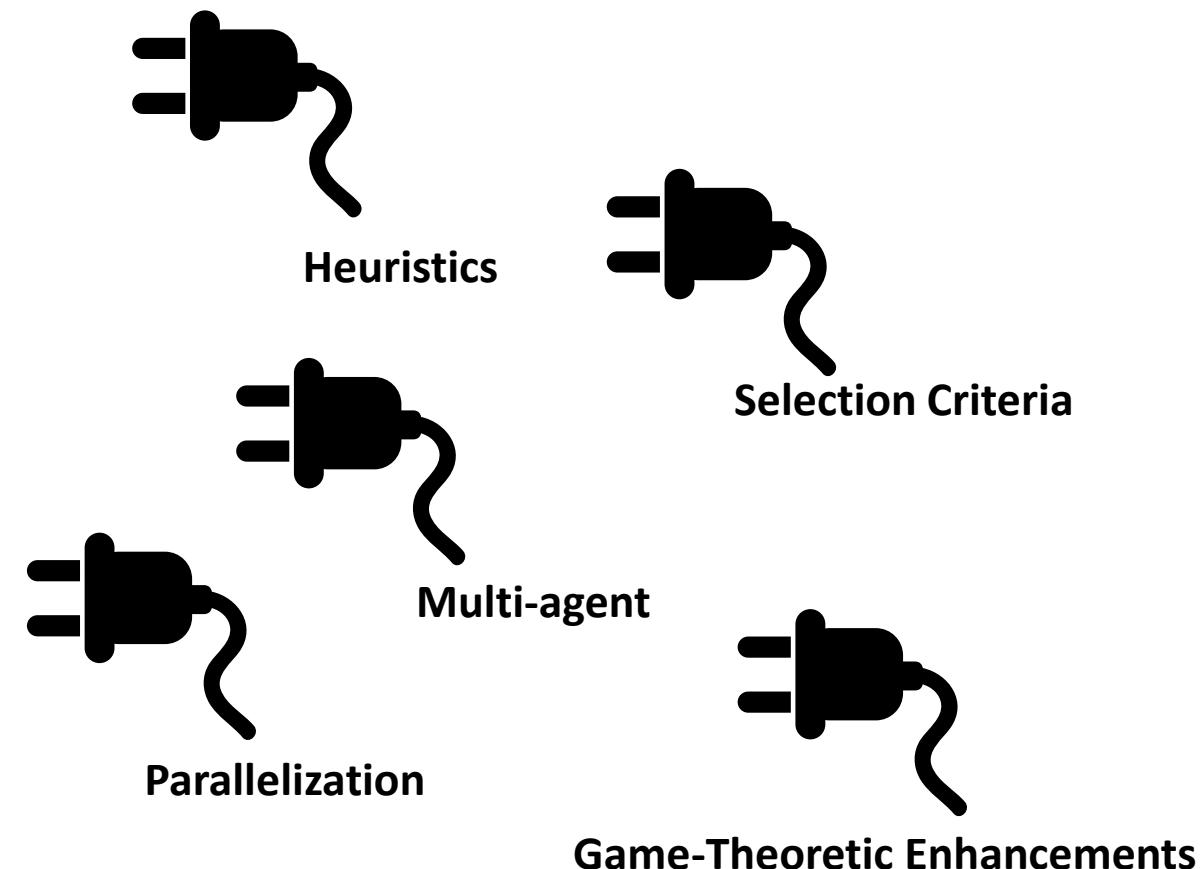
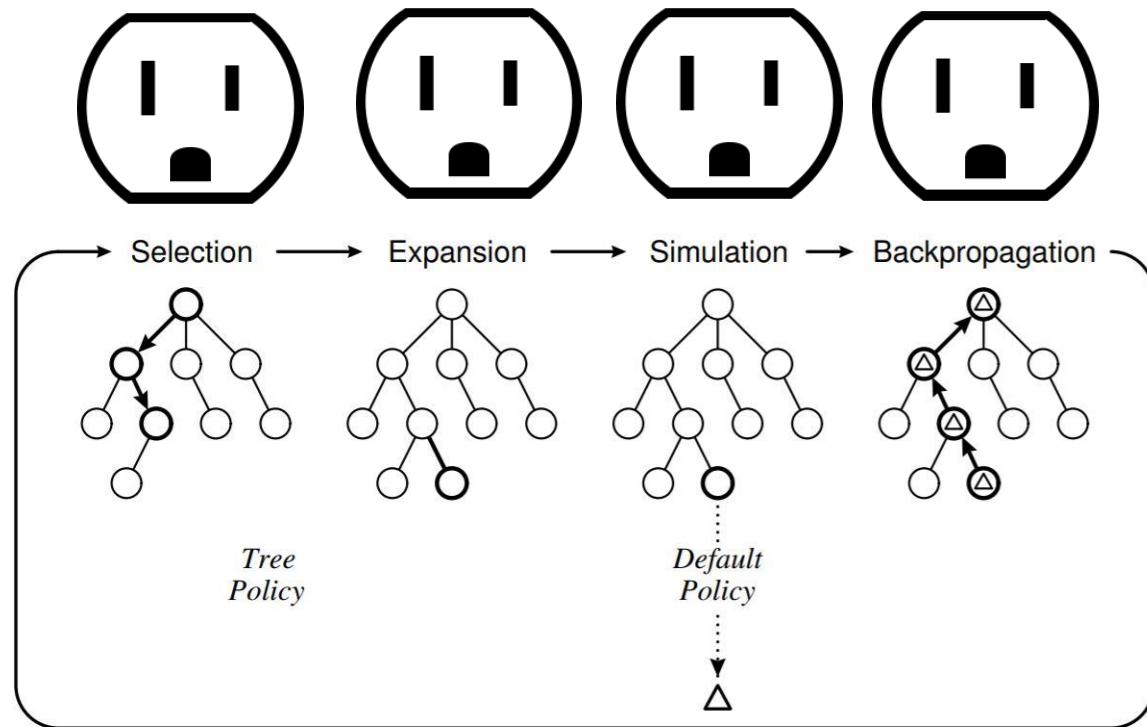
Default Policy

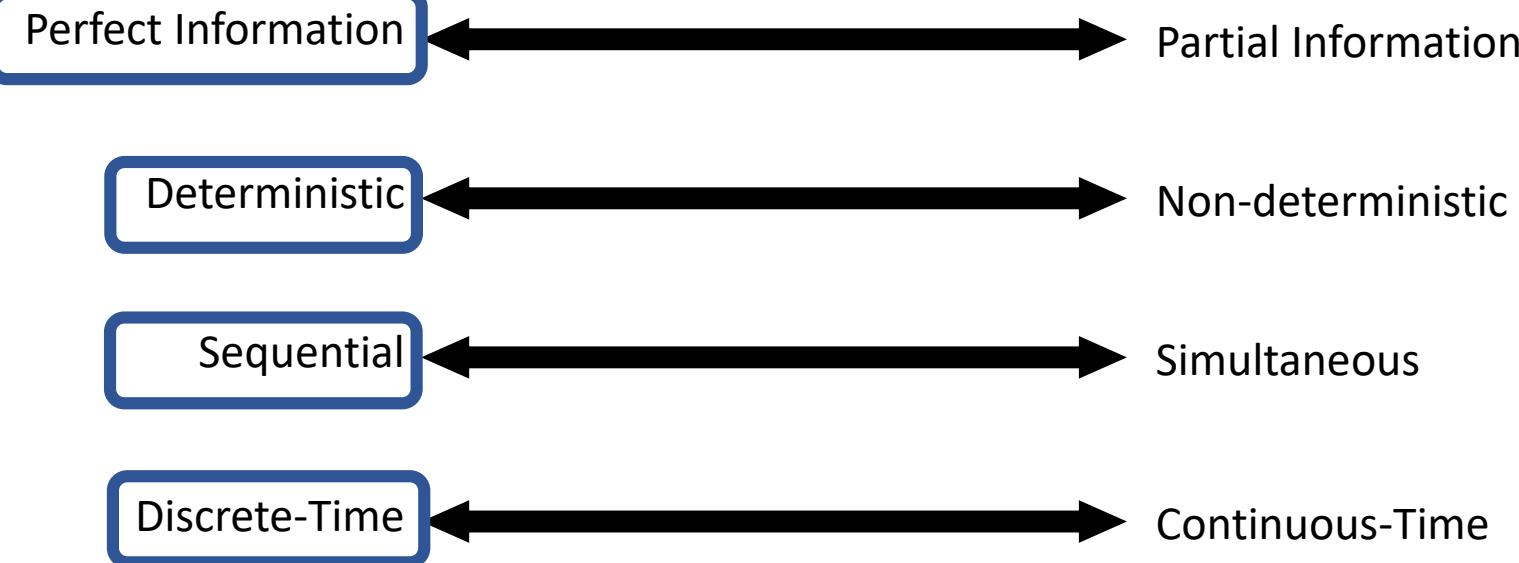
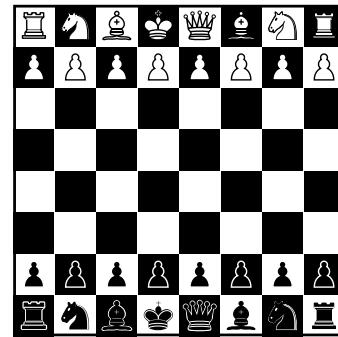
Simulation Phase
Example: Random Rollout

From *A Survey of Monte Carlo Tree Search Methods*

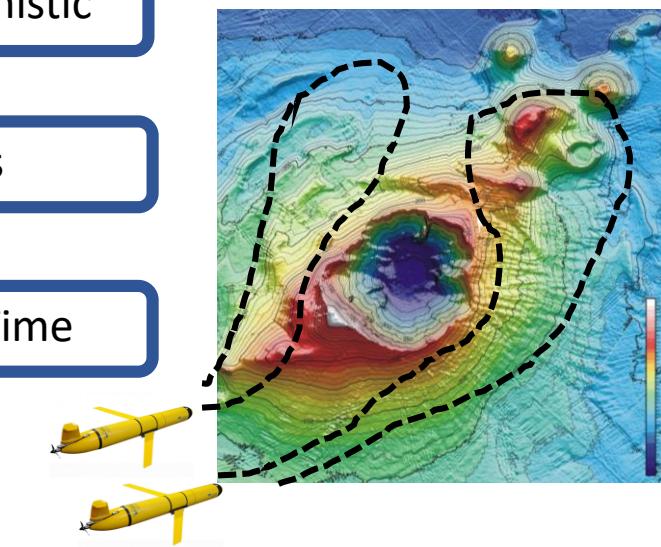
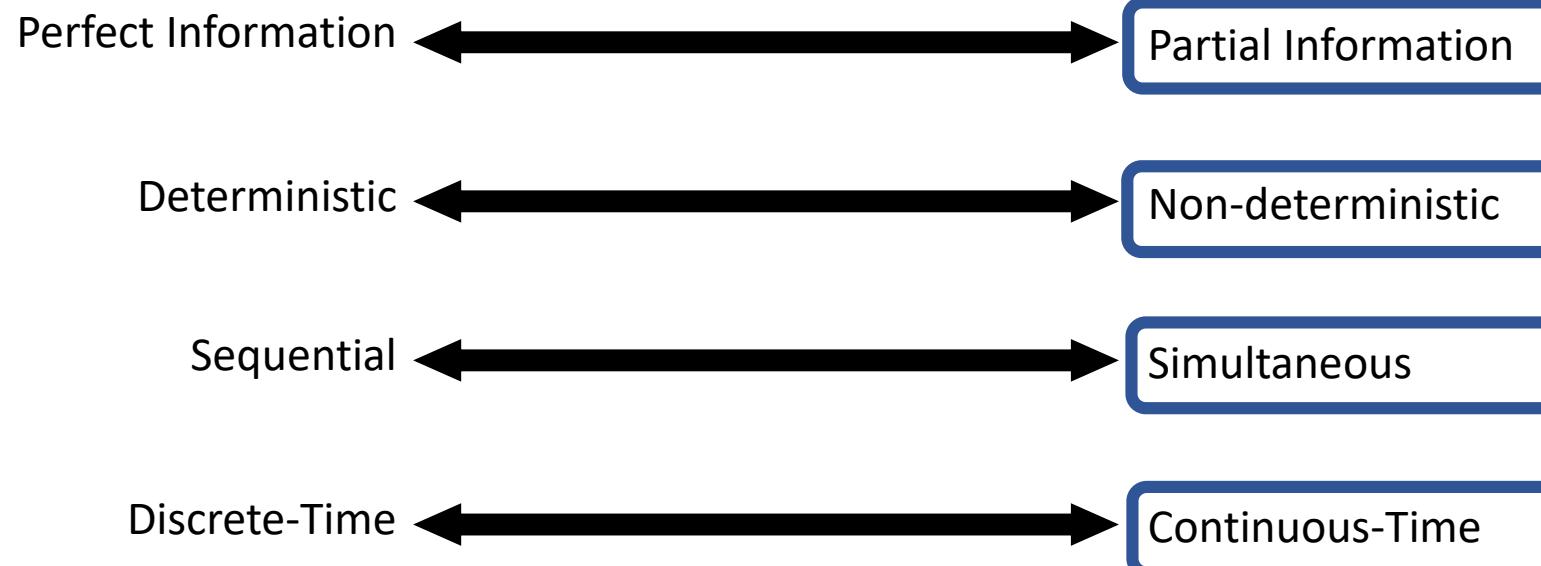


MCTS is Modular

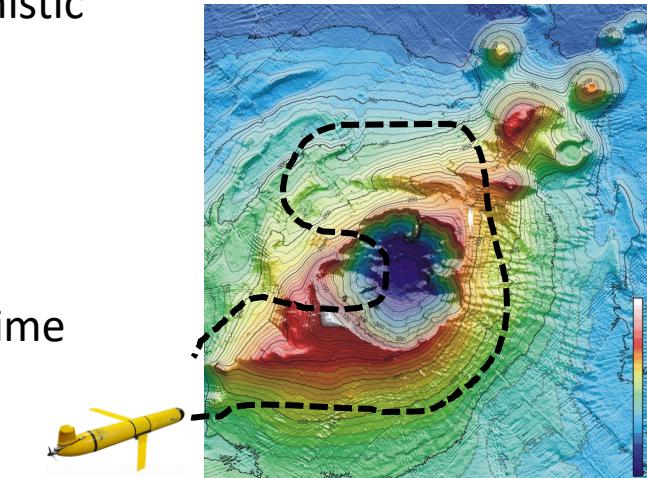
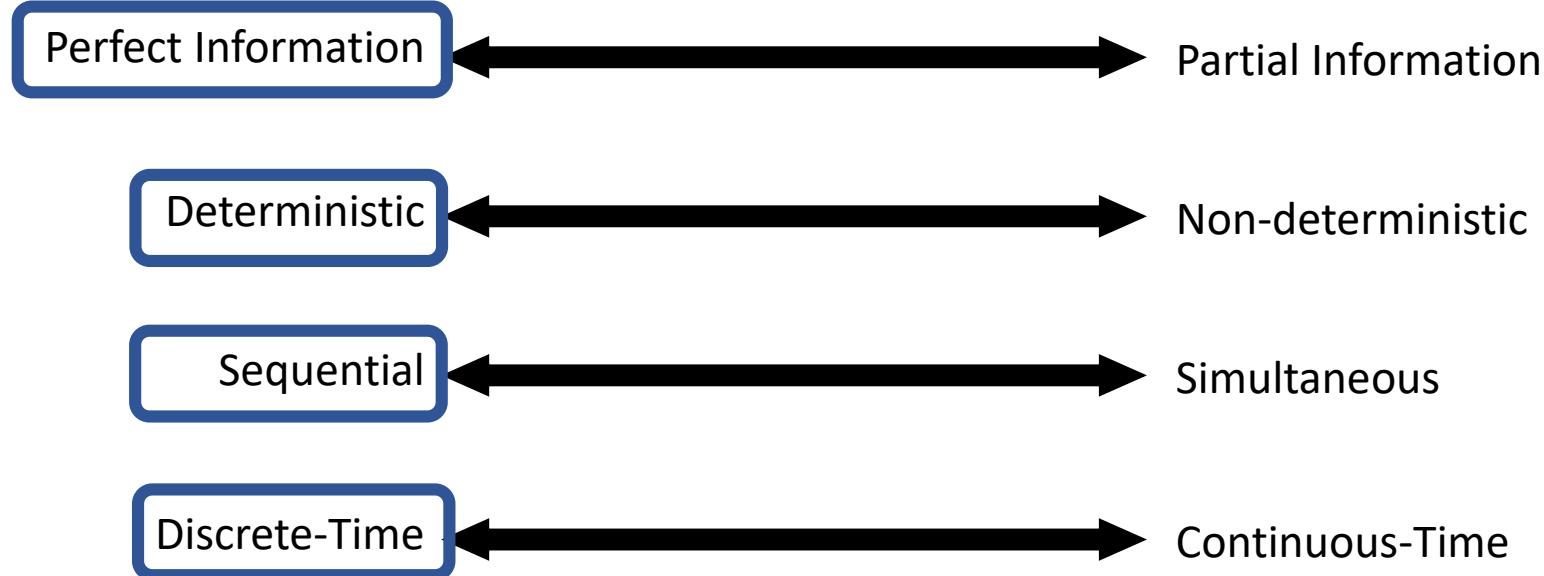




Our (Multi-Agent) Problem

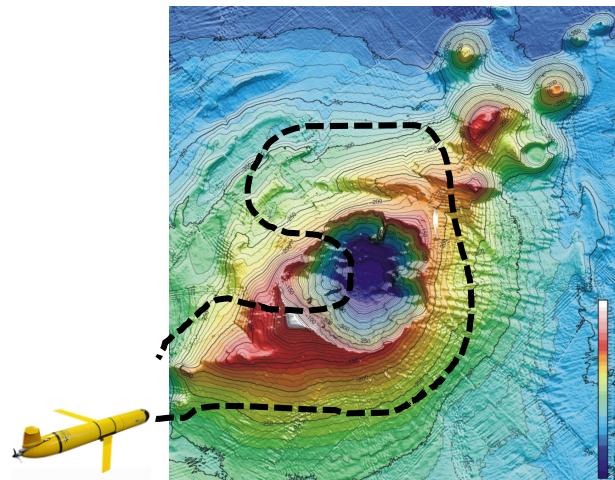


Our (Single-Agent) Problem



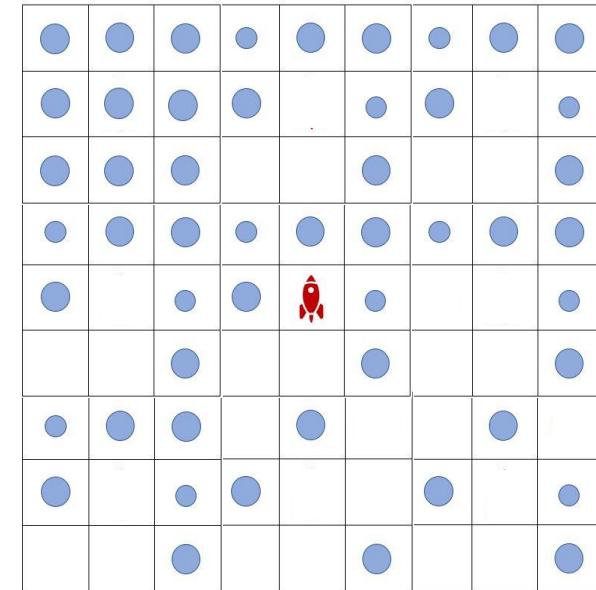
Our (Single-Agent) Problem

Investigating collaborative planning for AUVs in the Kolumbo crater.



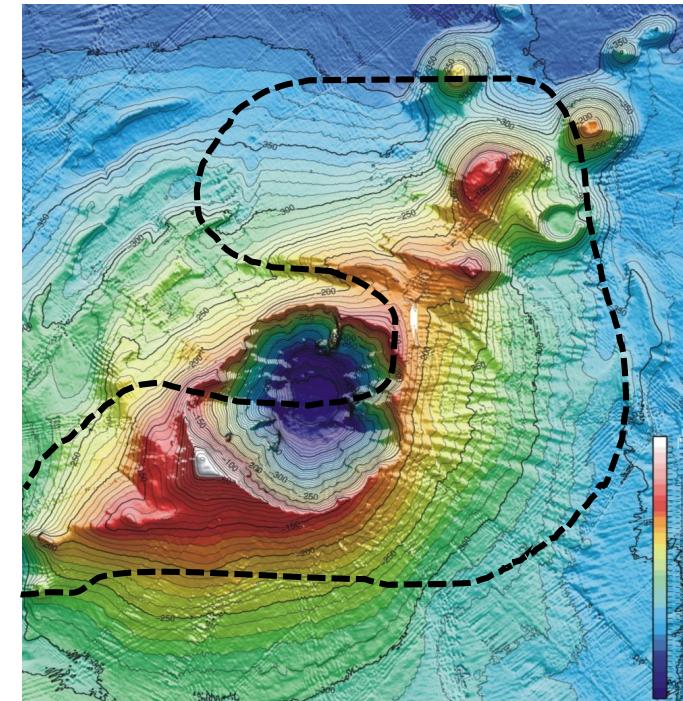
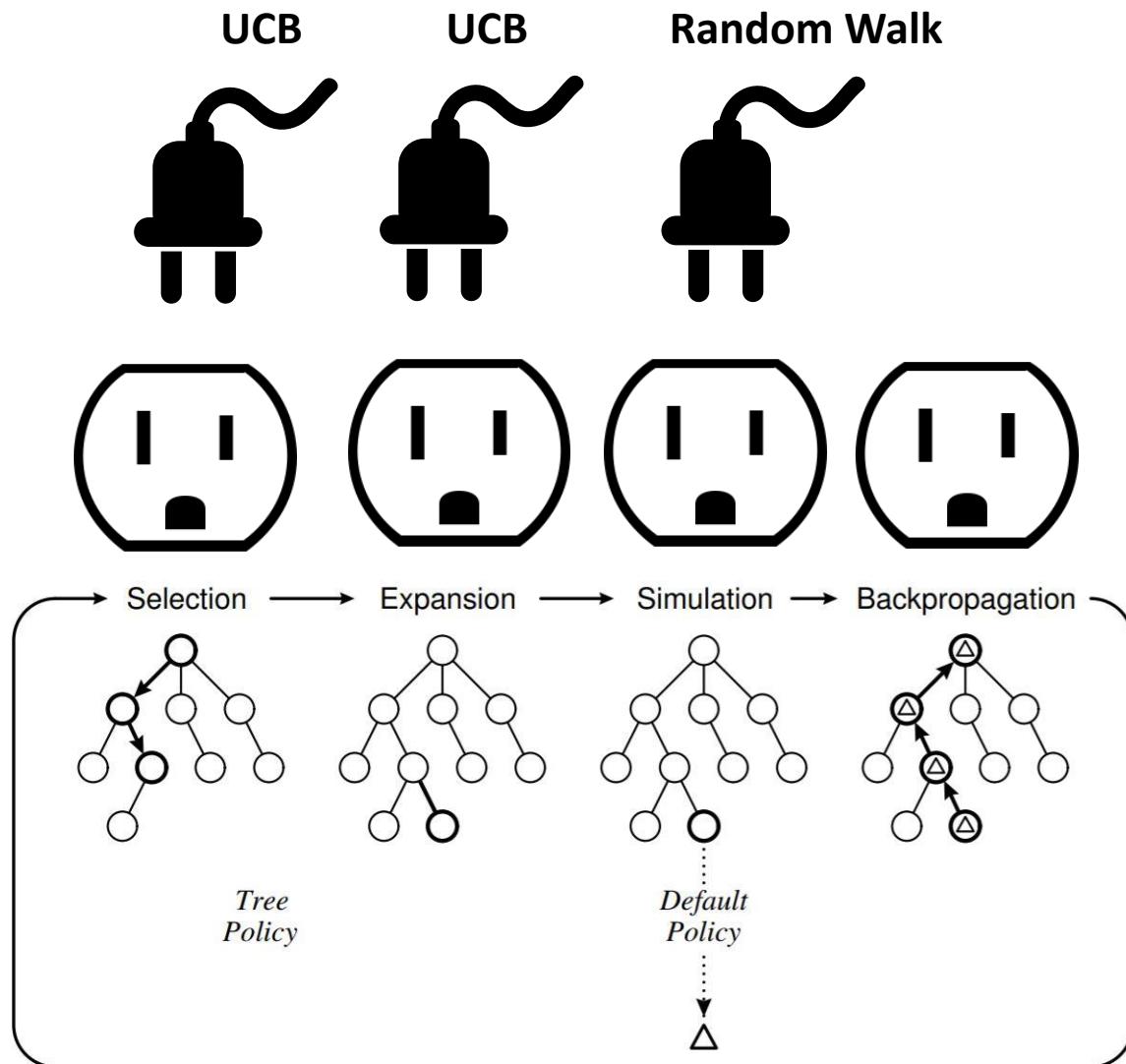
Simplified representation:

- Discretised map
- AUV can move (1 quadrant per time step)
- No dynamics
- Fixed mission duration
- **Single-Agent first: No collaboration**

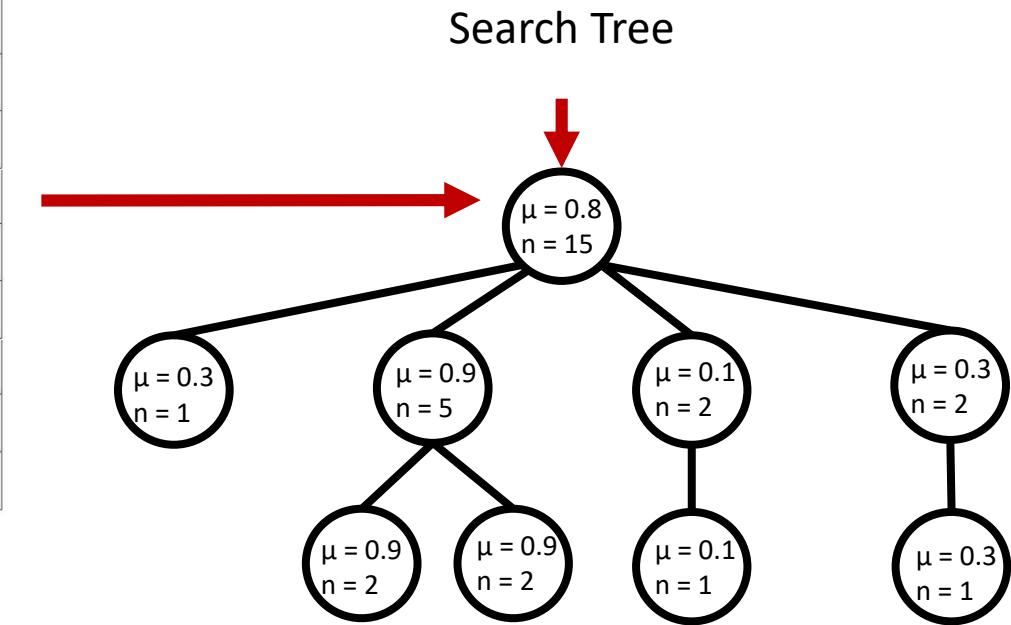
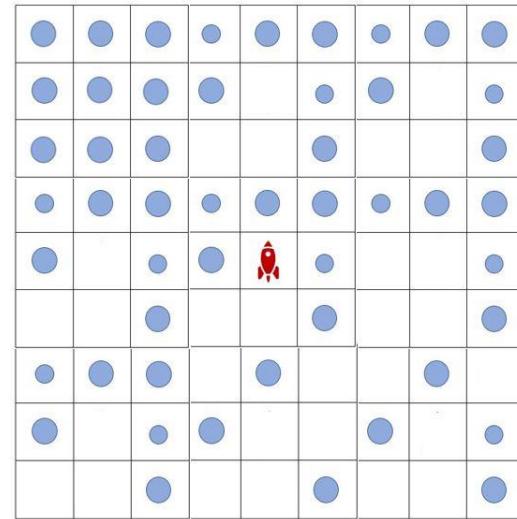
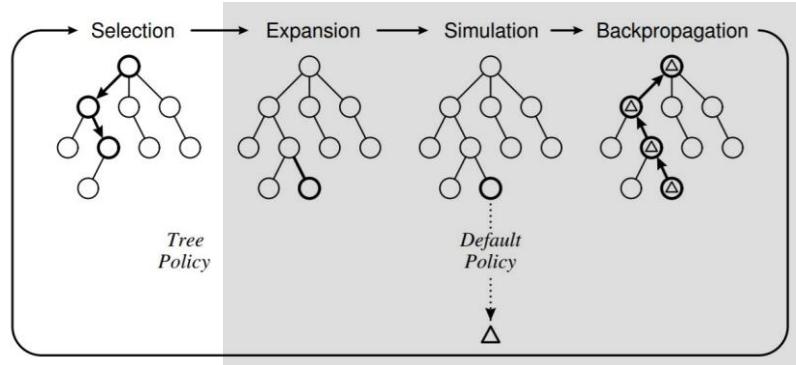


= 0.1

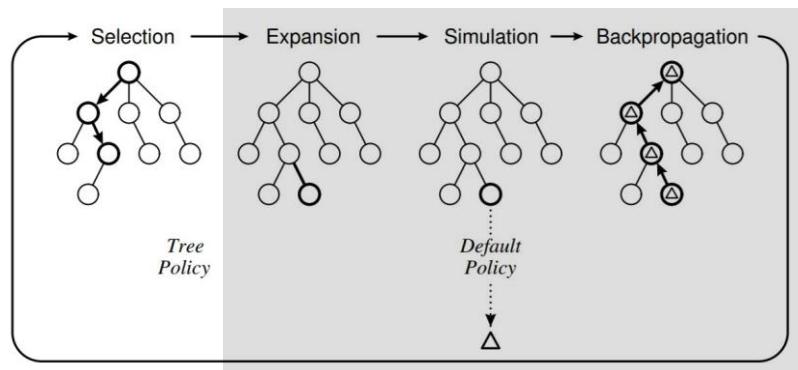
= 0.2



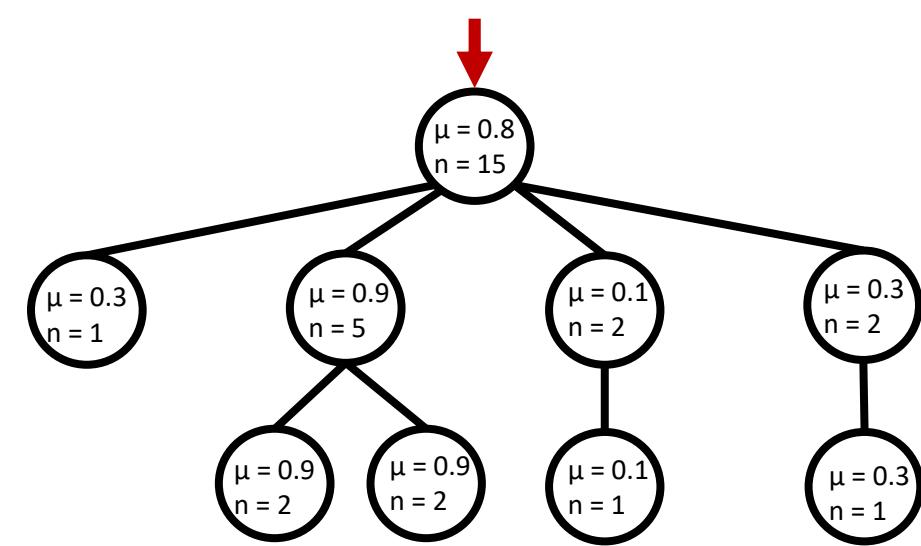
Selection



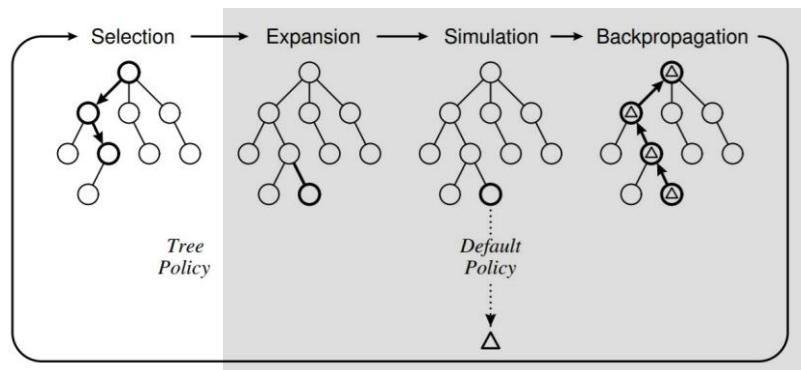
Selection



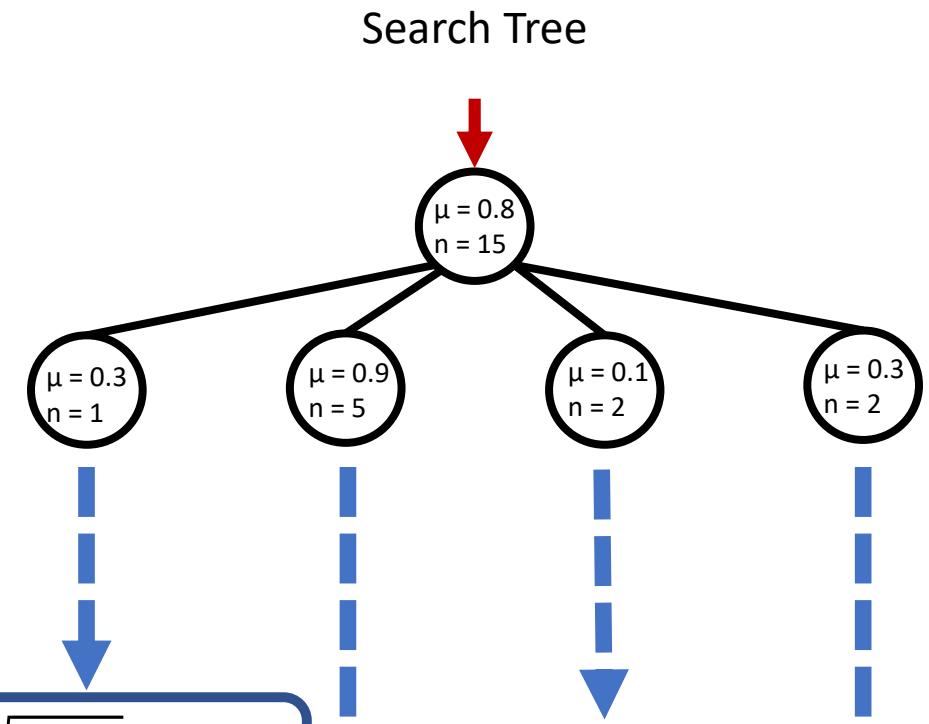
$$\text{UCB1} = \bar{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$



Selection



$$\text{UCB1} = \bar{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

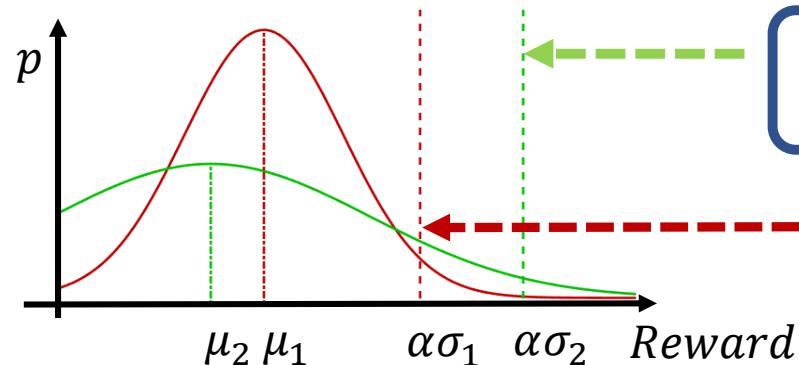


$$\text{UCB1} = 0.3 + \sqrt{\frac{2 \ln 15}{1}} = 2.63$$

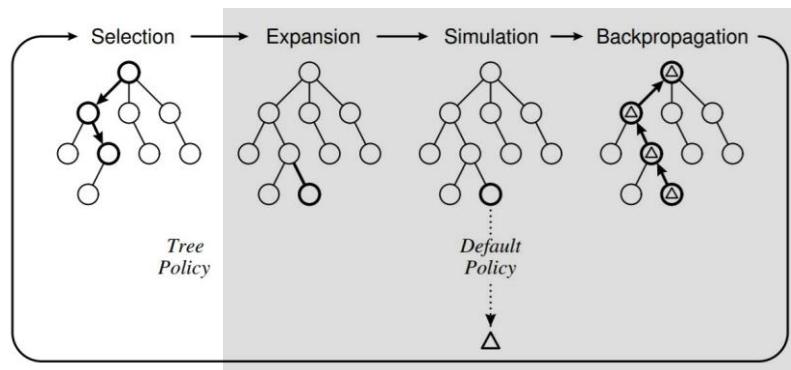
$$\text{UCB1} = 1.74$$

$$\text{UCB1} = 0.9 + \sqrt{\frac{2 \ln 15}{5}} = 1.94$$

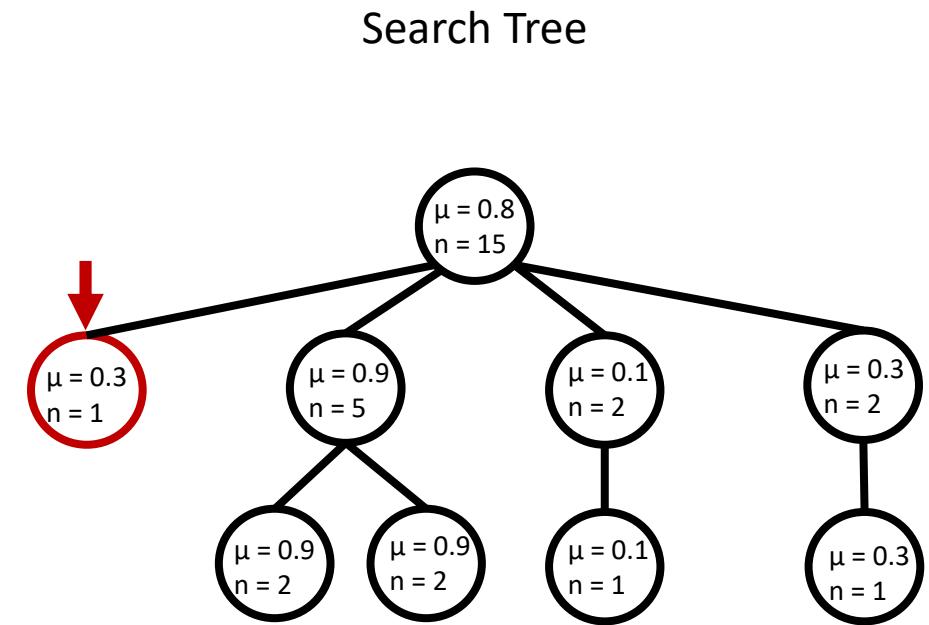
$$\text{UCB1} = 1.95$$



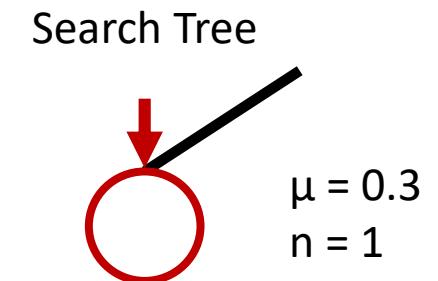
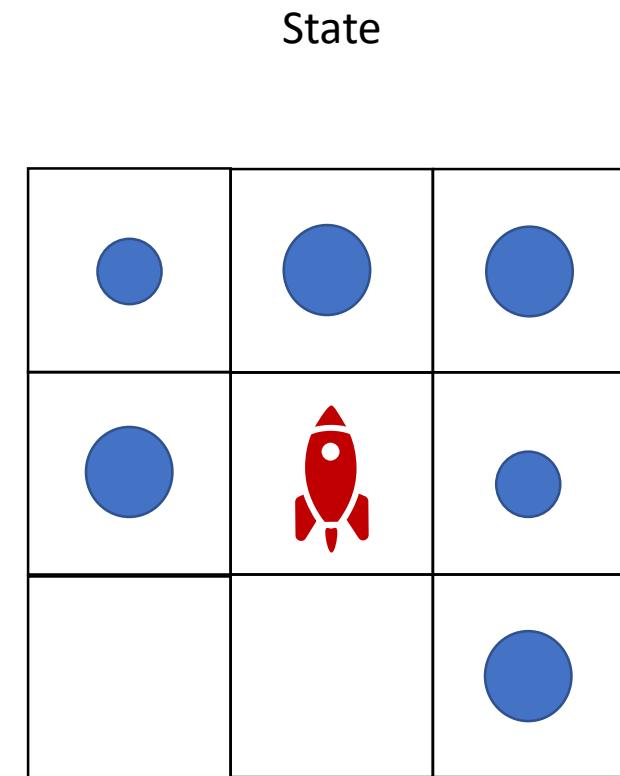
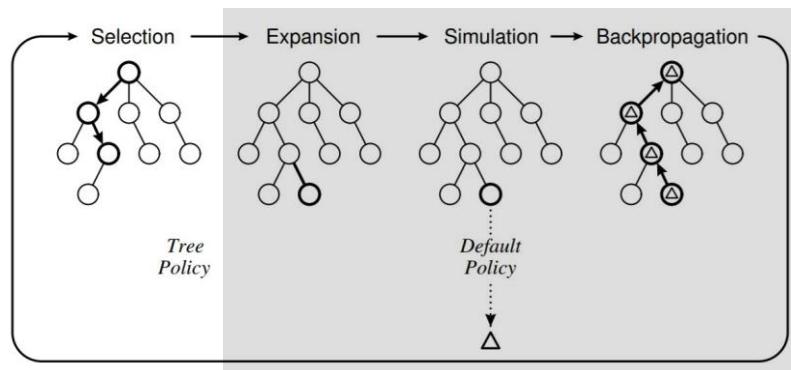
Selection



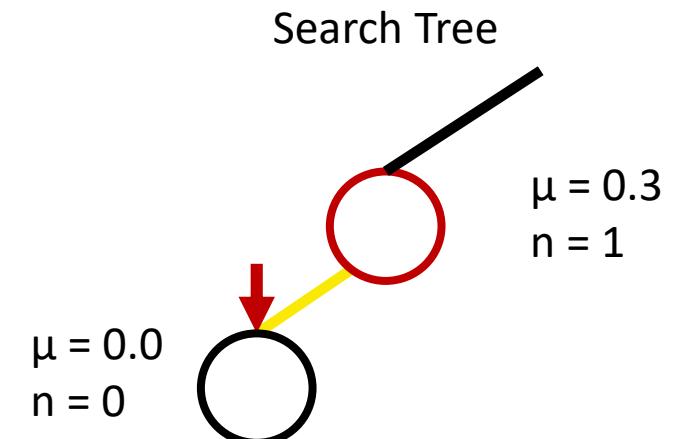
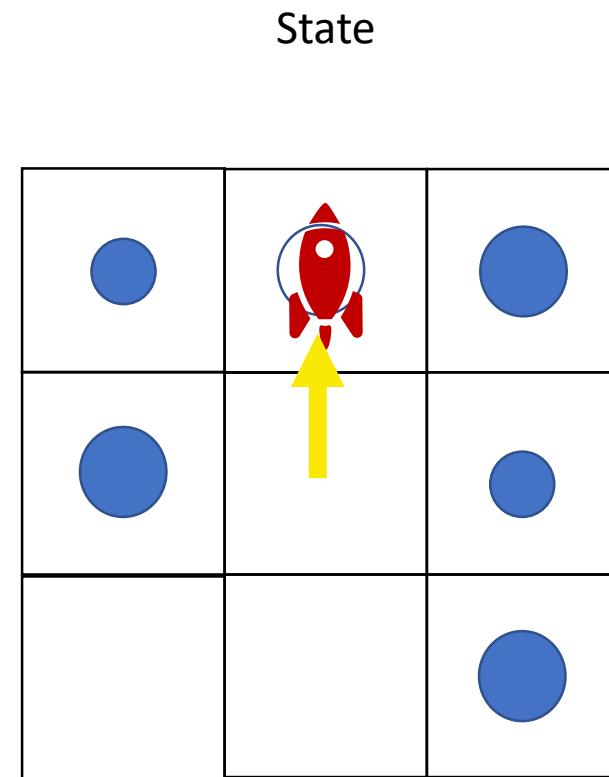
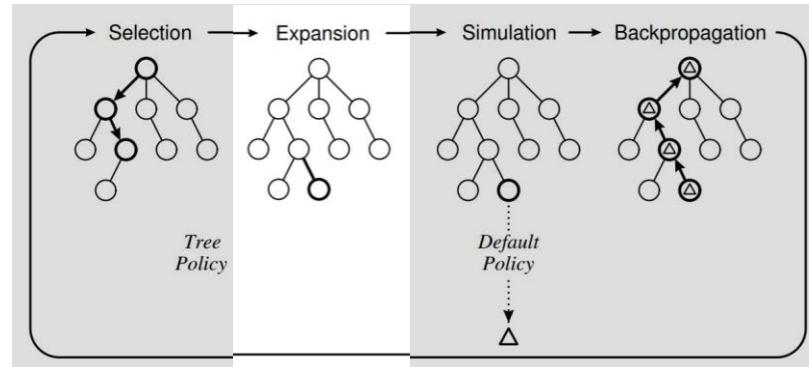
$$\text{UCB1} = \bar{X}_j + \sqrt{\frac{2 \ln n}{n_j}}$$

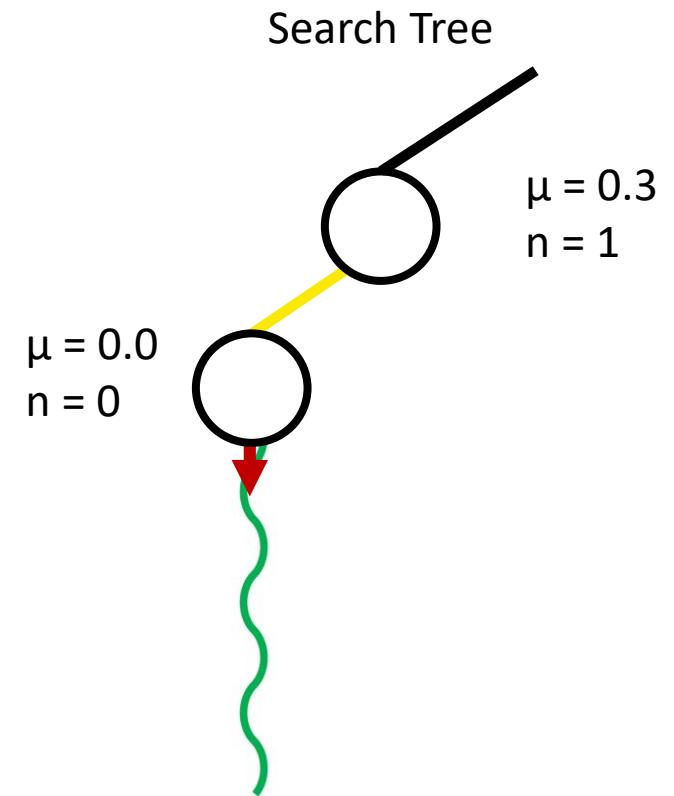
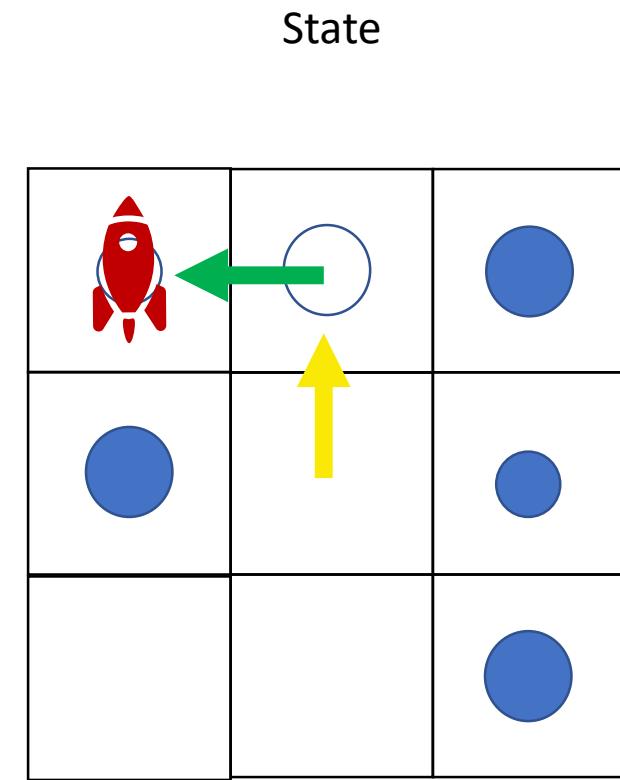
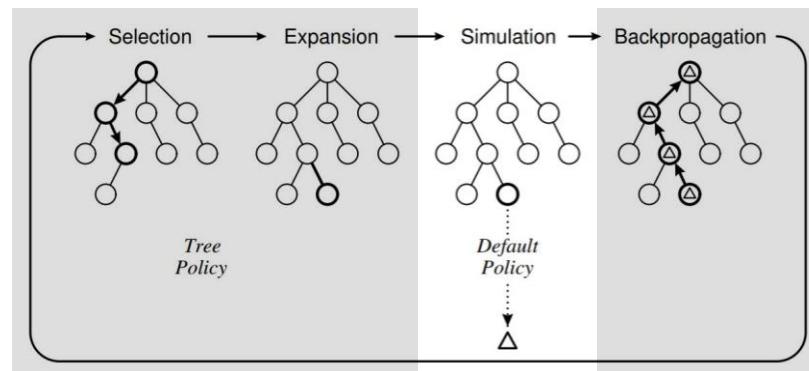


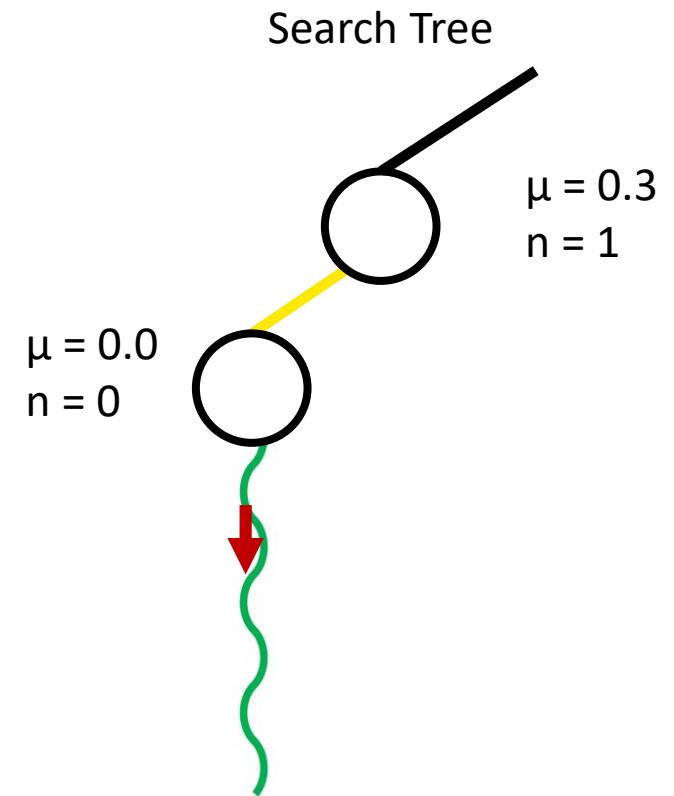
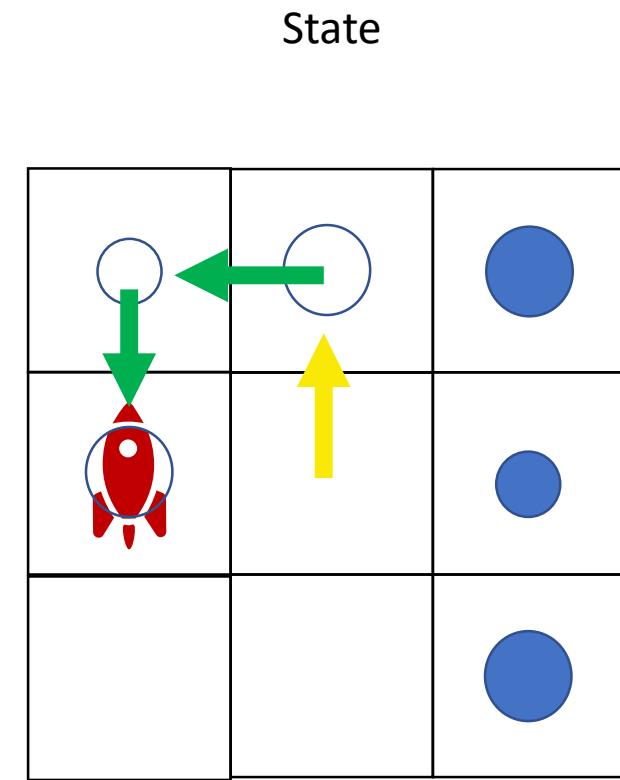
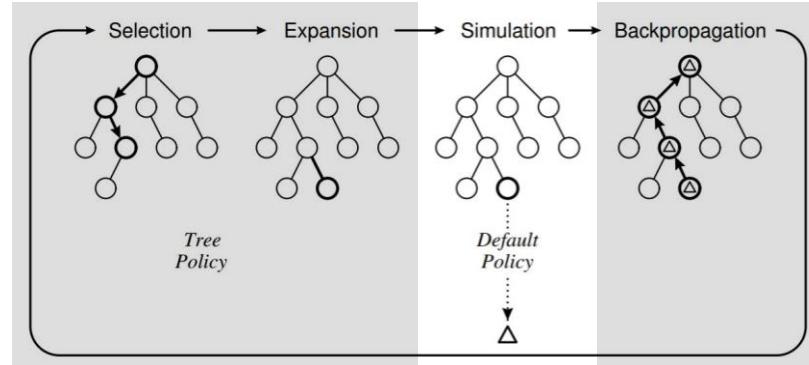
Selection



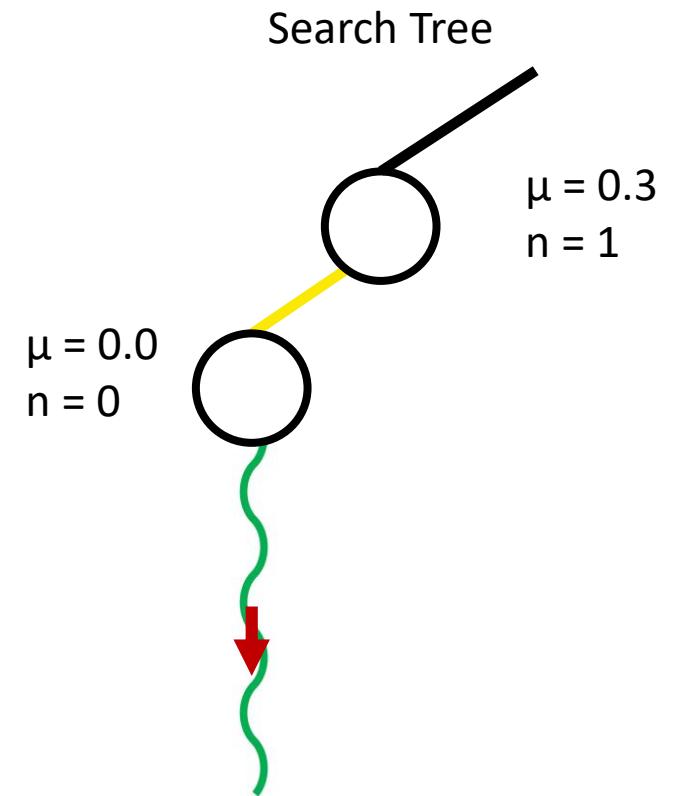
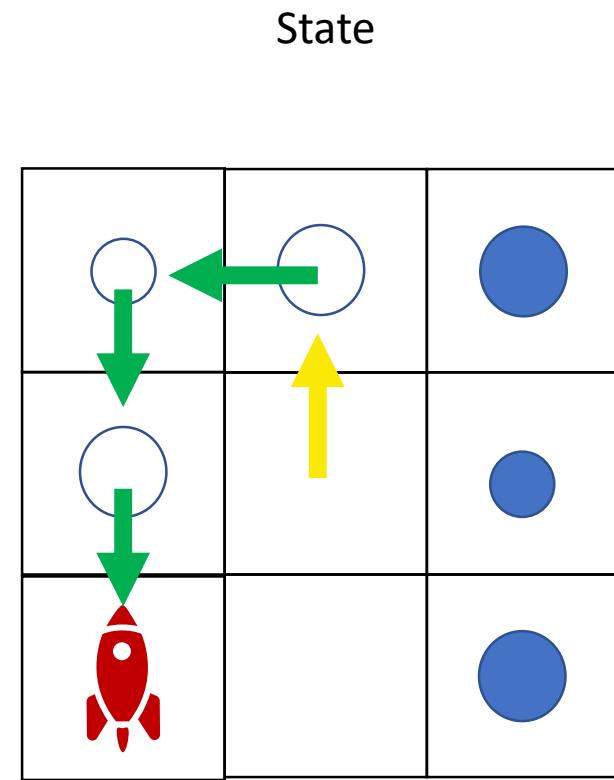
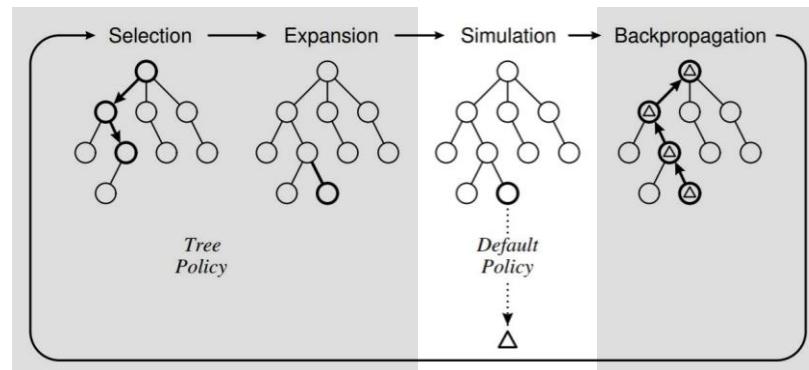
Expansion



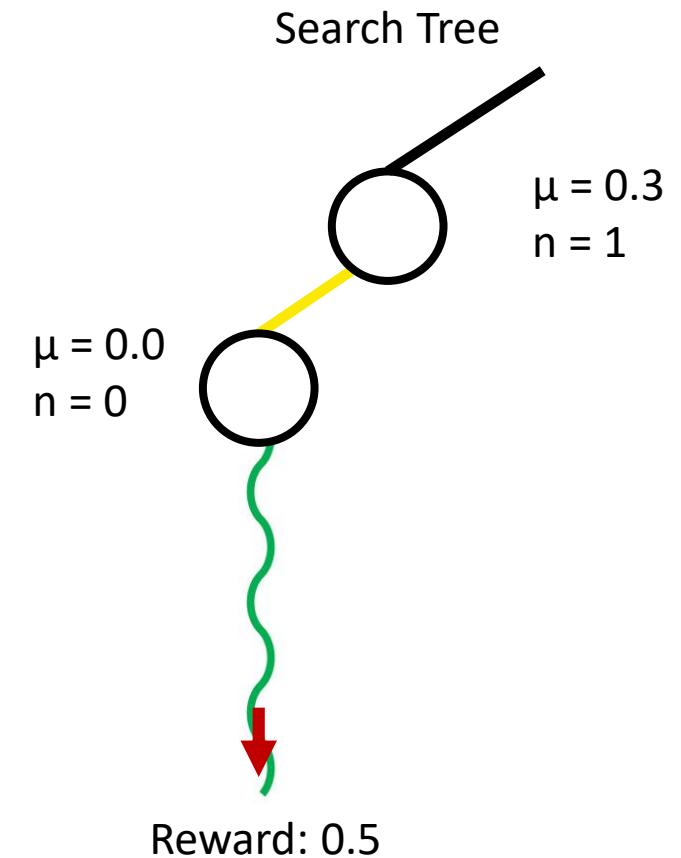
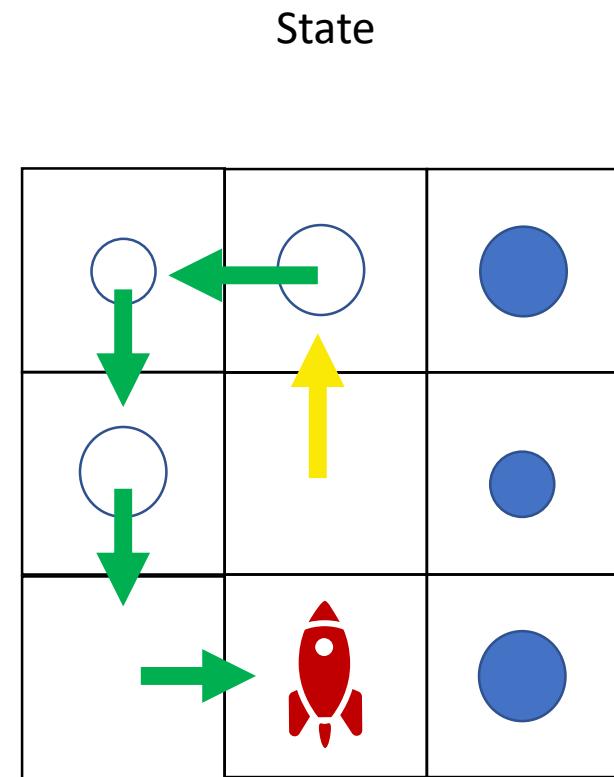
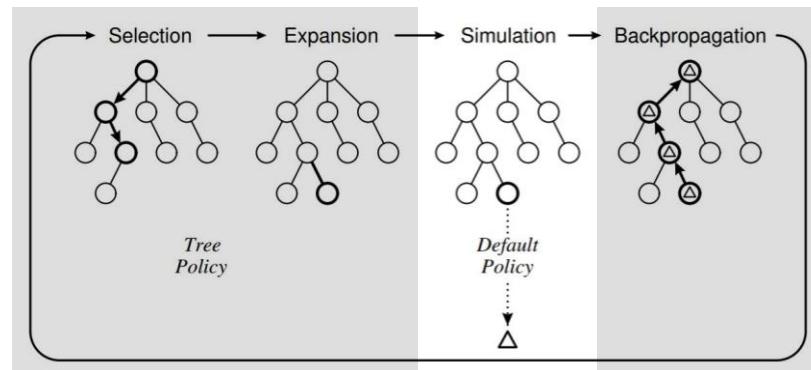




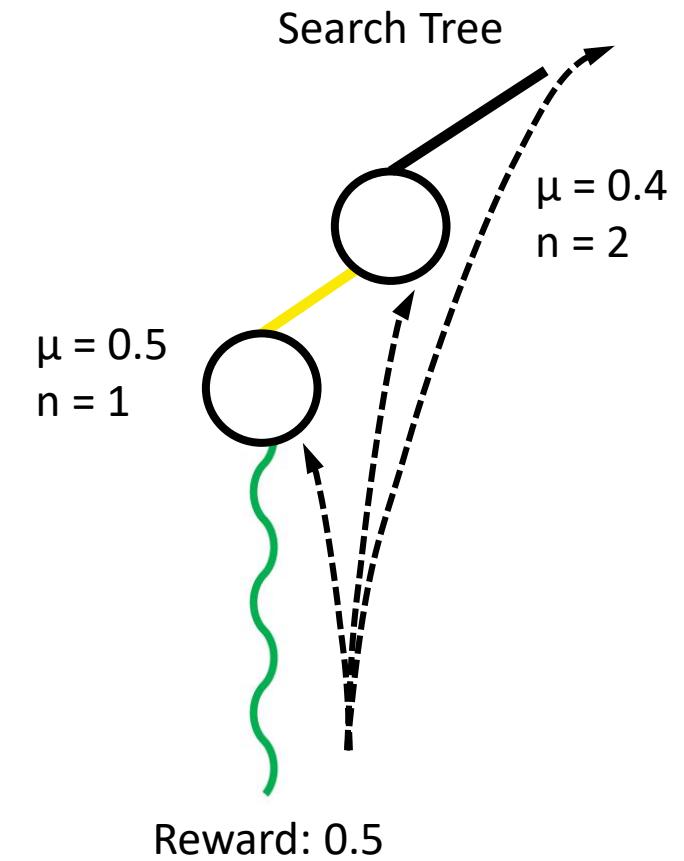
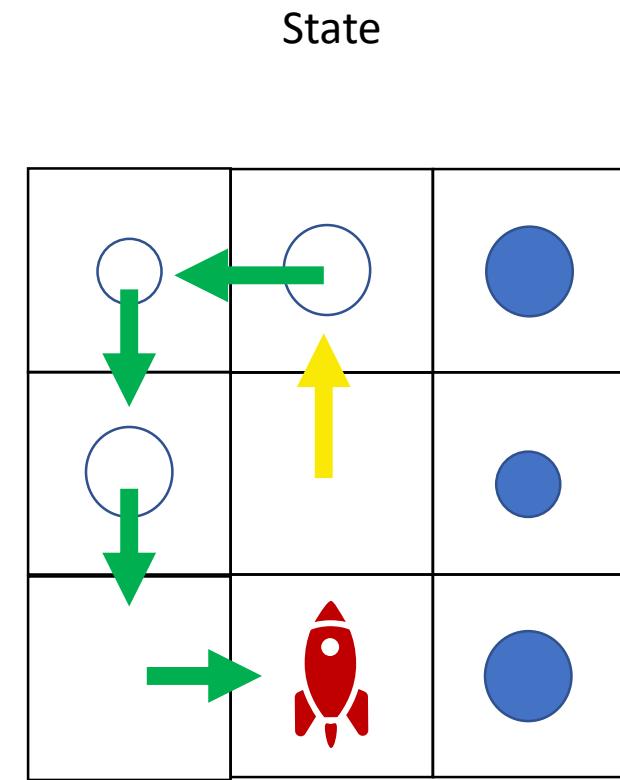
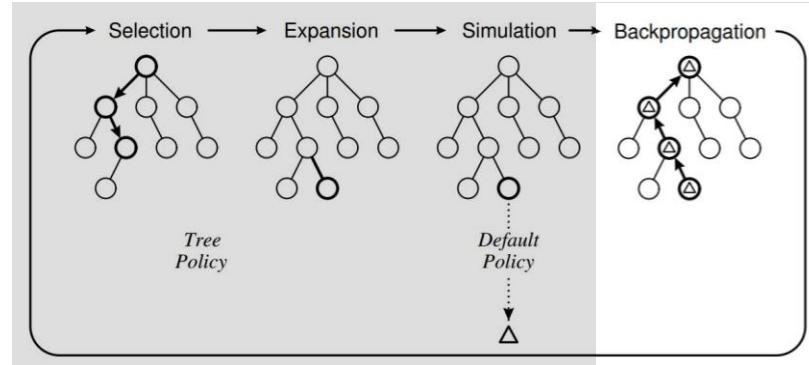
Simulation



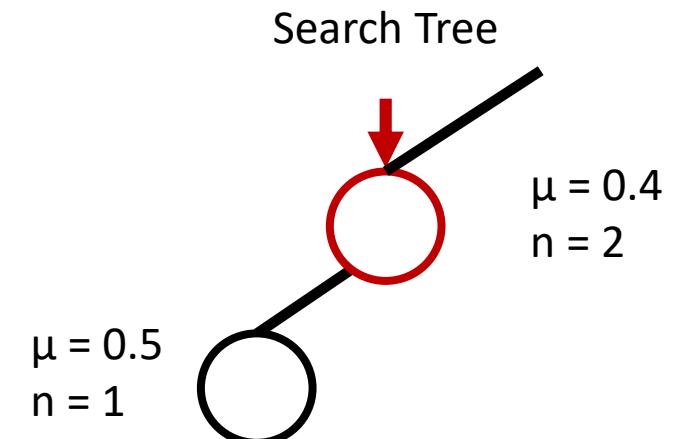
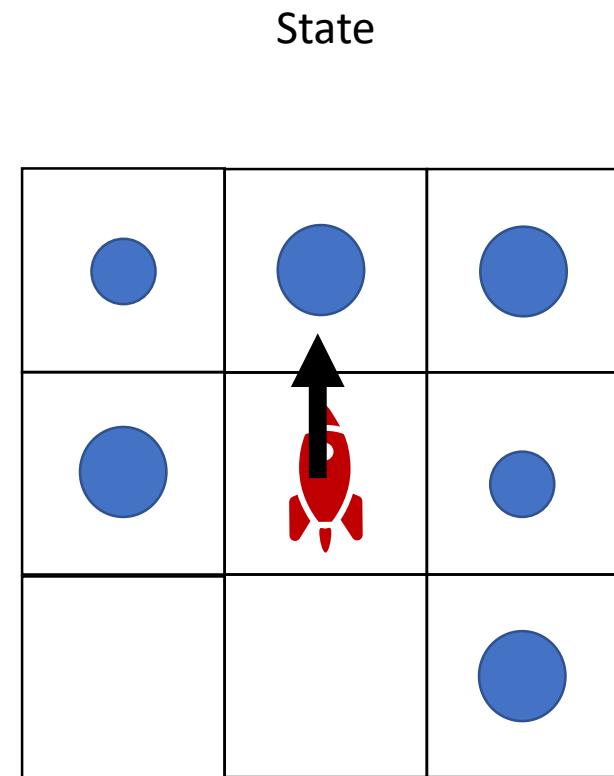
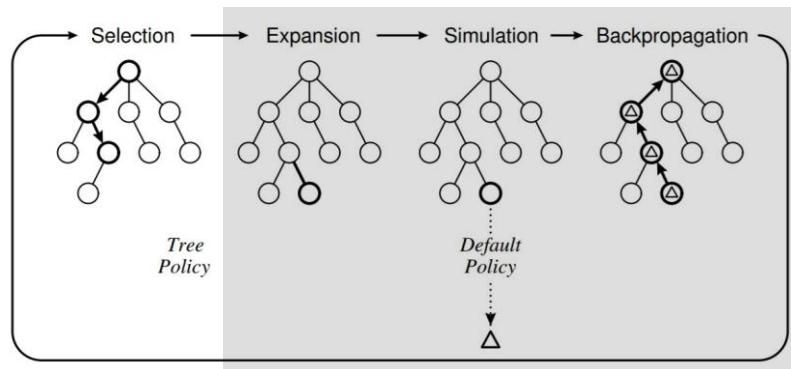
Simulation



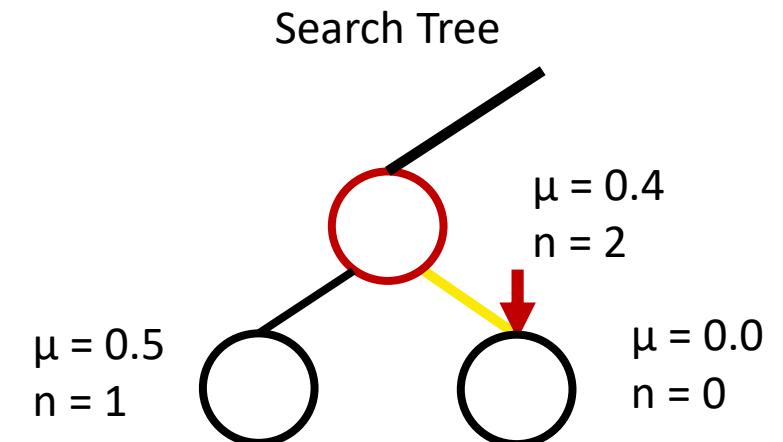
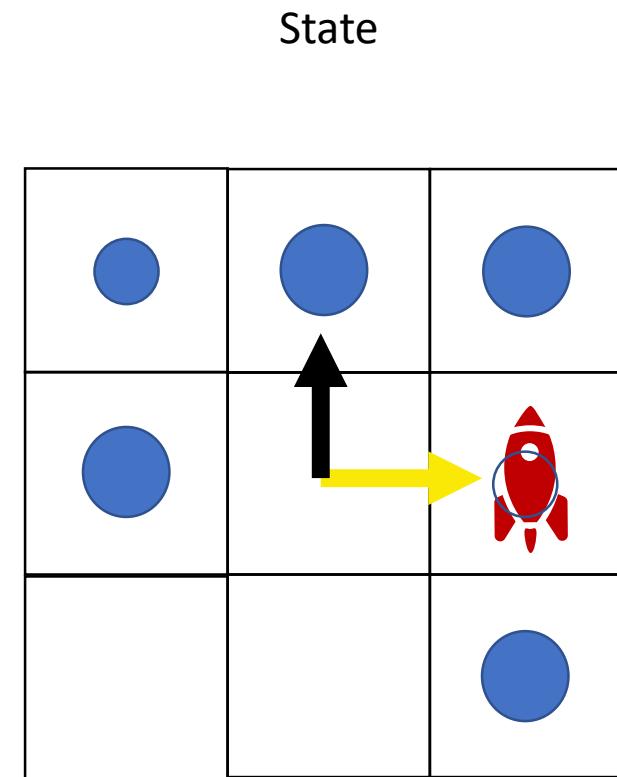
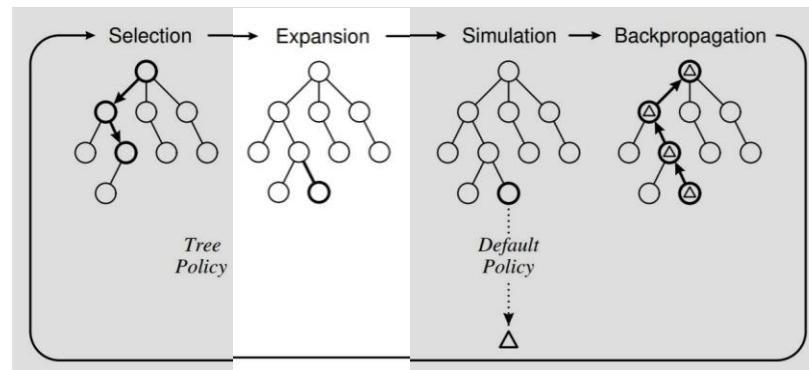
Backpropagation



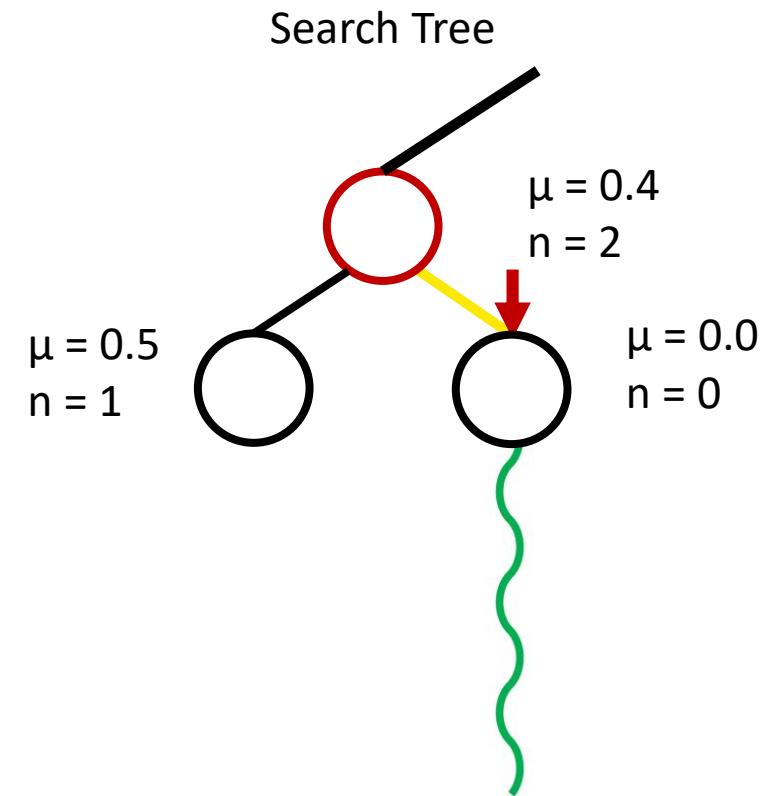
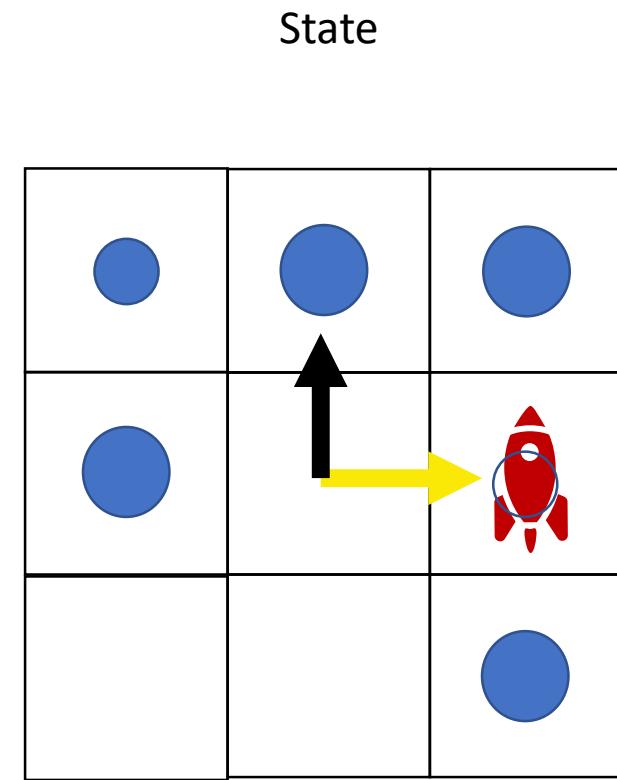
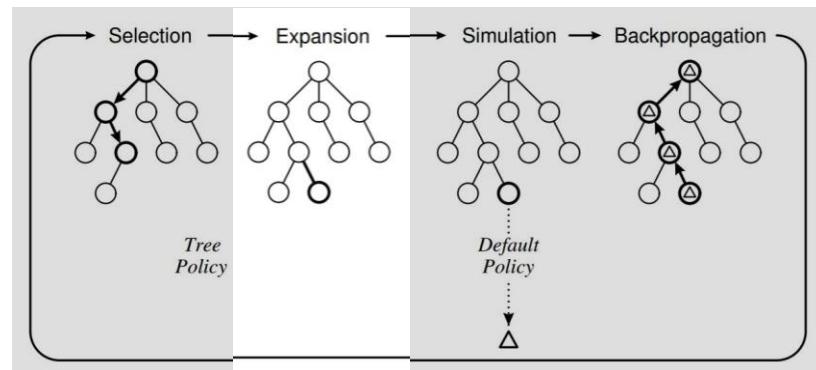
Selection



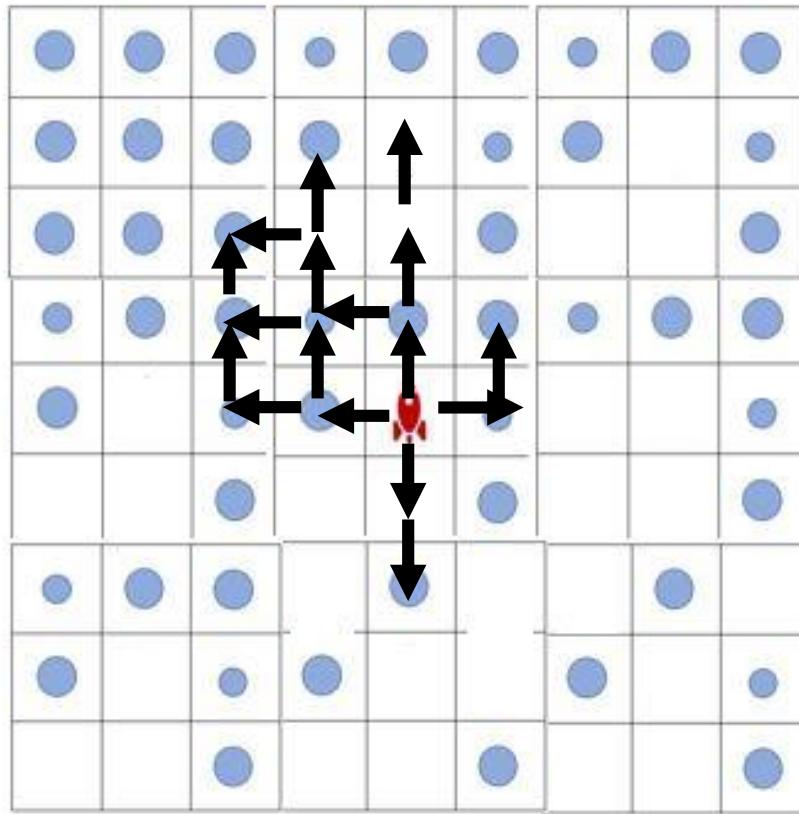
Expansion



Simulation

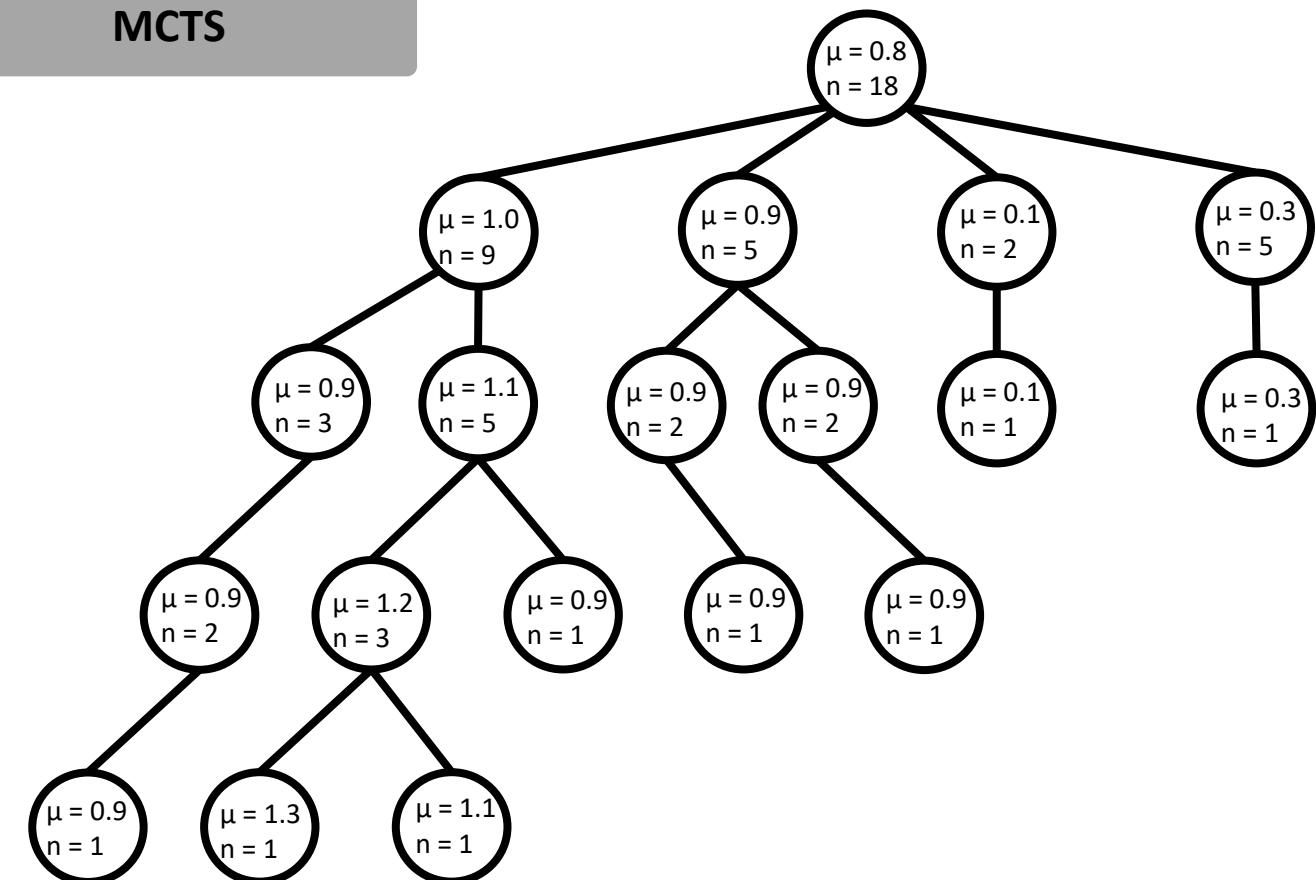


State

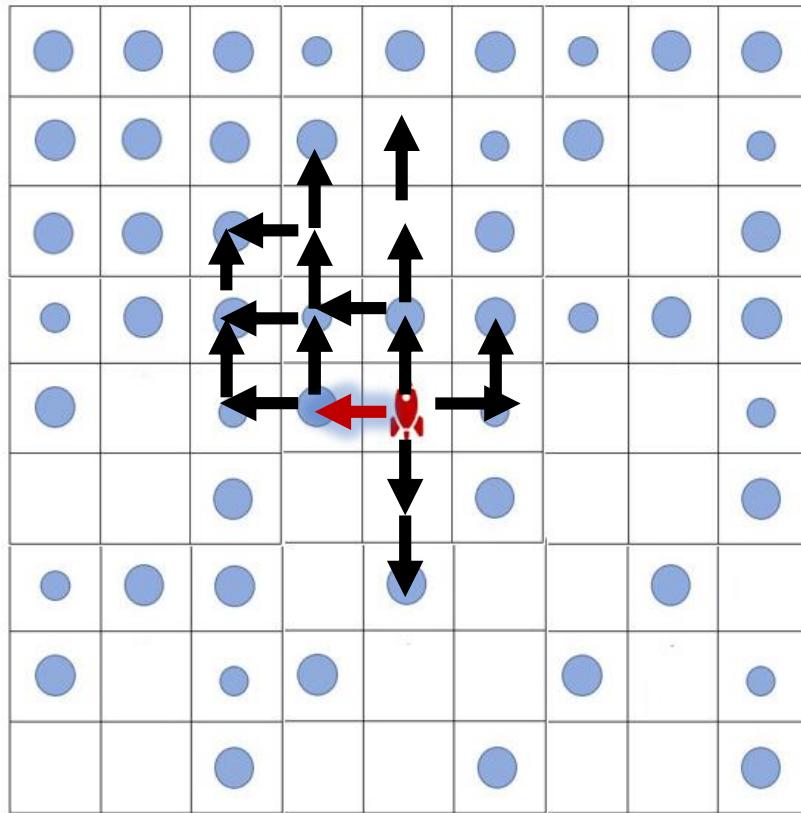


MCTS

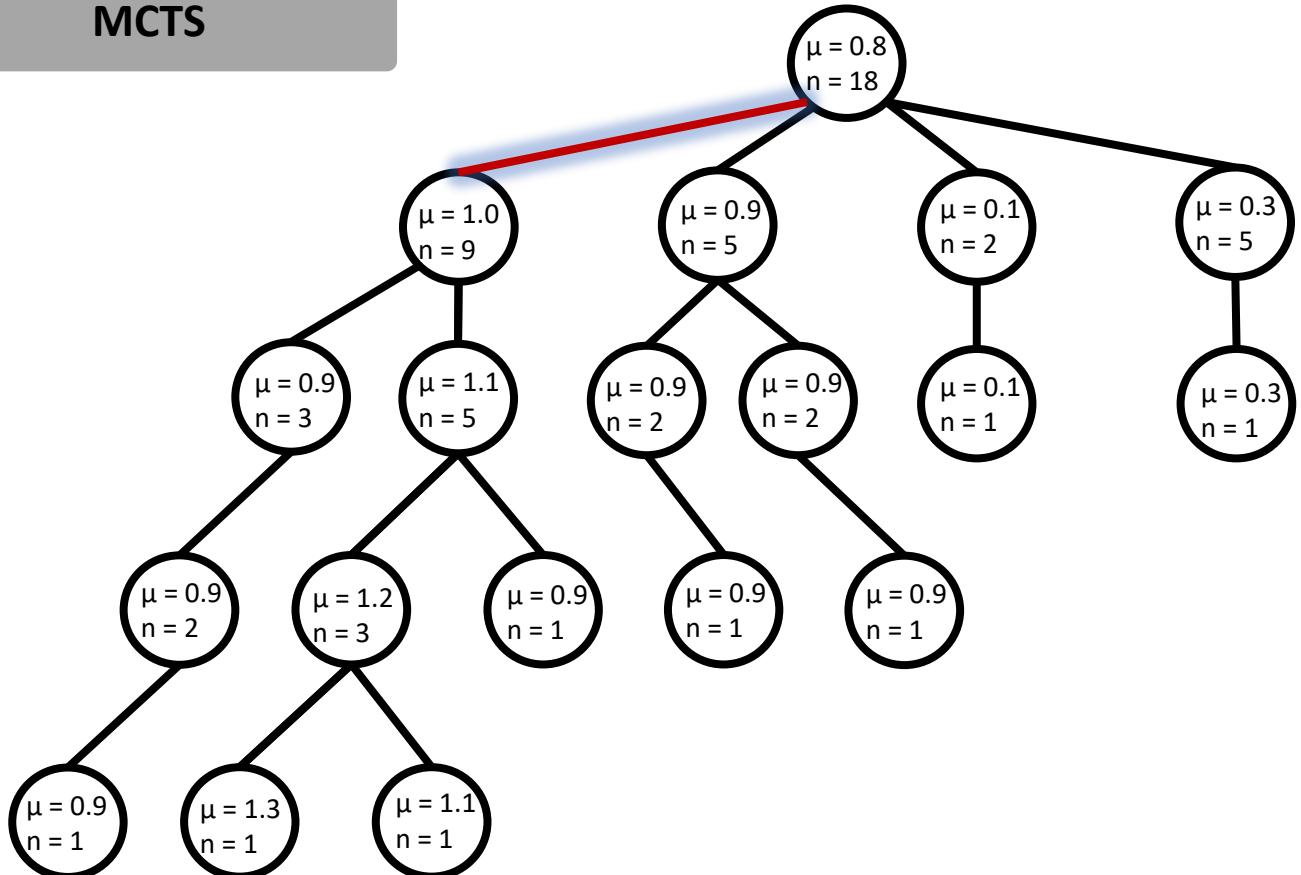
Search Tree



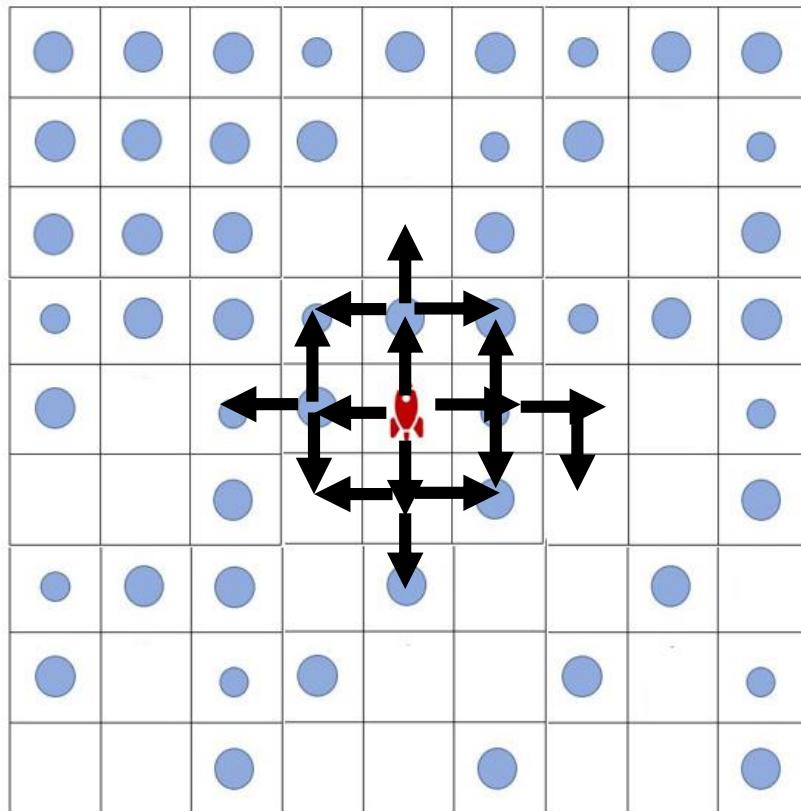
State



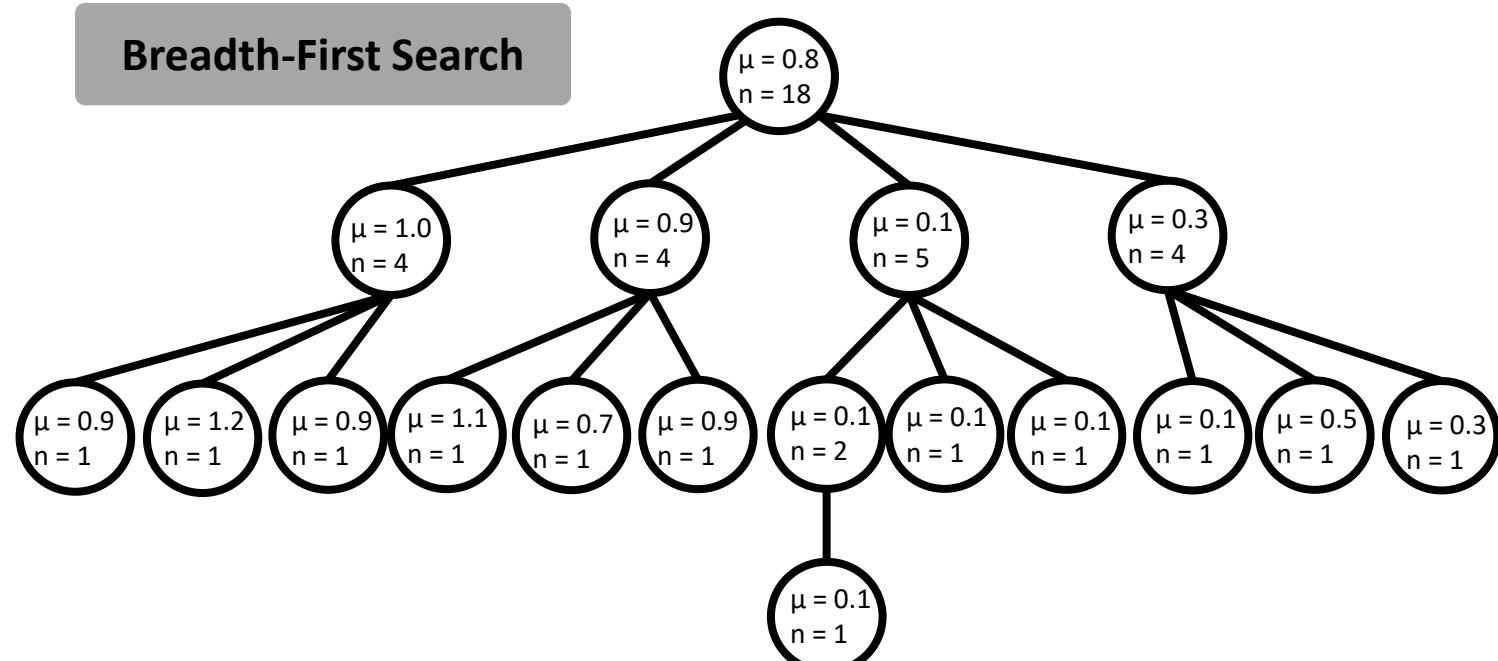
MCTS

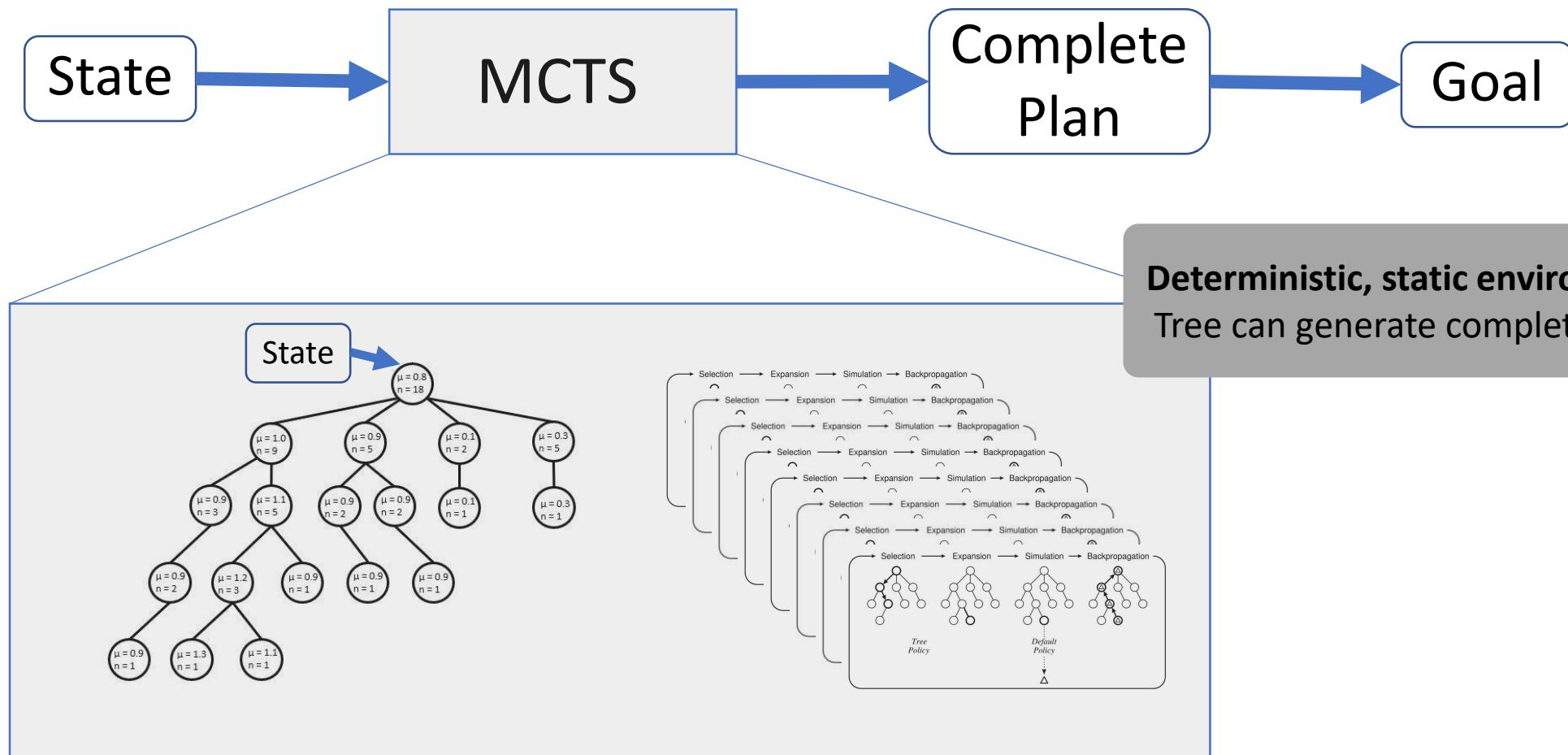


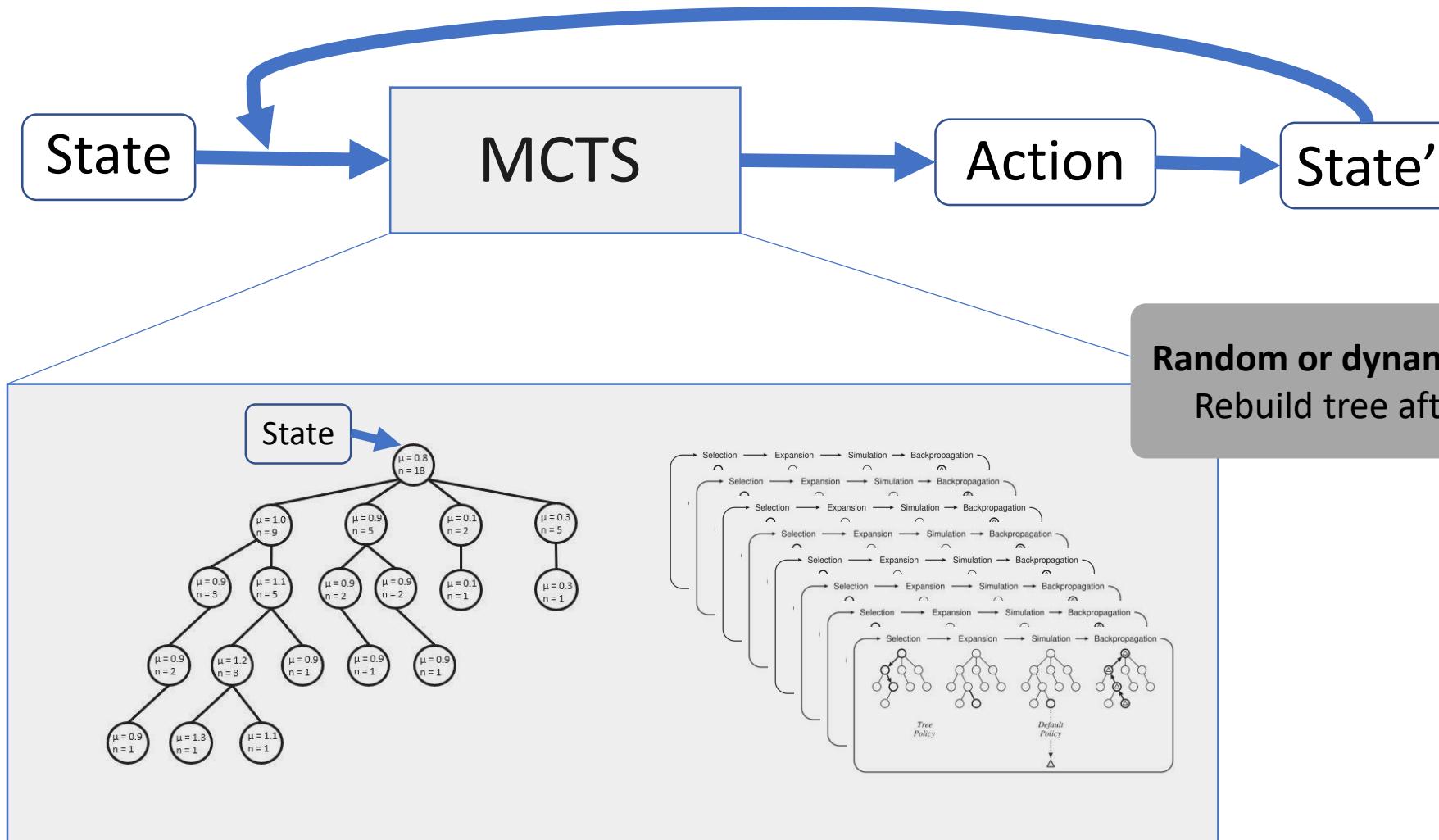
State



Search Tree







Aheuristic

A heuristic is learned rather than hand-defined

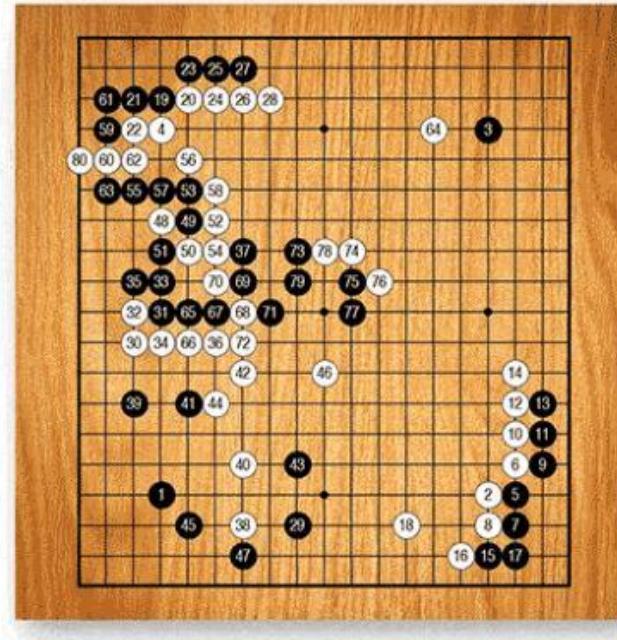
Asymmetric

Search tree focuses more on interesting areas

Anytime

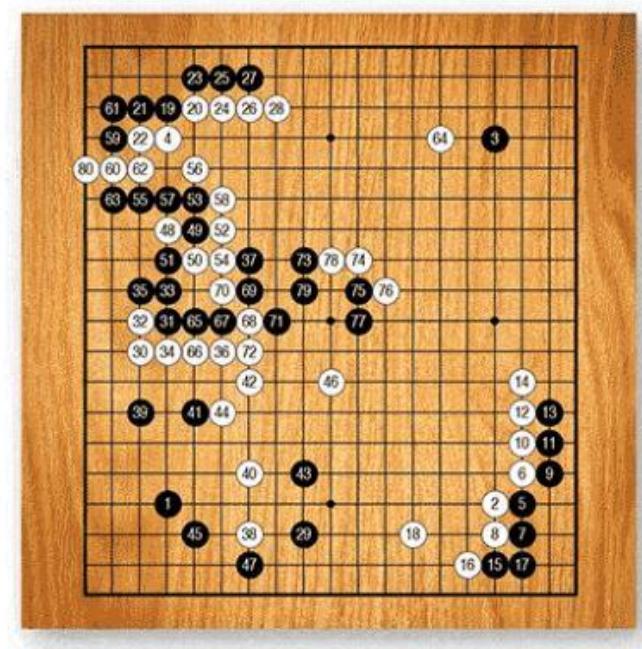
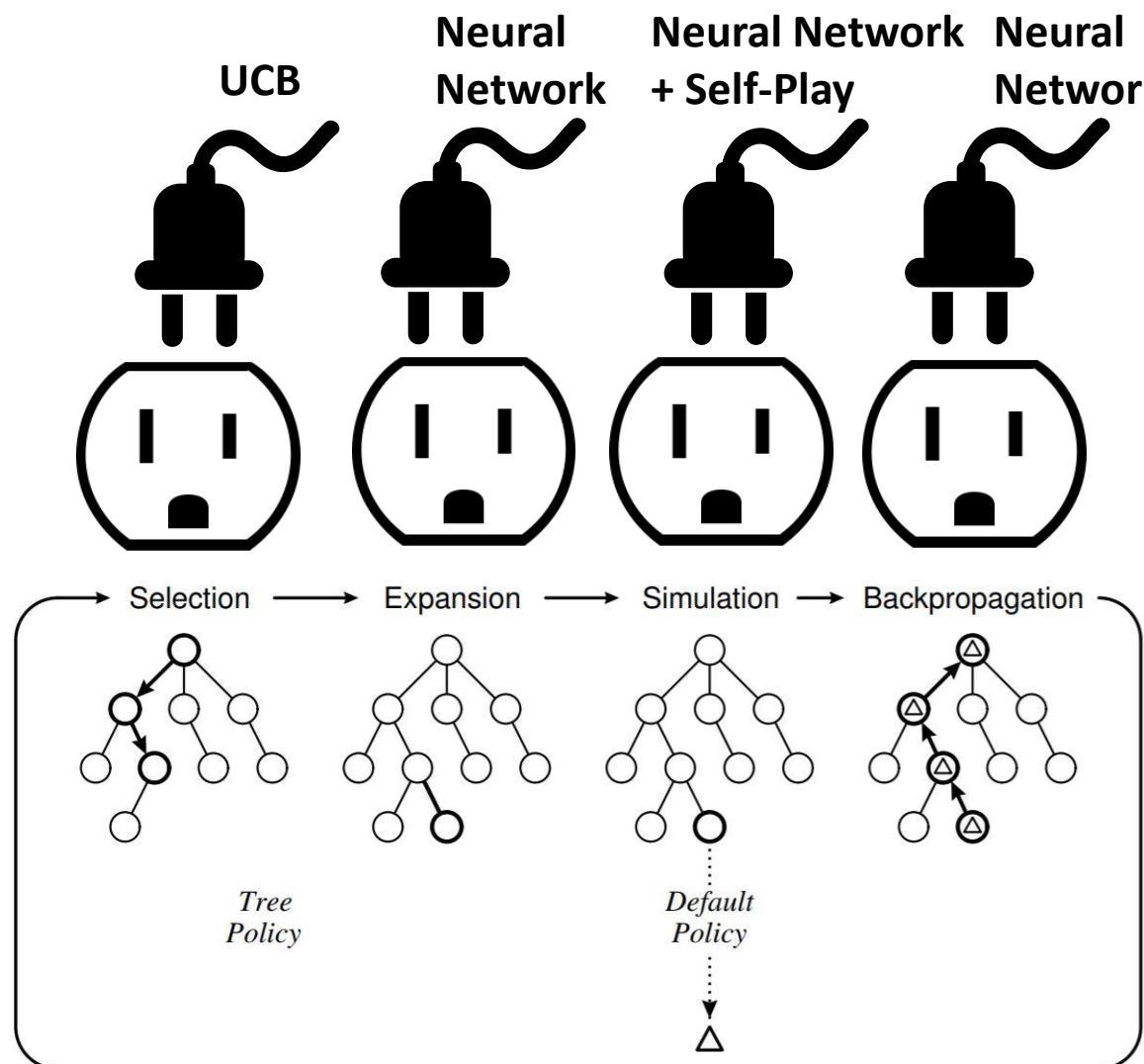
Proposes a complete policy at each step

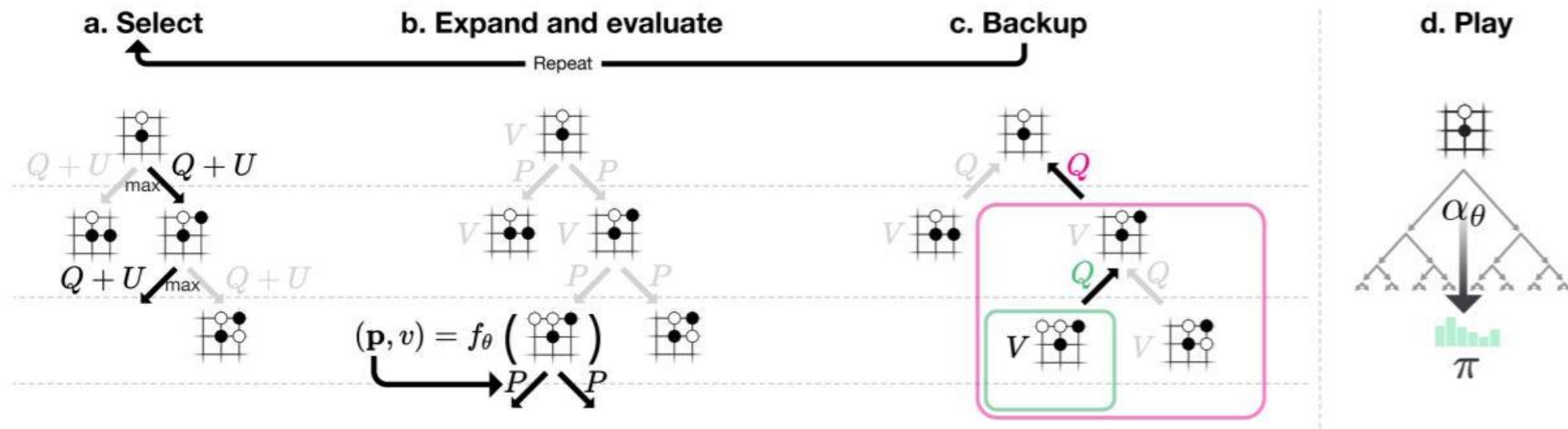
MCTS is modular and different configurations
may “tweak” these properties!



A game of Go.

Images from DeepMind website

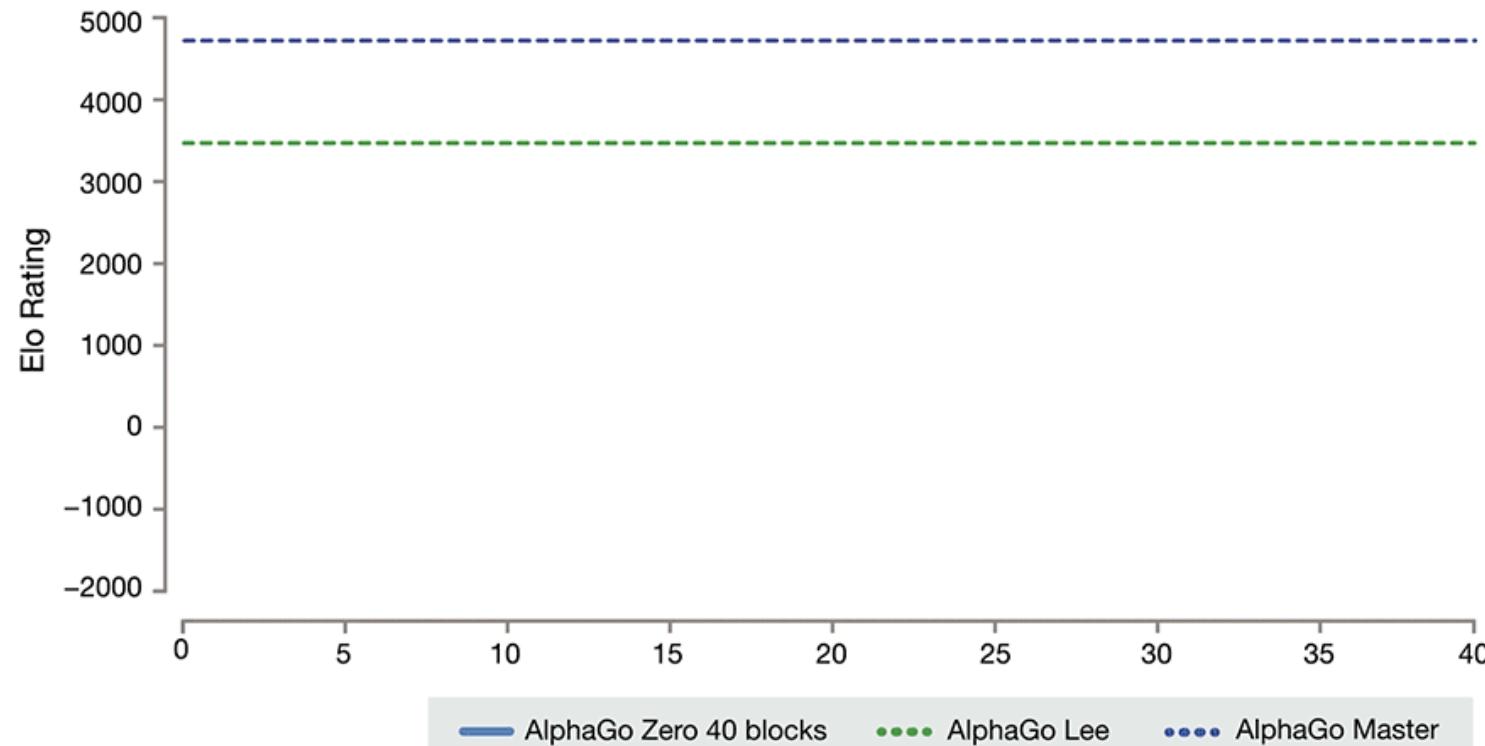




MCTS in AlphaGo Zero

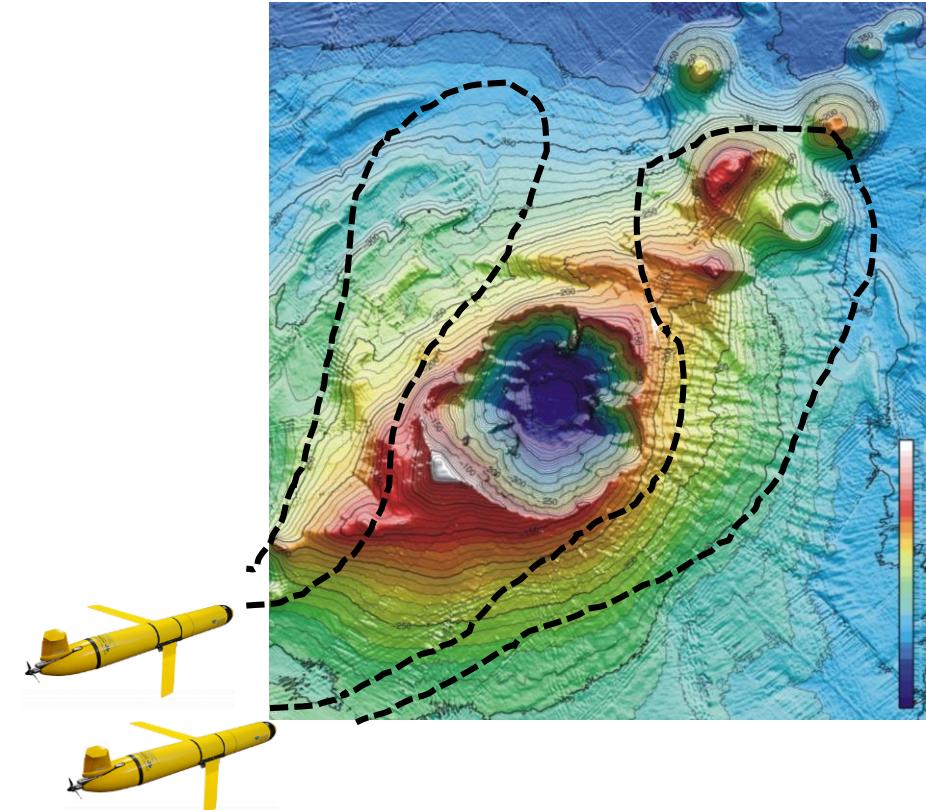
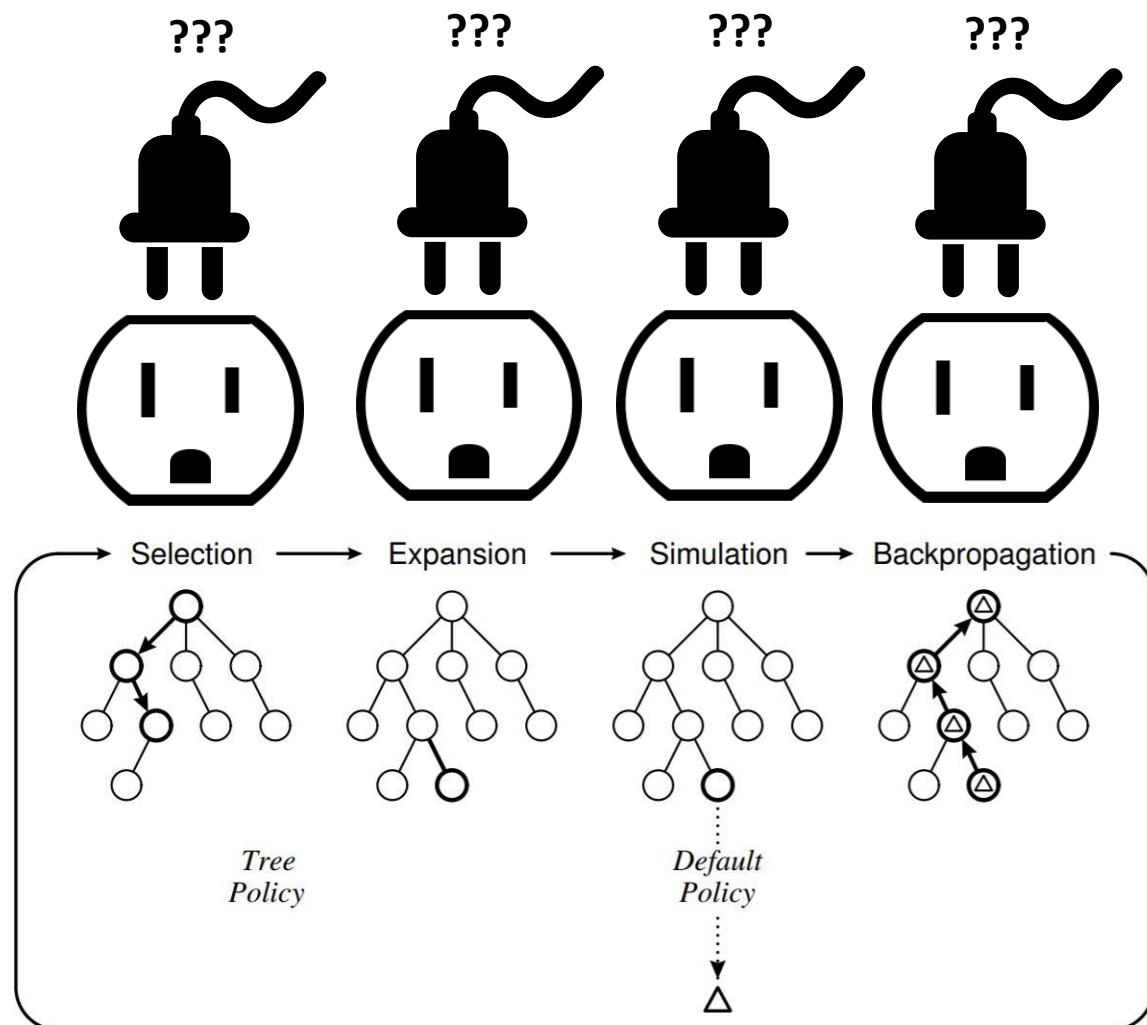
Images from DeepMind website

AlphaGo Zero Performance

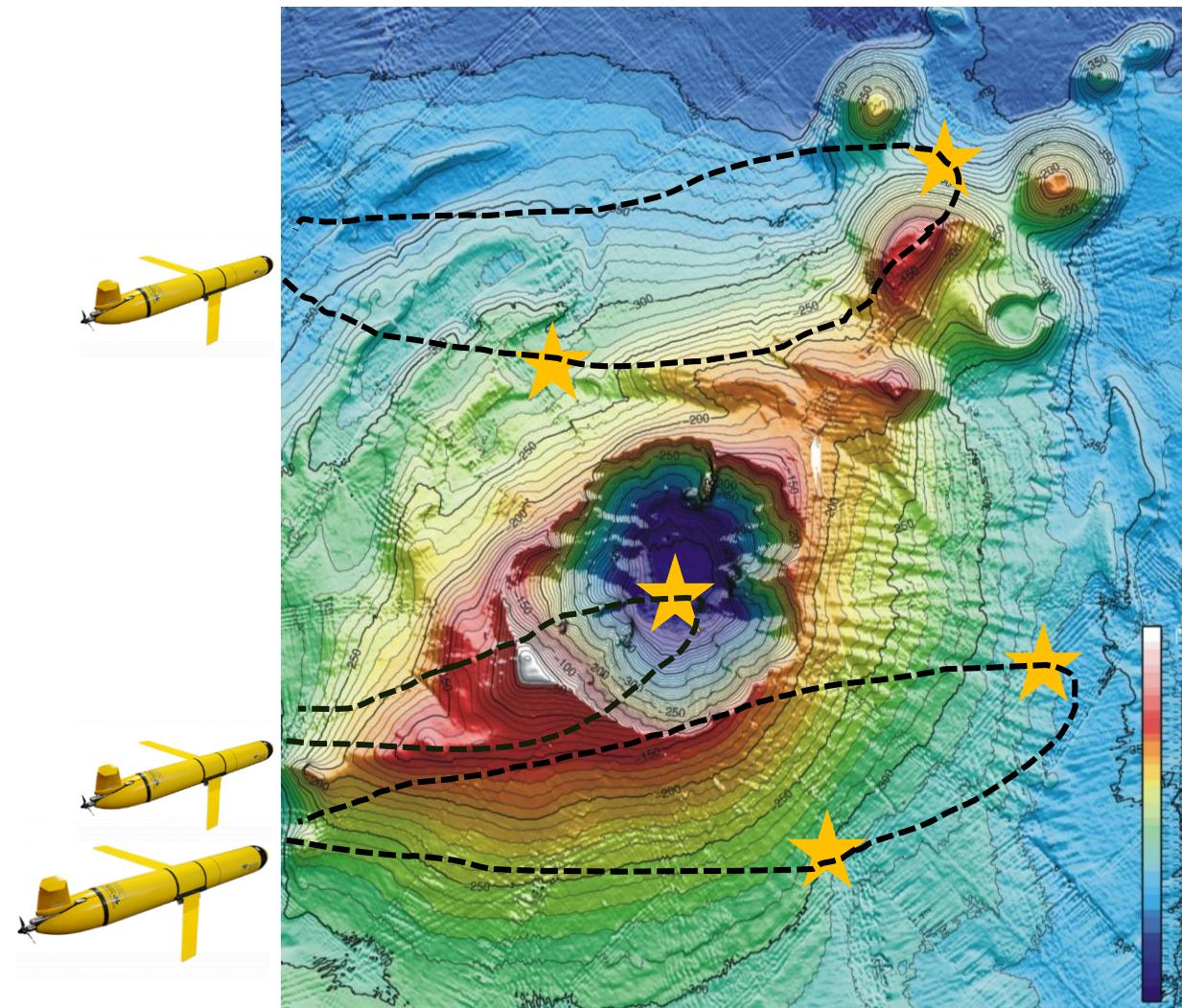


AlphaGo Zero performance.

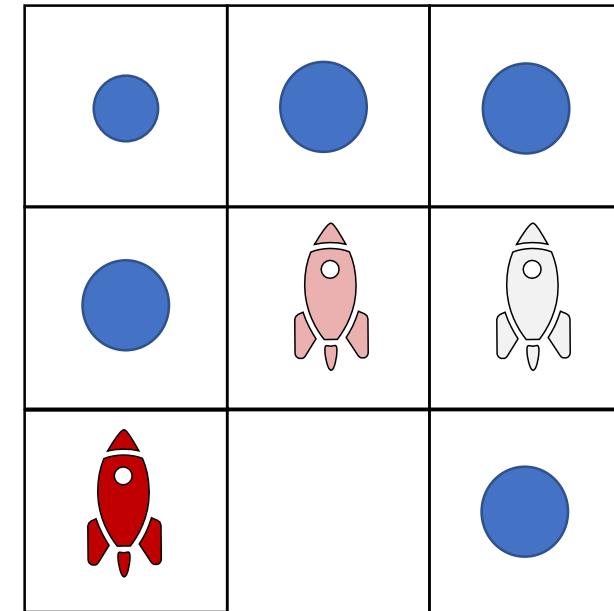
Images from DeepMind website



Extension to Collaboration

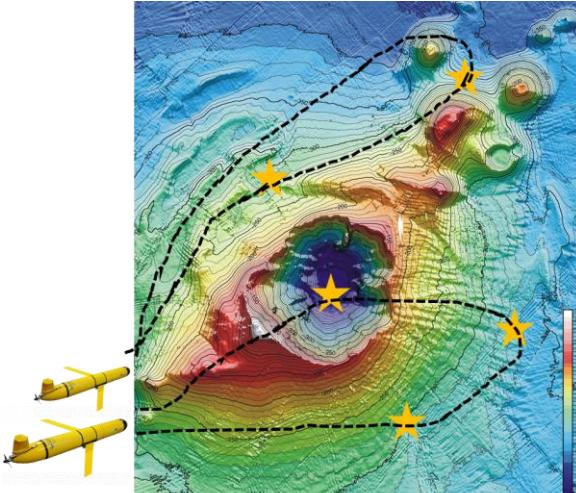
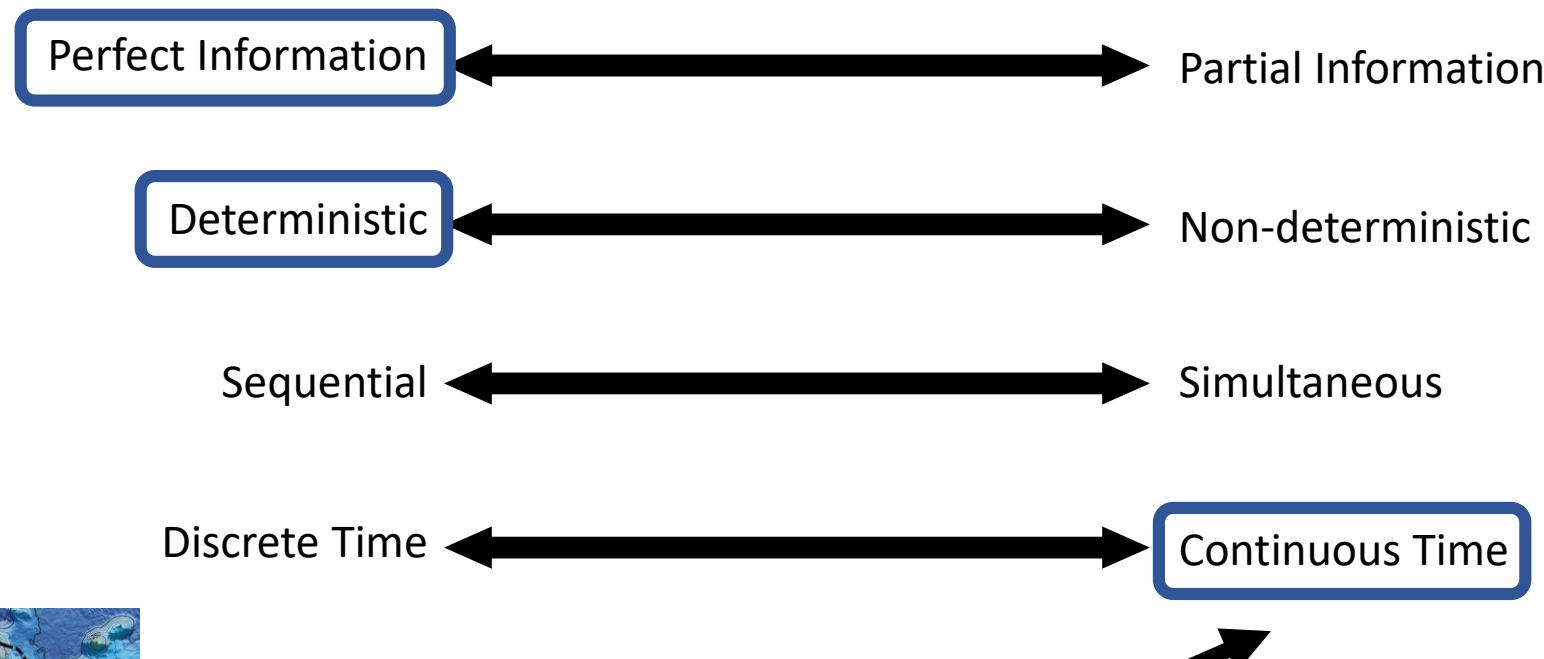


Discretized version with multiple gliders



Great. Solved, right?

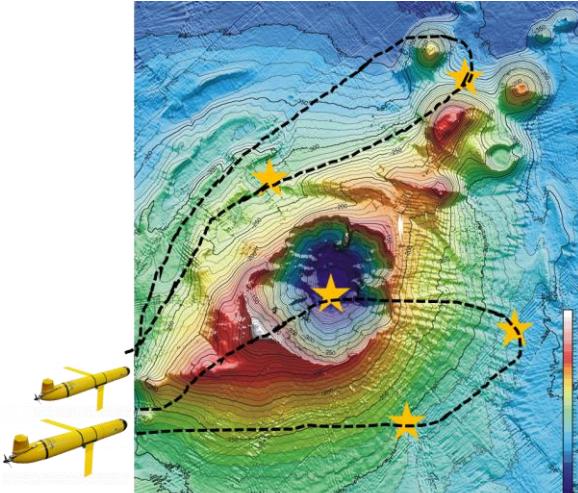
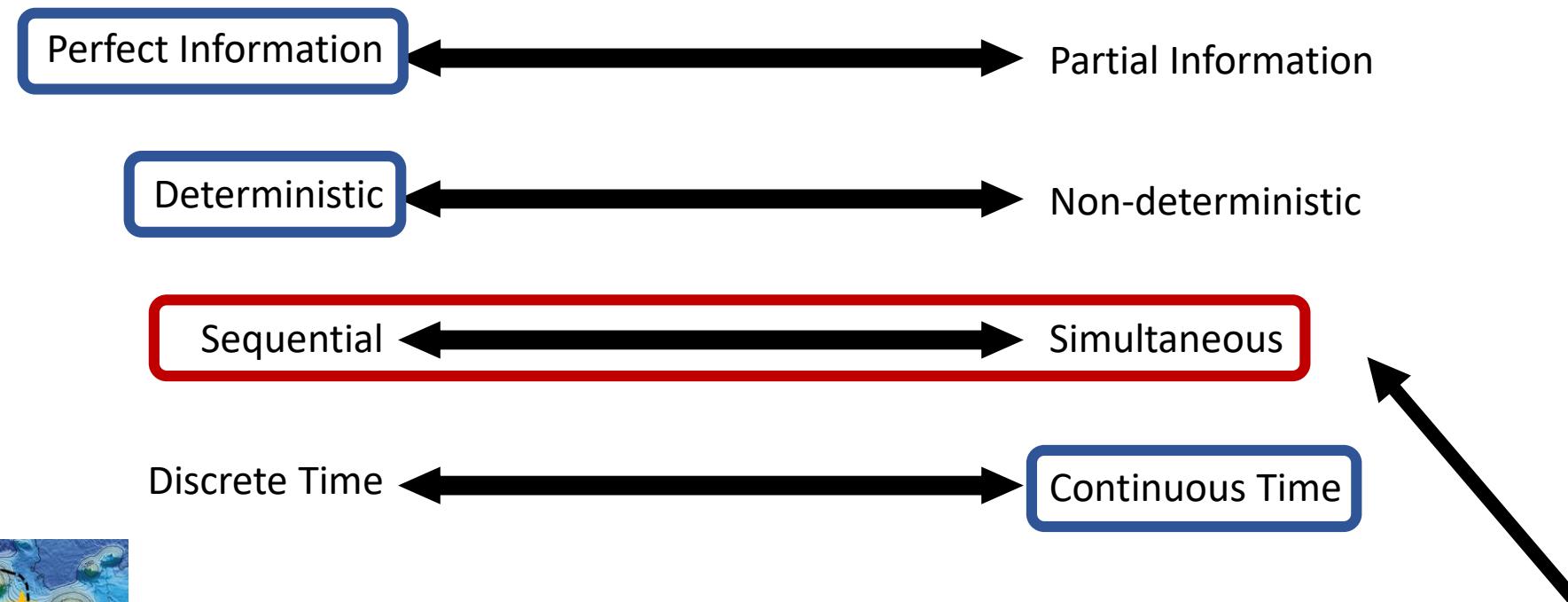
Extension to Collaboration



Because we have multiple agents acting, the world changes whether an individual agent takes a step or not. Thus our problem is continuous time. This results in:

- We must re-make the tree each time
- We have a limited amount of time to plan

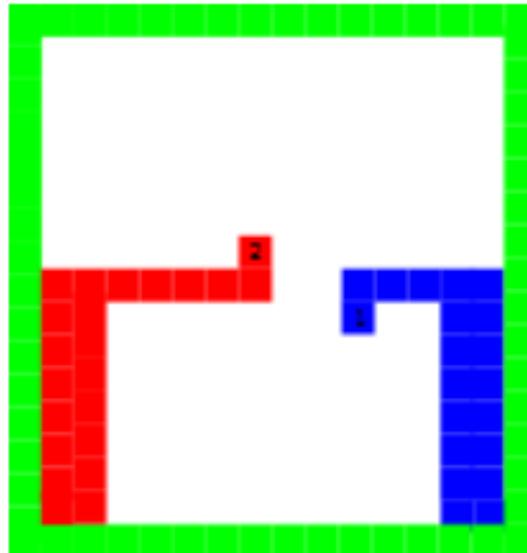
Extension to Collaboration



Unlike tic tac toe and the single agent application of MCTS to the search of the Kolumbo Volcano site, our multiple agents would **not act sequentially but rather in parallel.**

Does the non-sequential action of the agents matter?

Short answer: Probably not



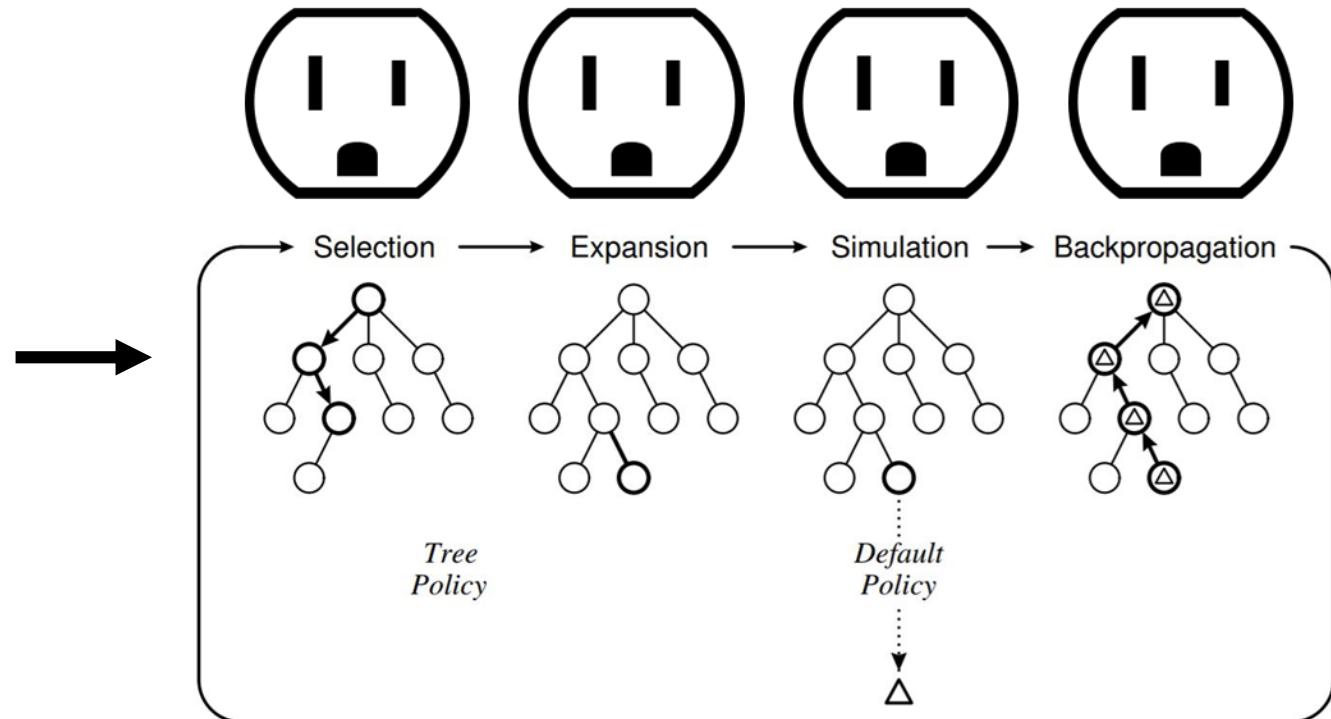
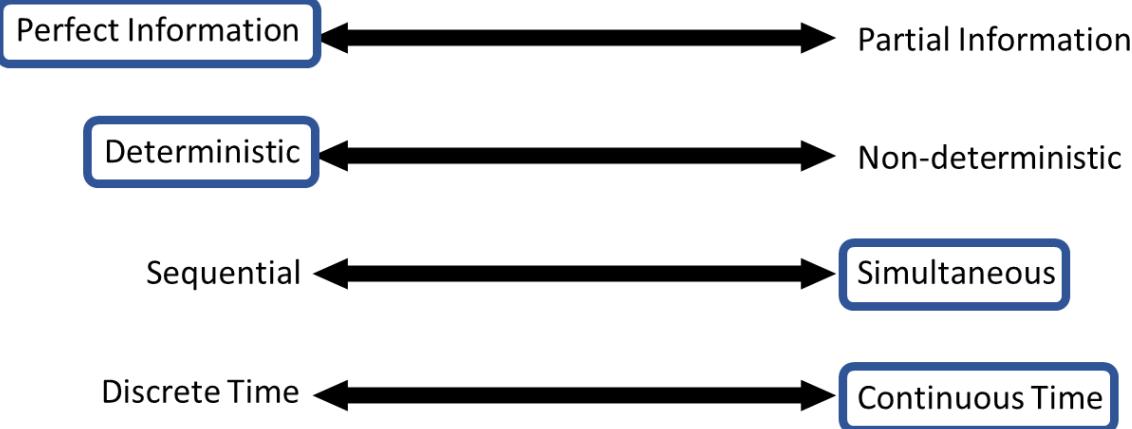
Samothrakis et al. applied MCTS to the game TRON in 2010. They stated:

“Although a simultaneous move game, it can be transformed into an alternate moves games by assuming that in each player’s tree node the current player plays always first”

They pointed to the University of Waterloo Tron Competition, where the best versions of MCTS—employing TRON-specific leaf heuristics—performed comparably to other approaches.

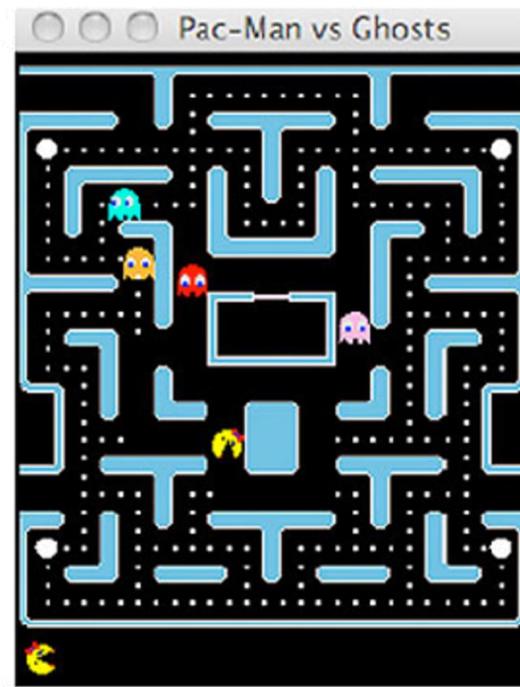
S. Samothrakis, D. Robles, and S. M. Lucas, “A uct agent for tron: Initial investigations,” in Proceedings of the 2010 IEEE Conference on Computational Intelligence and Games, 2010, pp. 365–371. DOI: 10.1109/CIG.2010.5560750

The guts of MCTS will remain the same

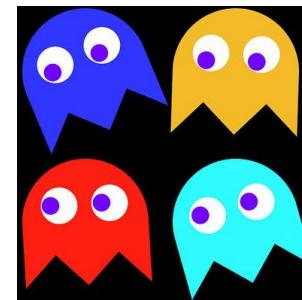


- Background
 - Monte Carlo Simulation
 - Upper Confidence Bounds
- Implementing MCTS with Upper Confidence Bounds for Trees
- Collaborative MCTS
 - Discuss Max-N, a variant of MCTS for multiple agents
 - Discuss how to decentralize MCTS
 - Discuss how to improve the default policy in decentralized MCTS

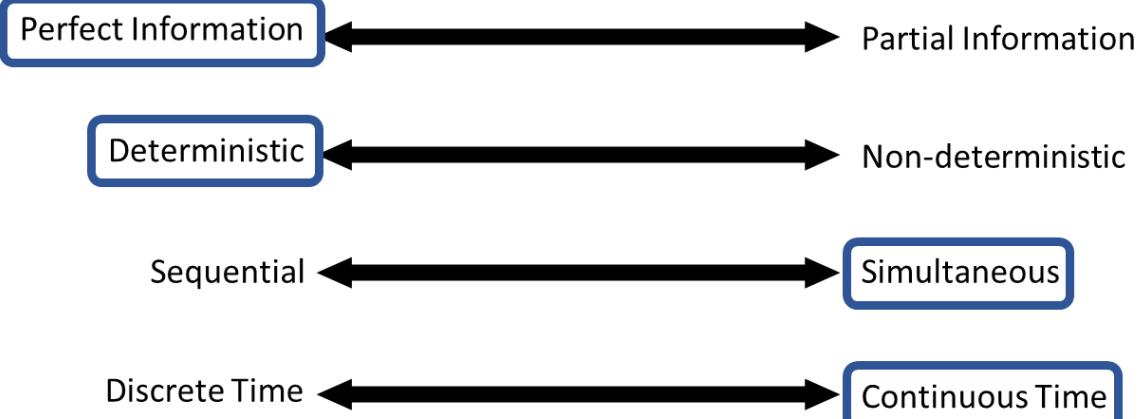
Ms. Pacman Competition (SM Lucas, 2007)



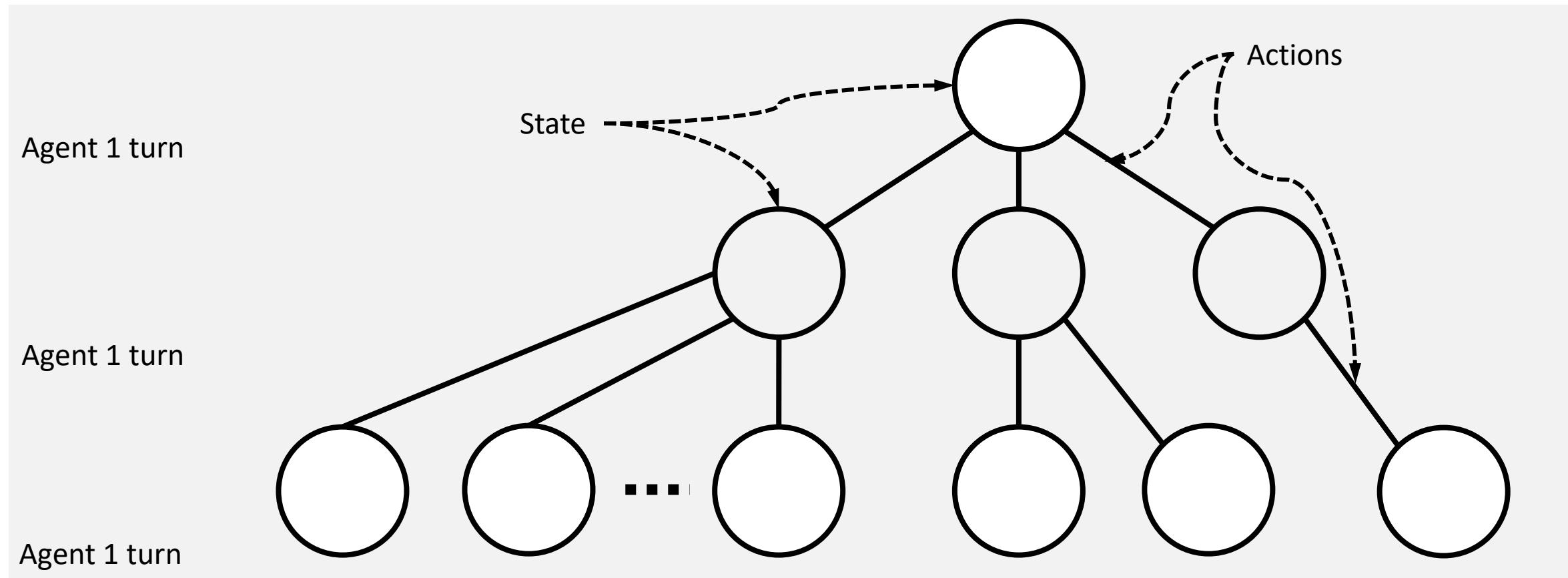
Team 1

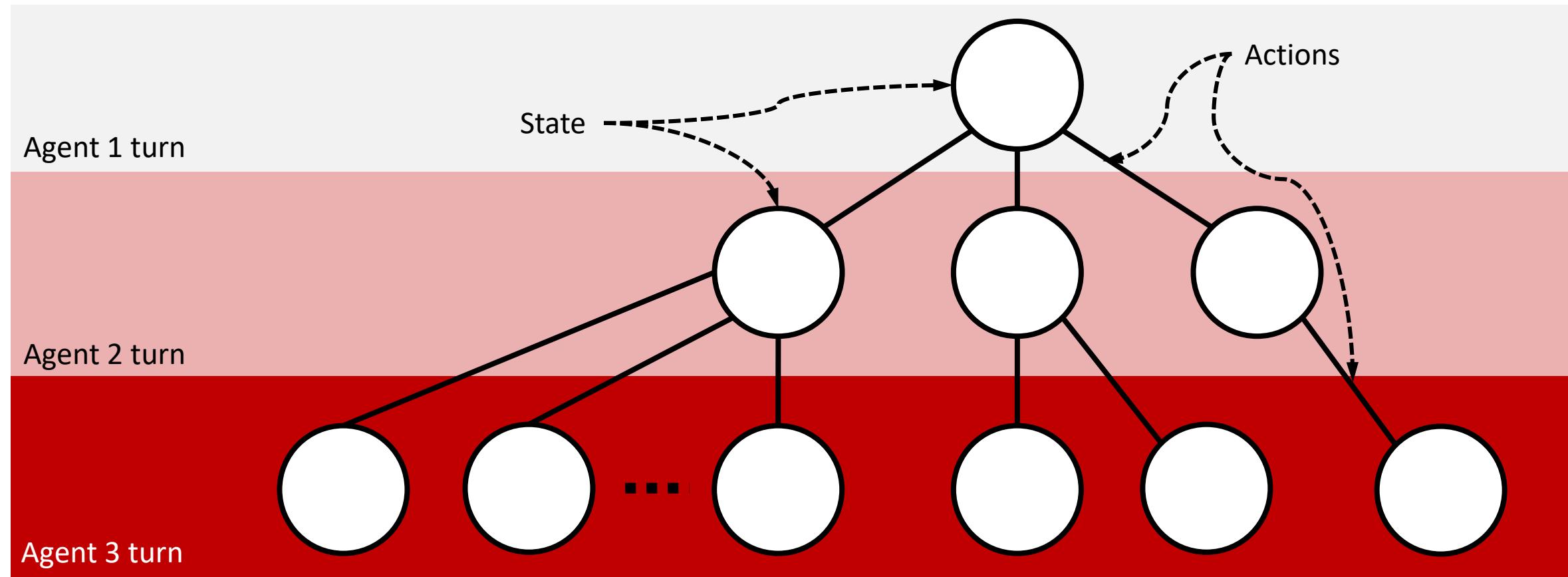


Team 2

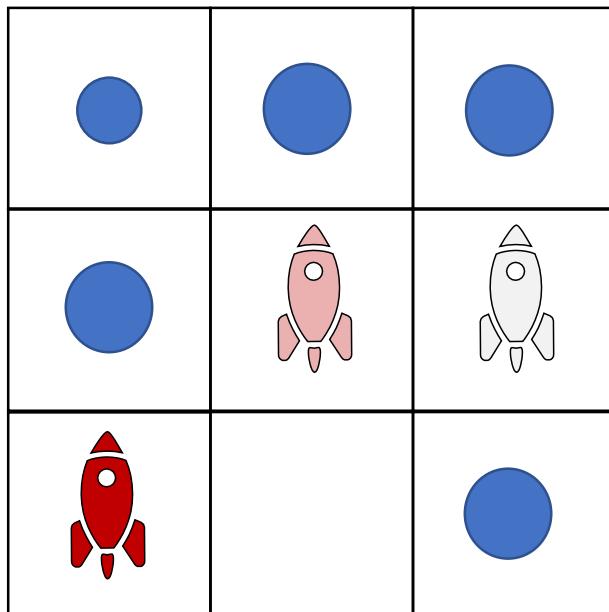


Ghosts work collaboratively
to eat Ms. Pacman

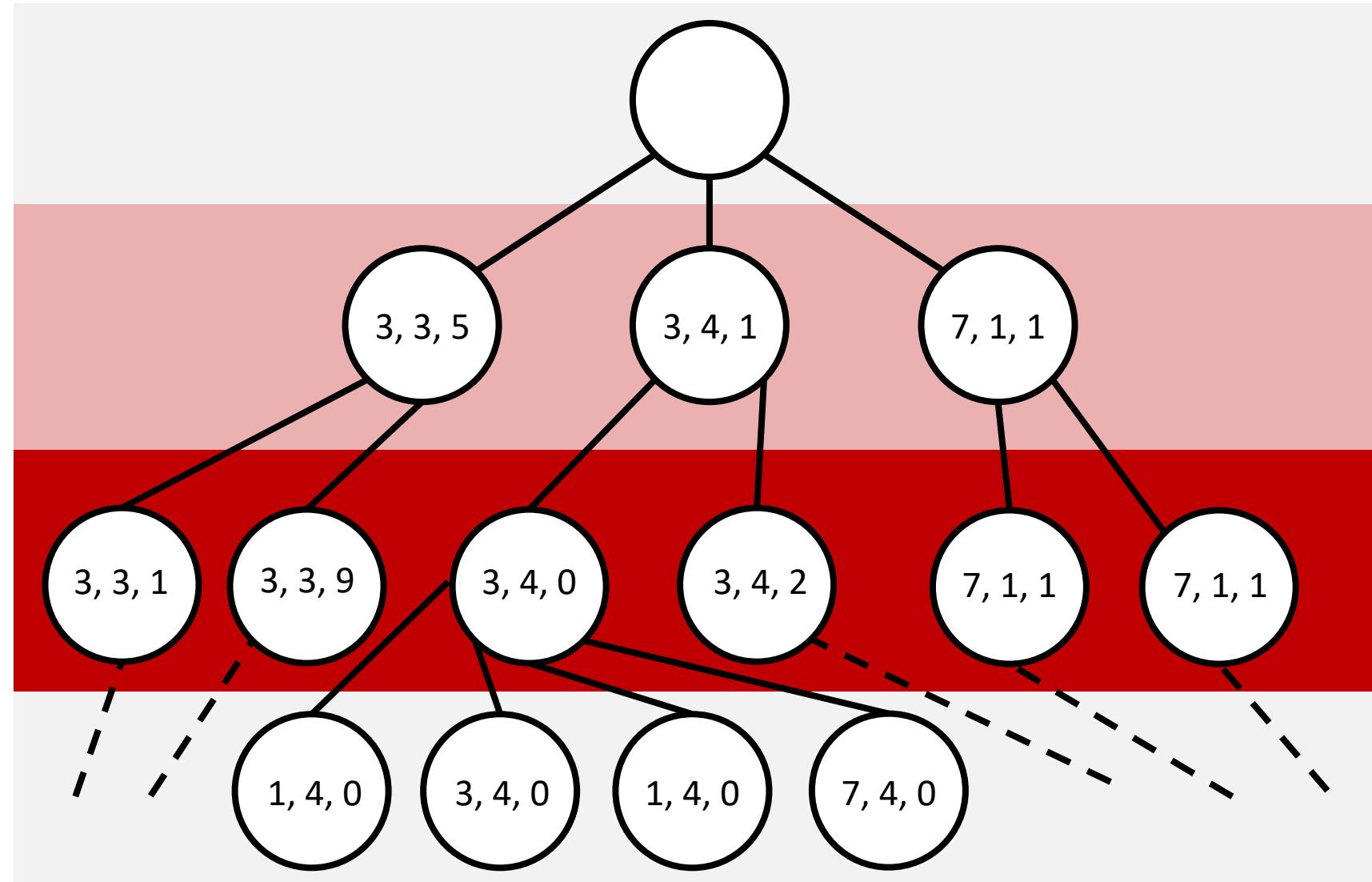


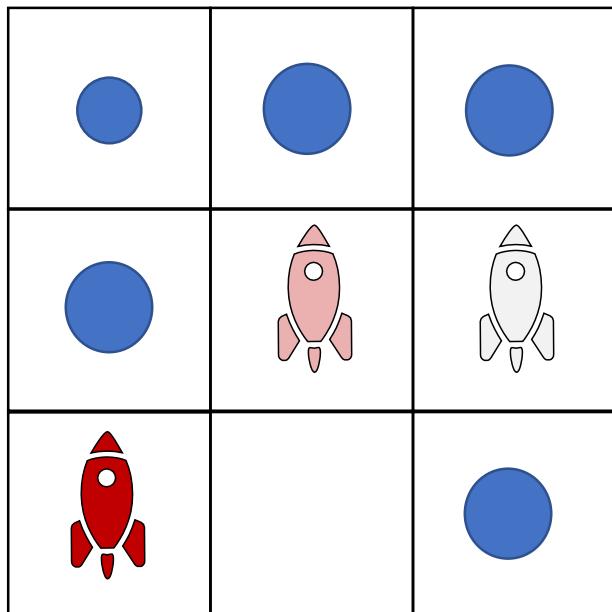


Leverage the finding from the TRON paper (Samothrakis 2010)
to assign the agents' moves sequentially

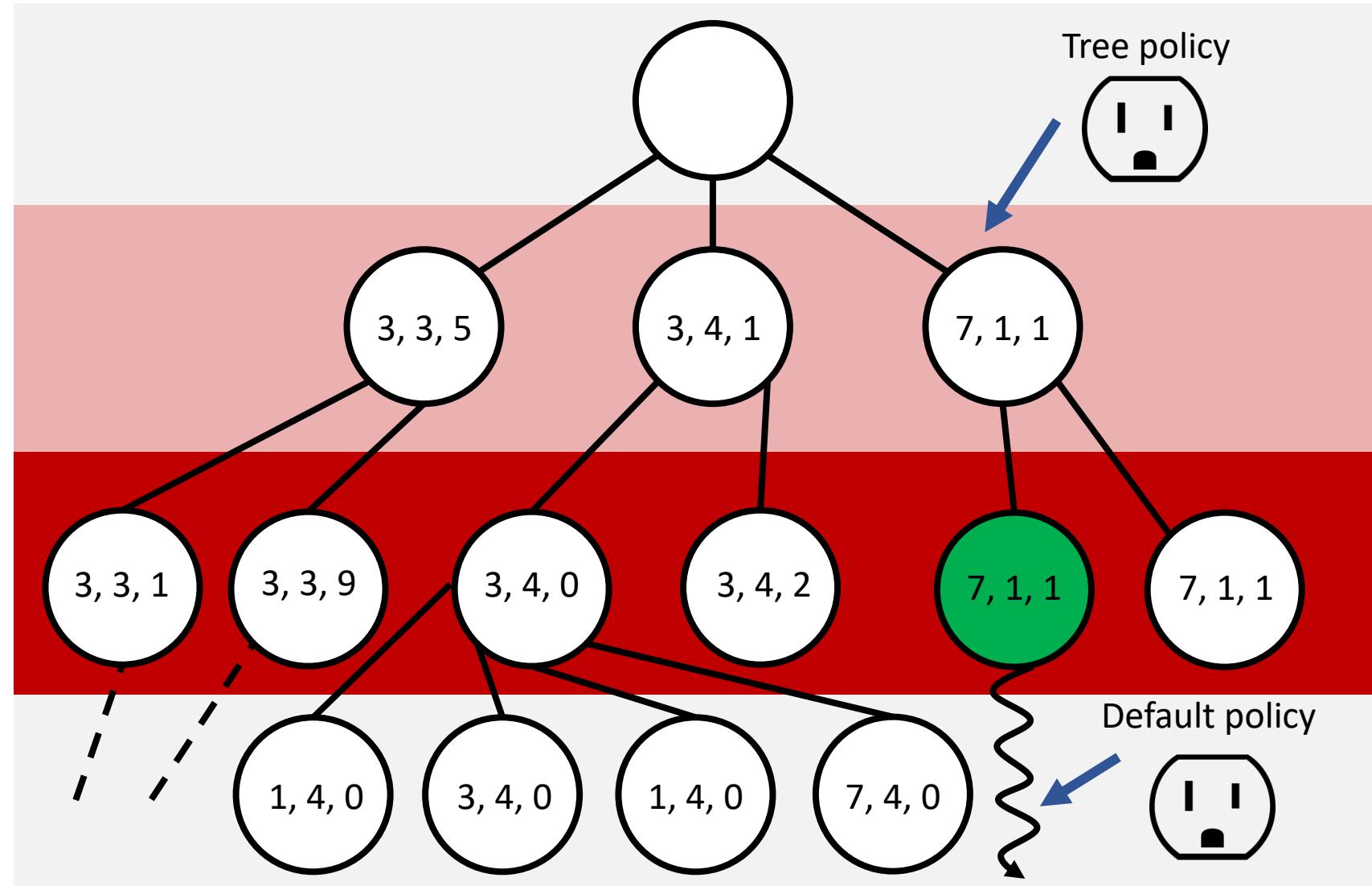


Instead of returning a single value for the playout, we return an N-tuple for N agents



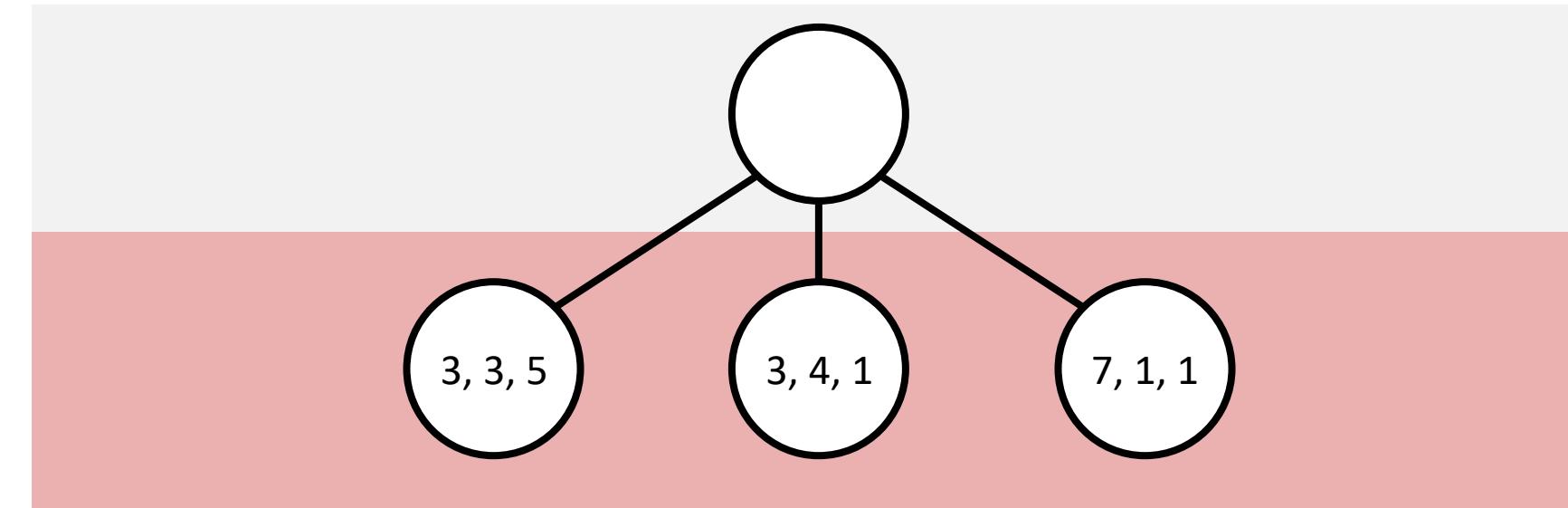


Instead of returning a single value for the playout, we return an N-tuple for N agents



Canonical Max-N:
Players are adversarial

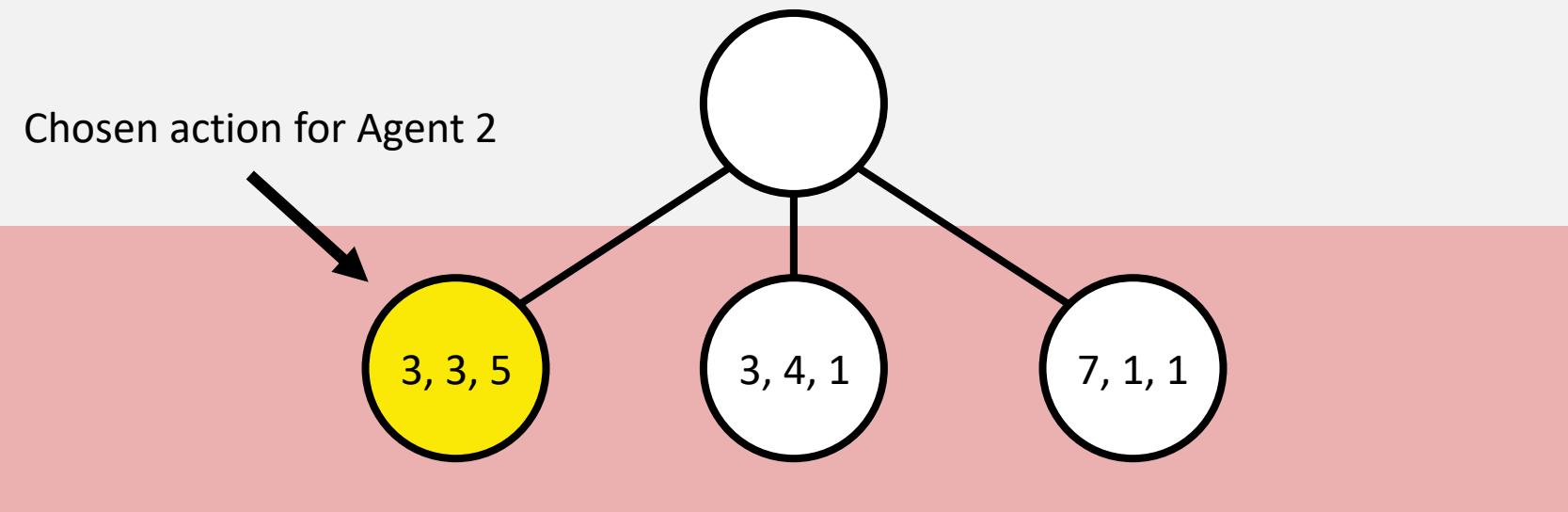
- All players against you
 - Paranoid UCT
- Fixed coalitions
 - UCT with Alliances
- Flexible coalitions
 - Confident UCT



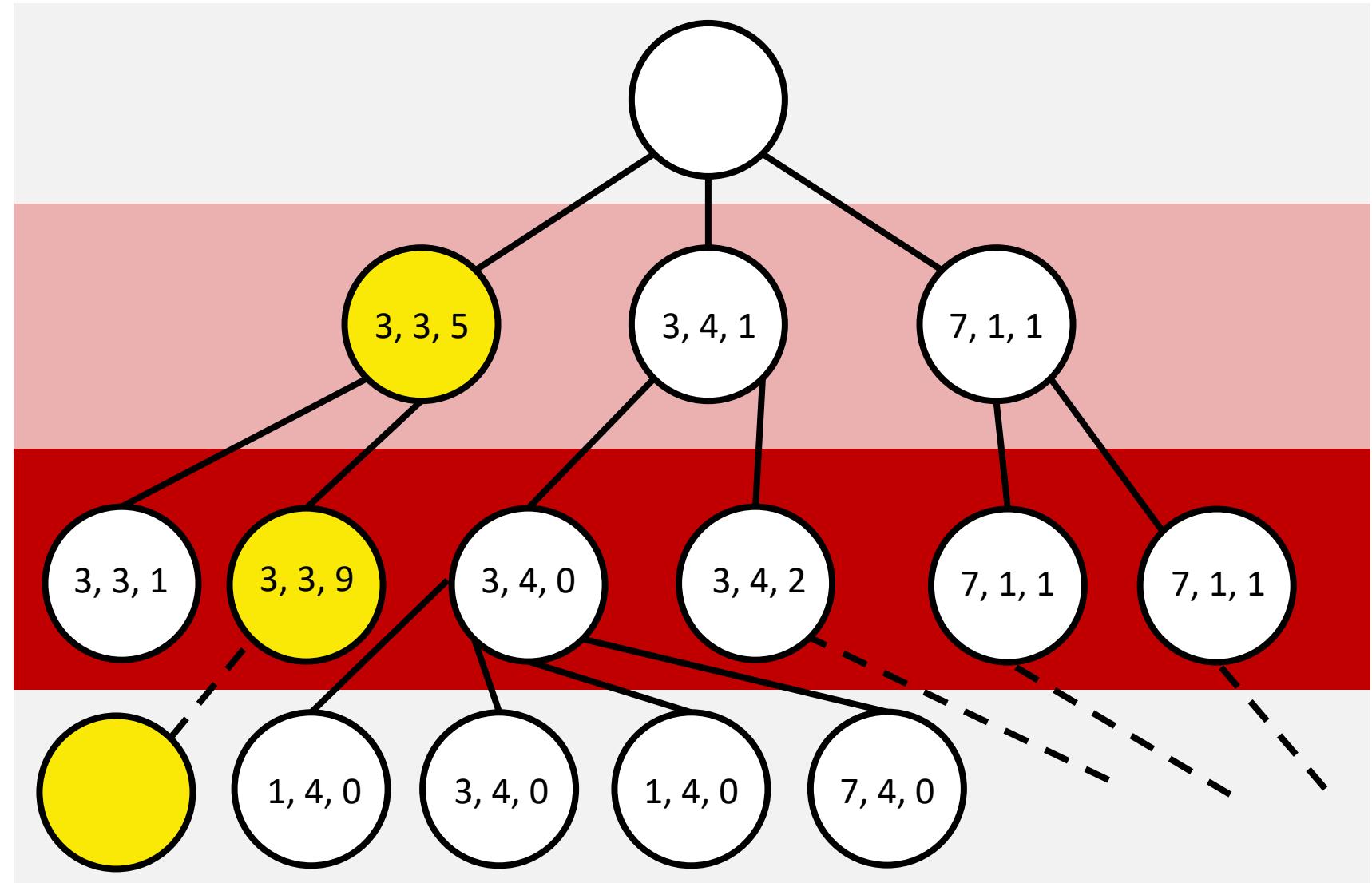
Because our agents are collaborative,
we can choose the action that
maximizes sum of the scores

Canonical Max-N:
Players are adversarial

- All players against you
 - Paranoid UCT
- Fixed coalitions
 - UCT with Alliances
- Flexible coalitions
 - Confident UCT



Because our agents are collaborative,
we can choose the action that
maximizes sum of the scores



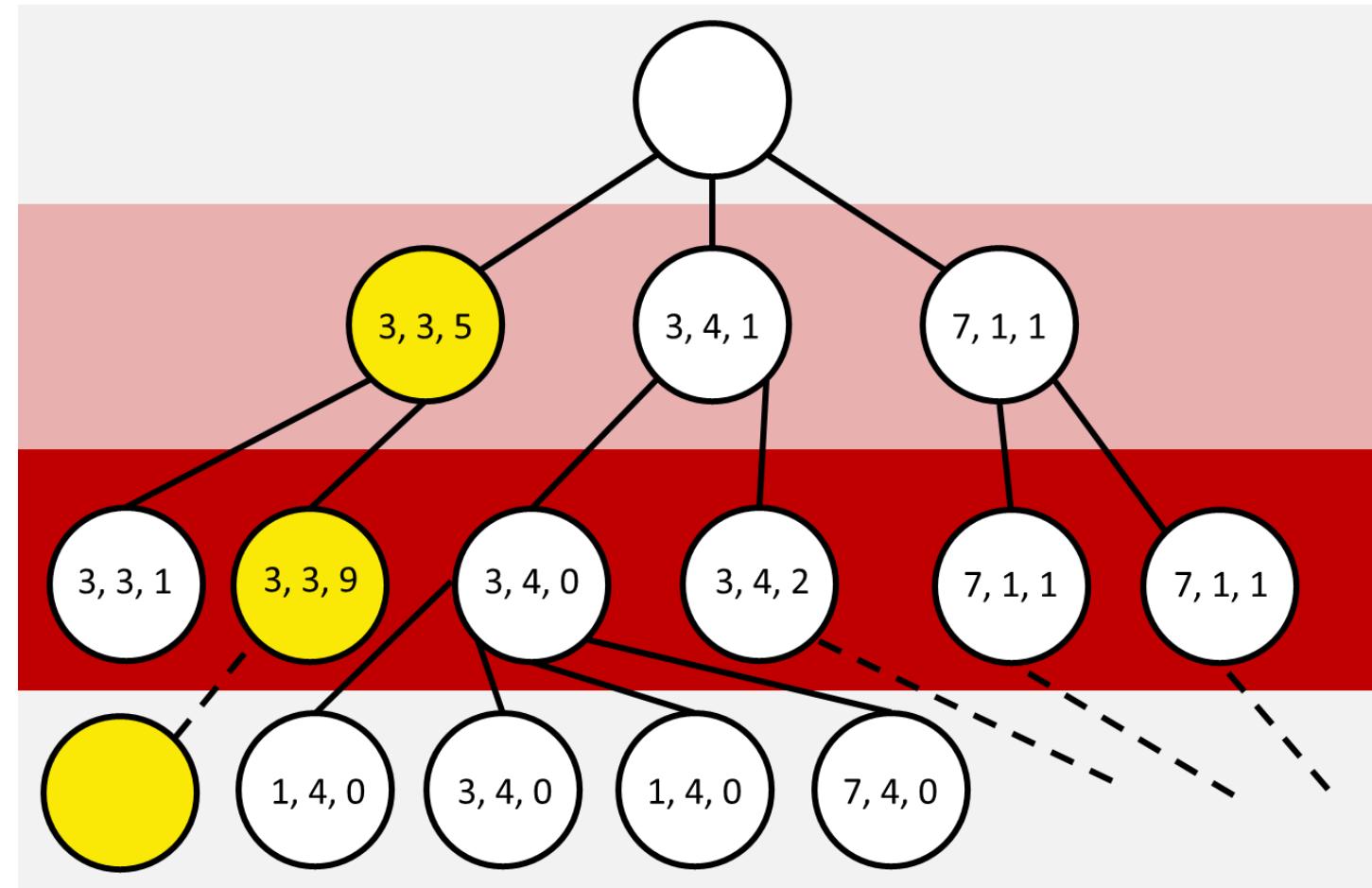
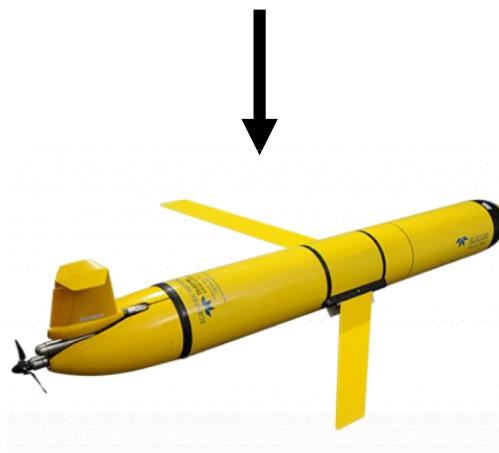
Bilal Kartal, Ernesto Nunes, Julio Godoy, and Maria Gini. Monte carlo tree search for multi-robot task allocation. In Thirtieth AAAI Conference on Artificial Intelligence, pages 4222–4223, 2016.

Benefits

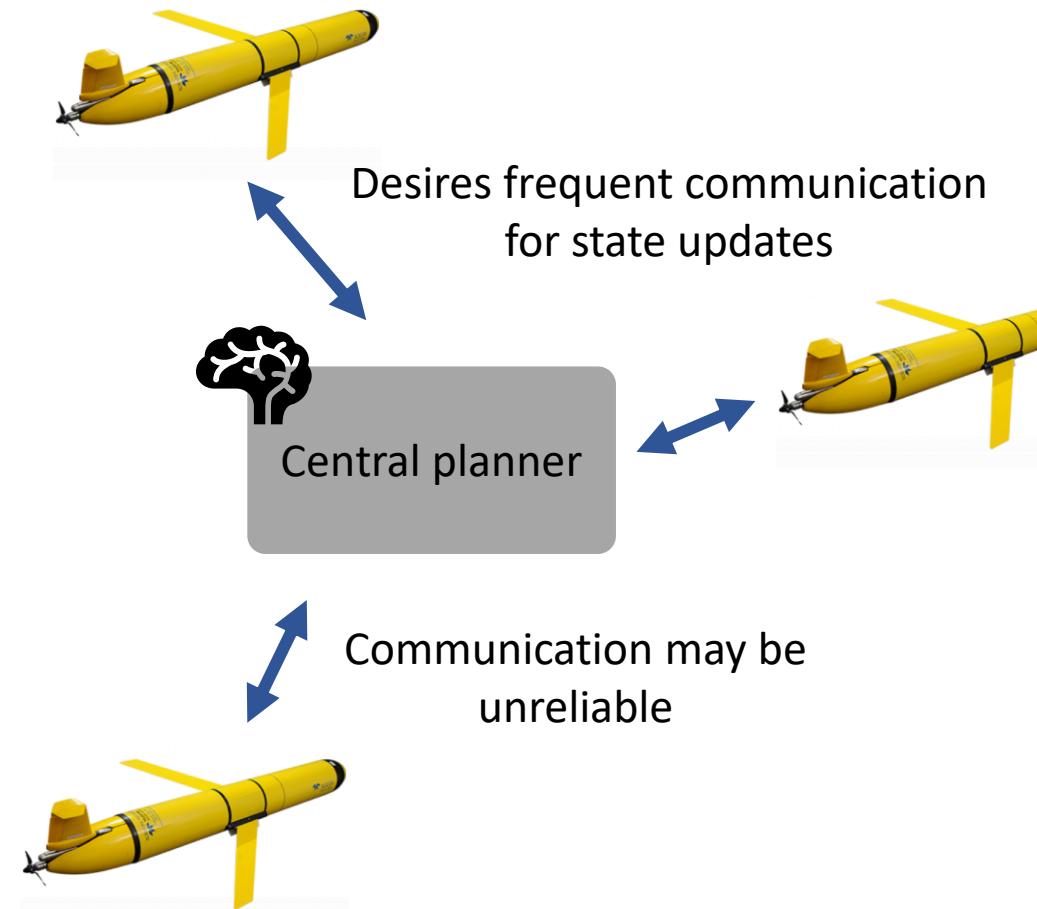
- Simple, modular

Disadvantages

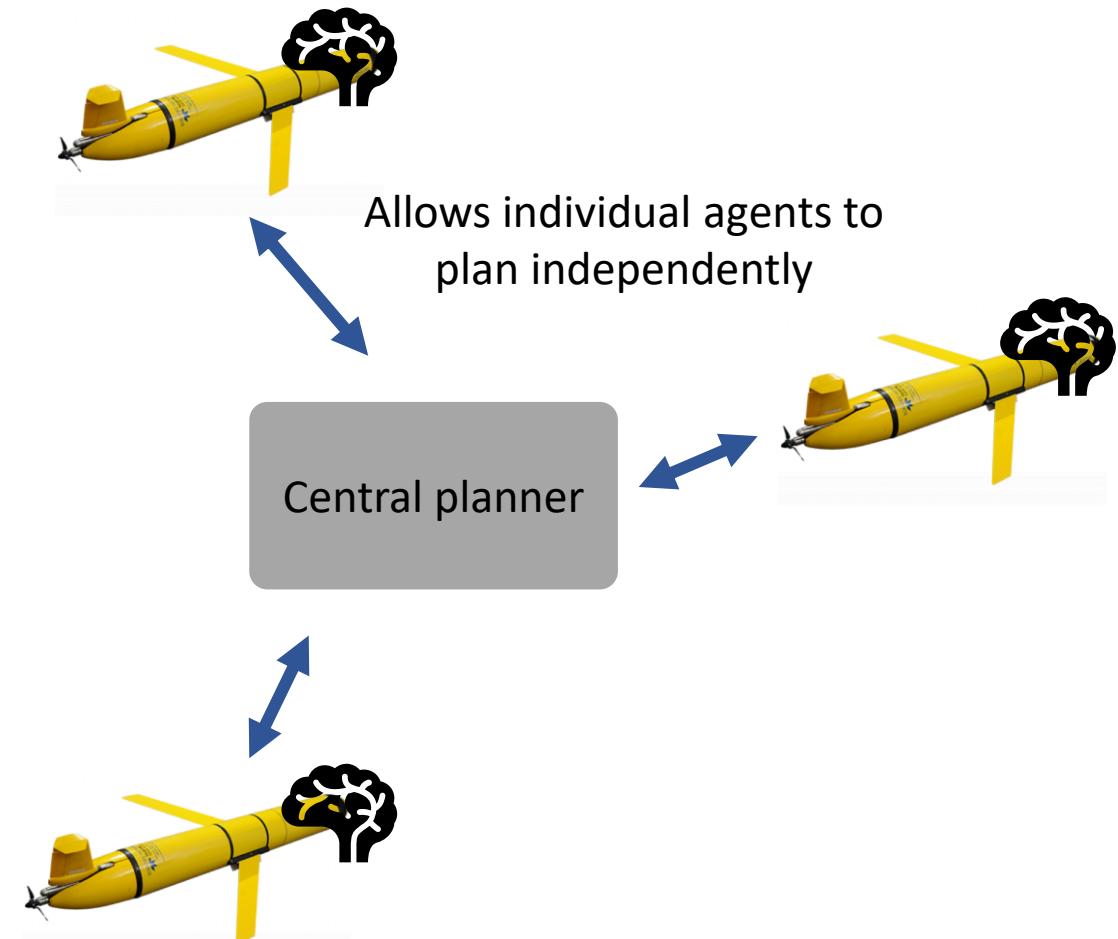
- Requires robust communication between agents and master



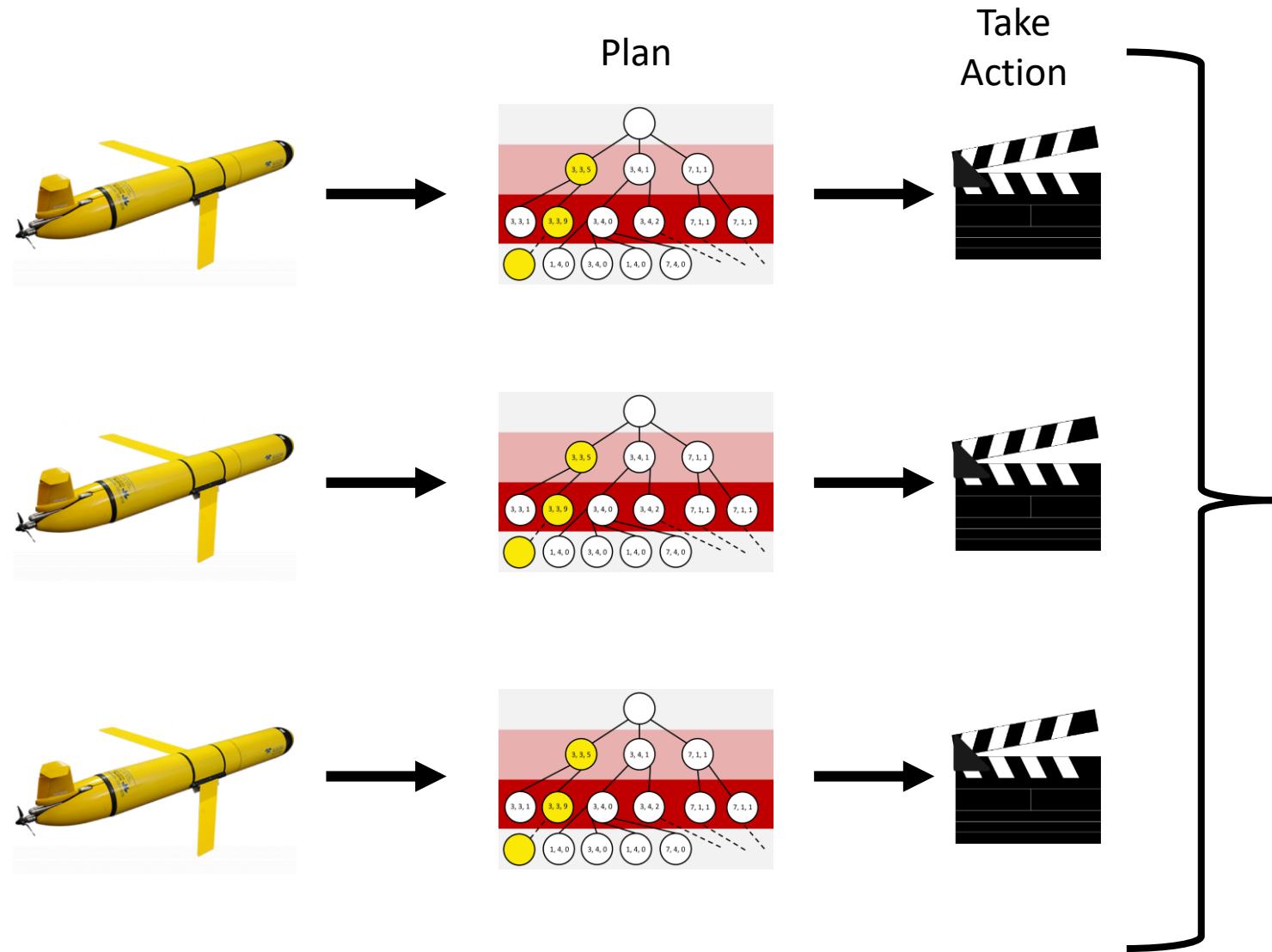
Central planner has authority



Hand authority over to agents



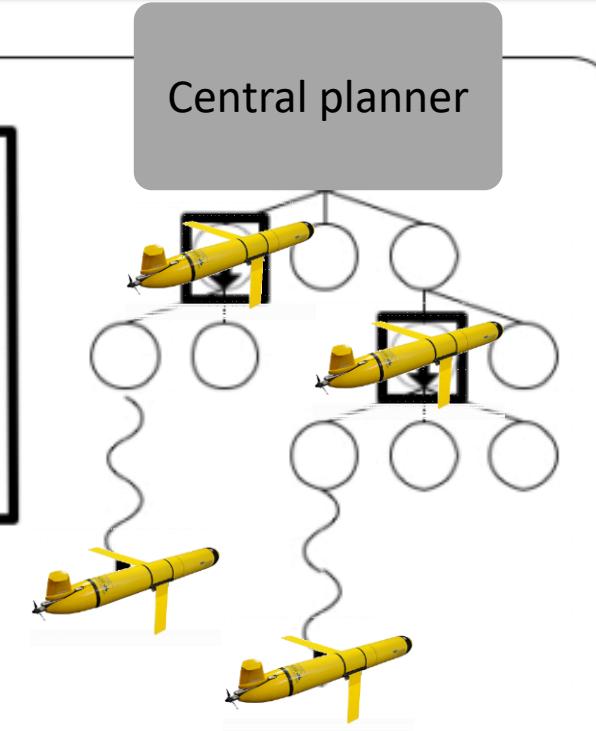
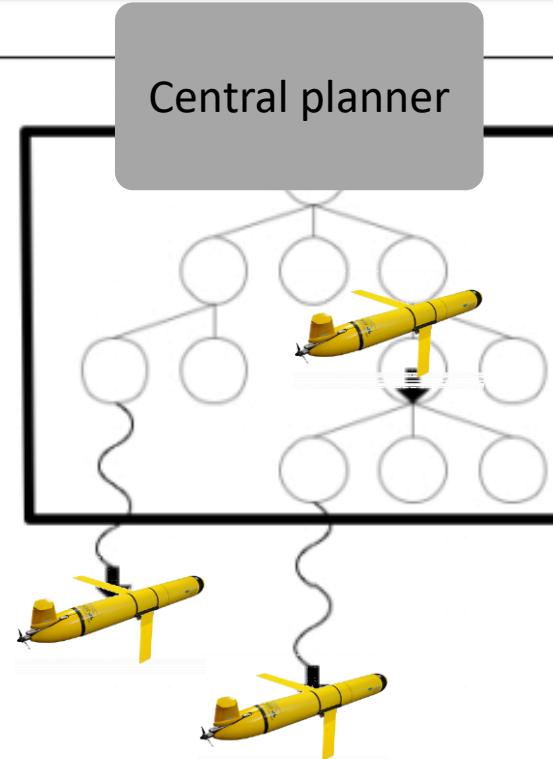
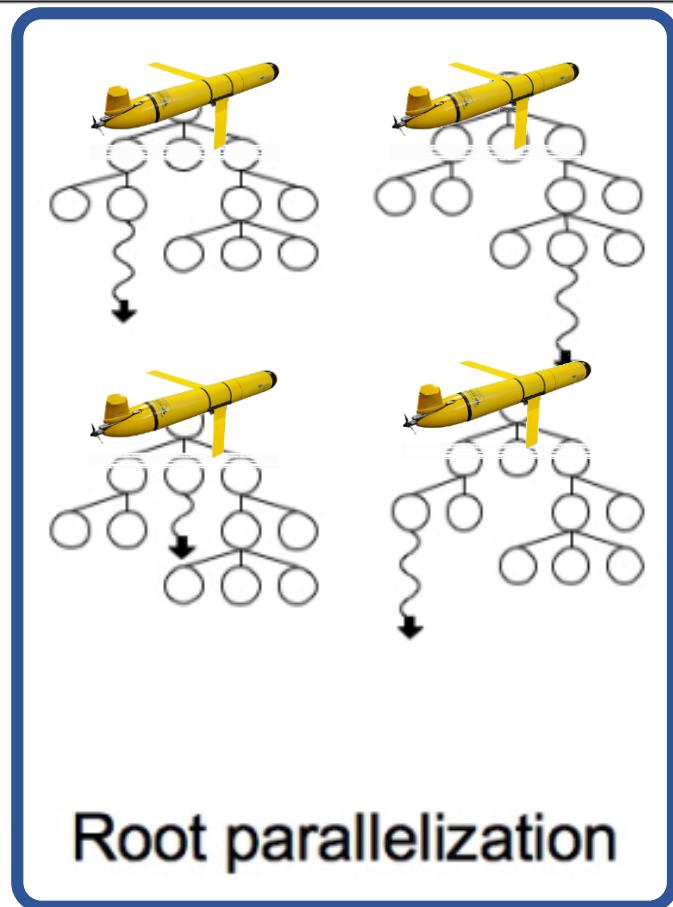
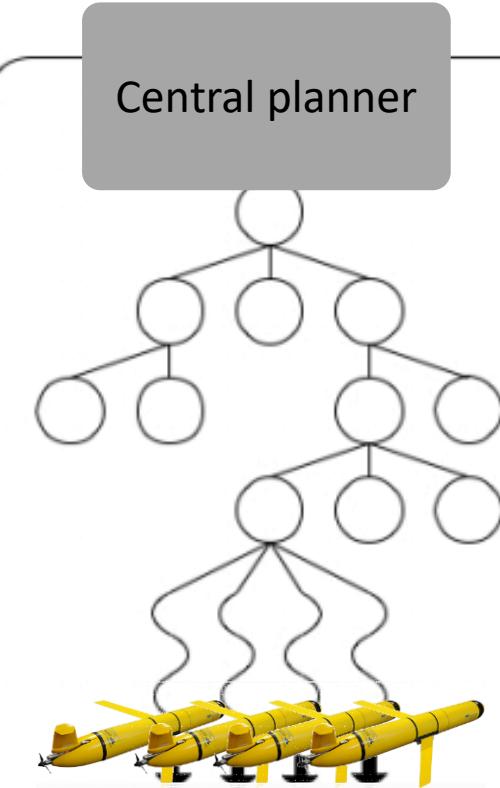
Handing over authority to the agents



- Each agent builds a multi-agent tree
- Keep tree policy as UCB
- Keep default policy as random simulation for each agent

What kind of
parallelization is this?

Brief aside on types of parallelization within MCTS



○ = Tree node

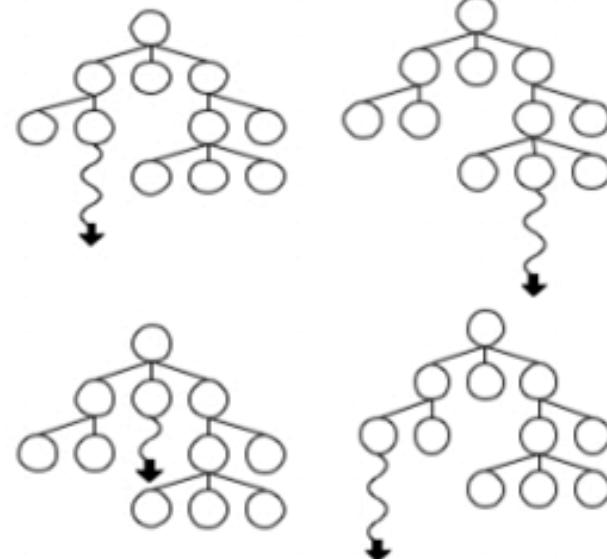
↓ = Thread location

□ = Locked memory section

⌘ = Tree-independent simulation

Browne, C. et al. A survey of Monte-Carlo tree search methods. *IEEE Trans. Comput. Intell. AI in Games* 4, 1–43 (2012)

Additional step in Root Parallelization

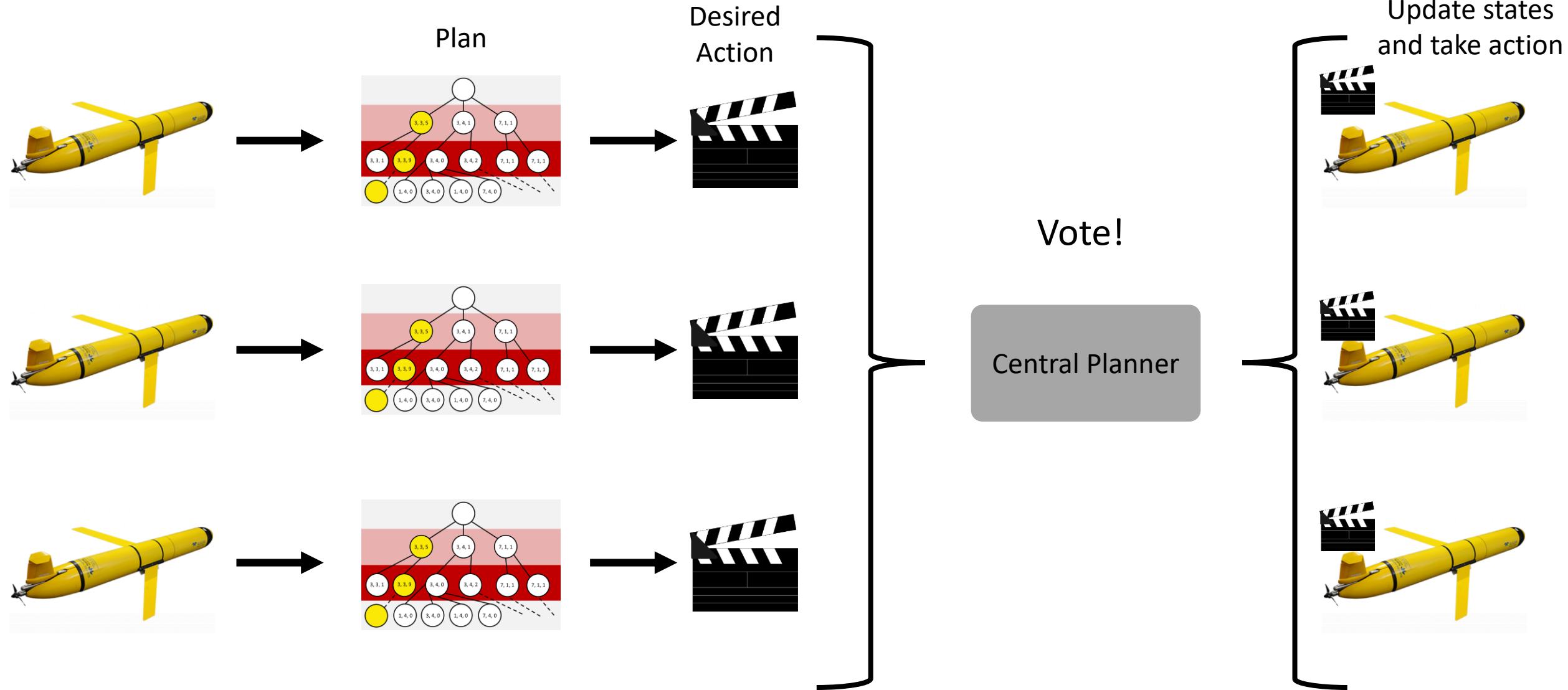


- Each tree votes on an action sequence
- Action sequence with largest number of votes is taken
- Soejima et al. showed evidence that a majority voting scheme gives better performance than conventional approaches



They hypothesized that this is because root parallelization, which utilizes unique random seeds for each tree, escapes local optima better.

Root parallelization for our agents



Browne, C. et al. A survey of Monte-Carlo tree search methods. *IEEE Trans. Comput. Intell. AI in Games* 4, 1–43 (2012)

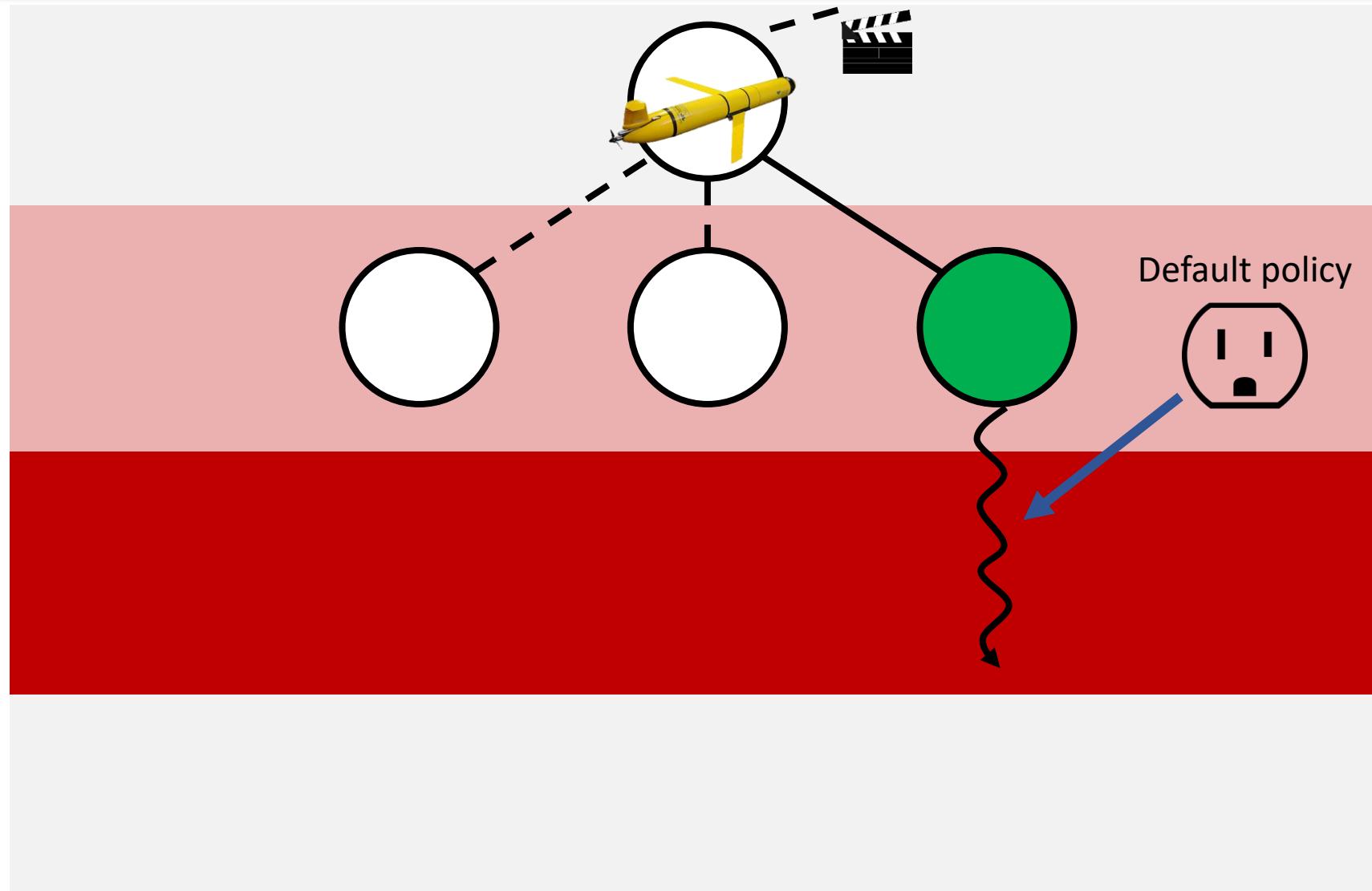
1. We are still uses random playouts for a default policy. We can do better.
2. We have seemingly reintroduced the need for robust communication between the agents.



To improve the default policy, we can pass action sequences from the other agents to each agent

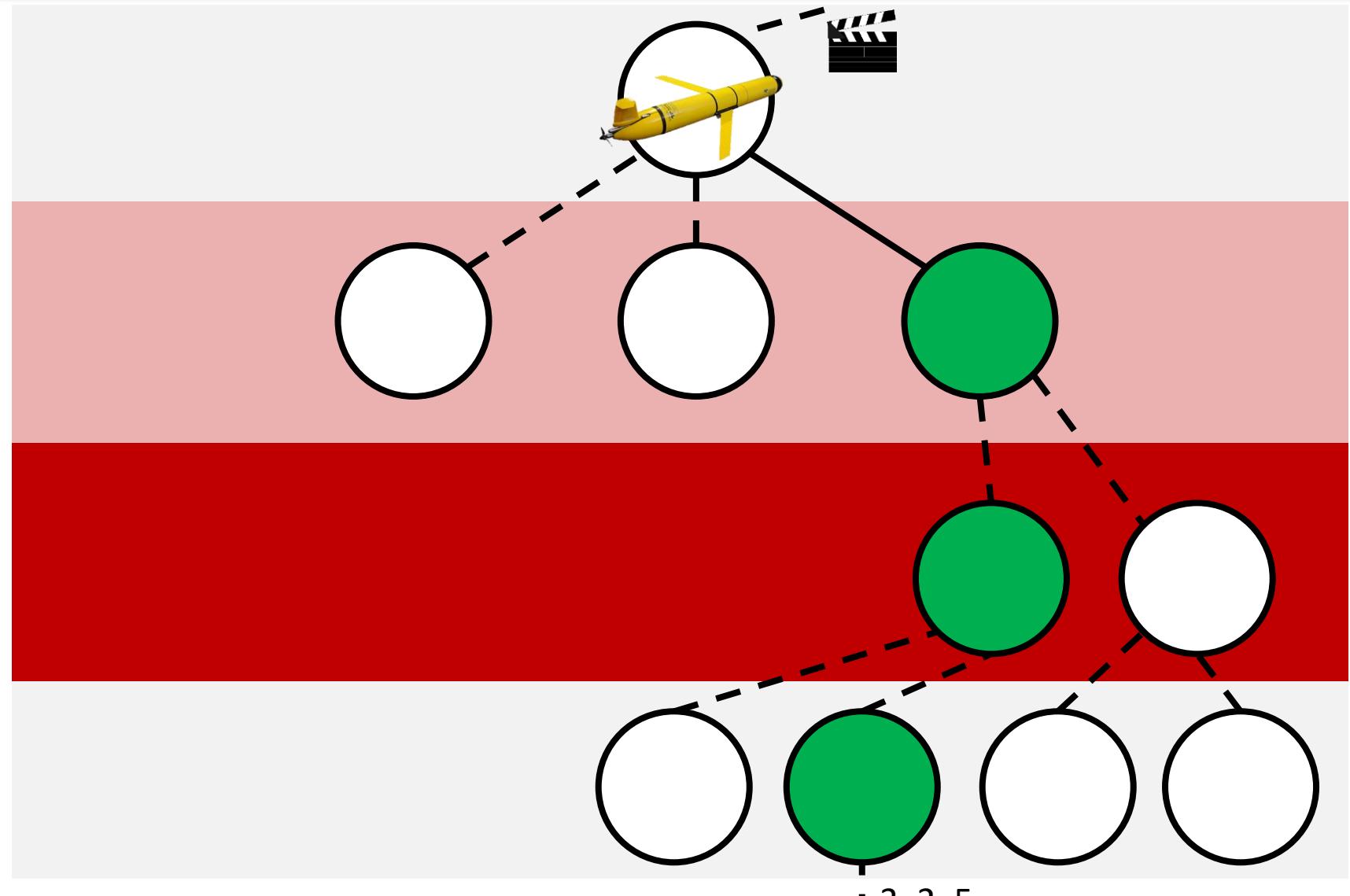
Passing action sequences detailed look

- Our original default policy was to simulate with random moves



Passing action sequences: Ms Pacman ghosts

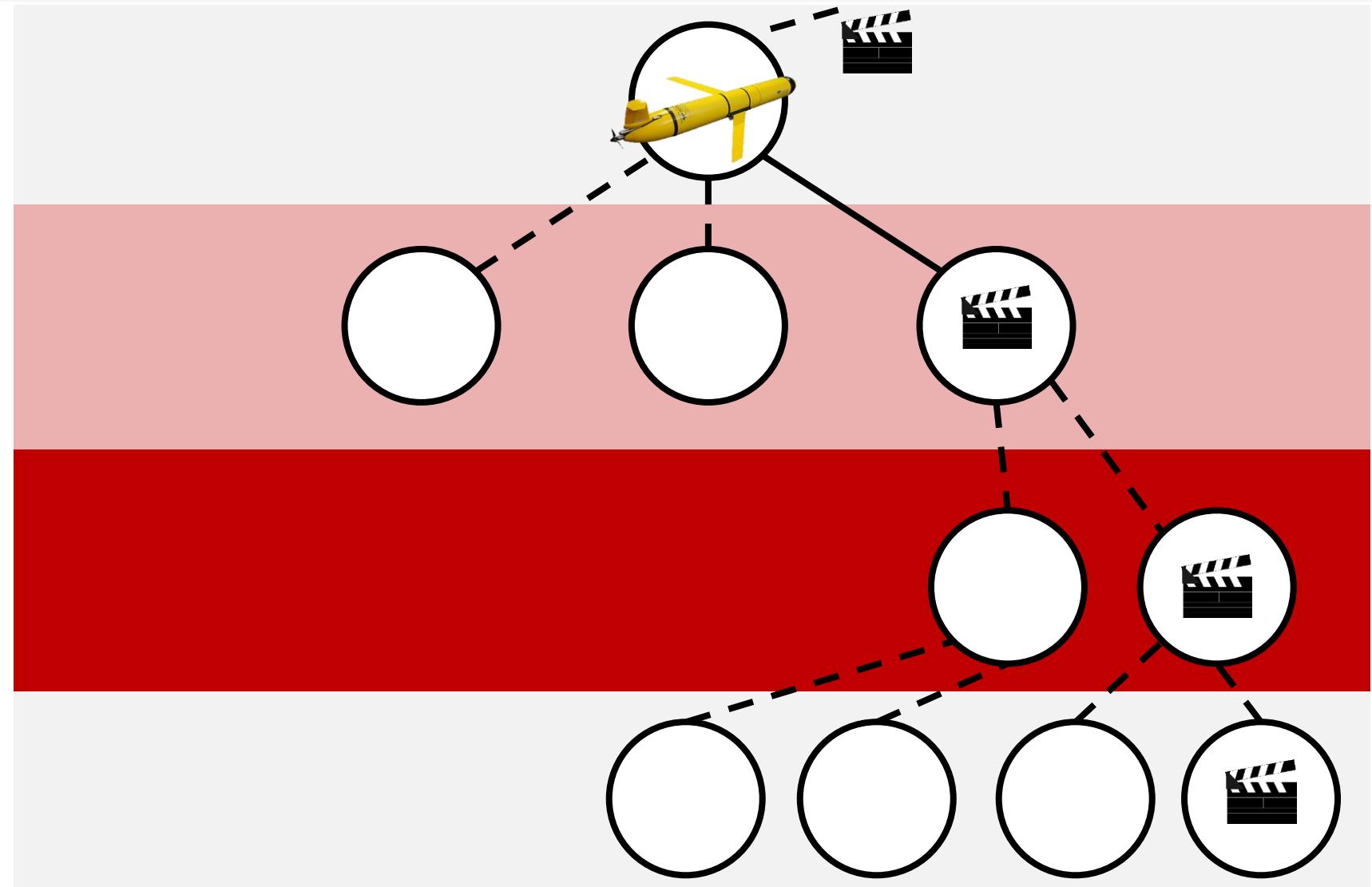
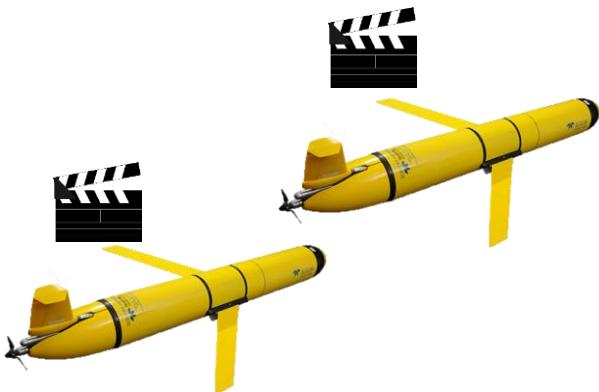
- Our original default policy was to simulate with random moves



| 3, 2, 5

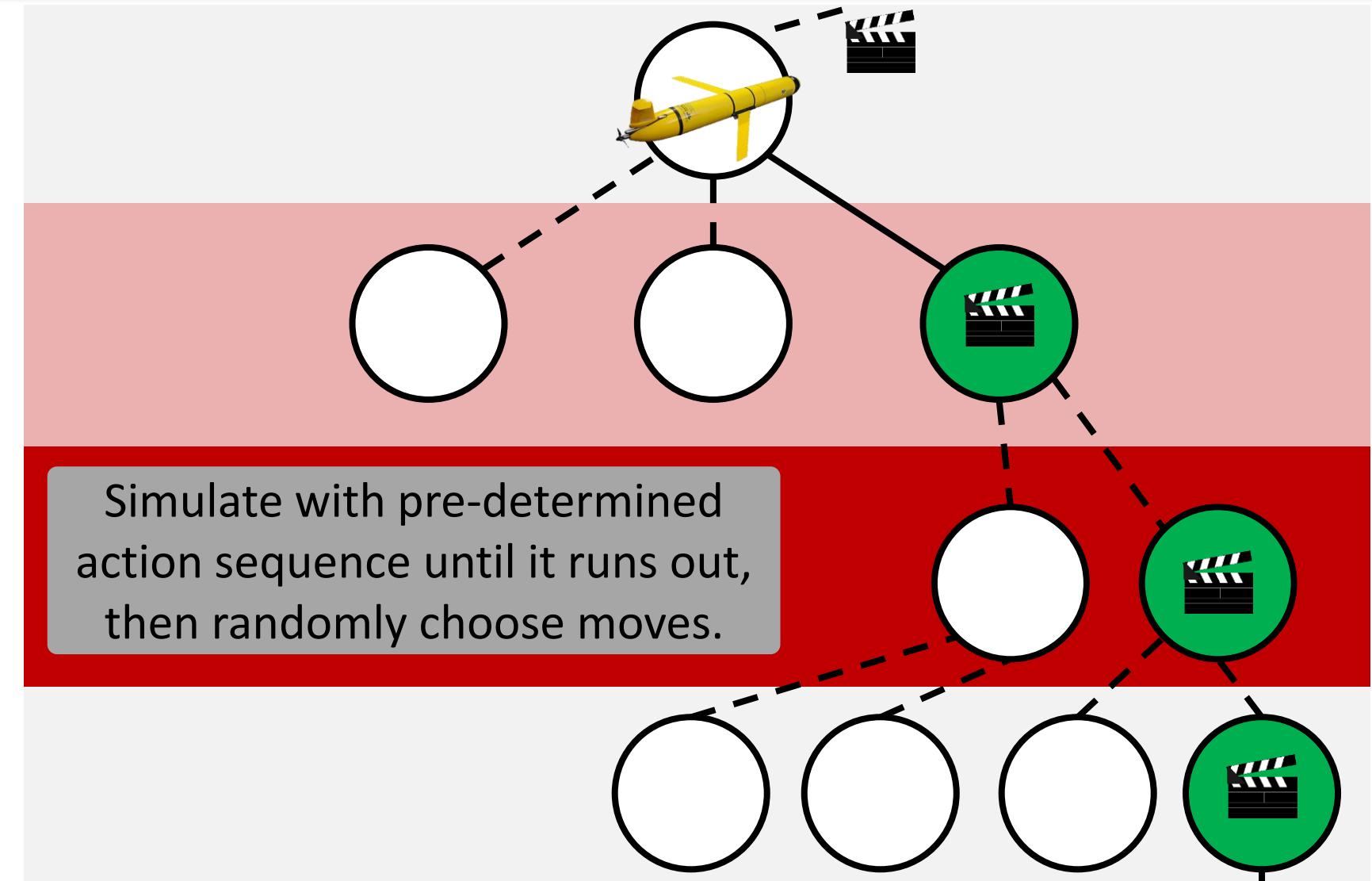
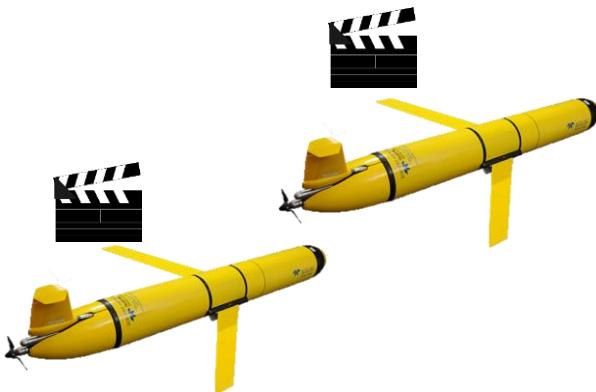
Passing action sequences: Ms Pacman ghosts

- Our original default policy was to simulate with random moves
- Instead, if our agent had a good guess on what all the other agents were going to do, we could improve our simulation



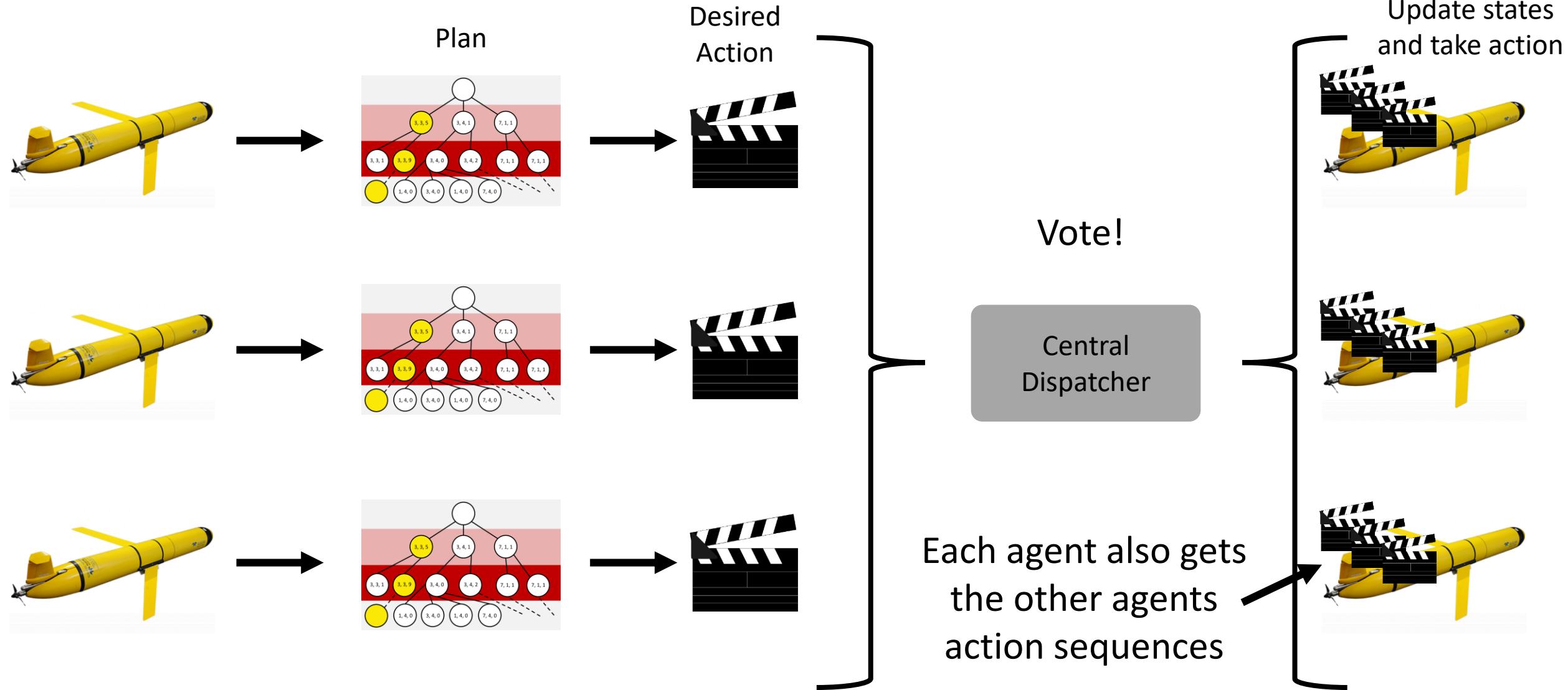
Passing action sequences: Ms Pacman ghosts

- Our original default policy was to simulate with random moves
- Instead, if our agent had a good guess on what all the other agents were going to do, we could improve our simulation



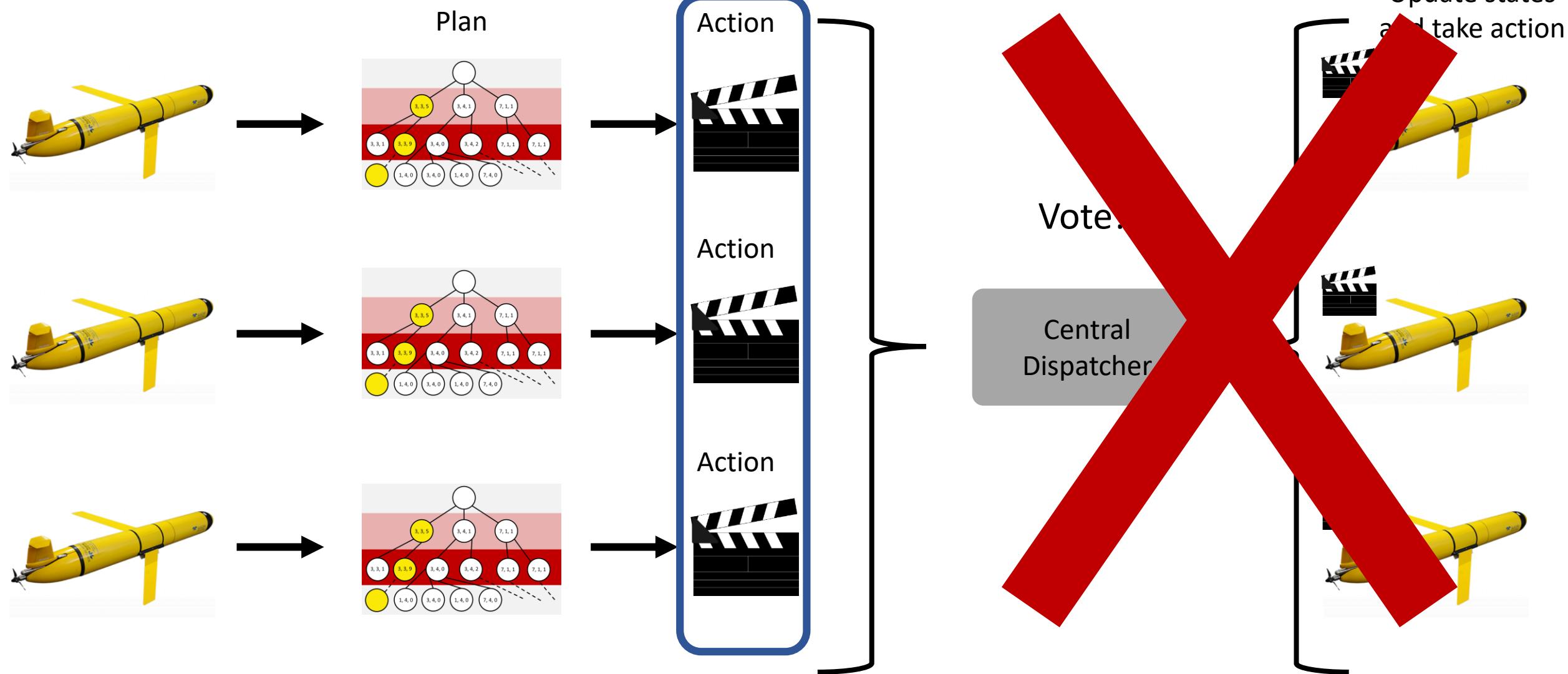
| 4, 6, 5

Root parallelization for our agents



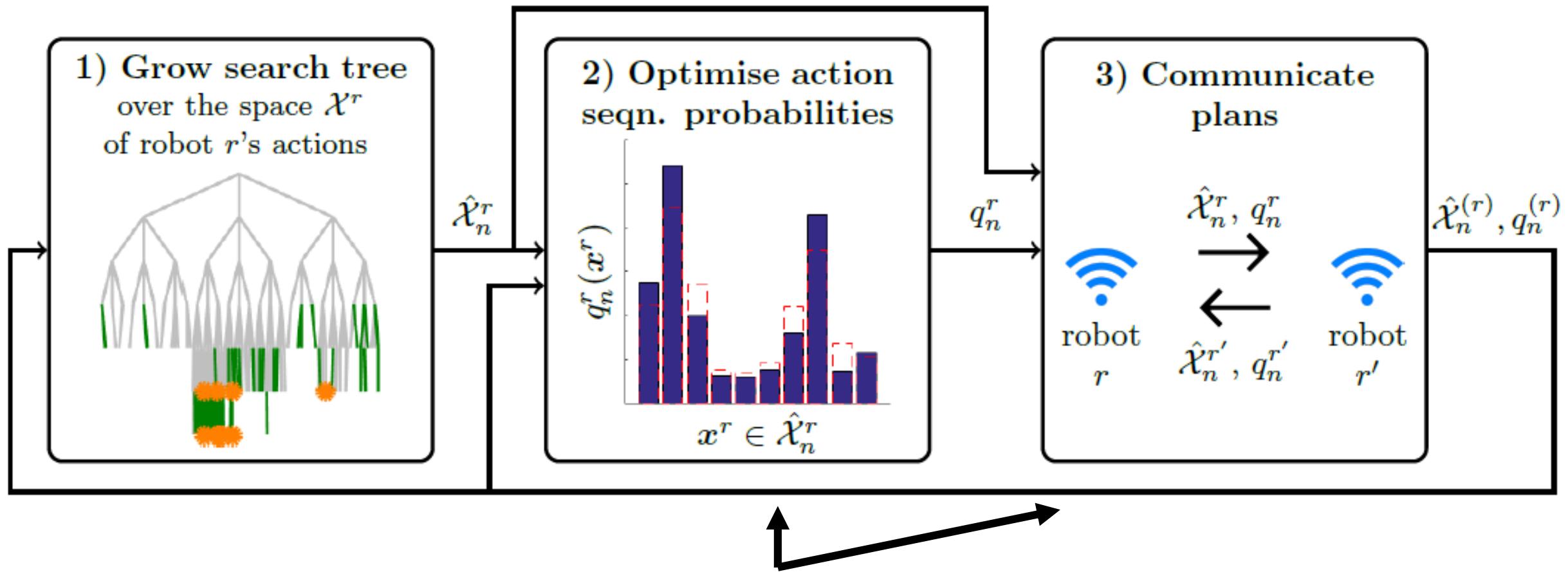
Browne, C. et al. A survey of Monte-Carlo tree search methods. *IEEE Trans. Comput. Intell. AI in Games* 4, 1–43 (2012)

What if/when communication fails?

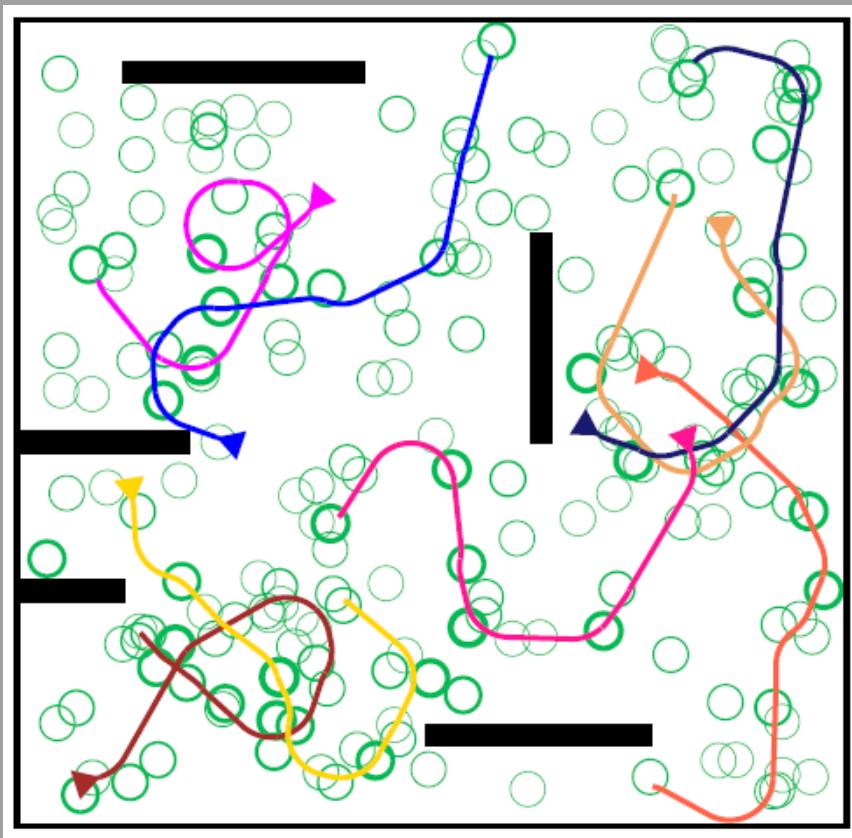


Browne, C. et al. A survey of Monte-Carlo tree search methods. *IEEE Trans. Comput. Intell. AI in Games* 4, 1–43 (2012)

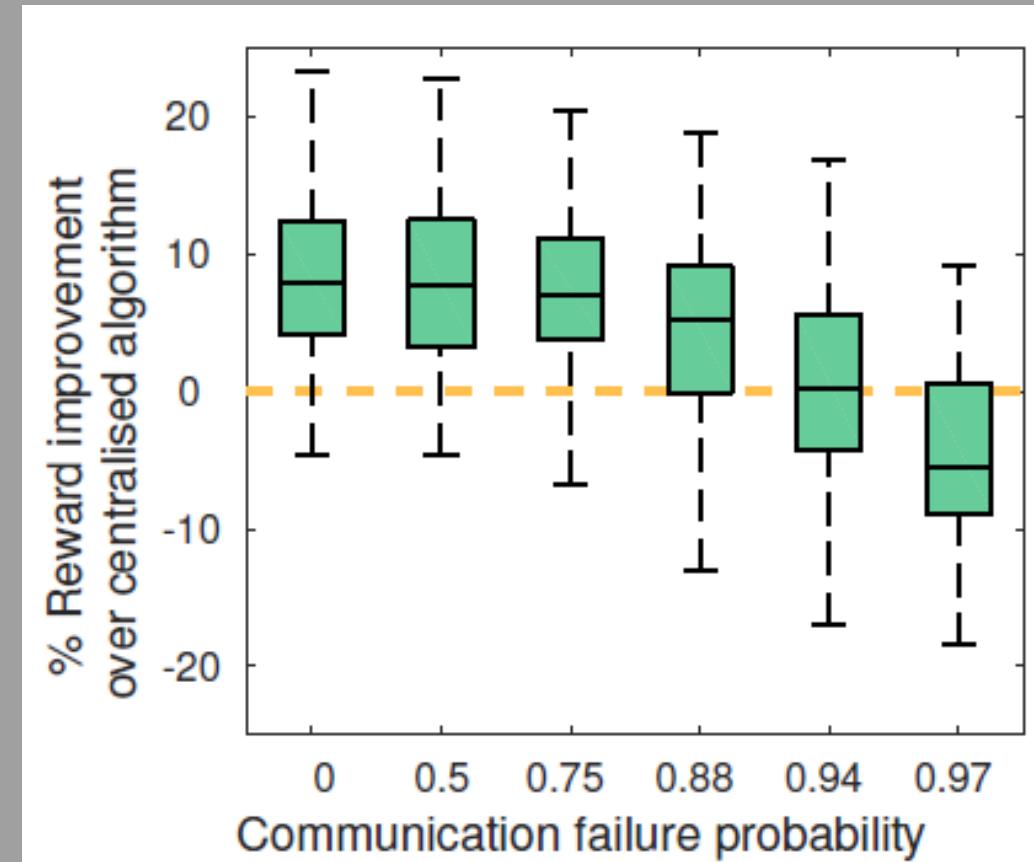
Best et al. (2018) investigated this with a spoofed up version of MCTS

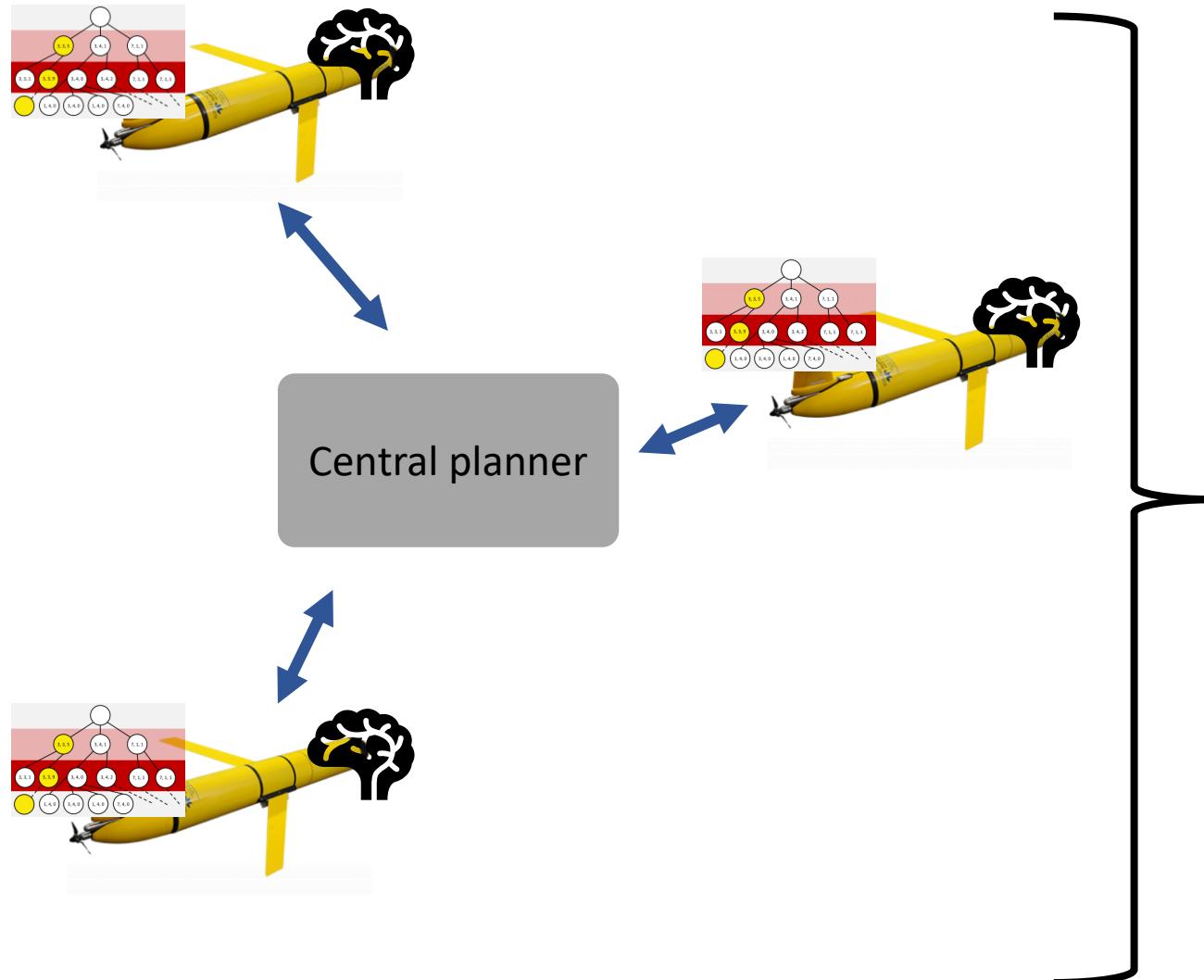


Problem environment



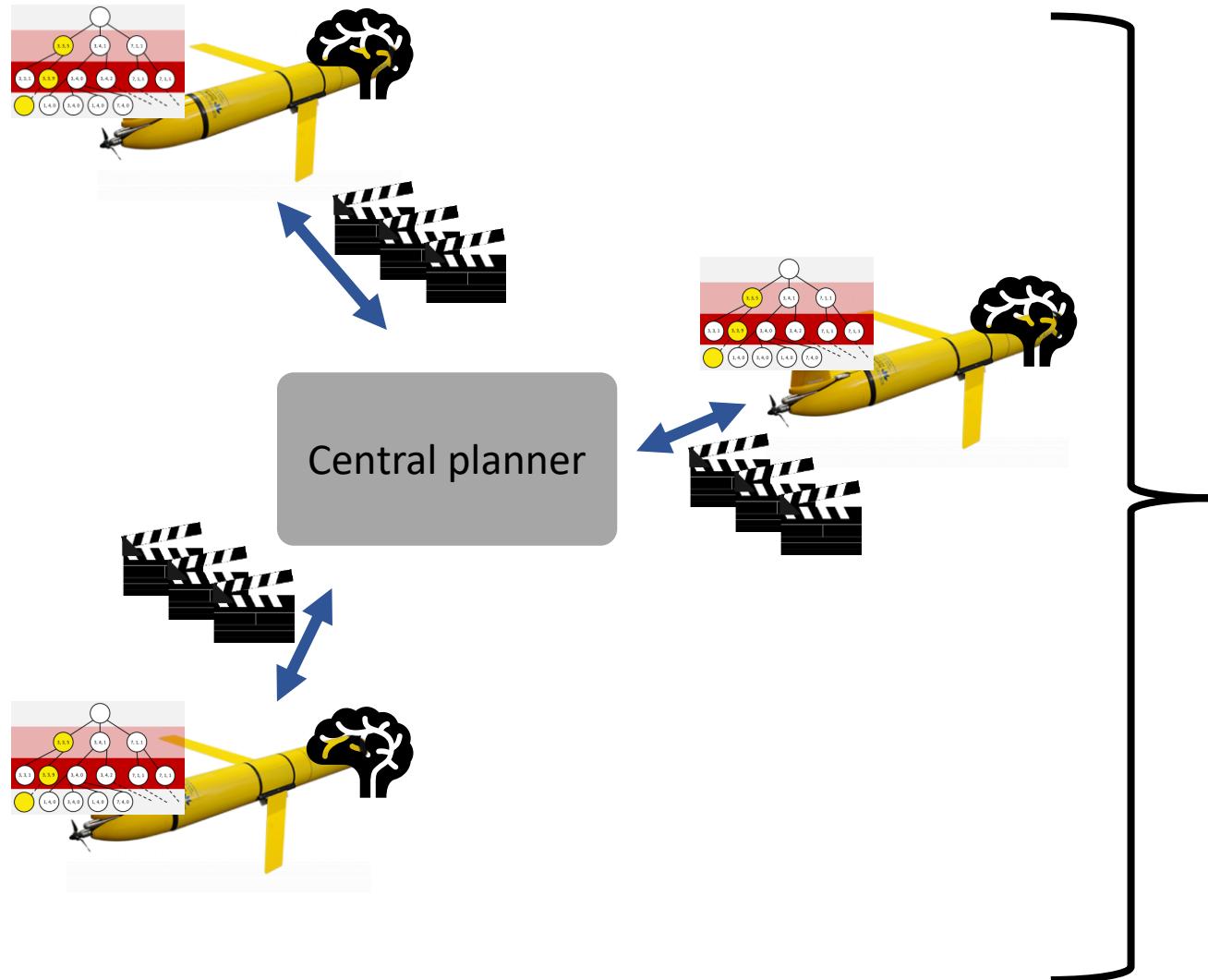
Communication Failure vs Reward Improvement





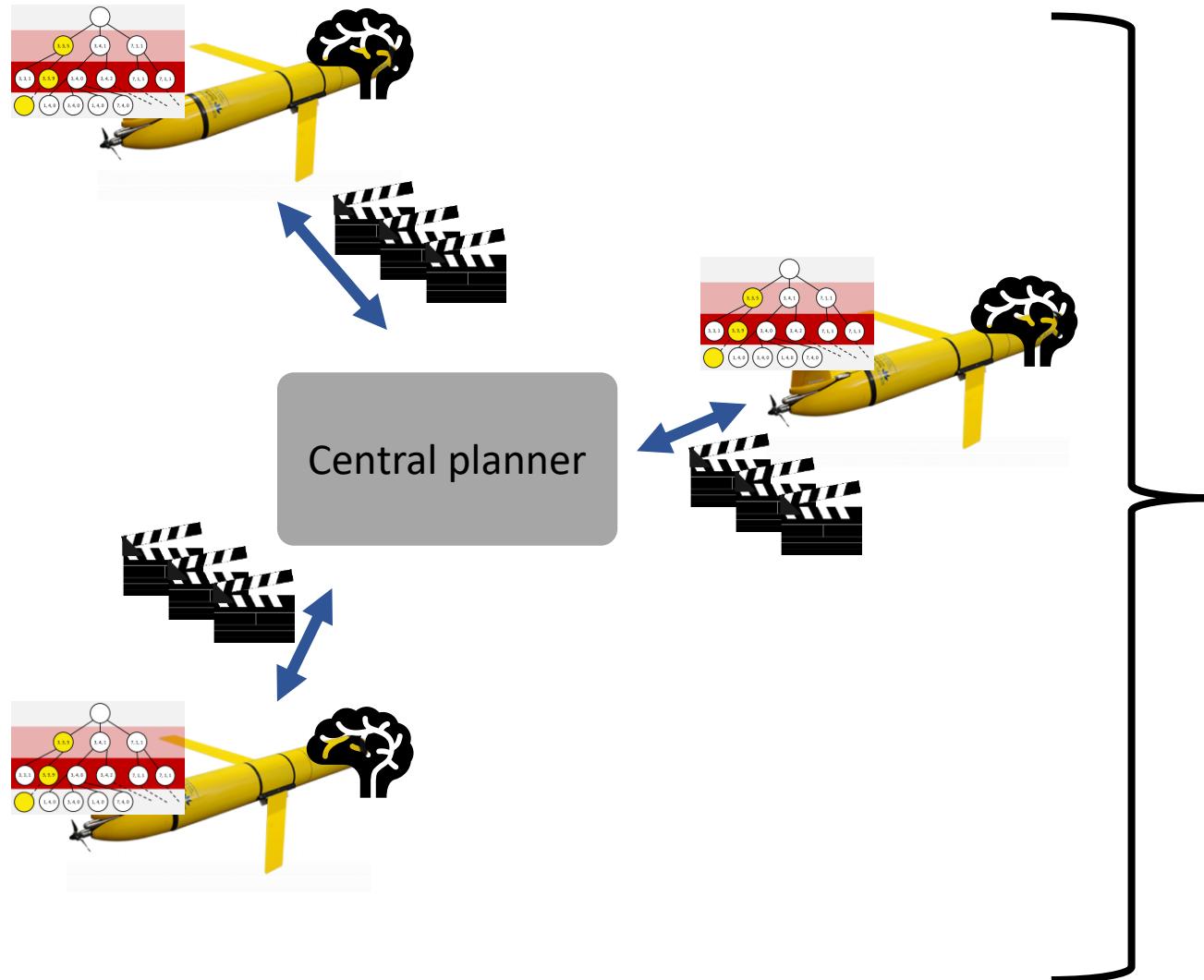
- Each agent runs own MCTS, incorporating heuristics and other performance boosting modifications

MCTS for Collaboration

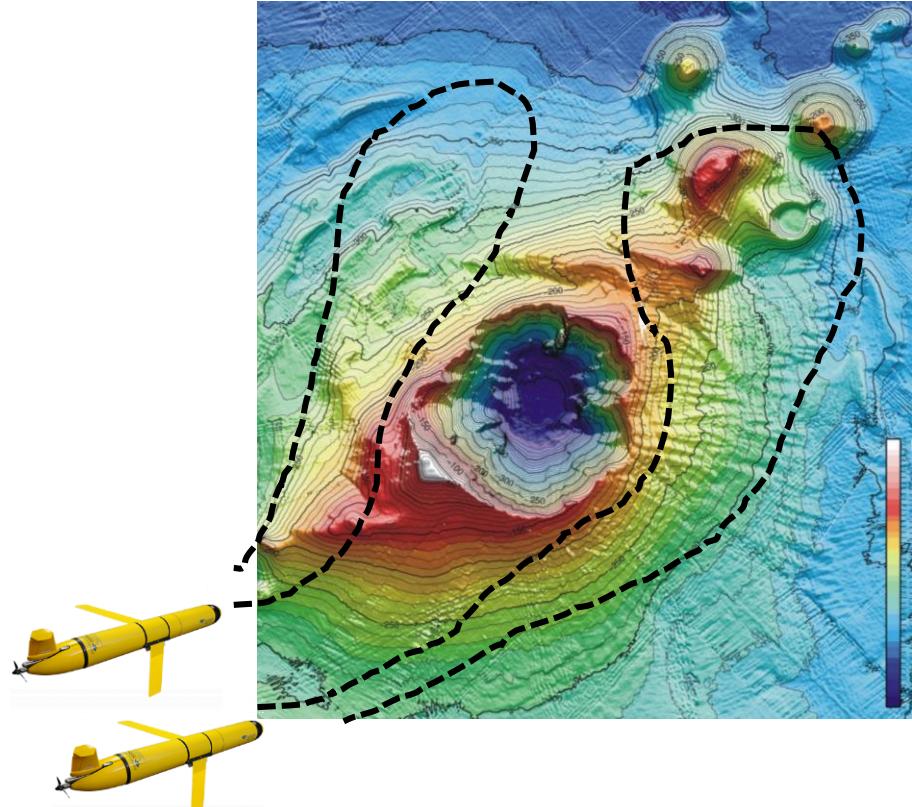


- Each agent runs own MCTS, incorporating heuristics and other performance boosting modifications
- When communication allows, action sequences are passed between the gliders and the central planner

MCTS for Collaboration



- Each agent runs own MCTS, incorporating heuristics and other performance boosting modifications
- When communication allows, action sequences are passed between the gliders and the central planner
- These action sequences are then incorporated into each local agents MCTS



Questions?