

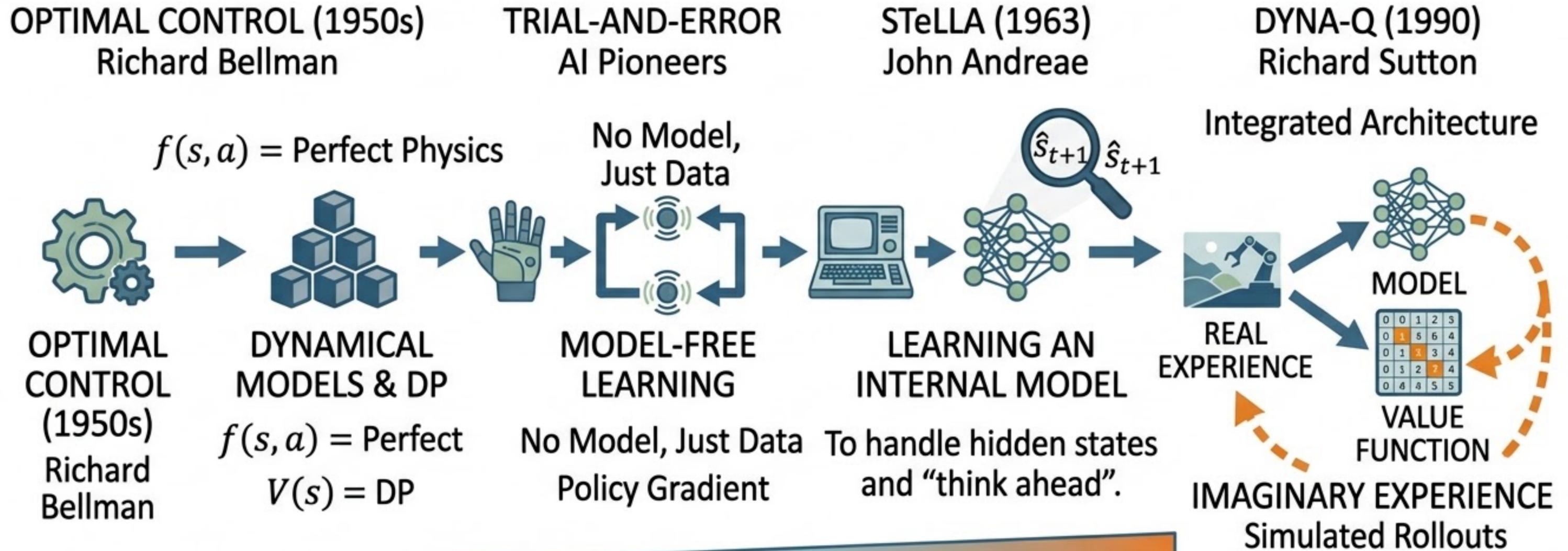
Deep Learning

Randall Balestriero

Lecture 25: Model-Based RL

Recap!

Origins of Model-Based RL: The Key Lineages



Key Driver: SAMPLE EFFICIENCY. Use models to “practice” and learn from imaginary data when real data is expensive or dangerous.

Notations

Notations

- On-policy and off-policy

Notations

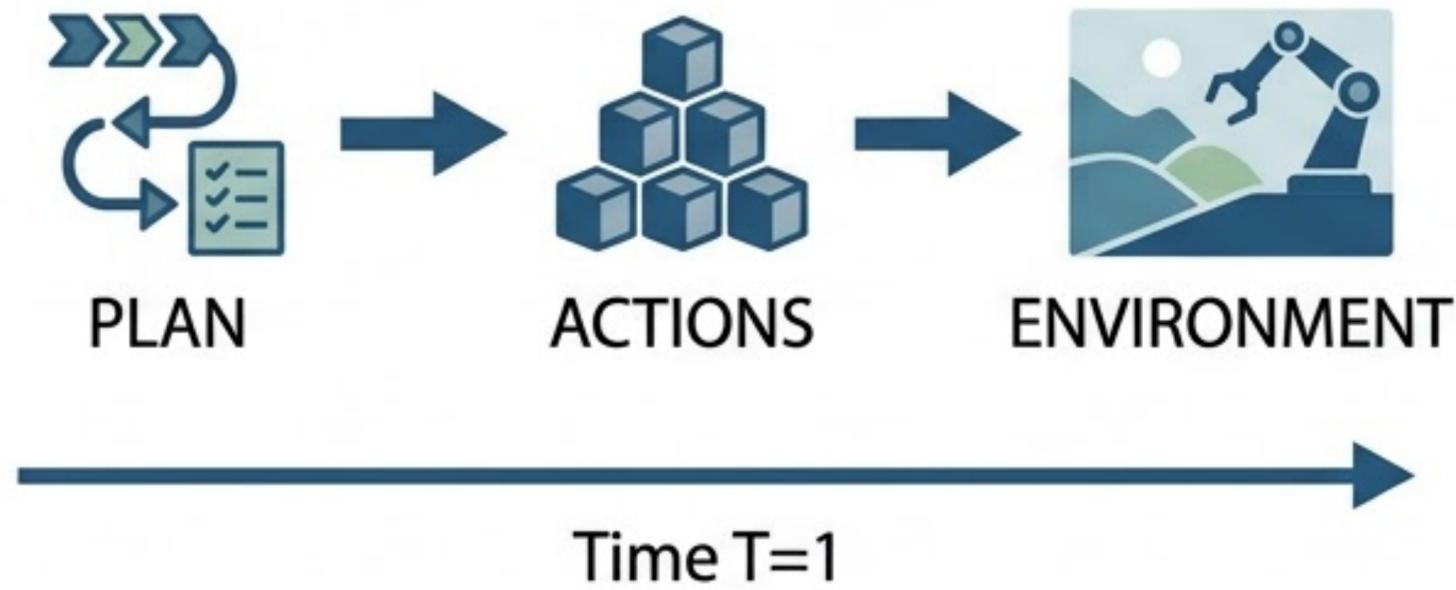
- On-policy and off-policy
- Model-free and model-based

Notations

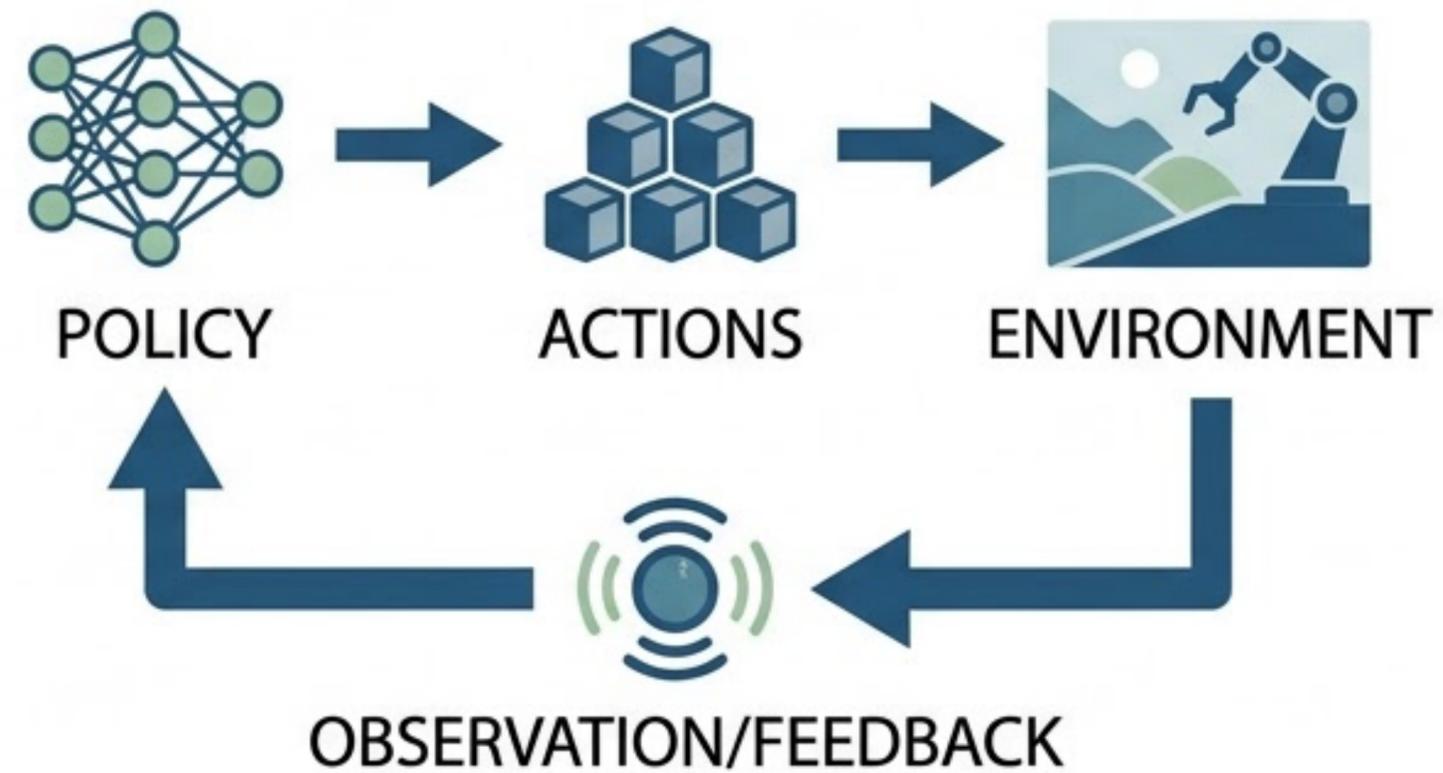
- On-policy and off-policy
- Model-free and model-based
- Open-loop and closed-loop

Open-Loop vs. Closed-Loop Control

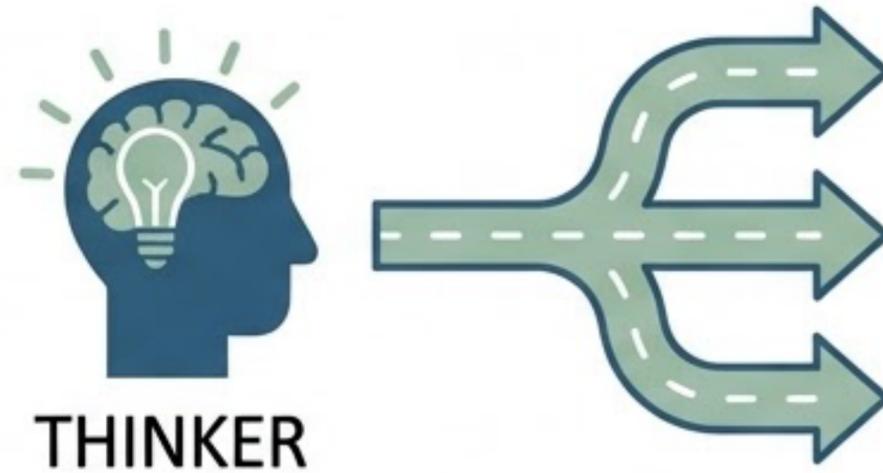
Open-Loop Control



Closed-Loop Control



Why Use Open-Loop Planning?

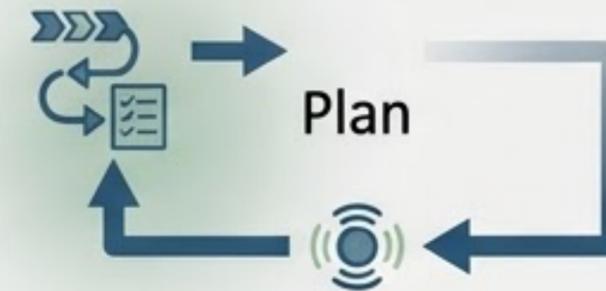


Computational Efficiency



Faster and cheaper to optimize a single trajectory of actions than a complex policy.

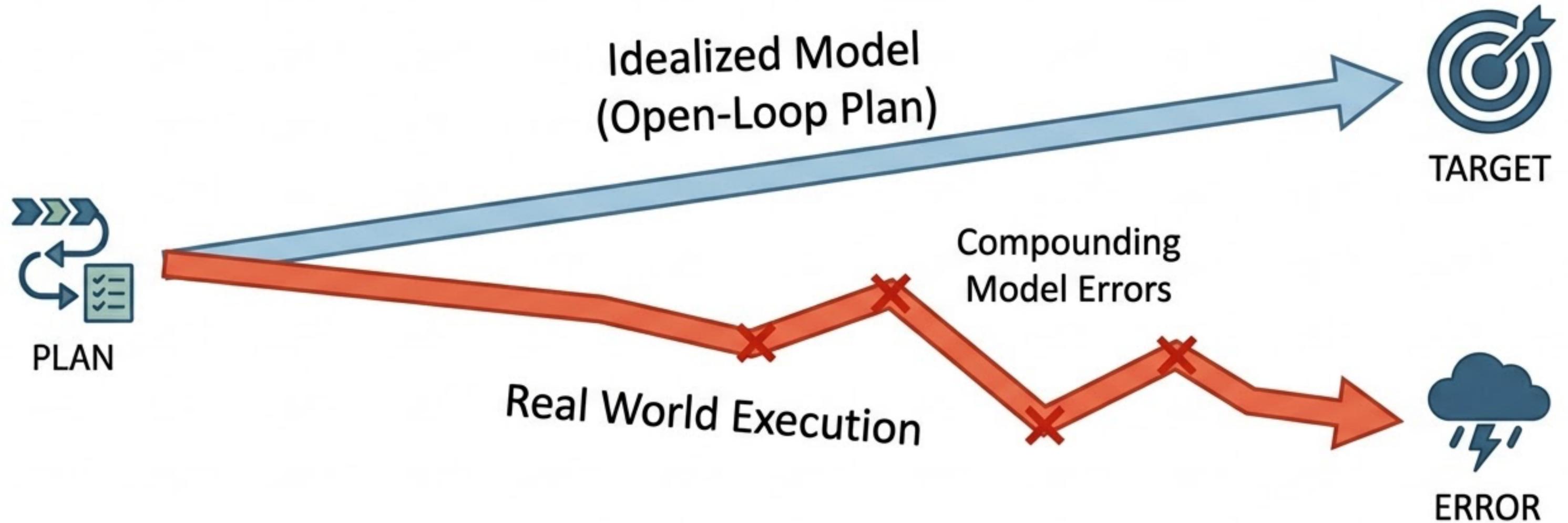
Foundation for Closed-Loop



Essential component for Model Predictive Control (MPC): We plan open-loop but execute closed-loop.

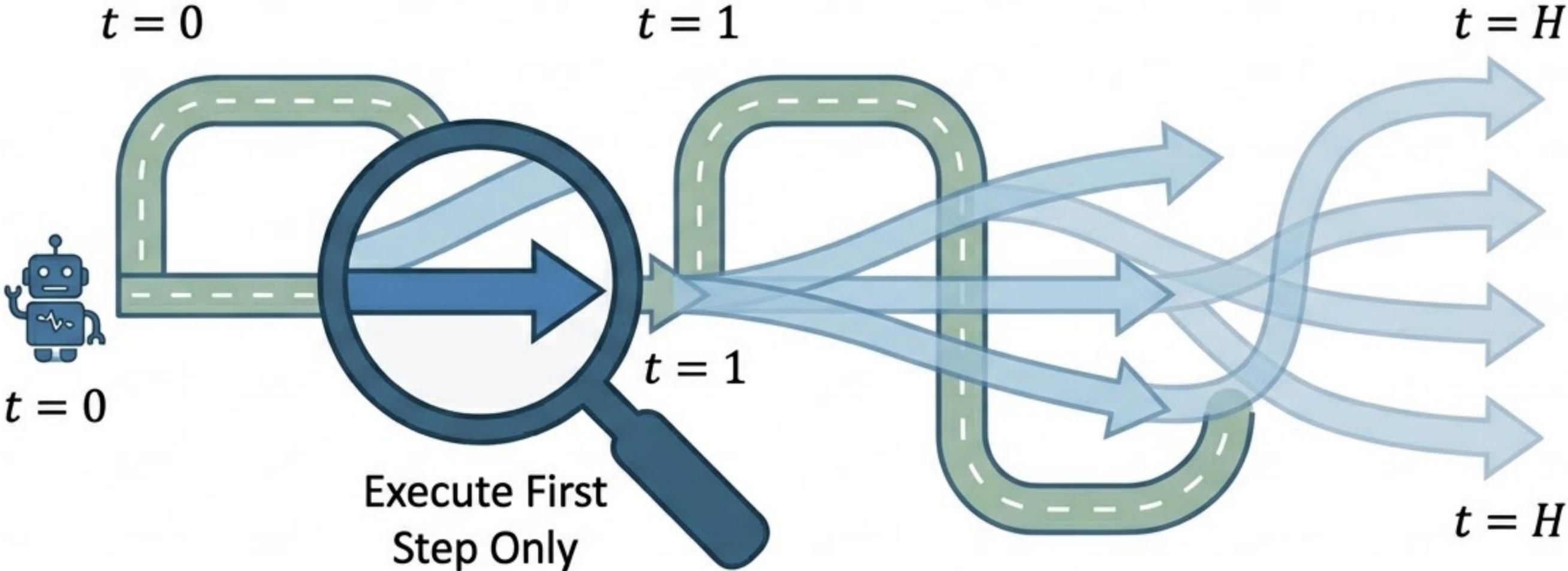
But still... why not open-loop?

The Problem: Model Drift & The Reality Gap



Small prediction errors at the start lead to massive failures over long horizons without feedback.

Model Predictive Control (MPC): The Receding Horizon

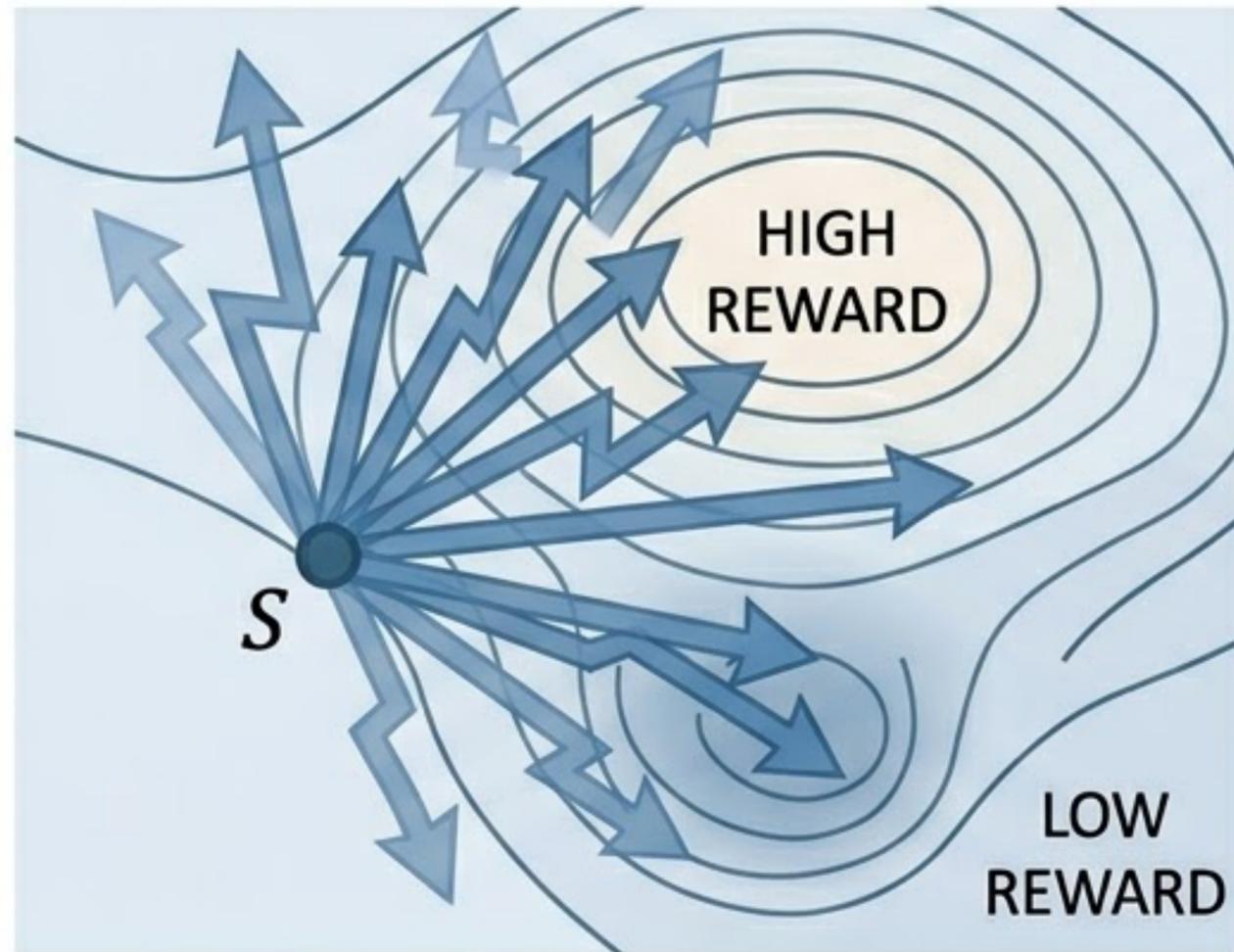


[Plan $t = 0 \dots H$] \rightarrow [Execute a_0] \rightarrow [Observe s_1] \rightarrow [Re-Plan $t = 1 \dots H + 1$]

**How to generate candidate
actions?**

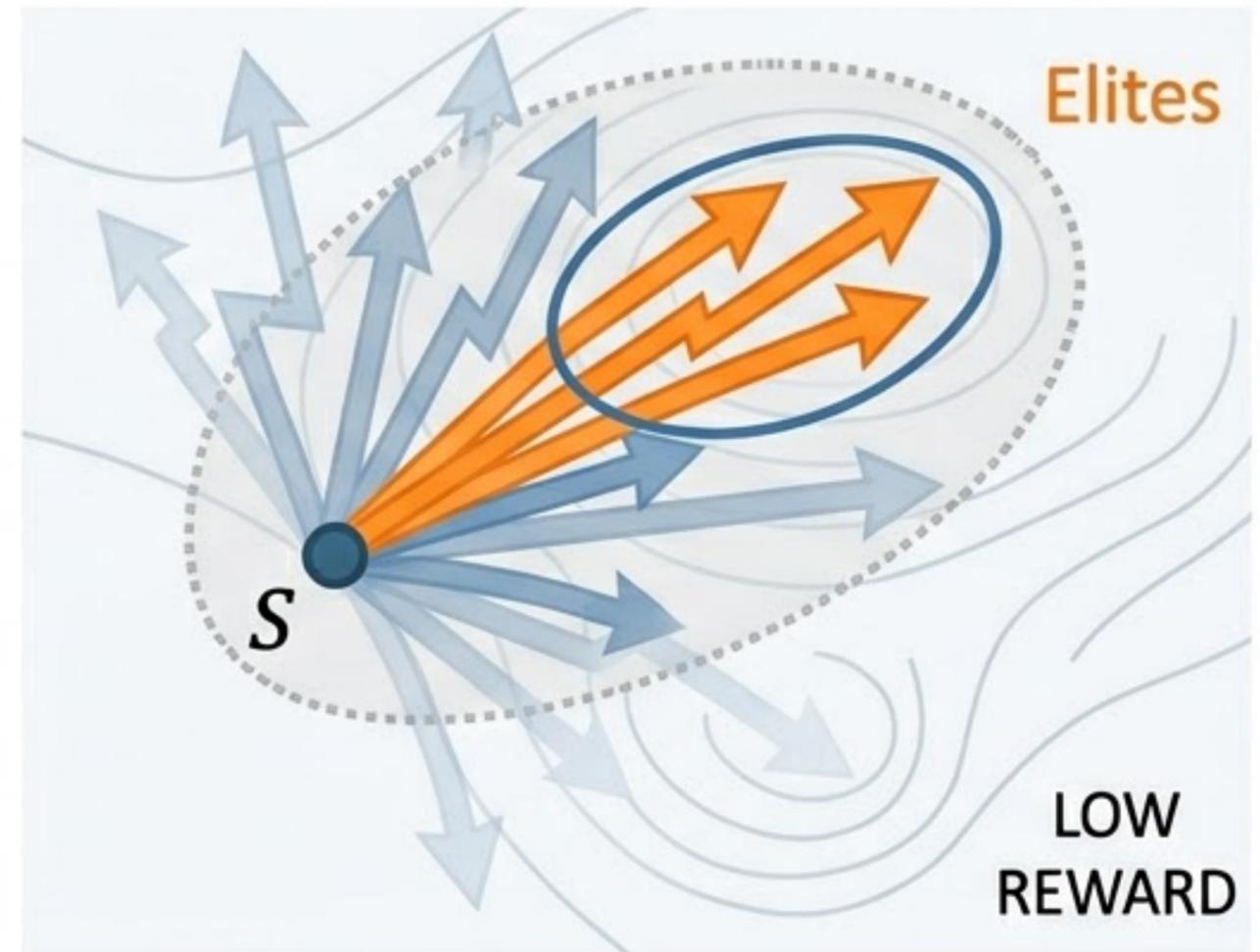
Optimization for Planning: Shooting Methods

Random Shooting



Sample N sequences, pick the best.

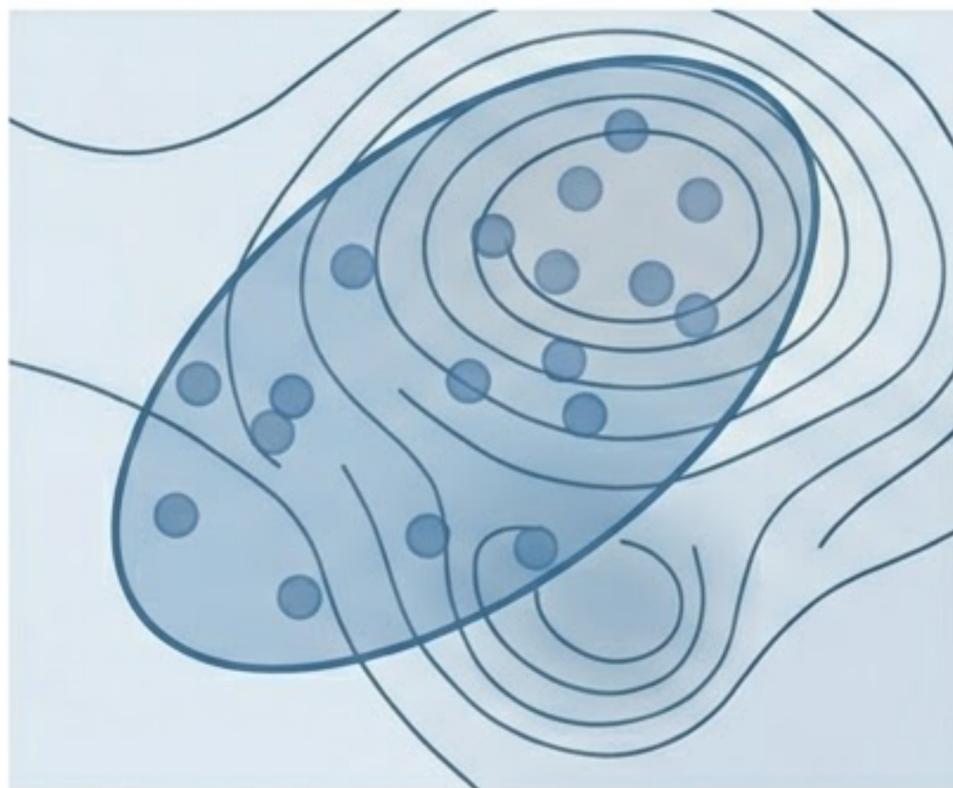
Cross-Entropy Method (CEM)



Iteratively refit distribution to elite samples.

The Cross-Entropy Method (CEM) Algorithm

1. Sample and Evaluate



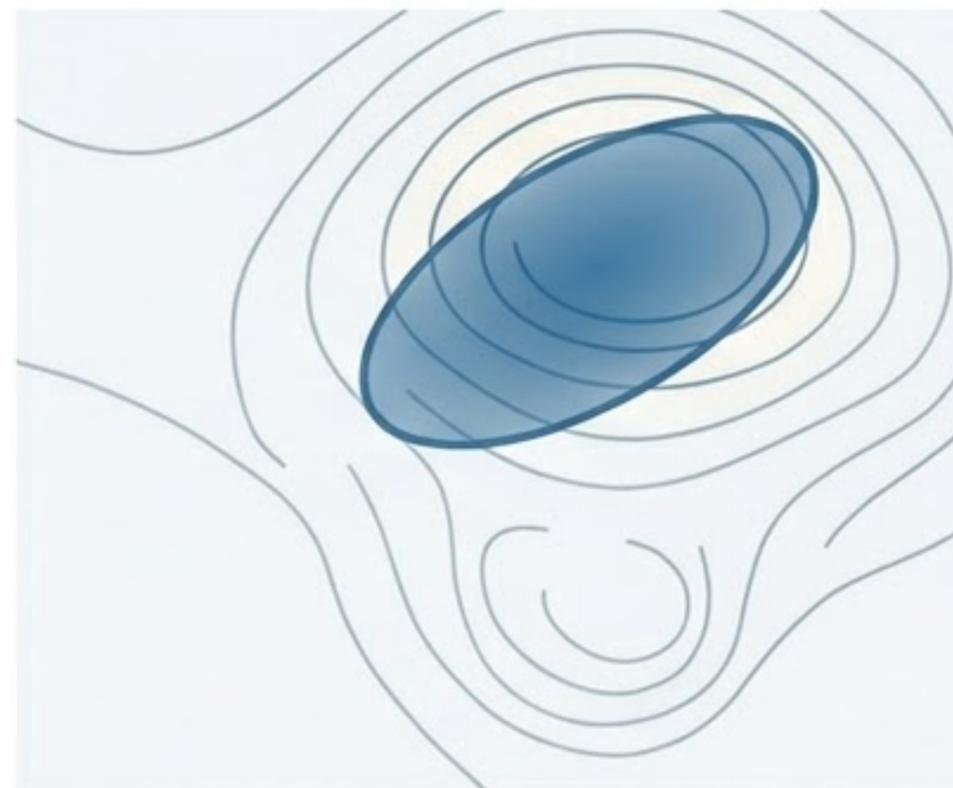
Sample N trajectories from initial distribution $P(A)$.

2. Select Elites



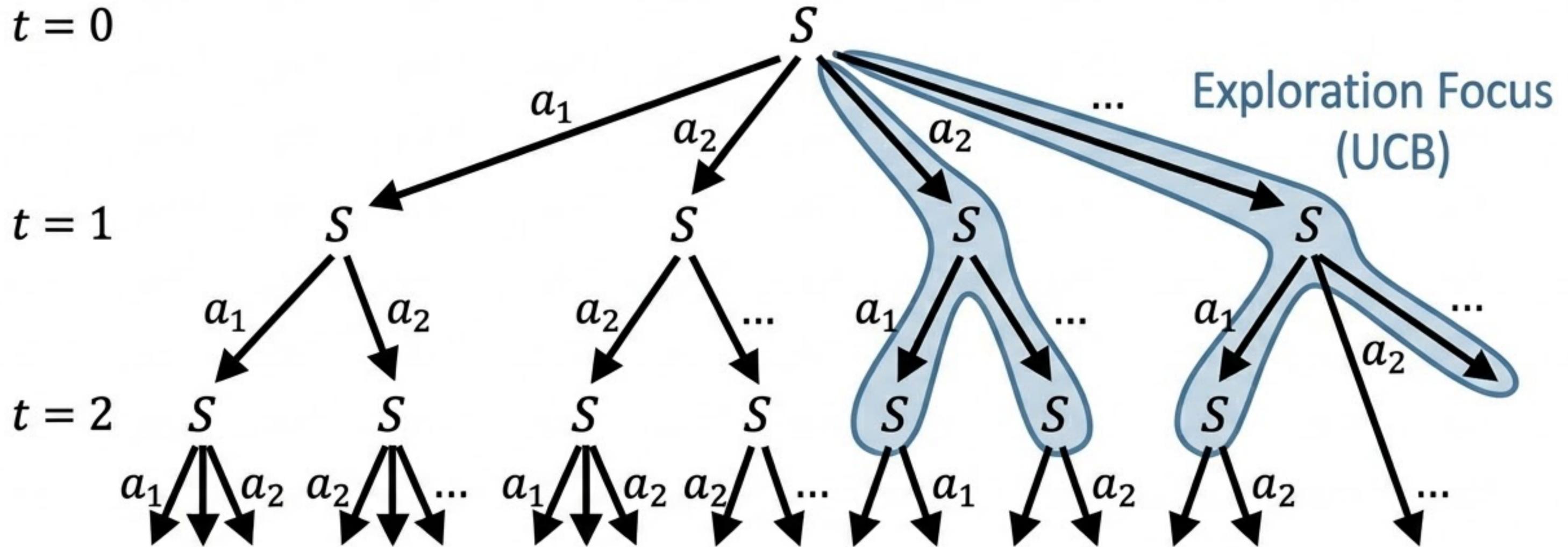
Keep top k (e.g., 10%) highest-reward samples.

3. Refit and Resample



Re-calculate parameters of $P(A)$ (e.g., mean and variance) to fit the Elites. Repeat.

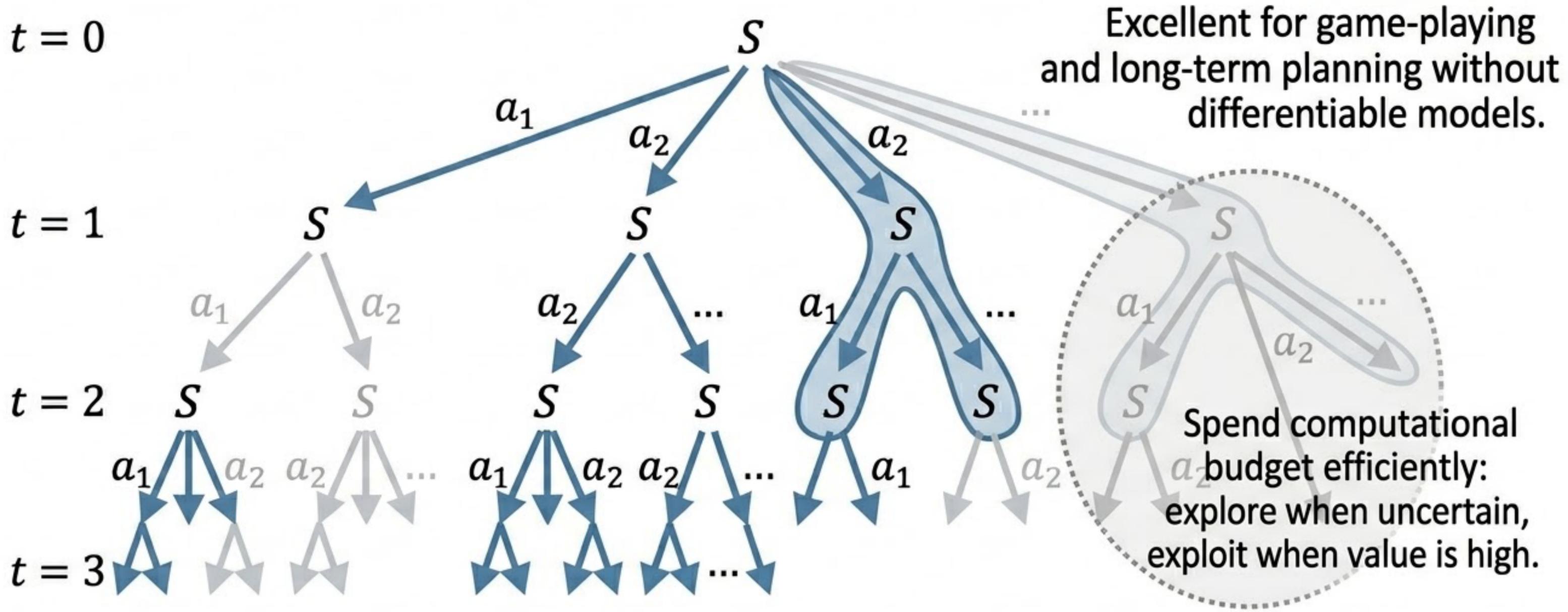
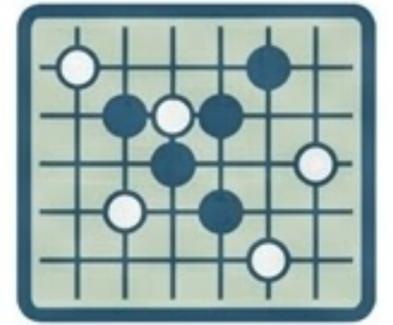
Monte Carlo Tree Search (MCTS): Planning with Discrete Actions



Best for discrete actions and long horizons.

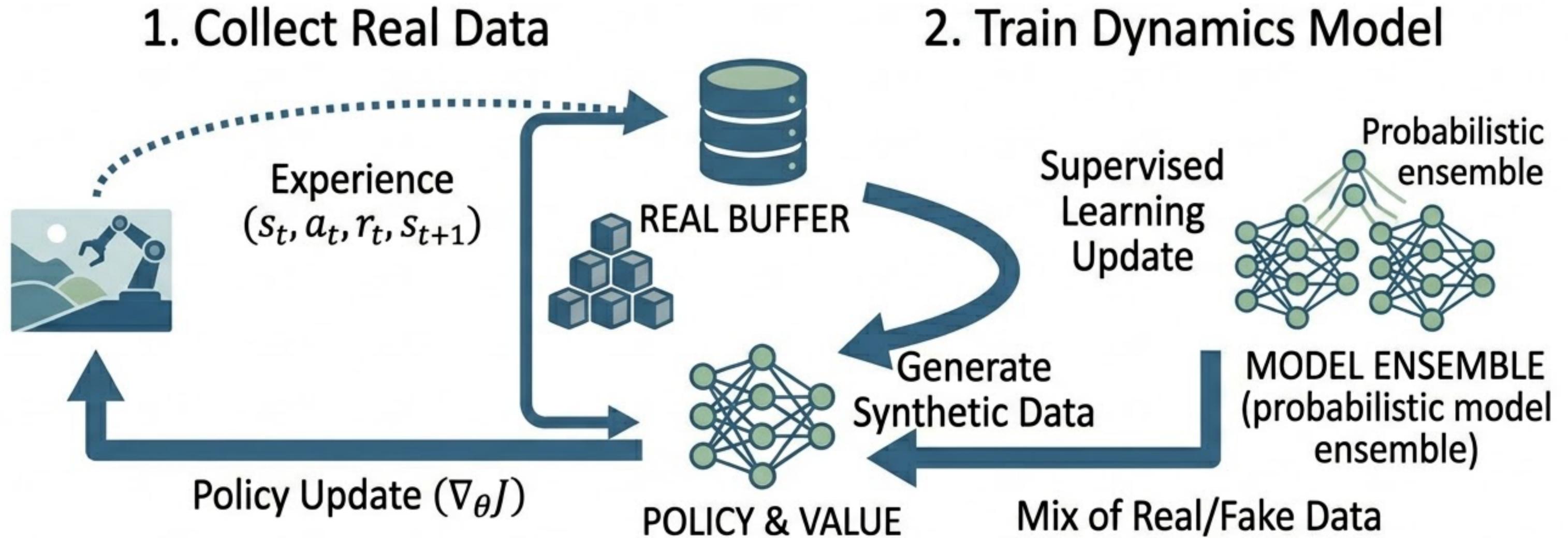
Dynamically builds an asymmetric search tree using a simulation model.

Why Monte Carlo Tree Search (MCTS)?



Asymmetric growth: focus search on promising paths.

The MBPO Algorithm Loop



Hybrid approach: Learn a model, hallucinate short rollouts from real states, and train policy on augmented data.

**Fundamental question:
What is s_t ?**



See you on Monday!