

# Lecture 8: Integrating Learning and Planning

David Silver

# Outline

## 1 Introduction

→ build model for env.

## 2 Model-Based Reinforcement Learning

→ combine all the ideas.

## 3 Integrated Architectures

→ imagine & learn by imagination.

## 4 Simulation-Based Search

# Outline

1 Introduction

2 Model-Based Reinforcement Learning

3 Integrated Architectures

4 Simulation-Based Search

# Model-Based Reinforcement Learning

- *Last lecture*: learn **policy** directly from experience
- *Previous lectures*: learn **value function** directly from experience
- *This lecture*: learn **model** directly from experience  
and use **planning** → *look ahead / imagine* to construct a value function or policy
- Integrate learning and planning into a single architecture

# Model-Based and Model-Free RL

## ■ Model-Free RL

- No model
- Learn value function (and/or policy) from experience

no model ↴ not represent transition dynamics or reward function that reward env operates on.

↳ weights

# Model-Based and Model-Free RL

- Model-Free RL

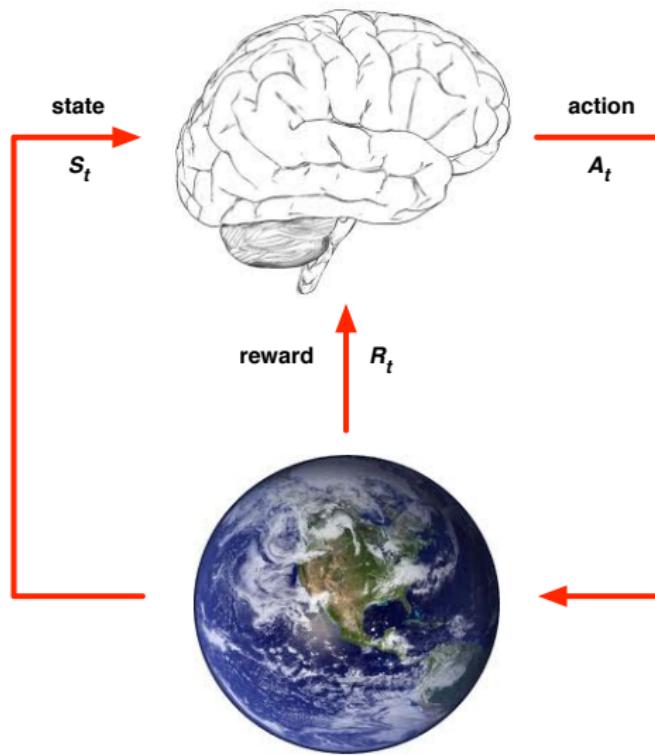
- No model
  - Learn value function (and/or policy) from experience

- Model-Based RL

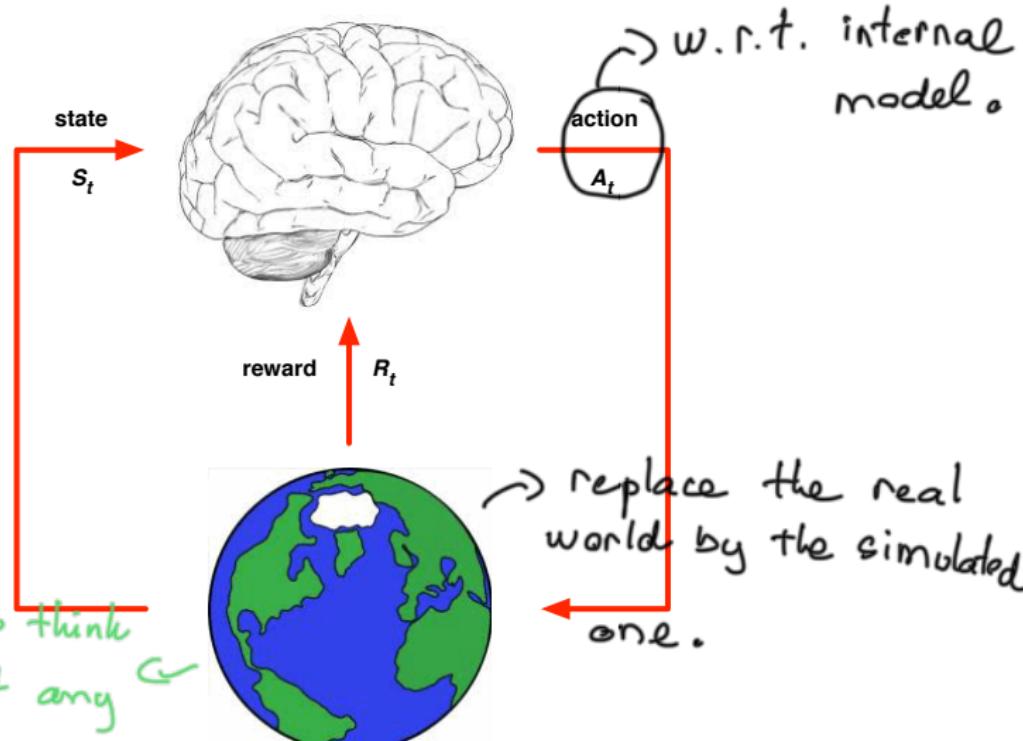
- Learn a model from experience
  - Plan value function (and/or policy) from model

construct + transition dynamic reward  
look ahead to select which right  
fune or right action are to select.

# Model-Free RL



# Model-Based RL



# Outline

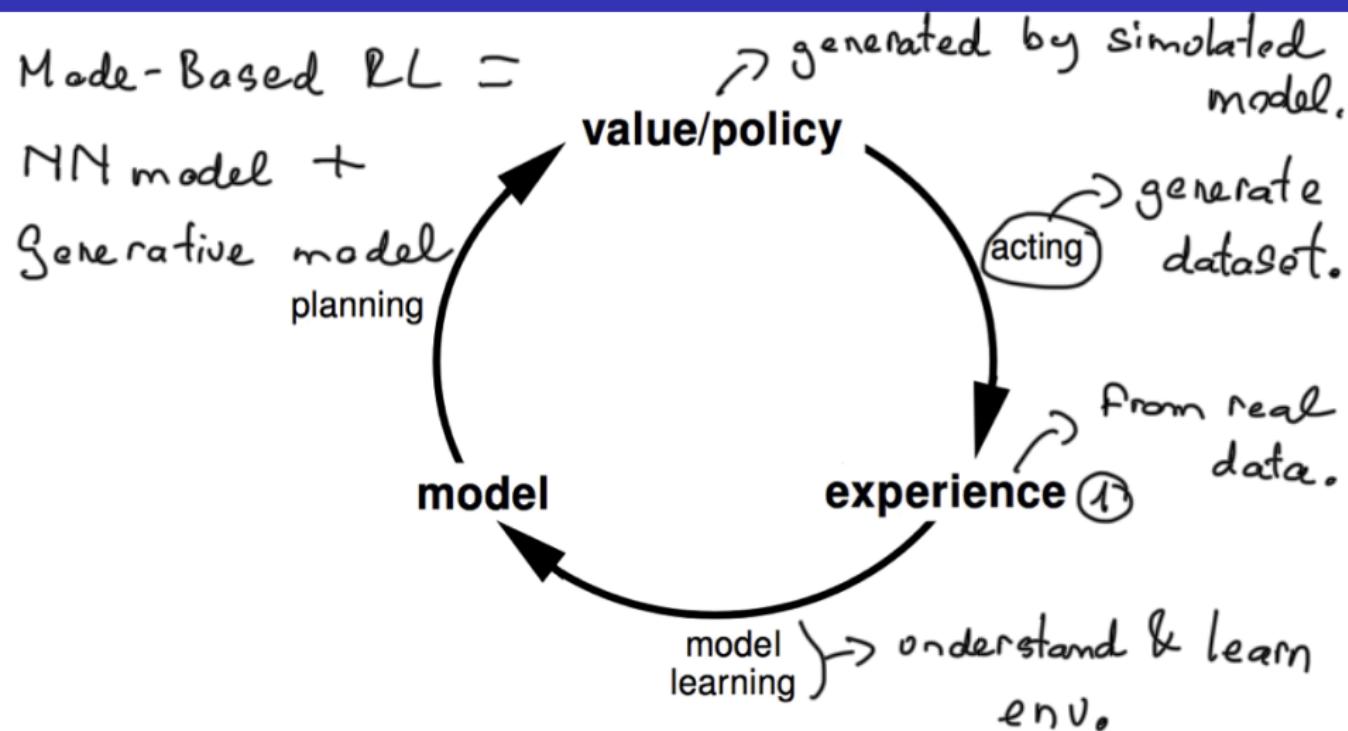
1 Introduction

2 Model-Based Reinforcement Learning

3 Integrated Architectures

4 Simulation-Based Search

# Model-Based RL



## Advantages of Model-Based RL

illustrate some MDP  
's really hard to optimize  
without planning.

↳ chess ↳ how to give reward for  
one movement?  
Am I likely to win from this position?

construct value func

Advantages:

- Can efficiently learn model by supervised learning methods
- Can reason about model uncertainty

↳ data + label

Disadvantages:

- First learn a model, then construct a value function  
⇒ two sources of approximation error

# What is a Model?

- A *model*  $\mathcal{M}$  is a representation of an MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R} \rangle$ , parametrized by  $\eta$
- We will assume state space  $\mathcal{S}$  and action space  $\mathcal{A}$  are known
- So a model  $\mathcal{M} = \langle \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$  represents state transitions  $\mathcal{P}_\eta \approx \mathcal{P}$  and rewards  $\mathcal{R}_\eta \approx \mathcal{R}$

*→ Transition*

$$S_{t+1} \sim \mathcal{P}_\eta(S_{t+1} | S_t, A_t)$$

$$R_{t+1} = \mathcal{R}_\eta(R_{t+1} | S_t, A_t) \quad \rightarrow \text{reward}$$

- Typically assume conditional independence between state transitions and rewards

$$\mathbb{P}[S_{t+1}, R_{t+1} | S_t, A_t] = \mathbb{P}[S_{t+1} | S_t, A_t] \mathbb{P}[R_{t+1} | S_t, A_t]$$

# Model Learning

- Goal: estimate model  $M_\eta$  from experience  $\{S_1, A_1, R_2, \dots, S_T\}$
- This is a supervised learning problem (transformed into)

state space &  
action space concerning

Input      Labels  
 $S_1, A_1 \rightarrow R_2, S_2$   
 $S_2, A_2 \rightarrow R_3, S_3$   
 $\vdots$   
 $S_{T-1}, A_{T-1} \rightarrow R_T, S_T$

- Learning  $s, a \rightarrow r$  is a regression problem
- Learning  $s, a \rightarrow s'$  is a density estimation problem
- Pick loss function, e.g. mean-squared error, KL divergence, ...
- Find parameters  $\eta$  that minimise empirical loss

## Examples of Models

- Table Lookup Model
  - Linear Expectation Model
  - Linear Gaussian Model  $\rightarrow E[\cdot]$
  - Gaussian Process Model
  - Deep Belief Network Model
  - ...
- $p(s, a) \rightarrow \text{not scale}!!$
- $w^T \phi(s, a)$

# Table Lookup Model

- Model is an explicit MDP,  $\hat{\mathcal{P}}, \hat{\mathcal{R}}$
- Count visits  $N(s, a)$  to each state action pair  
 $\hookrightarrow$  use empirical count to construct prob.

$$\hat{\mathcal{P}}_{s,s'}^a = \frac{1}{N(s, a)} \sum_{t=1}^T \mathbf{1}(S_t, A_t, S_{t+1} = s, a, s') \rightarrow \text{transition dynamics.}$$

$$\hat{\mathcal{R}}_s^a = \frac{1}{N(s, a)} \sum_{t=1}^T \mathbf{1}(S_t, A_t = s, a) R_t \rightarrow \text{average reward.}$$

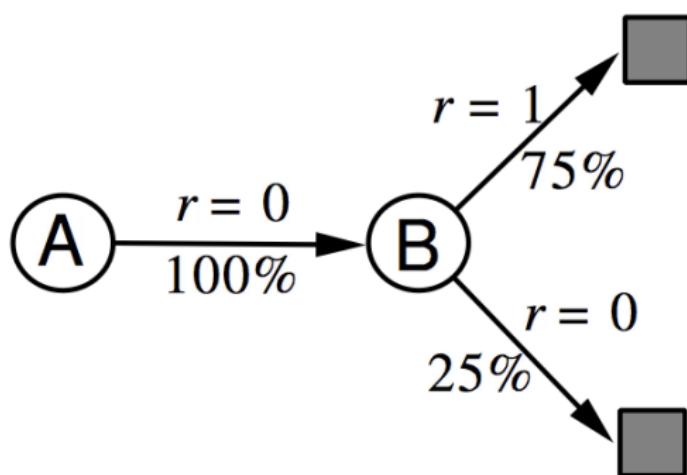
- Alternatively
  - At each time-step  $t$ , record experience tuple  $(S_t, A_t, R_{t+1}, S_{t+1}) \rightarrow$  sample uniformly by fixing  $(S_t, t)$ .
  - To sample model, randomly pick tuple matching  $(s, a, \cdot, \cdot)$

parametric VS. non-parametric

## AB Example

Two states  $A, B$ ; no discounting; 8 episodes of experience

ep <sup>1</sup>	A, 0, B, 0
ep <sup>2</sup>	B, 1
ep <sup>3</sup>	B, 1
ep <sup>4</sup>	B, 1
ep <sup>5</sup>	B, 1
ep <sup>6</sup>	B, 1
ep <sup>7</sup>	B, 1
ep <sup>8</sup>	B, 0



We have constructed a **table lookup model** from the experience

# Planning with a Model

→ estimated MDP

- Given a model  $M_\eta = \langle P_\eta, R_\eta \rangle$
  - Solve the MDP  $\langle S, A, P_\eta, R_\eta \rangle$  → try to find the best things to do.
  - Using favourite planning algorithm
    - Value iteration
    - Policy iteration
    - Tree search
    - ...
- DP with assumption of MDP being given.

# Sample-Based Planning

- A simple but powerful approach to planning
- Use the model **only** to generate samples
- Sample experience from model

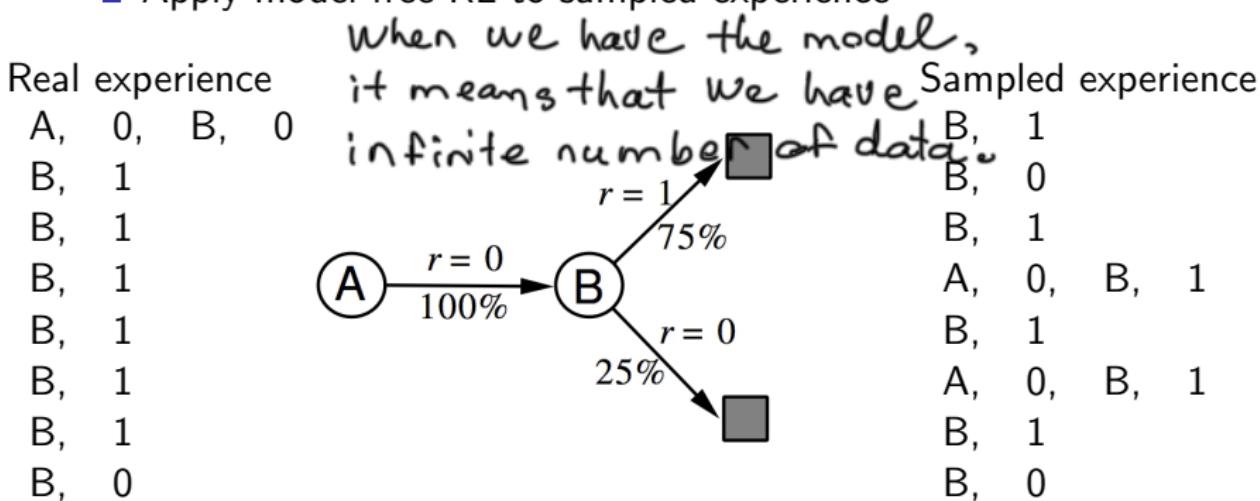
$$S_{t+1} \sim \mathcal{P}_\eta(S_{t+1} \mid S_t, A_t)$$

$$R_{t+1} = \mathcal{R}_\eta(R_{t+1} \mid S_t, A_t)$$

- Apply **model-free** RL to samples, e.g.:
  - Monte-Carlo control
  - Sarsa
  - Q-learning
- Sample-based planning methods are often more efficient

## Back to the AB Example

- Construct a table-lookup model from real experience
- Apply model-free RL to sampled experience



e.g. Monte-Carlo learning:  $V(A) = 1, V(B) = 0.75$

$$= \frac{1+1}{2} = 1 \quad = \frac{6}{8} = 0.75$$

Real exp  $\longrightarrow$  build model  $\rightarrow$  sample exp.

learn from sampled exp.  $\leftarrow$   
(MC, TD, ...)

How to evaluate our mdp model is good or  
bad?

## Planning with an Inaccurate Model

- Given an imperfect model  $\langle \mathcal{P}_\eta, \mathcal{R}_\eta \rangle \neq \langle \mathcal{P}, \mathcal{R} \rangle$
- Performance of model-based RL is limited to optimal policy for approximate MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_\eta, \mathcal{R}_\eta \rangle$
- i.e. Model-based RL is only as good as the estimated model
- When the model is inaccurate, planning process will compute a suboptimal policy
- Solution 1: when model is wrong, use model-free RL
- Solution 2: reason explicitly about model uncertainty

Model-based can be applied to continuous state-action space.

# Outline

1 Introduction

2 Model-Based Reinforcement Learning

3 Integrated Architectures

4 Simulation-Based Search

bring model-free  
&

model-based  
together

# Real and Simulated Experience

We consider two sources of experience

Real experience Sampled from environment (true MDP)

$$S' \sim \mathcal{P}_{s,s'}^a$$

$$R = \mathcal{R}_s^a$$

Simulated experience Sampled from model (approximate MDP)

$$S' \sim \mathcal{P}_\eta(S' | S, A)$$

$$R = \mathcal{R}_\eta(R | S, A)$$

How we learn to explore?

next lecture !!!

# Integrating Learning and Planning

- Model-Free RL
  - No model
  - **Learn** value function (and/or policy) from real experience

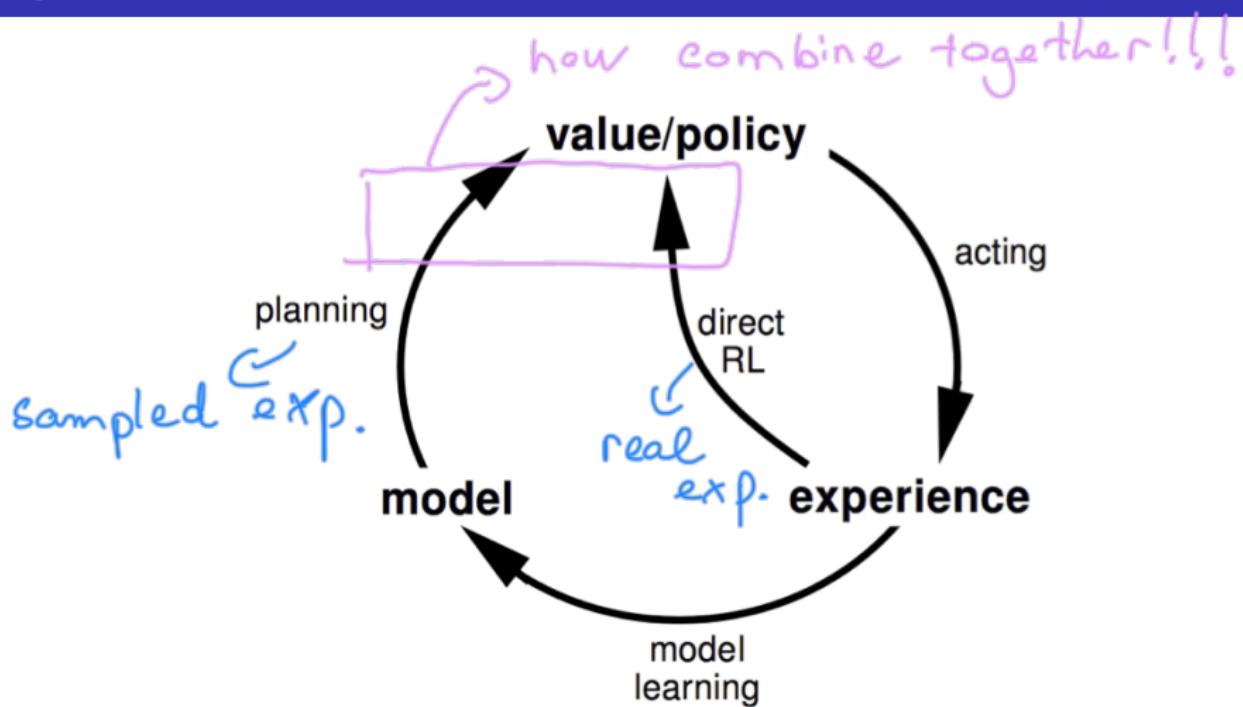
# Integrating Learning and Planning

- Model-Free RL
  - No model
  - Learn value function (and/or policy) from real experience
- Model-Based RL (using Sample-Based Planning)
  - Learn a model from real experience
  - Plan value function (and/or policy) from simulated experience

# Integrating Learning and Planning

- Model-Free RL
  - No model
  - **Learn** value function (and/or policy) from real experience
- Model-Based RL (using Sample-Based Planning)
  - Learn a model from real experience
  - **Plan** value function (and/or policy) from simulated experience
- Dyna
  - Learn a model from real experience
  - **Learn and plan** value function (and/or policy) from real and simulated experience

# Dyna Architecture



## Dyna-Q Algorithm

simply table lookup model

Initialize  $Q(s, a)$  and  $Model(s, a)$  for all  $s \in \mathcal{S}$  and  $a \in \mathcal{A}(s)$

Do forever:

(a)  $S \leftarrow$  current (nonterminal) state

(b)  $A \leftarrow \varepsilon\text{-greedy}(S, Q)$

(c) Execute action  $A$ ; observe resultant reward,  $R$ , and state,  $S'$

(d)  $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

(e)  $Model(S, A) \leftarrow R, S'$  (assuming deterministic environment)

(f) Repeat  $n$  times:  $\rightarrow$  how much thinking do I have to do?

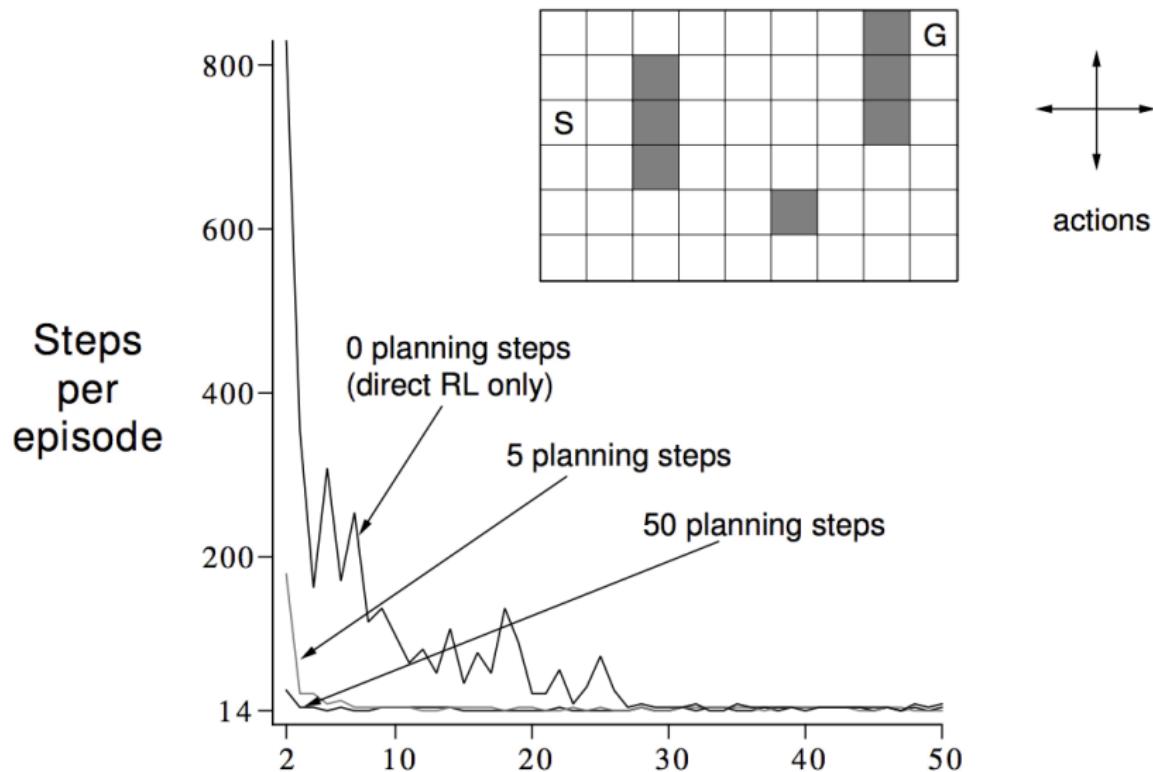
$S \leftarrow$  random previously observed state

$A \leftarrow$  random action previously taken in  $S$

$R, S' \leftarrow Model(S, A) \rightarrow$  imagine transitions

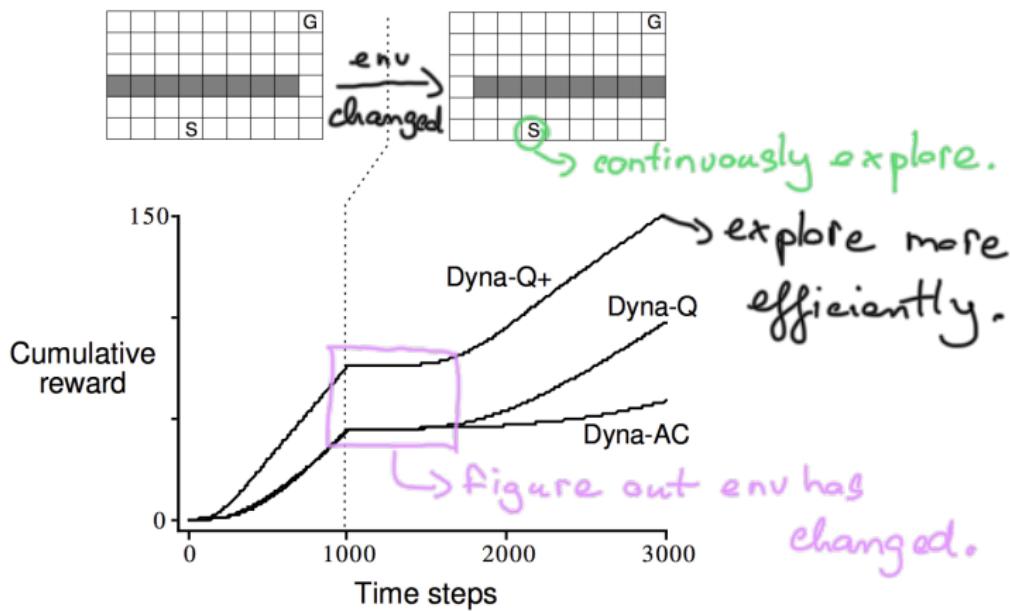
$Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a) - Q(S, A)]$

## Dyna-Q on a Simple Maze



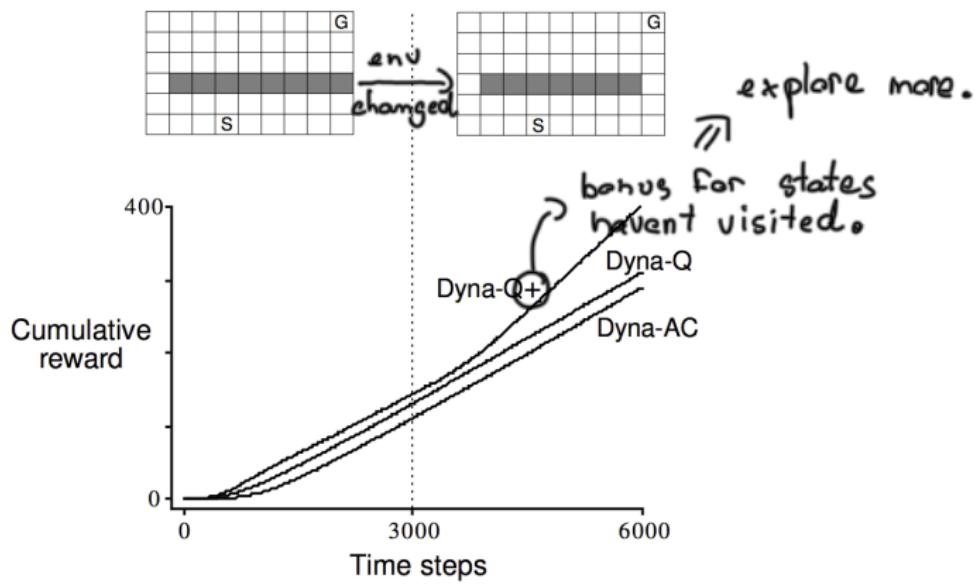
## Dyna-Q with an Inaccurate Model

- The changed environment is **harder**



## Dyna-Q with an Inaccurate Model (2)

- The changed environment is easier



# Outline

1 Introduction

2 Model-Based Reinforcement Learning

3 Integrated Architectures

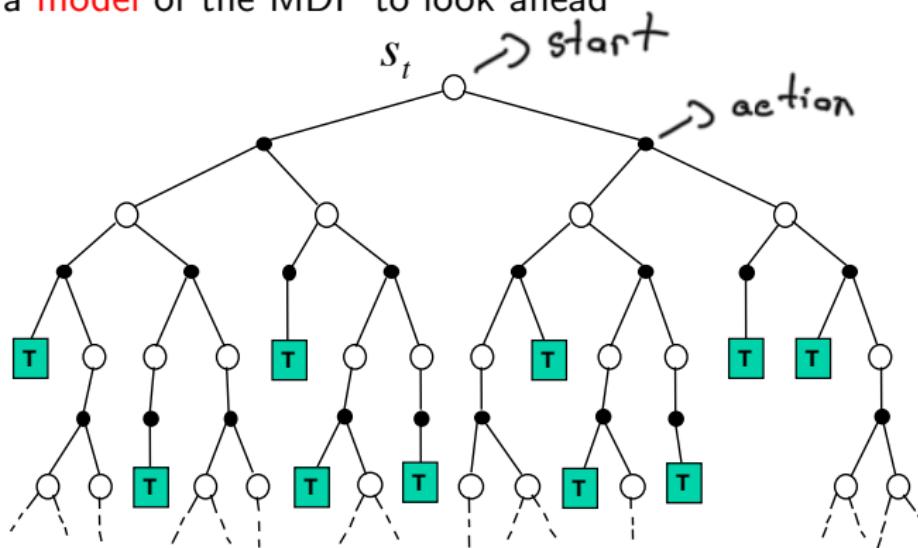
4 Simulation-Based Search

→ focus on planning part.  
→ how to plan effectively.

## Forward Search

↳ don't explore the entire space.

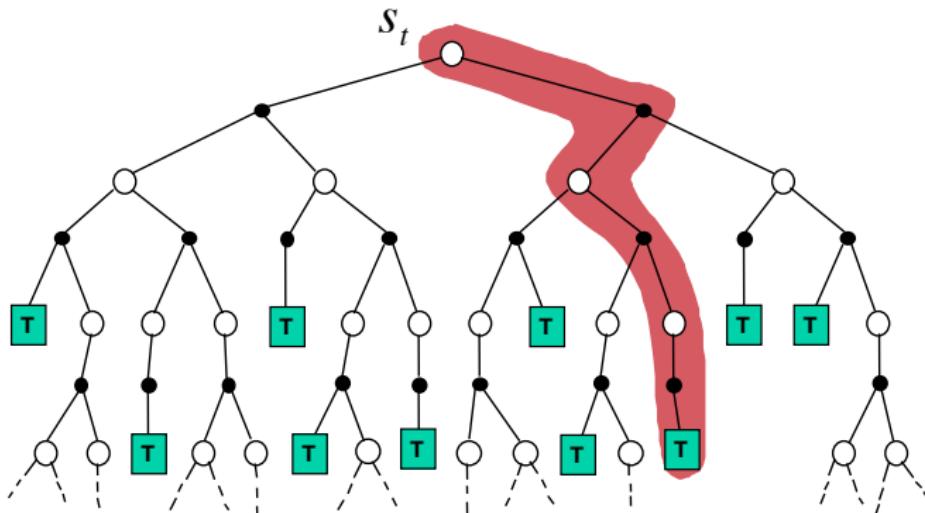
- Forward search algorithms select the best action by lookahead
- They build a search tree with the current state  $s_t$  at the root
- Using a model of the MDP to look ahead



- No need to solve whole MDP, just sub-MDP starting from now

# Simulation-Based Search

- Forward search paradigm using sample-based planning
- Simulate episodes of experience from now with the model
- Apply model-free RL to simulated episodes



## Simulation-Based Search (2)

- Simulate episodes of experience from now with the model

$$\{s_t^k, A_t^k, R_{t+1}^k, \dots, S_T^k\}_{k=1}^K \xrightarrow{\text{simulation-} t} \mathcal{M}_v$$

- Apply model-free RL to simulated episodes

- Monte-Carlo control → Monte-Carlo search
- Sarsa → TD search

# Simple Monte-Carlo Search

- Given a model  $\mathcal{M}_\nu$  and a simulation policy  $\pi$
- For each action  $a \in \mathcal{A}$ 
  - Simulate  $K$  episodes from current (real) state  $s_t$

$$\{\mathbf{s}_t, \mathbf{a}, R_{t+1}^k, S_{t+1}^k, A_{t+1}^k, \dots, S_T^k\}_{k=1}^K \sim \mathcal{M}_\nu, \pi$$

- Evaluate actions by mean return (Monte-Carlo evaluation)

$$Q(\mathbf{s}_t, \mathbf{a}) = \frac{1}{K} \sum_{k=1}^K G_t \xrightarrow{P} q_\pi(s_t, a)$$

- Select current (real) action with maximum value

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$$

## Monte-Carlo Tree Search (Evaluation)

- Given a model  $\mathcal{M}_\nu$ ,
- Simulate  $K$  episodes from current state  $s_t$  using current simulation policy  $\pi$

$$\{\textcolor{red}{s_t}, A_t^k, R_{t+1}^k, S_{t+1}^k, \dots, S_T^k\}_{k=1}^K \sim \mathcal{M}_\nu, \pi$$

- Build a search tree containing visited states and actions
- Evaluate states  $Q(s, a)$  by mean return of episodes from  $s, a$

$$Q(\textcolor{red}{s}, \textcolor{red}{a}) = \frac{1}{N(s, a)} \sum_{k=1}^K \sum_{u=t}^T \mathbf{1}(S_u, A_u = s, a) G_u \xrightarrow{P} q_\pi(s, a)$$

- After search is finished, select current (real) action with maximum value in search tree

$$a_t = \operatorname{argmax}_{a \in \mathcal{A}} Q(s_t, a)$$

# Monte-Carlo Tree Search (Simulation)

- In MCTS, the simulation policy  $\pi$  improves
- Each simulation consists of two phases (in-tree, out-of-tree)
  - Tree policy (improves): pick actions to maximise  $Q(S, A)_{\text{in-tree}}$ .
  - Default policy (fixed): pick actions randomly  $\rightarrow \text{out-of-tree}$ .
- Repeat (each simulation)
  - Evaluate states  $Q(S, A)$  by Monte-Carlo evaluation
  - Improve tree policy, e.g. by  $\epsilon$ -greedy( $Q$ )  
*starting at root state.*
- Monte-Carlo control applied to simulated experience
- Converges on the optimal search tree,  $Q(S, A) \rightarrow q_*(S, A)$

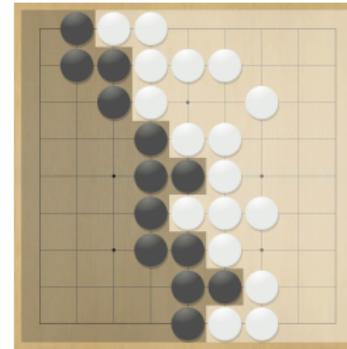
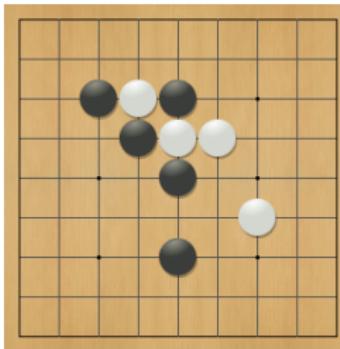
## Case Study: the Game of Go

- The ancient oriental game of Go is 2500 years old
- Considered to be the hardest classic board game
- Considered a grand challenge task for AI  
*(John McCarthy)*
- Traditional game-tree search has failed in Go



# Rules of Go

- Usually played on 19x19, also 13x13 or 9x9 board
- Simple rules, complex strategy
- Black and white place down stones alternately
- Surrounded stones are captured and removed
- The player with more territory wins the game



# Position Evaluation in Go

- How good is a position  $s$ ?
- Reward function (undiscounted):

$R_t = 0$  for all non-terminal steps  $t < T$

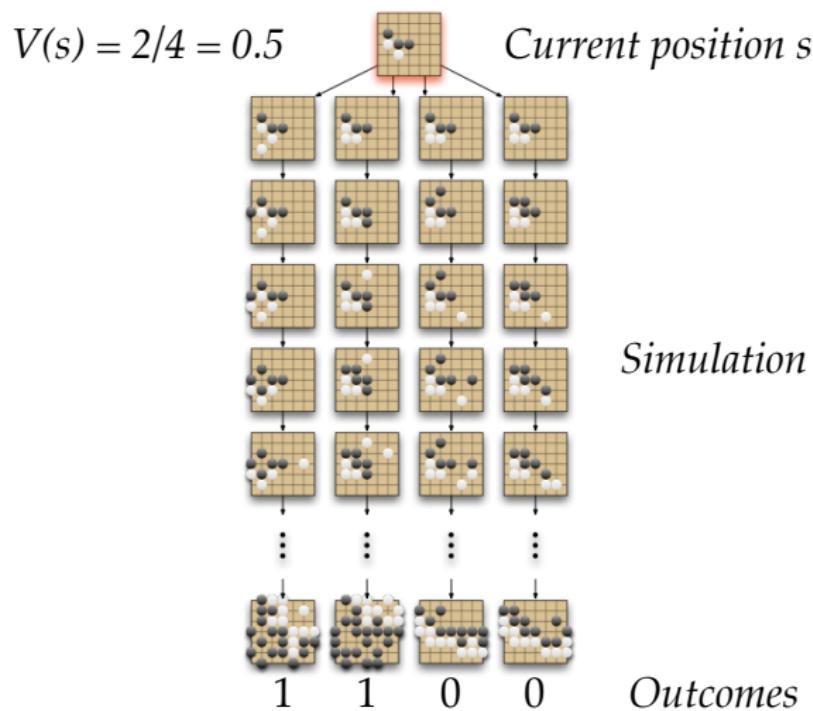
$$R_T = \begin{cases} 1 & \text{if Black wins} \\ 0 & \text{if White wins} \end{cases}$$

- Policy  $\pi = \langle \pi_B, \pi_W \rangle$  selects moves for both players
- Value function (how good is position  $s$ ):

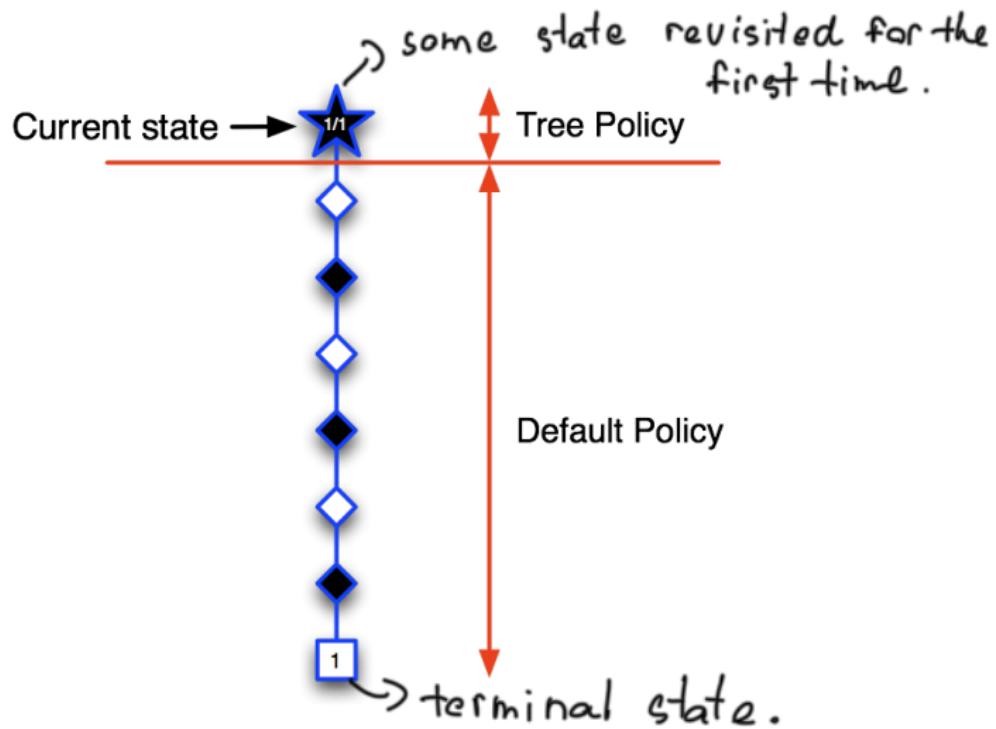
$$v_\pi(s) = \mathbb{E}_\pi [R_T \mid S = s] = \mathbb{P} [\text{Black wins} \mid S = s]$$

$$v_*(s) = \max_{\pi_B} \min_{\pi_W} v_\pi(s)$$

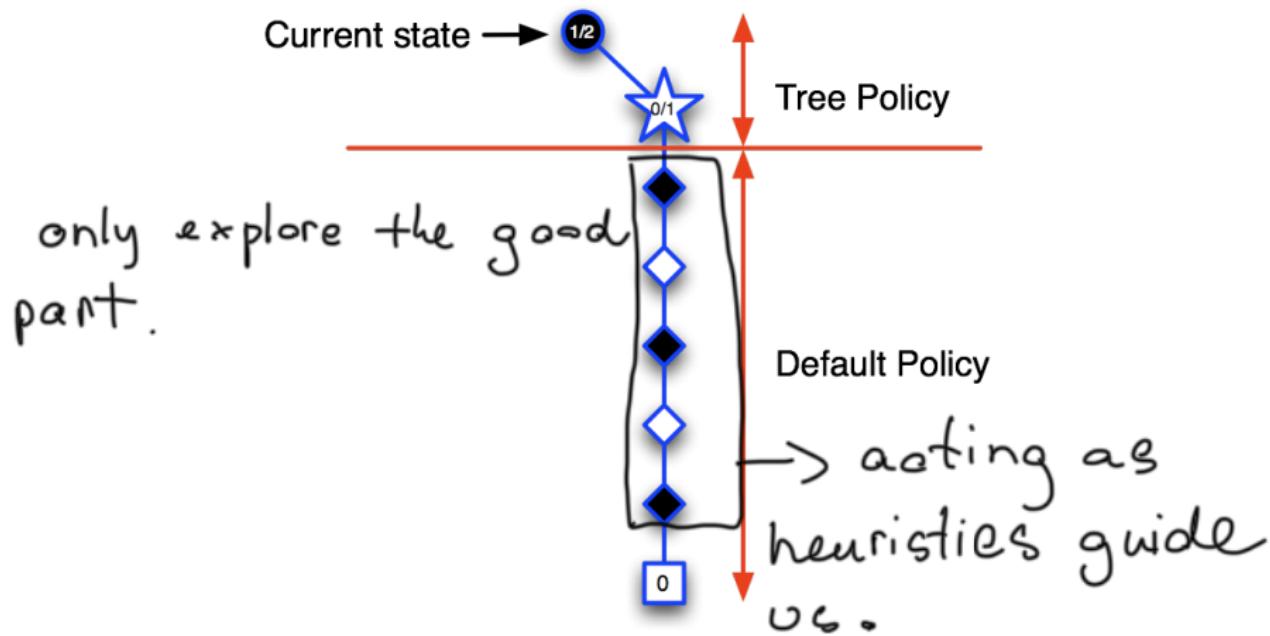
# Monte-Carlo Evaluation in Go



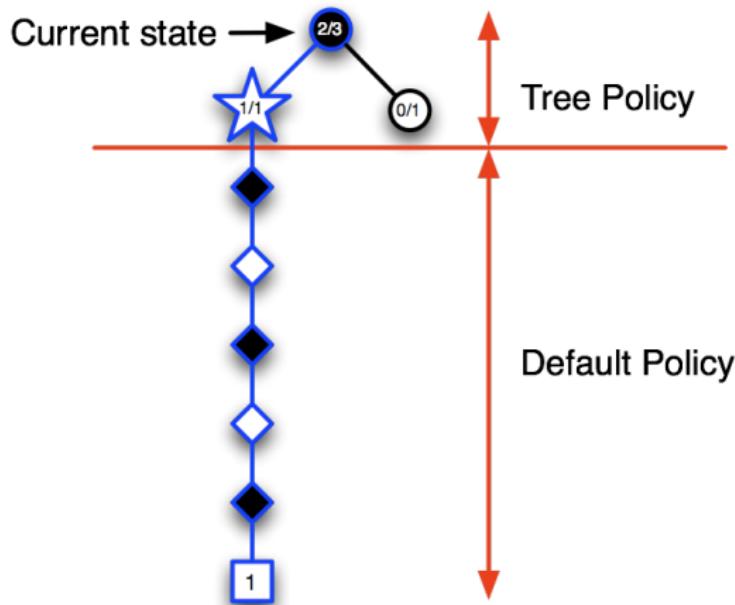
# Applying Monte-Carlo Tree Search (1)



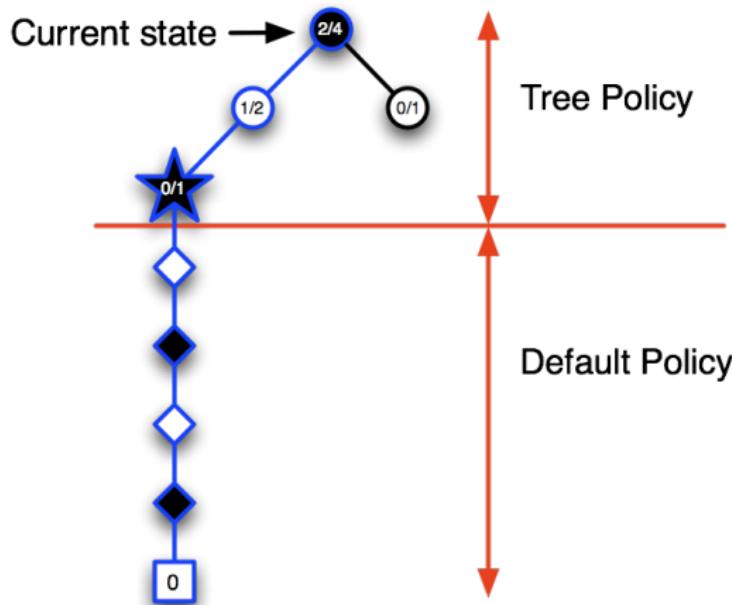
## Applying Monte-Carlo Tree Search (2)



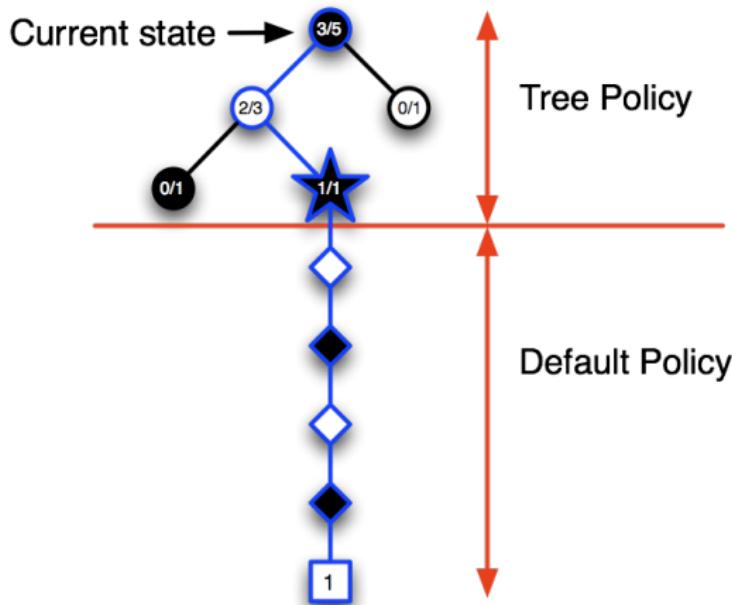
# Applying Monte-Carlo Tree Search (3)



# Applying Monte-Carlo Tree Search (4)



# Applying Monte-Carlo Tree Search (5)

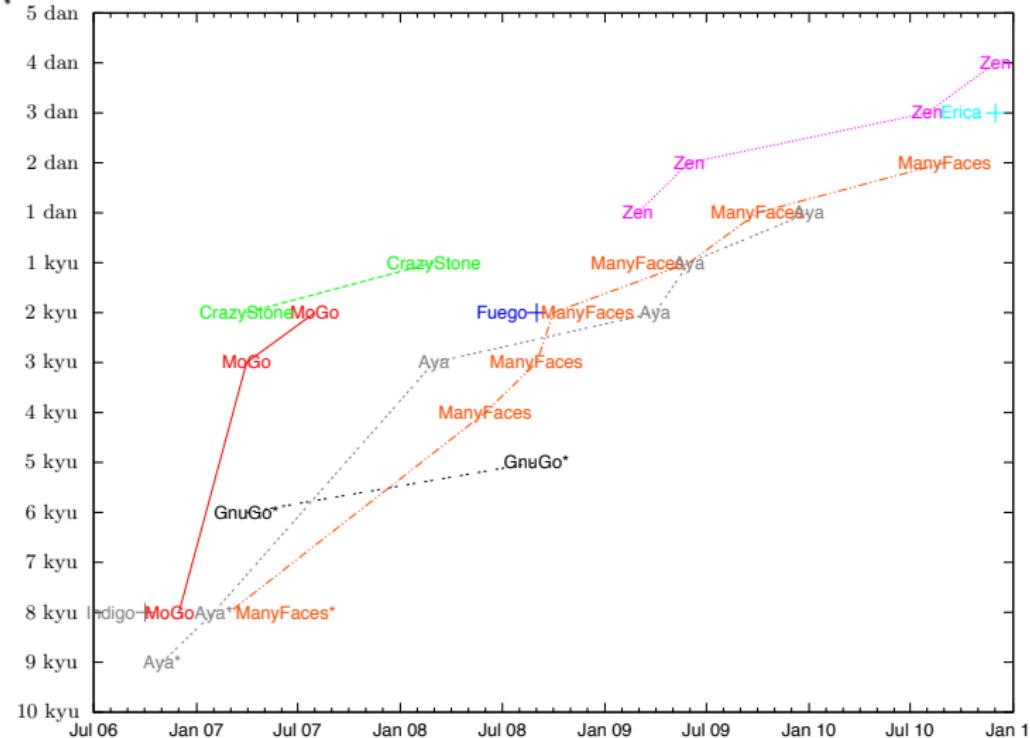


## Advantages of MC Tree Search

- Highly selective best-first search *don't have*
- Evaluates states *dynamically* (unlike e.g. DP) *↗ to consider all possible things.*
- Uses sampling to break curse of dimensionality *all possible*
- Works for “black-box” models (only requires samples) *+ things.*
- Computationally efficient, anytime, parallelisable

# Example: MC Tree Search in Computer Go

strength of model



# Temporal-Difference Search

- Simulation-based search
- Using TD instead of MC (bootstrapping)
- MC tree search applies MC control to sub-MDP from now
- TD search applies Sarsa to sub-MDP from now

## MC vs. TD search

- For model-free reinforcement learning, bootstrapping is helpful
  - TD learning reduces variance but increases bias
  - TD learning is usually more efficient than MC
  - $\text{TD}(\lambda)$  can be much more efficient than MC
- For simulation-based search, bootstrapping is also helpful
  - TD search reduces variance but increases bias
  - TD search is usually more efficient than MC search
  - $\text{TD}(\lambda)$  search can be much more efficient than MC search

# TD Search

- Simulate episodes from the current (real) state  $s_t$
- Estimate action-value function  $Q(s, a)$
- For each step of simulation, update action-values by Sarsa

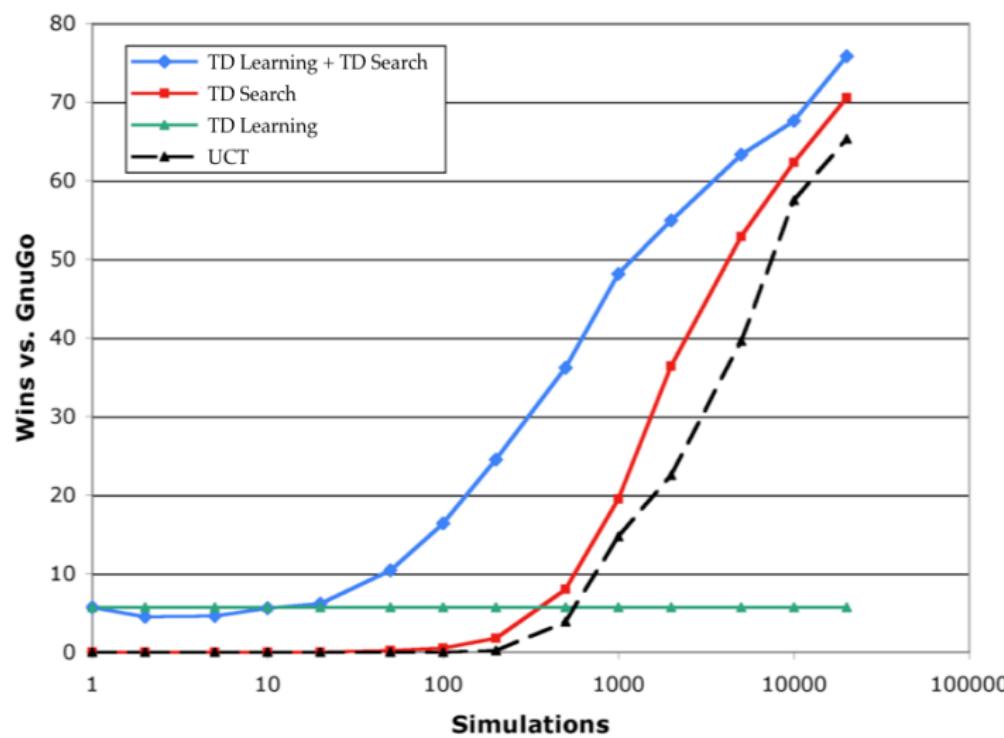
$$\Delta Q(S, A) = \alpha(R + \gamma Q(S', A') - Q(S, A))$$

- Select actions based on action-values  $Q(s, a)$ 
  - e.g.  $\epsilon$ -greedy
- May also use function approximation for  $Q$

## Dyna-2

- In Dyna-2, the agent stores two sets of feature weights
  - a value function | ■ Long-term memory
  - Short-term (working) memory
- Long-term memory is updated from real experience using TD learning
  - General domain knowledge that applies to any episode
- Short-term memory is updated from simulated experience using TD search
  - Specific local knowledge about the current situation
- Overall value function is sum of long and short-term memories

## Results of TD search in Go



# Questions?