

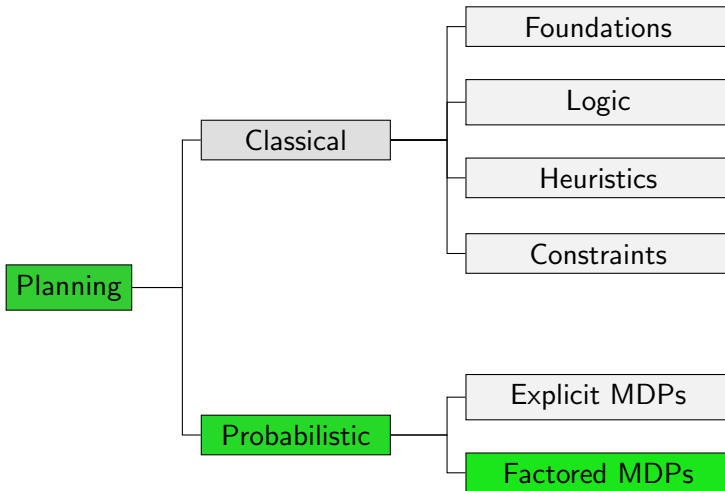
# Planning and Optimization

## G4. Monte-Carlo Tree Search: Framework

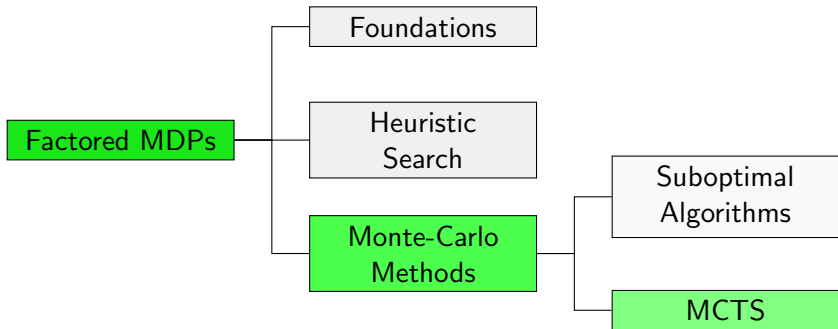
Malte Helmert and Gabriele Röger

Universität Basel

# Content of this Course



# Content of this Course: Factored MDPs



# Motivation

# Motivation

Previously discussed Monte-Carlo methods:

- Hindsight Optimization suffers from **assumption of clairvoyance**
- Policy Simulation overcomes assumption of clairvoyance by **sampling the execution of a policy**
- Policy Simulation is suboptimal **due to inability of policy to improve**
- Sparse Sampling achieves **near-optimality without considering all outcomes**
- Sparse Sampling wastes time in **non-promising parts of state space**

# Monte-Carlo Tree Search

**Monte-Carlo Tree Search** (MCTS) has several similarities with algorithms we have already seen:

- Like (L)RTDP, MCTS performs **trials** (also called **rollouts**)
- Like Policy Simulation, trials **simulate execution** of a policy

# Monte-Carlo Tree Search

**Monte-Carlo Tree Search** (MCTS) has several similarities with algorithms we have already seen:

- Like (L)RTDP, MCTS performs **trials** (also called **rollouts**)
- Like Policy Simulation, trials **simulate execution** of a policy
- Like other Monte-Carlo methods, **Monte-Carlo backups** are performed

# Monte-Carlo Tree Search

**Monte-Carlo Tree Search** (MCTS) has several similarities with algorithms we have already seen:

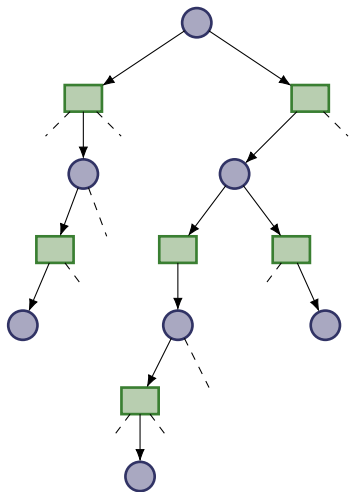
- Like (L)RTDP, MCTS performs **trials** (also called **rollouts**)
- Like Policy Simulation, trials **simulate execution** of a policy
- Like other Monte-Carlo methods, **Monte-Carlo backups** are performed
- Like Sparse Sampling, an outcome is only explicated if it is **sampled** in a trial



# MCTS Tree

# MCTS Tree

- Unlike previous methods, the SSP is **explicitated as a tree**
- **Duplicates** (also: **transpositions**) possible, i.e., multiple **search nodes** with identical associated state
- Search tree can (and often will) have **unbounded** depth





# MCTS Tree

## Definition (MCTS Tree)

An **MCTS tree** is given by a tuple  $\mathcal{G} = \langle d_0, D, C, E \rangle$ , where

- $D$  and  $C$  are disjoint sets of **decision** and **chance** nodes (simply **search node** if the type does not matter)
- $d_0 \in D$  is the **root node**
- $E \subseteq (D \times C) \cup (C \times D)$  is the set of **edges** such that the graph  $\langle D \cup C, E \rangle$  is a tree

Note: can be regarded as an AND/OR tree

# Search Node Annotations

## Definition (Search Node Annotations)

Let  $\mathcal{G} = \langle d_0, D, C, E \rangle$  be an MCTS Tree.

- Each search node  $n \in D \cup C$  is annotated with
  - a visit counter  $N(n)$
  - a state  $s(n)$
- Each decision node  $d \in D$  is annotated with
  - a state-value estimate  $\hat{V}(d)$
  - a probability  $p(d)$
- Each chance node  $c \in C$  is annotated with
  - an action-value estimate (or Q-value estimate)  $\hat{Q}(c)$
  - an action  $a(c)$

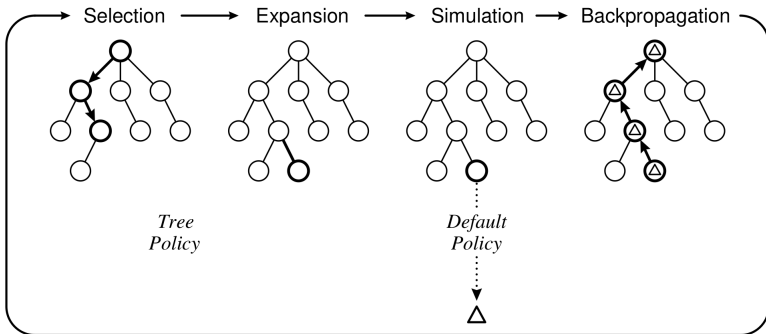
**Note:** some annotations can be computed on the fly to save memory

# Framework

# Trials

- The MCTS tree is built in **trials**
- Trials are performed as long as resources (deliberation time, memory) allow
- Initially, the MCTS tree consists of only the **root node** for the initial state
- Trials (may) **add search nodes** to the tree
- MCTS tree at the end of the  $i$ -th trial is denoted with  $\mathcal{G}^i$
- Use same superscript for annotations of search nodes

# Trials



Taken from Browne et al., "A Survey of Monte Carlo Tree Search Methods", 2012



# Phases of Trials

Each trial consists of (up to) four **phases**:

- **Selection**: traverse the tree by **sampling** the execution of the **tree policy** until
  - ① an action is applicable that is not explicated, or
  - ② an outcome is sampled that is not explicated, or
  - ③ a goal state is reached (jump to backpropagation)
- **Expansion**: **create search nodes** for the applicable action and a sampled outcome (case 1) or just the outcome (case 2)
- **Simulation**: simulate **default policy** until a goal is reached
- **Backpropagation**: update visited nodes **in reverse order** by
  - increasing visit counter by 1
  - performing Monte-Carlo backup of state-/action-value estimate

## Monte-Carlo Backups in MCTS Tree

- let  $d_0, c_0, \dots, c_{n-1}, d_n$  be the decision and chance nodes that were visited in a trial of MCTS (including explicated ones),
- let  $h$  be the cost incurred by the simulation of the default policy until a goal state is reached
- each decision node  $d_j$  for  $0 \leq j \leq n$  is updated by

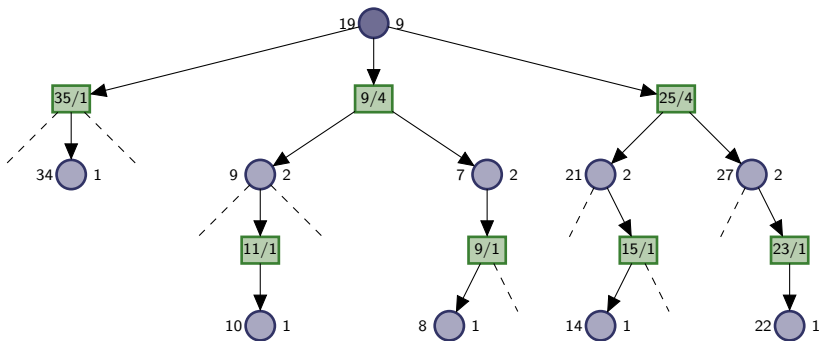
$$\hat{V}^i(d_j) := \hat{V}^{i-1}(d_j) + \frac{1}{N^i(d_j)} \left( \sum_{k=j}^{n-1} \text{cost}(a(c_k)) + h - \hat{V}^{i-1}(d_j) \right)$$

- each chance node  $c_j$  for  $0 \leq j < n$  is updated by

$$\hat{Q}^i(c_j) := \hat{Q}^{i-1}(c_j) + \frac{1}{N^i(c_j)} \left( \sum_{k=j}^{n-1} \text{cost}(a(c_k)) + h - \hat{Q}^{i-1}(c_j) \right)$$

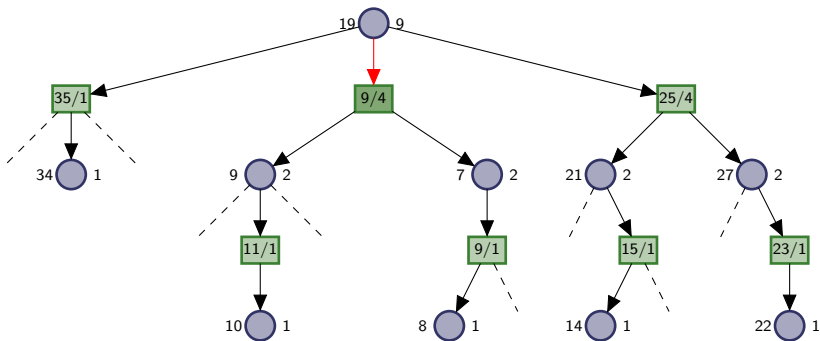
# MCTS: (Unit-cost) Example

**Selection phase:** apply tree policy to traverse tree



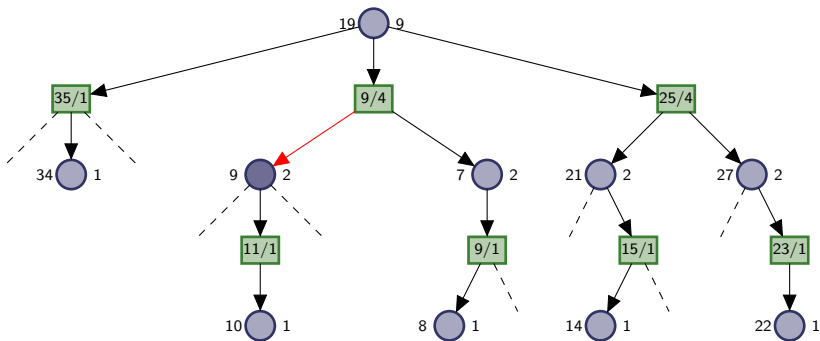
# MCTS: (Unit-cost) Example

**Selection phase:** apply tree policy to traverse tree



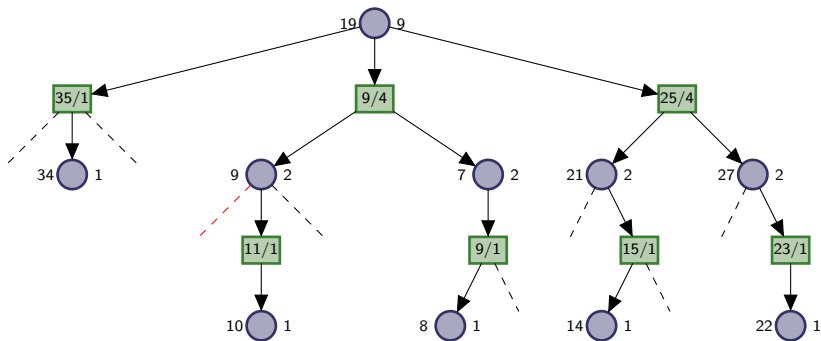
# MCTS: (Unit-cost) Example

**Selection phase:** apply tree policy to traverse tree



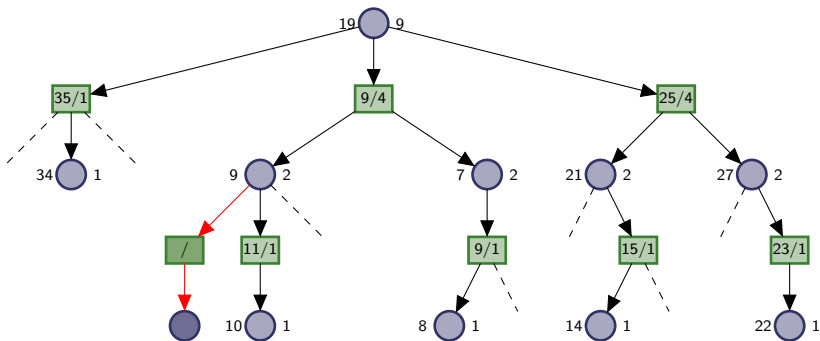
# MCTS: (Unit-cost) Example

**Selection phase:** apply tree policy to traverse tree



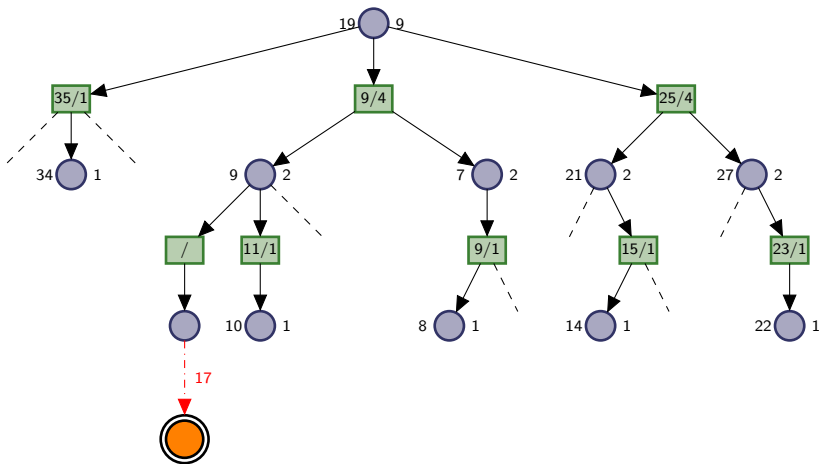
# MCTS: (Unit-cost) Example

Expansion phase: create search nodes



# MCTS: (Unit-cost) Example

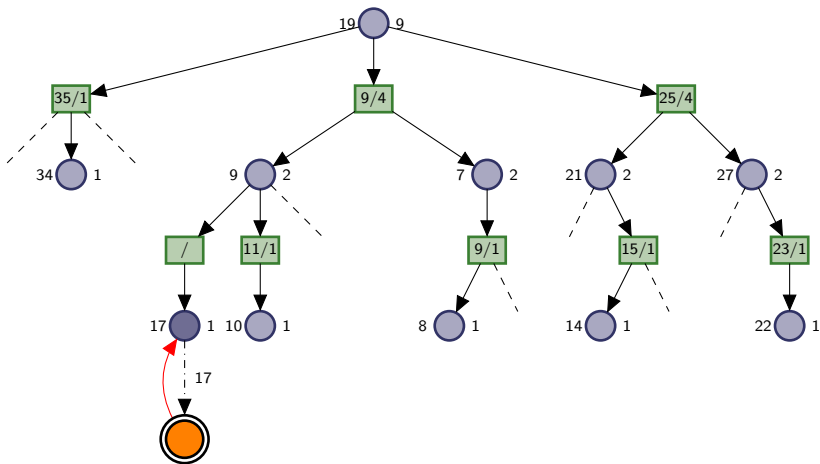
**Simulation phase:** apply default policy until goal





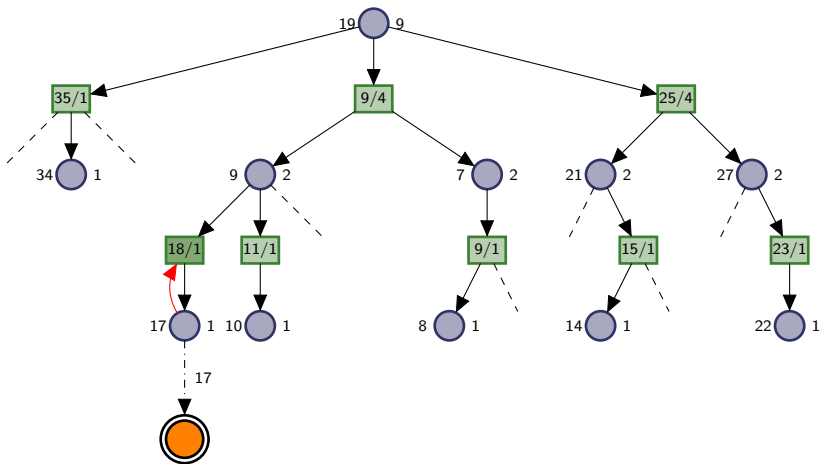
# MCTS: (Unit-cost) Example

**Backpropagation phase:** update visited nodes



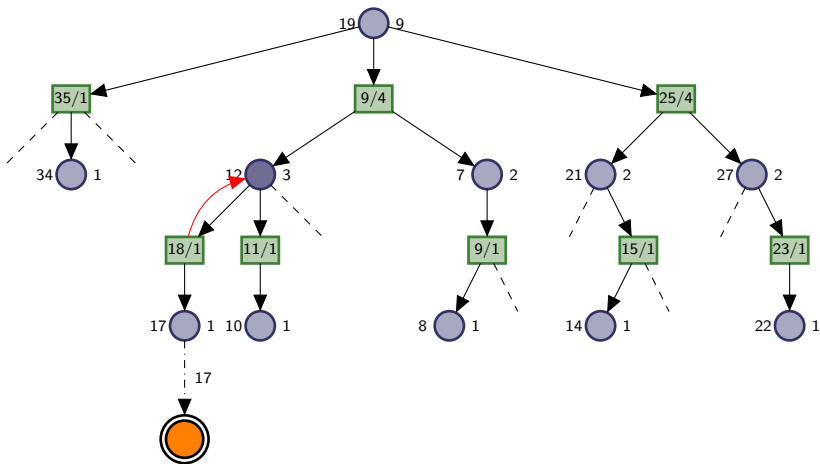
# MCTS: (Unit-cost) Example

Backpropagation phase: update visited nodes



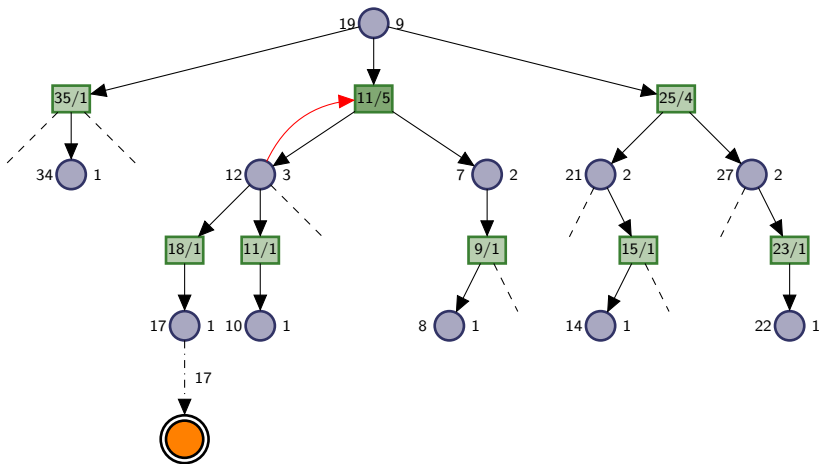
# MCTS: (Unit-cost) Example

**Backpropagation phase:** update visited nodes



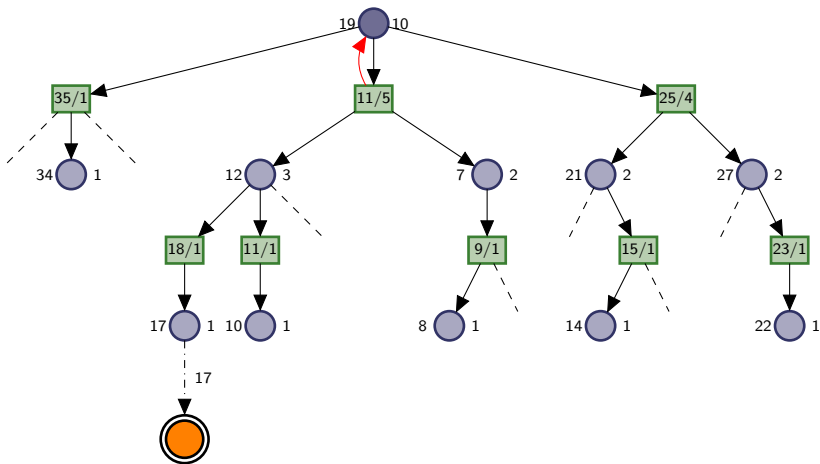
# MCTS: (Unit-cost) Example

Backpropagation phase: update visited nodes



# MCTS: (Unit-cost) Example

Backpropagation phase: update visited nodes



# MCTS Framework

Member of MCTS **framework** are specified in terms of:

- Tree policy
- Default policy

# MCTS Tree Policy

## Definition (Tree Policy)

Let  $\mathcal{T}$  be an SSP. An **MCTS tree policy** is a probability distribution  $\pi(a \mid d)$  over all  $a \in A(s(d))$  for each decision node  $d$ .

**Note:** The tree policy may take information annotated in the current tree into account.

# MCTS Default Policy

## Definition (Default Policy)

Let  $\mathcal{T}$  be an SSP. An **MCTS default policy** is a probability distribution  $\pi(a | s)$  over actions  $a \in A(s)$  for each state  $s$ .

**Note:** The default policy is independent of the MCTS tree.



# Monte-Carlo Tree Search

MCTS for SSP  $\mathcal{T} = \langle S, A, c, T, s_0, S_* \rangle$

$d_0$  = create root node associated with  $s_0$

**while** time allows:

    visit\_decision\_node( $d_0, \mathcal{T}$ )

**return**  $a(\arg \min_{c \in \text{children}(d_0)} \hat{Q}(c))$

## MCTS: Visit a Decision Node

visit\_decision\_node for decision node  $d$ , SSP

$\mathcal{T} = \langle S, A, c, T, s_0, S_* \rangle$

**if**  $s(d) \in S_*$  **then return** 0

**if** there is  $a \in A(s(d))$  s.t.  $a(c) \neq a$  for all  $c \in \text{children}(d)$ :

select such an  $a$  and add node  $c$  with  $a(c) = a$  to  $\text{children}(d)$

**else:**

$c = \text{tree\_policy}(d)$

cost = visit\_chance\_node( $c, \mathcal{T}$ )

$N(d) := N(d) + 1$

$\hat{V}(d) := \hat{V}(d) + \frac{1}{N(d)} \cdot (\text{cost} - \hat{V}(d))$

**return** cost

# MCTS: Visit a Chance Node

visit\_chance\_node for chance node  $c$ , SSP  $\mathcal{T} = \langle S, L, c, T, s_0, S_\star \rangle$

$s' \sim \text{succ}(s(c), a(c))$

let  $d$  be the node in  $\text{children}(c)$  with  $s(d) = s'$

**if** there is no such node:

    add node  $d$  with  $s(d) = s'$  to  $\text{children}(c)$

    cost = sample\_default\_policy( $s'$ )

$N(d) := 1, \hat{V}(d) := \text{cost}$

**else:**

    cost = visit\_decision\_node( $d, \mathcal{T}$ )

cost = cost + cost( $s(c), a(c)$ )

$N(c) := N(c) + 1$

$\hat{Q}(c) := \hat{Q}(c) + \frac{1}{N(c)} \cdot (\text{cost} - \hat{Q}(c))$

**return** cost

# Summary

# Summary

- Monte-Carlo Tree Search is a **framework** for algorithms
- MCTS algorithms perform trials
- Each trial consists of (up to) 4 phases
- MCTS algorithms are specified by two policies:
  - a **tree policy** that describes behavior “in” tree
  - and a **default policy** that describes behavior “outside” of tree