

# Monte-Carlo Tree Search by Best Arm Identification

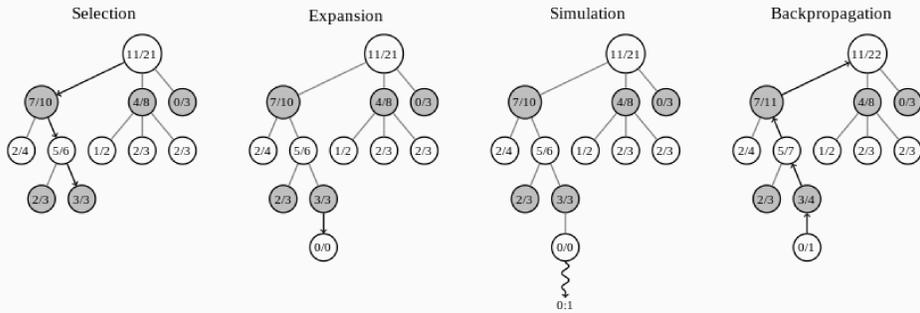
Emilie Kaufmann Wouter M. Koolen



CWI

Centrum Wiskunde & Informatica

## Monte Carlo Tree Search Successful ...



source: Wikipedia

## Goal

A MCTS algorithm should find the **best move at the root**:

$$V_s = \begin{cases} \mu_s & \text{if } s \in \mathcal{L}, \\ \max_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MAX node,} \\ \min_{c \in \mathcal{C}(s)} V_c & \text{if } s \text{ is a MIN node.} \end{cases} \quad s^* = \operatorname{argmax}_{s \in \mathcal{C}(s_0)} V_s$$

## A PAC learning framework

MCTS algorithm:  $(L_t, \tau, \hat{s}_\tau)$ , where

- $L_t$  is the **sampling rule**
- $\tau$  is the **stopping rule**
- $\hat{s}_\tau \in \mathcal{C}(s_0)$  is the **recommendation rule**

is  $(\epsilon, \delta)$ -PAC if  $\mathbb{P}(V_{\hat{s}_\tau} \geq V_{s^*} - \epsilon) \geq 1 - \delta$ .

Goal:  $(\epsilon, \delta)$ -PAC algorithm with a small **sample complexity**  $\tau$ .

## BAI-MCTS Architecture

**Input:** a BAI algorithm

**Initialization:**  $t = 1$ .

**while not** BAIStop  $(\{s \in \mathcal{C}(s_0)\})$  **do**

$R_t = \text{BAIStep}(\{s \in \mathcal{C}(s_0)\})$

    Sample the **representative leaf**  $L_t = \ell_{R_t}(t)$

    Update conf. intervals and representative leaves;  $t = t + 1$ .

**end**

**Output:** BAIReco  $(\{s \in \mathcal{C}(s_0)\})$

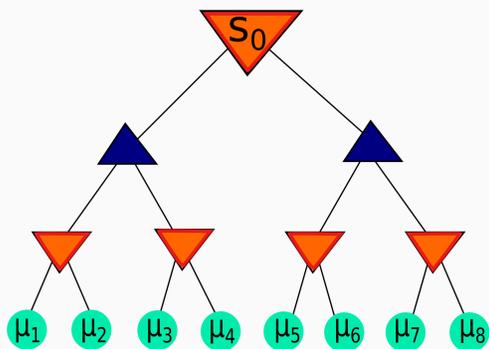
## ... But Theory Missing!

We introduce an idealized model:

- fixed** maximin tree
- i.i.d.** payouts starting from each leaf

and propose **new algorithms** with **sample complexity guarantees** based on **Best Arm Identification** methods.

## Stylised Model



A fixed MAXMIN game tree  $\mathcal{T}$ , with leaves  $\mathcal{L}$ .

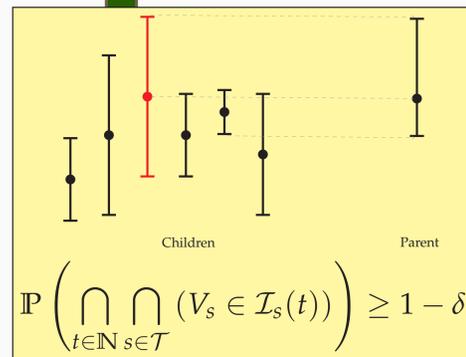
▲ MAX node (= your move)

▲ MIN node (= adversary move)

● Leaf  $\ell$ : stochastic oracle  $\mathcal{O}_\ell$  that evaluates the position

## Confidence Intervals at Depth One

$\ell_s(t)$ : **representative leaf** of internal node  $s \in \mathcal{T}$ .



**Idea:** alternate optimistic/pessimistic moves starting from  $s$

## Theoretical Results

UGapE-MCTS is  $(\epsilon, \delta)$ -PAC for confidence intervals of the form

$$\hat{\mu}_\ell(t) \pm \sqrt{\frac{\beta(N_\ell(t), \delta)}{2N_\ell(t)}}$$

where  $\beta(s, \delta) = \log(|\mathcal{L}|/\delta) + 3 \log \log(|\mathcal{L}|/\delta) + (3/2) \log(\log s + 1)$ .

**Sample complexity:**  $\tau = O\left(\sum_{\ell \in \mathcal{L}} \frac{1}{\Delta_\ell^2 V \Delta_*^2 v \epsilon^2} \log\left(\frac{1}{\delta}\right)\right)$  w.p.  $\geq 1 - \delta$ ,

where

$$\Delta_* := V(s^*) - V(s_2^*),$$

$$\Delta_\ell := \max_{s \in \text{Ancestors}(\ell) \setminus \{s_0\}} |V_{\text{Parent}(s)} - V_s|.$$

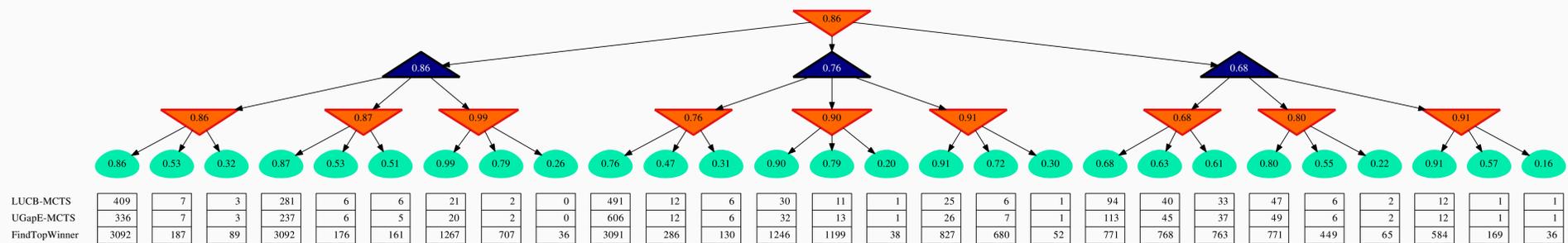
## Learning Protocol

At round  $t$  a MCTS algorithm:

- picks a path down to a leaf  $L_t$
- get an evaluation of this leaf  $X_t \sim \mathcal{O}_{L_t}$

Assumption: i.i.d. successive evaluations,  $\mathbb{E}_{X \sim \mathcal{O}_\ell}[X] = \mu_\ell$

## Numerical Results



Our benchmark 3-way tree of depth 3. Shown below the leaves are the numbers of pulls of 3 algorithms: LUCB-MCTS (0.72% errors, 1551 samples), UGapE-MCTS (0.75%, 1584), and FindTopWinner (0%, 20730). Numbers are averages over 10K repetitions with  $\epsilon = 0$  and  $\delta = 0.1 \cdot 27$ .