

# Extreme Value Monte Carlo Tree Search

Masataro Asai<sup>1</sup> and Stephen Wissow<sup>2</sup>

<sup>1</sup>MIT-IBM Watson AI Lab, USA, <sup>2</sup>University of New Hampshire, USA

**MIT-IBM**  
Watson AI Lab

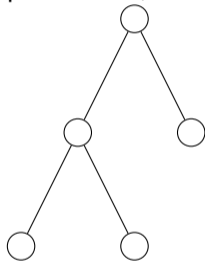


**University of  
New Hampshire**

# Monte Carlo Tree Search (MCTS/UCT)

Auer, Cesa-Bianchi, and Fischer (2002, UCB1), Kocsis and Szepesvári (2006, UCT), Schulte and Keller (2014, THTS)

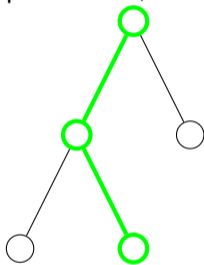
partial tree, mid-search:



# Monte Carlo Tree Search (MCTS/UCT)

Auer, Cesa-Bianchi, and Fischer (2002, UCB1), Kocsis and Szepesvári (2006, UCT), Schulte and Keller (2014, THTS)

partial tree, mid-search:



1. **select** a leaf node by

$$LCB1_i = \hat{\mu}_i - c\sqrt{\frac{2 \log T}{t_i}}$$

$\hat{\mu}_i$  : mean of child  $i$

$t_i$  : visit count of  $i$

$T$  : parent visit count

MCTS

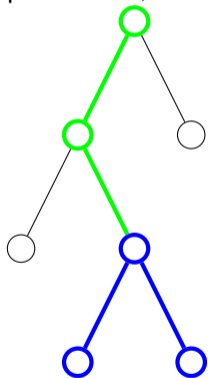
Extreme Value  
Theory

Results

# Monte Carlo Tree Search (MCTS/UCT)

Auer, Cesa-Bianchi, and Fischer (2002, UCB1), Kocsis and Szepesvári (2006, UCT), Schulte and Keller (2014, THTS)

partial tree, mid-search:



1. **select** a leaf node by  
$$LCB1_i = \hat{\mu}_i - c\sqrt{\frac{2 \log T}{t_i}}$$

$\hat{\mu}_i$  : mean of child  $i$   
 $t_i$  : visit count of  $i$   
 $T$  : parent visit count
2. **expand** a leaf node

MCTS

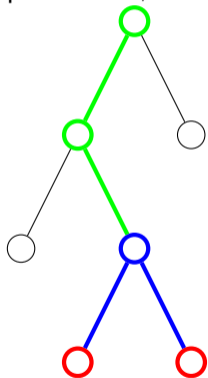
Extreme Value  
Theory

Results

# Monte Carlo Tree Search (MCTS/UCT)

Auer, Cesa-Bianchi, and Fischer (2002, UCB1), Kocsis and Szepesvári (2006, UCT), Schulte and Keller (2014, THTS)

partial tree, mid-search:



1. **select** a leaf node by  
$$LCB1_i = \hat{\mu}_i - c\sqrt{\frac{2 \log T}{t_i}}$$
$$\hat{\mu}_i$$
 : mean of child  $i$ 
$$t_i$$
 : visit count of  $i$ 
$$T$$
 : parent visit count
2. **expand** a leaf node
3. **evaluate** heuristics of children

MCTS

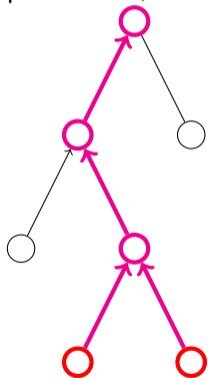
Extreme Value  
Theory

Results

# Monte Carlo Tree Search (MCTS/UCT)

Auer, Cesa-Bianchi, and Fischer (2002, UCB1), Kocsis and Szepesvári (2006, UCT), Schulte and Keller (2014, THTS)

partial tree, mid-search:



- select** a leaf node by  

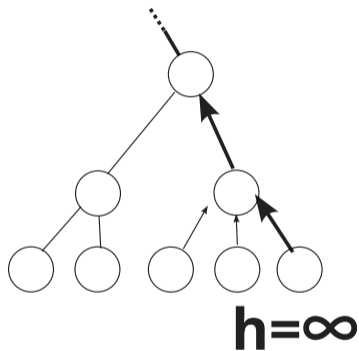
$$\text{LCB1}_i = \hat{\mu}_i - c\sqrt{\frac{2 \log T}{t_i}}$$
 $\hat{\mu}_i$  : mean of child  $i$   
 $t_i$  : visit count of  $i$   
 $T$  : parent visit count
- expand** a leaf node
- evaluate** heuristics of children
- backup** information to ancestors  
 $t = \sum_i t_i$  : Sum of children  
 $\hat{\mu} = \frac{\sum_i t_i \hat{\mu}_i}{\sum_i t_i}$  : Weighted avg

MCTS

Extreme Value  
Theory

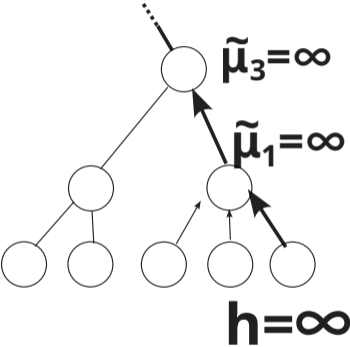
Results

# Average is Weird in Planning: What happens at a dead-end?



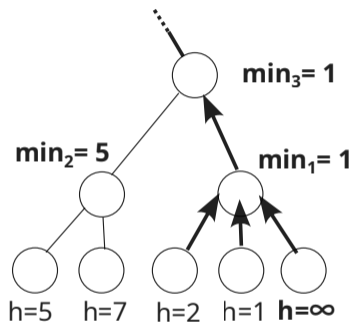
# Average is Weird in Planning: What happens at a dead-end?

► A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$



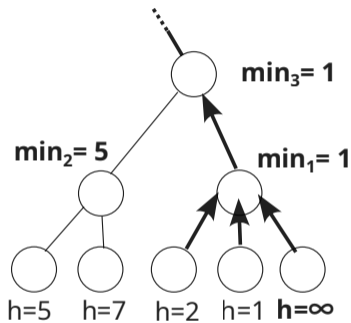


# Average is Weird in Planning: What happens at a dead-end?



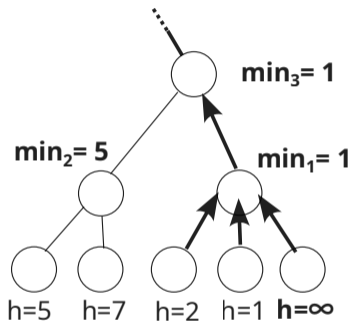
- ▶ A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$
- ▶ min has no such problem

# Average is Weird in Planning: What happens at a dead-end?



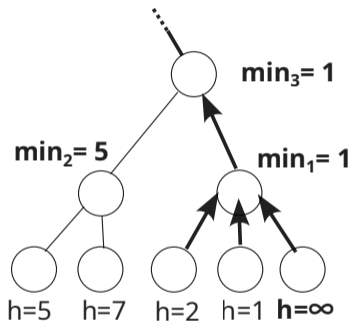
- ▶ **A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$**
- ▶ min has no such problem
- ▶ backup min  $h$ , select min  $h =$  GBFS  
 backup avg  $h$ , select UCB1 = GUCT  
 backup min  $h$ , select UCB1 = GUCT\*  
 (Schulte and Keller 2014, THTS)  
 However...

# Average is Weird in Planning: What happens at a dead-end?



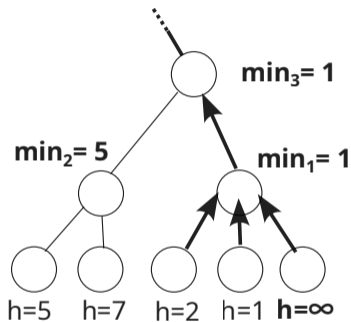
- ▶ **A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$**
- ▶ min has no such problem
- ▶ backup min  $h$ , select min  $h =$  GBFS  
 backup avg  $h$ , select UCB1 = GUCT  
 backup min  $h$ , select UCB1 = GUCT\*  
 (Schulte and Keller 2014, THTS)  
 However...
- ▶ **min lacks statistical interpretation**

# Average is Weird in Planning: What happens at a dead-end?



- ▶ **A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$**
- ▶ min has no such problem
- ▶ backup min  $h$ , select min  $h =$  GBFS  
 backup avg  $h$ , select UCB1 = GUCT  
 backup min  $h$ , select UCB1 = GUCT\*  
 (Schulte and Keller 2014, THTS)  
 However...
- ▶ **min lacks statistical interpretation**
- ▶ GUCT/\* also remove  $\infty$  from the tree

# Average is Weird in Planning: What happens at a dead-end?



- ▶ **A single  $h = \infty \rightarrow$  all  $\hat{\mu} = \infty$**
- ▶ min has no such problem
- ▶ backup min  $h$ , select min  $h =$  GBFS  
 backup avg  $h$ , select UCB1 = GUCT  
 backup min  $h$ , select UCB1 = GUCT\*  
 (Schulte and Keller 2014, THTS)  
 However...
- ▶ **min lacks statistical interpretation**
- ▶ GUCT/\* also remove  $\infty$  from the tree
- ▶ **Lacks statistical interpretation**

# The key slide™

The average is weird, but how to use the minimum with statistical rigor?

The average is weird, but how to use the minimum with statistical rigor?

The statistical theory of the average:  
is **Central Limit Theorem (CLT)**.

The average is weird, but how to use the minimum with statistical rigor?

The statistical theory of the average:  
is **Central Limit Theorem (CLT)**.

The statistical theory of the minimum (or maximum)?



The average is weird, but how to use the minimum with statistical rigor?

The statistical theory of the average:  
is **Central Limit Theorem (CLT)**.

The statistical theory of the minimum (or maximum)?  
It's **Extreme Value Theory (EVT)**!

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

- ▶ **Monthly average water level**

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

- ▶ **Monthly average water level**  
→ Gaussian distribution

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

- ▶ **Monthly average water level**  
→ Gaussian distribution
- ▶ **Exceedance over the safety limit**

**Safety-critical** applications: e.g. **Maximum water level in rivers**

- ▶ **Monthly average water level**  
→ Gaussian distribution
- ▶ **Exceedance over the safety limit**  
→ **Generalized Pareto (GP) distribution**

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

▶ **Monthly average water level**

→ Gaussian distribution

▶ **Exceedance over the safety limit**

→ **Generalized Pareto (GP) distribution**

$$\text{GP}(x \mid \theta, \sigma, \xi) = \begin{array}{l} \text{Unimportant,} \\ \text{complicated math.} \end{array} (x > \theta)^* \text{ for threshold } \theta$$

# Extreme Value Theory (EVT)

**Safety-critical** applications: e.g. **Maximum water level in rivers**

▶ **Monthly average water level**

→ Gaussian distribution

▶ **Exceedance over the safety limit**

→ **Generalized Pareto (GP) distribution**

$$\text{GP}(x \mid \theta, \sigma, \xi) = \begin{array}{l} \text{Unimportant,} \\ \text{complicated math.} \end{array} \quad (x > \theta)^* \text{ for threshold } \theta$$

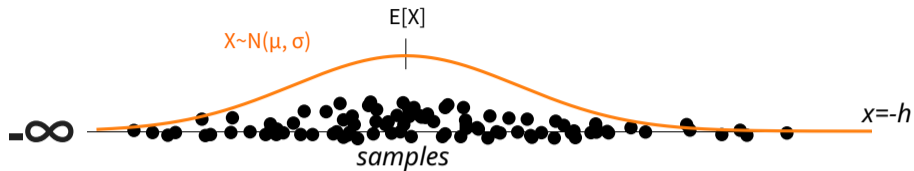
**\*justifies removing dead-ends:**  $x = -h = -\infty$  (minimize  $h$ ; maximize  $x$ )



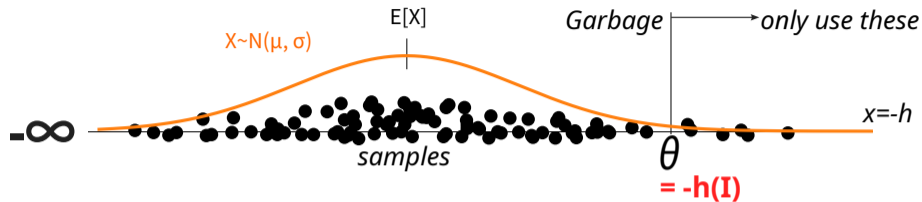
# Backup = Fitting a distribution



# Backup = Fitting a distribution

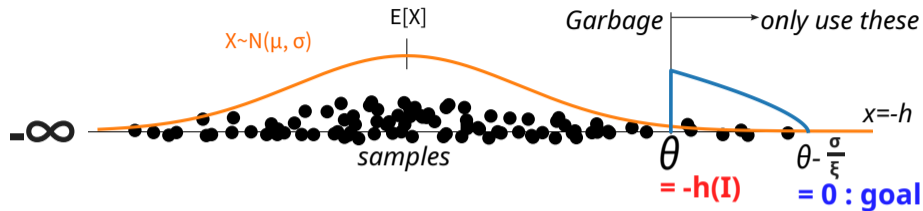


# Backup = Fitting a distribution



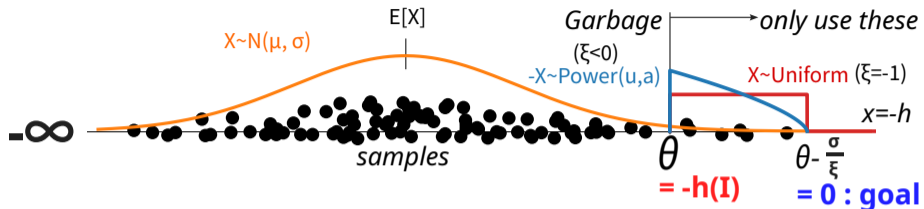
► initial state  $I$

# Backup = Fitting a distribution



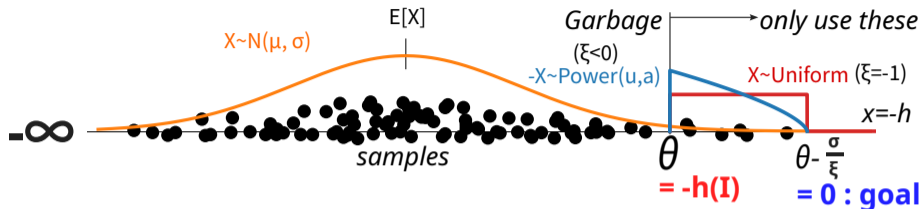
- ▶ **initial state  $I$**
- ▶ We fit **exceedance over  $\theta = -h(I)$  to  $GP(\theta, \sigma, \xi)$**

# Backup = Fitting a distribution



- ▶ **initial state  $I$**
- ▶ We fit **exceedance over  $\theta = -h(I)$  to  $GP(\theta, \sigma, \xi)$**
- ▶ We use special cases of GP: Power ( $\xi < 0$ ) and Uniform ( $\xi = -1$ )

# Backup = Fitting a distribution



- ▶ **initial state  $I$**
- ▶ We fit **exceedance over  $\theta = -h(I)$  to  $GP(\theta, \sigma, \xi)$**
- ▶ We use special cases of GP: Power ( $\xi < 0$ ) and Uniform ( $\xi = -1$ )
- ▶ We define two Bandits:

$$\begin{aligned} \text{LCB1-Uniform}_i &= \frac{\hat{u}_i + \hat{l}_i}{2} - (\hat{u}_i - \hat{l}_i) \sqrt{6t_i \log T} \\ \text{LCB1-Power}_i &= \frac{\hat{u}_i \hat{a}_i}{\hat{a}_i + 1} - \hat{u}_i \sqrt{6t_i \log T} \end{aligned}$$

Num. solved on 24 IPC domains w/  $10^4$  evaluations

$h =$	$h^{\text{FF}}$	$h^{\text{add}}$	$h^{\text{max}}$	$h^{\text{GC}}$	$h^{\text{FF}+\text{PO}}$	$h^{\text{FF}+\text{DE}}$	$h^{\text{FF}+\text{DE}+\text{PO}}$
GBFS	538	518	224	354	-	489	-
Softmin-Type(h)	576	542.6	297.2	357.6	-	578	-
GUCT <b>Uses UCB1</b>	412	397.8	228.4	285.2	454	389.2	439.4
-Normal <small>Uses UCB1-Normal</small>	283.4	265	212	233.4	372.4	289	381.6
*-Normal <small>+ backprop min</small>	318.8	300	215.2	246.2	378.05	304.4	386.7
-Normal2	581.8	535.8	316.6	379	621	518	578
*-Normal2	567.2	533.8	263	341	618	511.4	567.8
-Power	<b>596</b>	541.8	<b>450.6</b>	463.2	<b>623.4</b>	413.6	<b>583</b>
-Uniform	<b>594.8</b>	<b>543.8</b>	<b>450.6</b>	<b>463.8</b>	<b>626.4</b>	416.4	<b>583</b>

FD/C++ implementation is on the way and showing promising results

Find our full paper: <https://arxiv.org/abs/2405.18248>

# References I

Auer, P.; Cesa-Bianchi, N.; and Fischer, P. 2002. Finite-Time Analysis of the Multiarmed Bandit Problem. *Machine Learning*, 47(2-3): 235–256.

Kocsis, L.; and Szepesvári, C. 2006. Bandit Based Monte-Carlo Planning. In *Proc. of ECML*, 282–293. Springer.

Schulte, T.; and Keller, T. 2014. Balancing Exploration and Exploitation in Classical Planning. In *Proc. of SOCS*.