# MCTS BASED ON SIMPLE REGRET AAAI (SOCS) 2012

# MULTI-ARMED BANDITS

- A set of *K* arms.
- Each arm can be pulled multiple times.



• When the ith arm is pulled, a random reward  $X_i$  is encountered.

## Regret minimization:

- Simple regret (SR): the reward of the last pull only is collected.
- Cumulative regret (CR): all rewards are accumulated.

# **UCB** AND **UCT**

- **UCB**(c) pulls arm i that maximizes upper confidence bound  $b_i$  on the reward:  $b_i = \overline{X}_i + \sqrt{\frac{c \log(n)}{n_i}}$
- UCB is nearly optimal in minimizing the *cumulative regret*.
- **UCT** extends UCB to MCTS by invoking UCB at every node of a rollout.

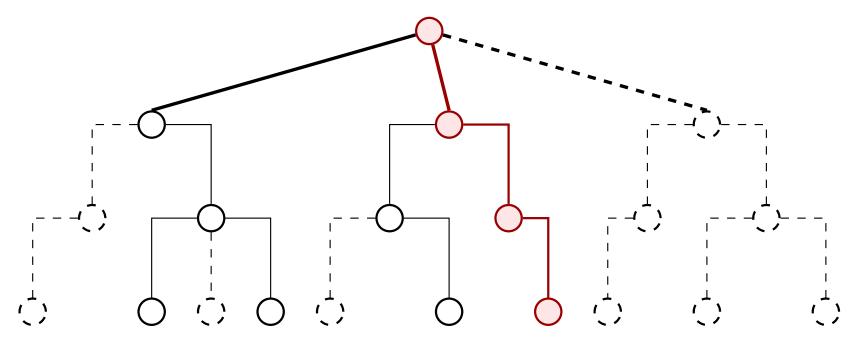
## **M**ETAREASONING

- A problem-solving agent can perform base-level actions from a known set  $\{A_i\}$ .
- Before committing to an action, the agent may perform a sequence of meta-level deliberation actions from a set  $\{S_j\}$ .
- At any given time there is a baselevel action  $A_{\alpha}$  that maximizes the agent's *expected utility*.
- The value of information  $VOI_j$  is the expected difference between the expected utilities of the new and the old selected base-level action after meta-level action  $S_j$  is taken.
- The agent selects a meta-level action that maximizes the VOI, or  $A_{\alpha}$  if no meta-level action has positive VOI.

## **A**CKNOWLEDGMENTS

- IMG4 Consortium under the MAGNET program of the Israeli Ministry of Trade and Industry
- Israel Science Foundation grant 305/09
- Lynne and William Frankel Center for Computer Sciences
- Paul Ivanier Center for Robotics Research and Production Management

# Monte-Carlo Sampling in Trees



- MCTS performs multiple *rollouts* to partially explore the search space.
- At the current root node, the sampling is aimed at finding the first move to perform: minimizing the simple regret is more appropriate at the root node.
- Deeper in the tree, minimizing cumulative regret results in a better estimate of the value of the state.
- An improvement over UCT can be achieved by combining different sampling schemes on the first step and during the rest of a rollout.

# MAIN RESULTS

## The SR+CR MCTS Scheme

- Selects an action at **the current root** suitable for minimizing the simple regret (FirstAction).
- **Deeper down,** selects actions according to UCB, that approximately minimizes the cumulative regret (NextAction).

```
Rollout (node, depth=1)
     if IsLEAF(node, depth)
 3
       return 0
 4
     else
       if depth=1 then action ← FirstAction(node)
 5
 6
       else action \leftarrow NextAction(node)
       next \leftarrow NextState(node, action)
 8
       reward \leftarrow Reward (node, action, next)
 9
                   + Rollout(next, depth+1)
10
      UpdateStats(node, action, reward)
      return reward
```

# **Sampling for Simple Regret**

- 1.  $\varepsilon$ -greedy sampling  $\left(\varepsilon = \frac{1}{2}\right)$ .
- 2. Modified version of **UCB** (optimized for *simple regret*).
- 3. **VOI-aware** sampling:

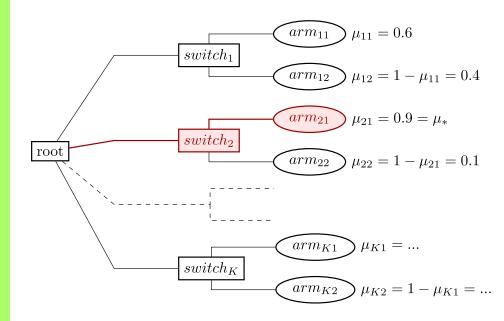
$$VOI_lpha pprox rac{\overline{X}_eta}{n_lpha+1} \exp\Bigl(-2(\overline{X}_lpha-\overline{X}_eta)^2 n_lpha\Bigr)$$

$$VOI_i pprox rac{1-\overline{X}_lpha}{n_i+1} \exp\Bigl(-2(\overline{X}_lpha-\overline{X}_i)^2 n_i\Bigr), \ i 
eq lpha$$

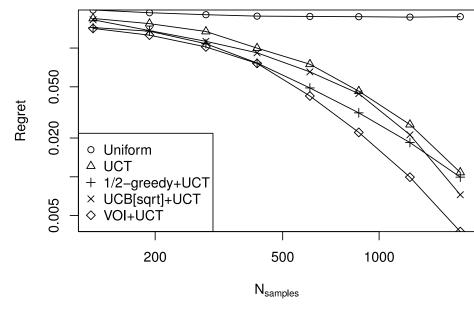
#### **EXPERIMENTS**

- SR+CR outperforms UCT.
- SR+UCT(*c*) is less dependent on tuning of the exploration factor *c*.

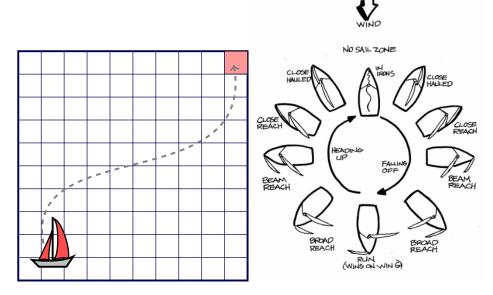
## **Random Trees**



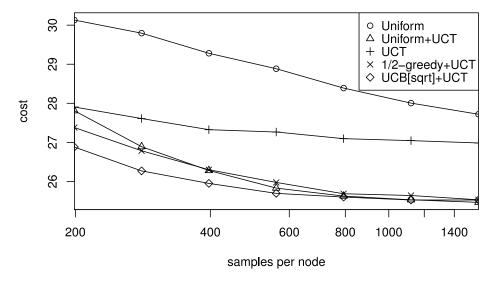
# Regret vs. number of samples:



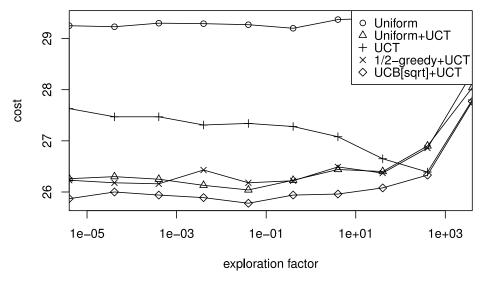
# **Sailing Domain**



## Path cost vs. number of samples:



## Path cost vs. exploration factor:



## **C**ONTRIBUTIONS

- Improved MCTS scheme SR+CR.
- SR+CR performs better than unmodified UCT.
- VOI-aware sampling for minimizing *simple regret*.

## FUTURE WORK

- Rational metareasoning in MCTS: theory and VOI estimates.
- Better sampling for non-root nodes.
- Application to Computer Go and other complex domains.

