

# AAAI 2012

*Ben Gurion University of the Negev,  
Beer Sheva, Israel*

## MULTI-ARMED BANDITS

- 

## UCB AND UCT

- $$b_i = \overline{X}_i + \sqrt{\frac{c \log(n)}{n_i}}$$

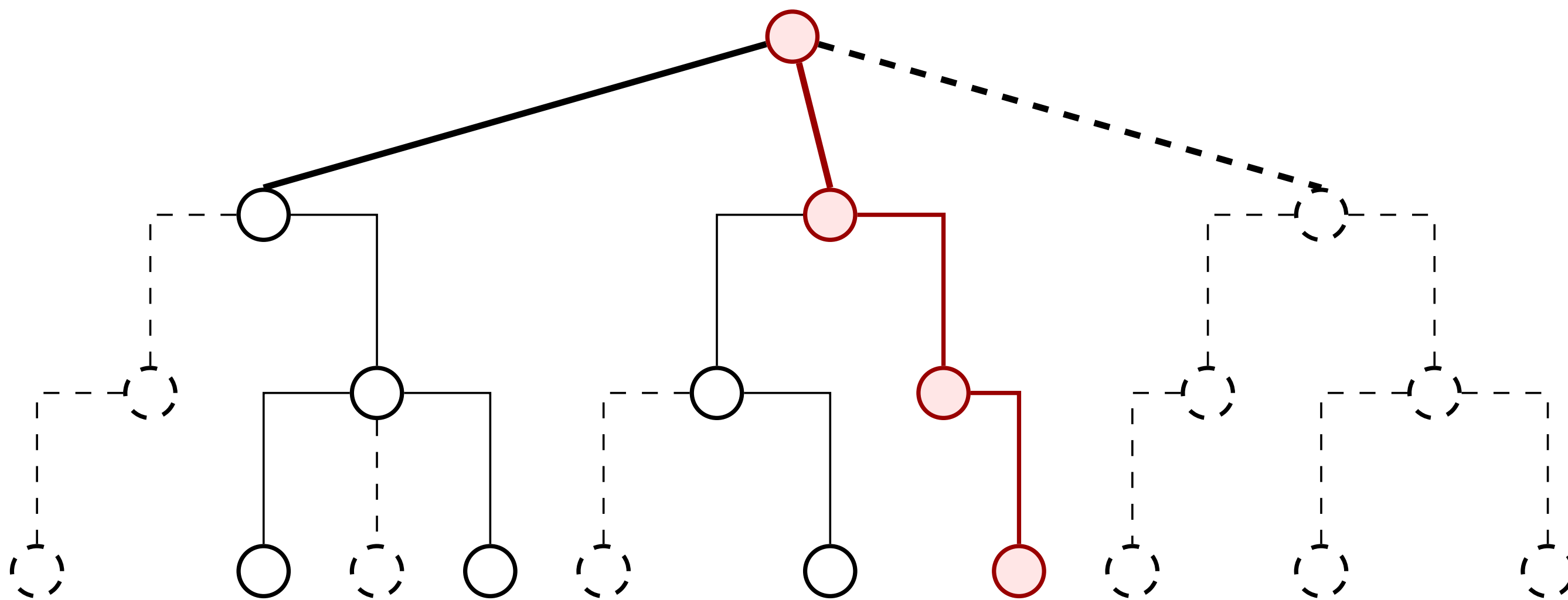
## METAREASONING

- 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.

## ACKNOWLEDGMENTS

- 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 more precise 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.
- Then selects actions according to UCB, that approximately minimizes the cumulative regret.

```

ROLLOUT(node, depth=1)
  if ISLEAF(node, depth)
    return 0
  else
    if depth=1 then action  $\leftarrow$  FIRSTACTION(node)
    else action  $\leftarrow$  NEXTACTION(node)
    next  $\leftarrow$  NEXTSTATE(node, action)
    reward  $\leftarrow$  REWARD(node, action, next)
               + ROLLOUT(next, depth+1)
    UPDATESTATS(node, action, reward)
  return reward

```

## Sampling for Simple Regret

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

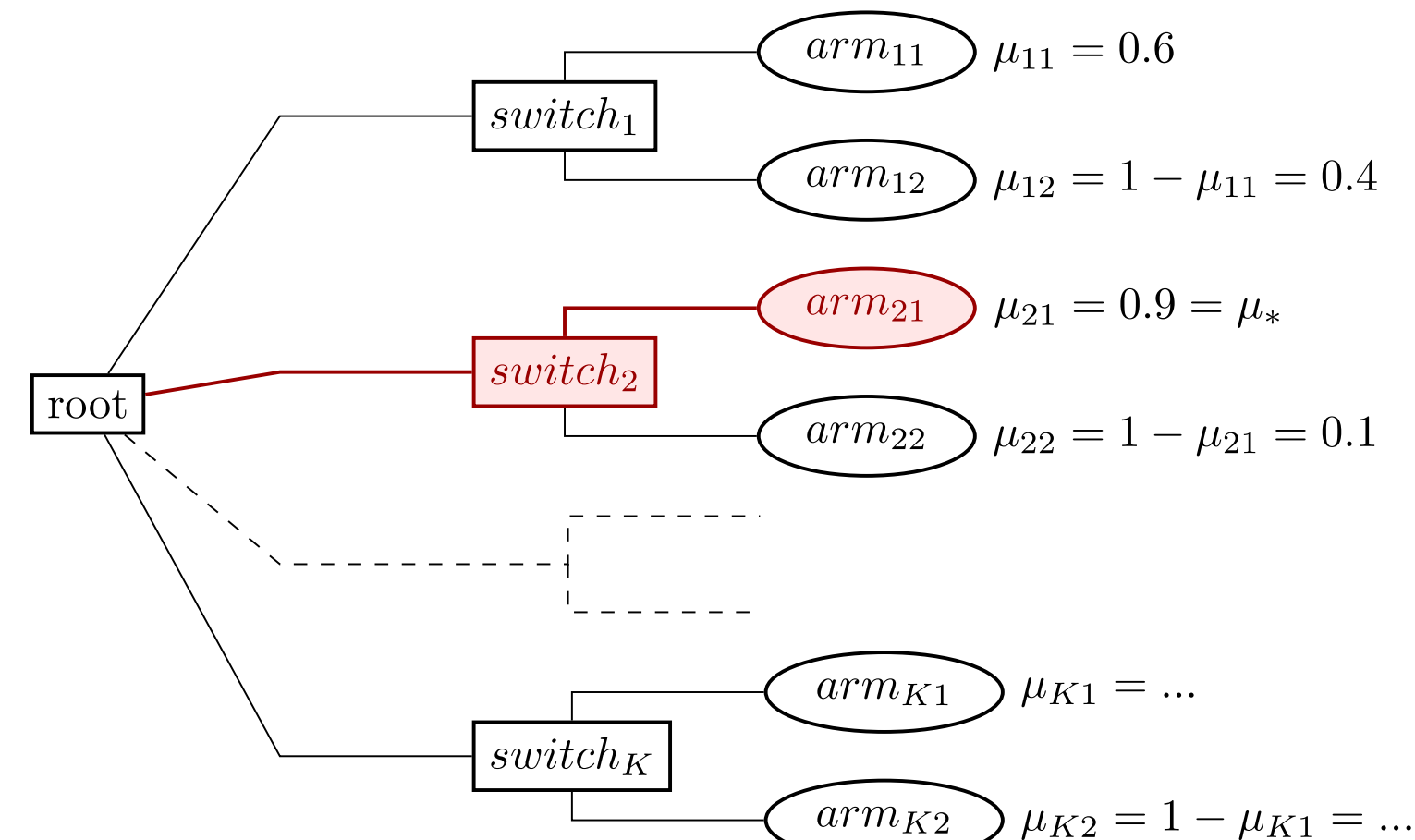
$$VOI_{\alpha} \approx \frac{\bar{X}_{\beta}}{n_{\alpha} + 1} \exp\left(-2(\bar{X}_{\alpha} - \bar{X}_{\beta})^2 n_{\alpha}\right)$$

$$VOI_i \approx \frac{1 - \bar{X}_\alpha}{n_i + 1} \exp\left(-2(\bar{X}_\alpha - \bar{X}_i)^2 n_i\right), i \neq \alpha$$

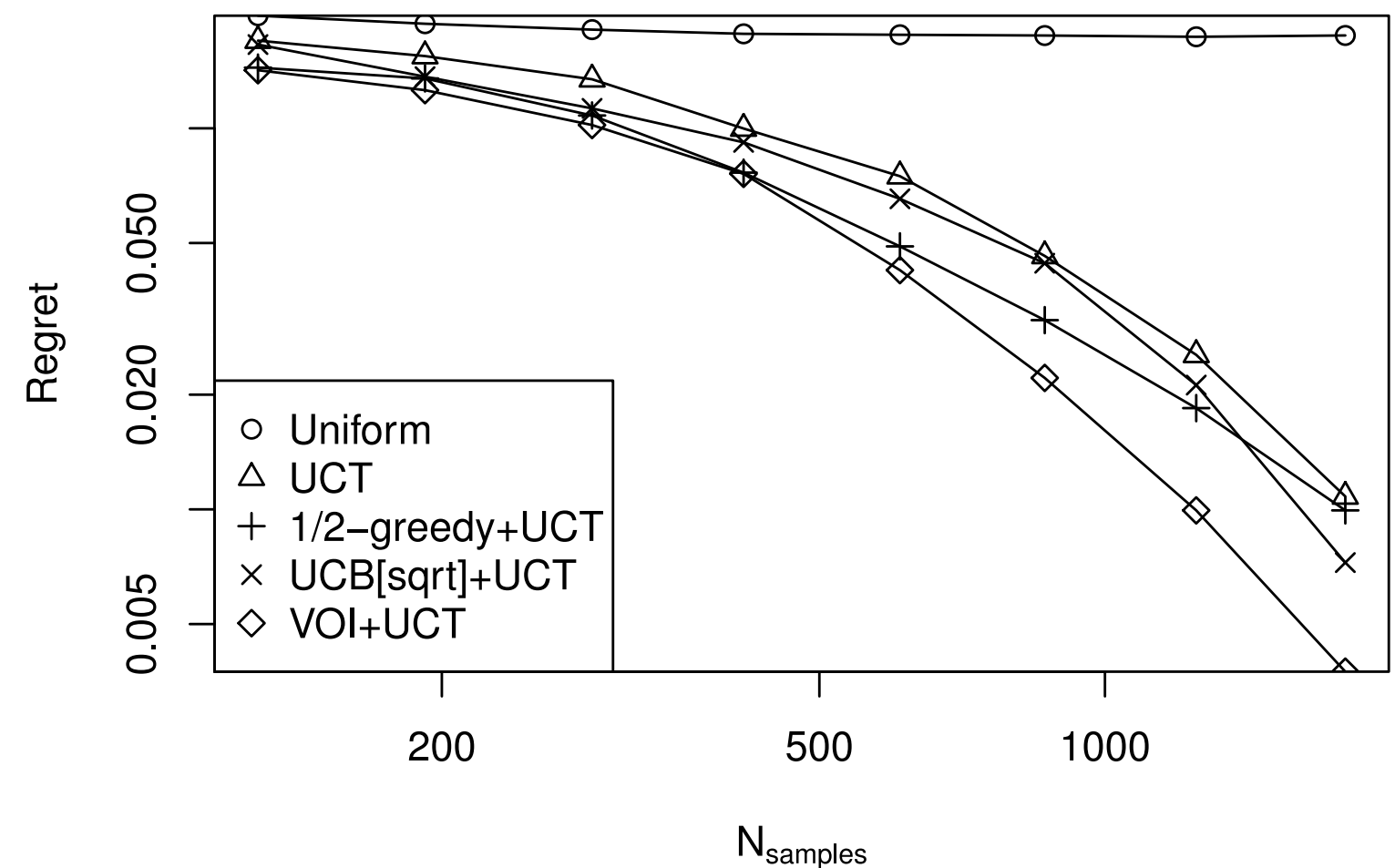
## EXPERIMENTS

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

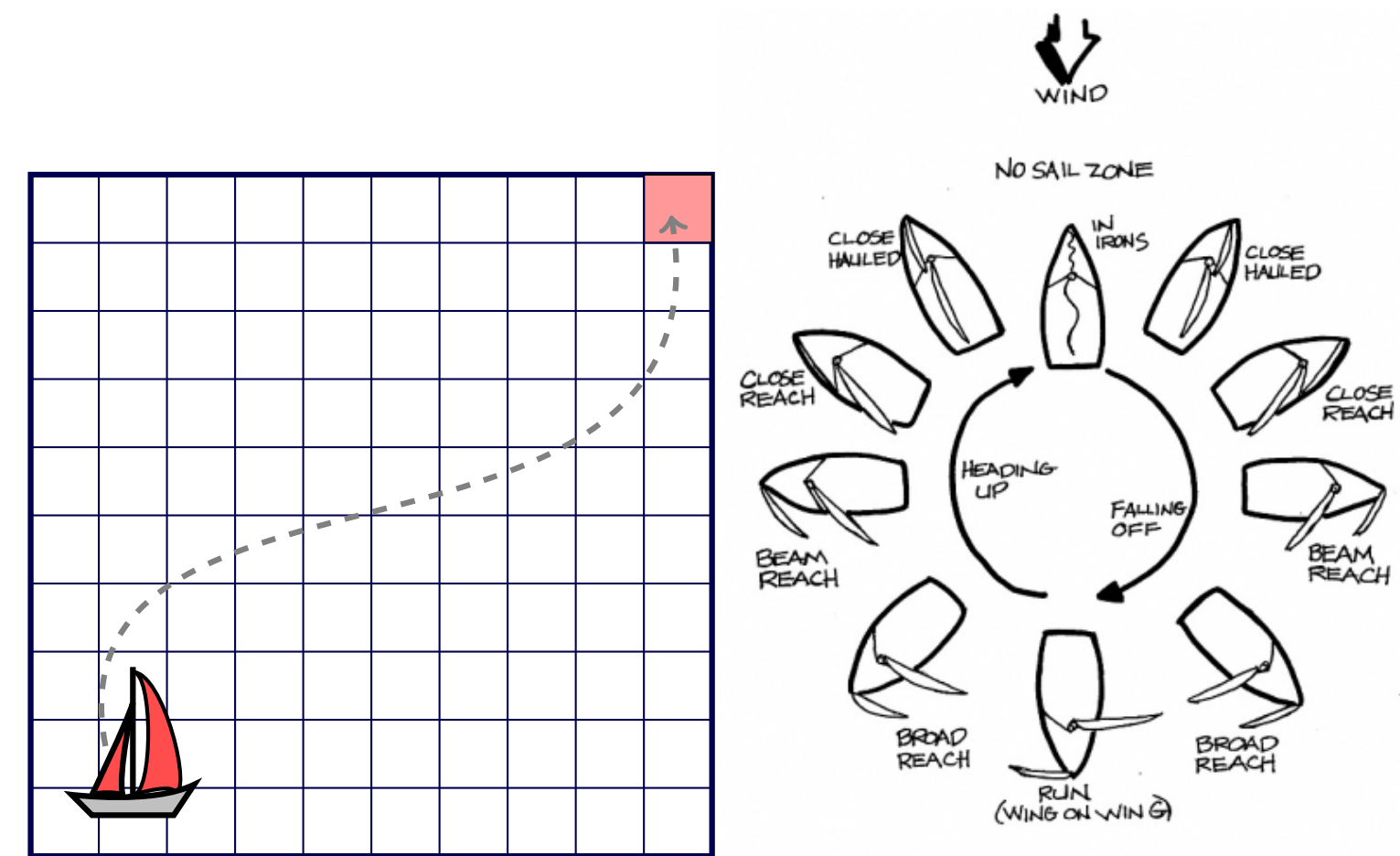
# Random Trees



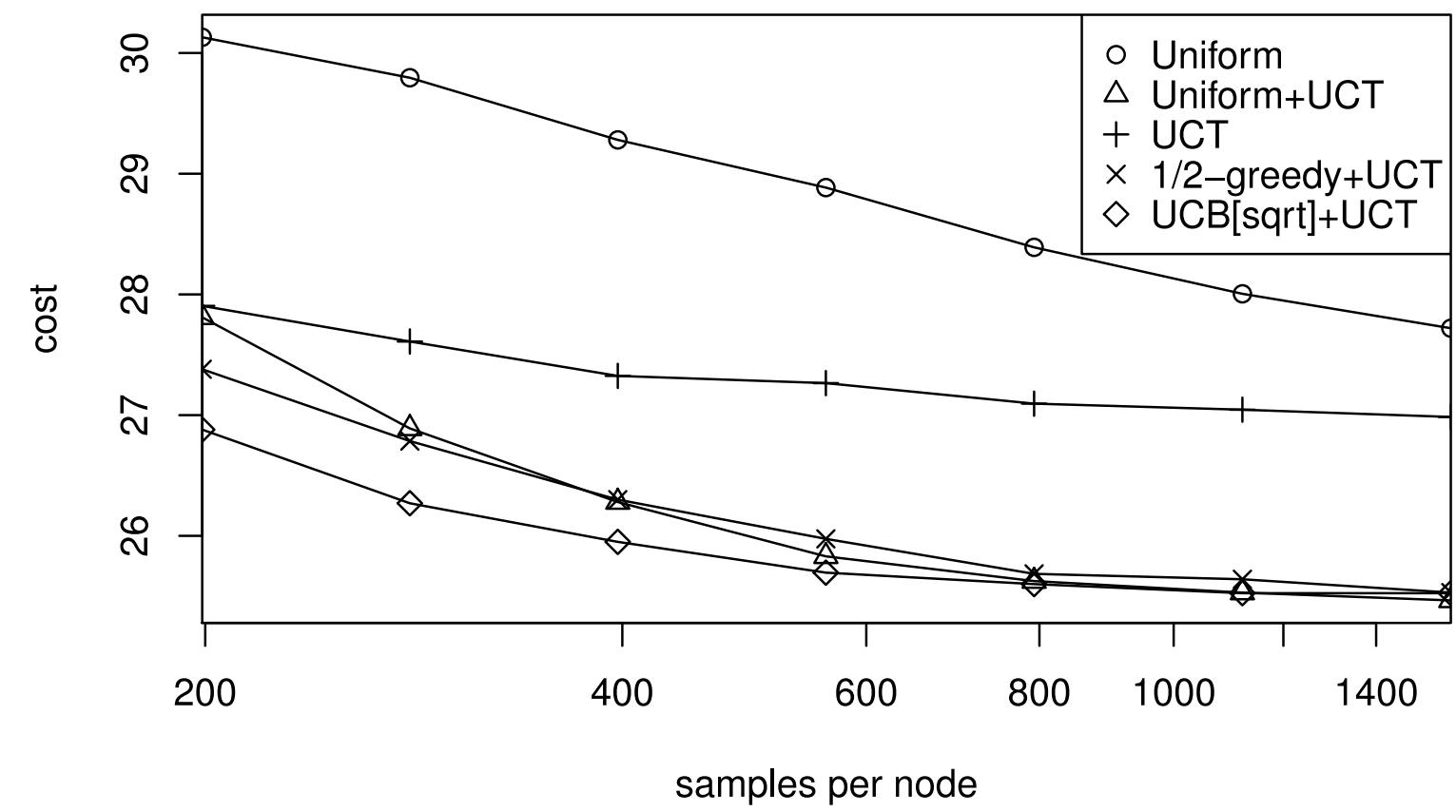
### Regret vs. number of samples:



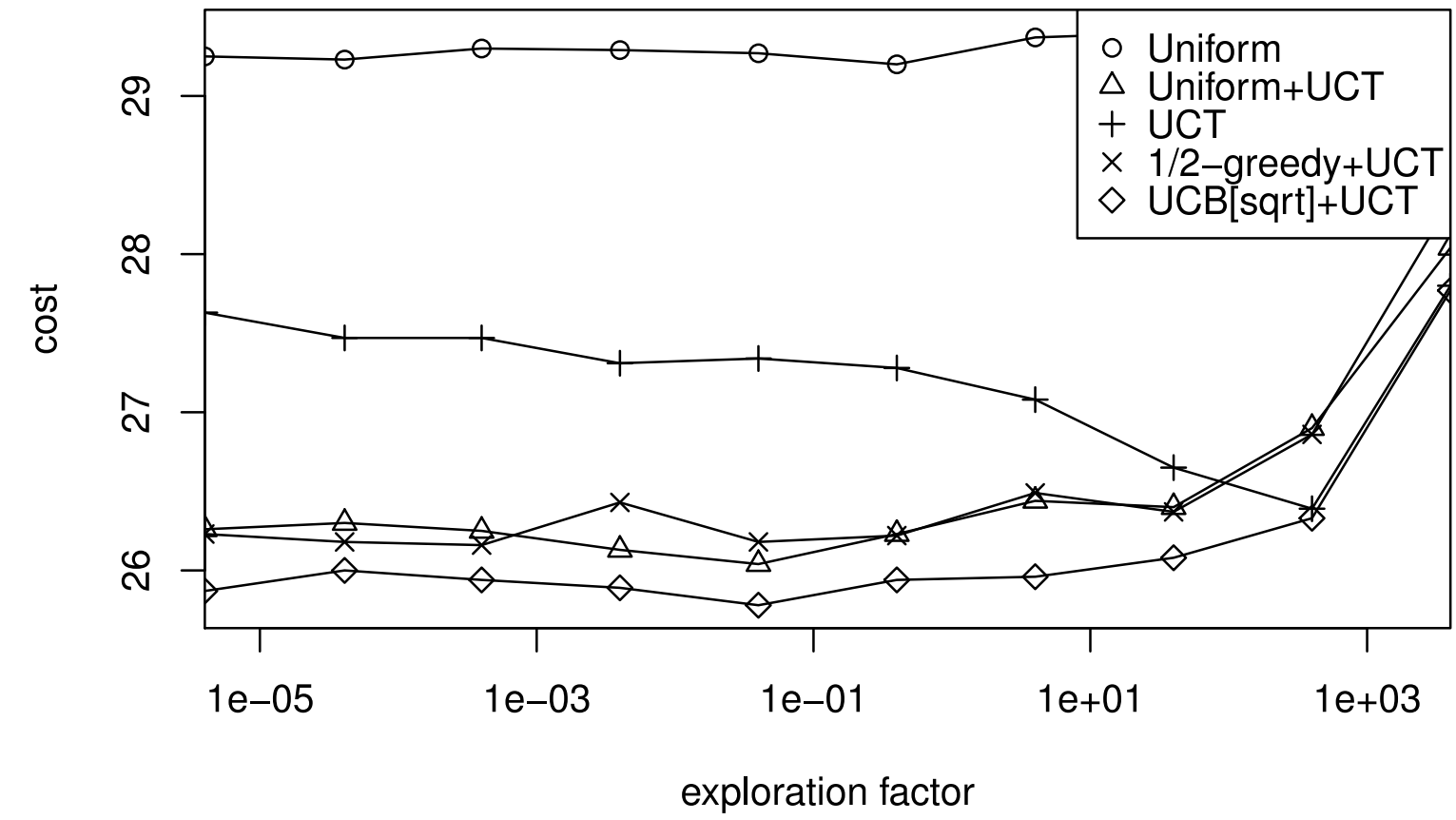
## Sailing Domain



### Path cost vs. number of samples:



### Path cost vs. exploration factor:



## CONTRIBUTIONS

- 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.

