

MONTE-CARLO (TREE) SEARCH

Stefan Edelkamp

PUI - CTU



STEFAN EDELKAMP

**AI CENTER
FEE CTU**

MONTE-CARLO TREE SEARCH

- General randomized heuristic search technique
- Less problem-specific knowledge to be added

Protagonists

- UCB/UCT, best in many games like Go, Amazons, GGP, etc. (Auer et al. 2002, Kocsis & Szepesvari, 2006, Coulom 2006)
- Nested-Monte-Carlo Tree Search (Cazenave, 2009)
- **Nested Rollout with Policy Adaptation (Rosin, 2011)**

Input: iteration width (exploit), nestedness level (explore)

Policy: (city-to-city) Mapping $N \times N \rightarrow IR$ to be learnt

MONTE-CARLO SEARCH

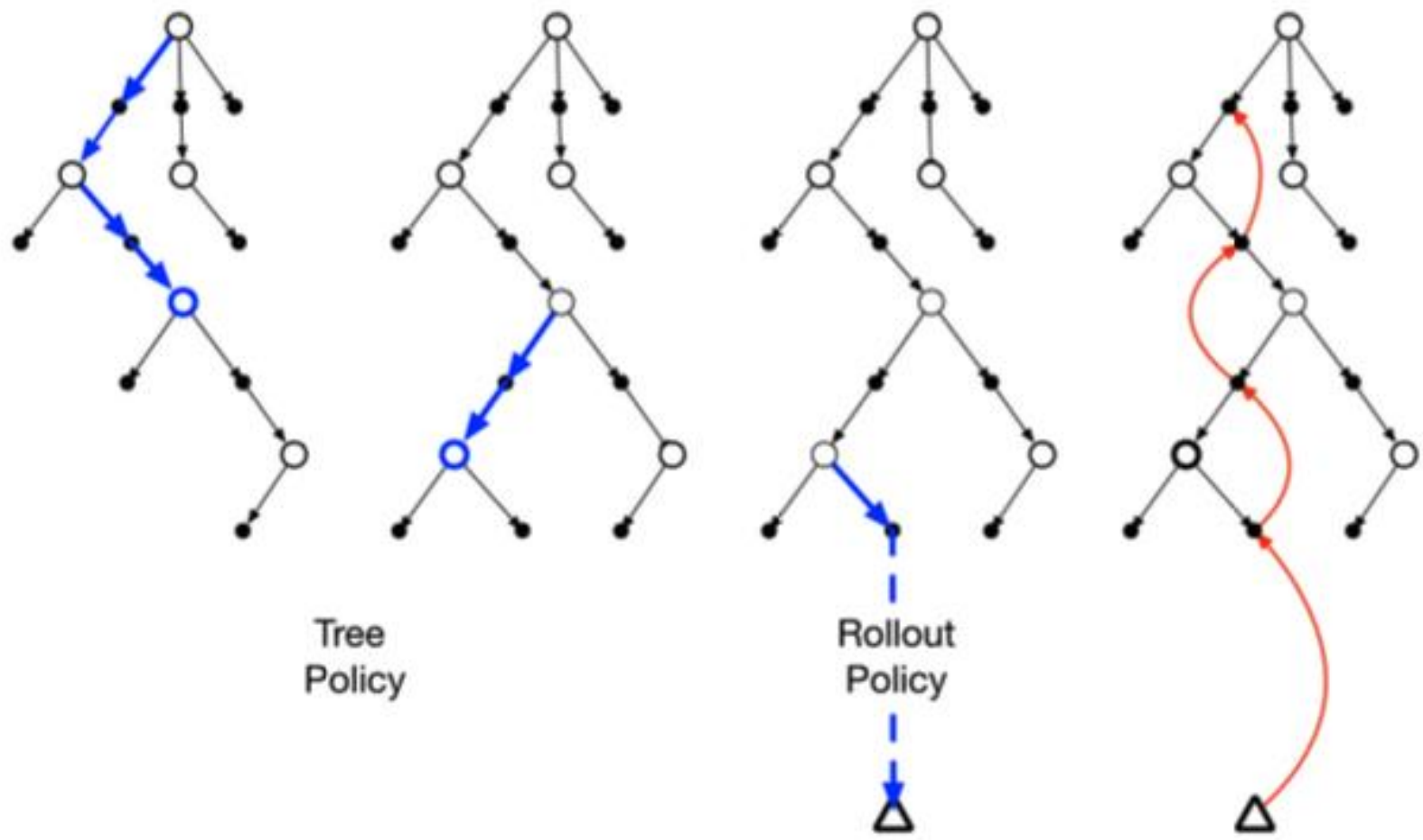
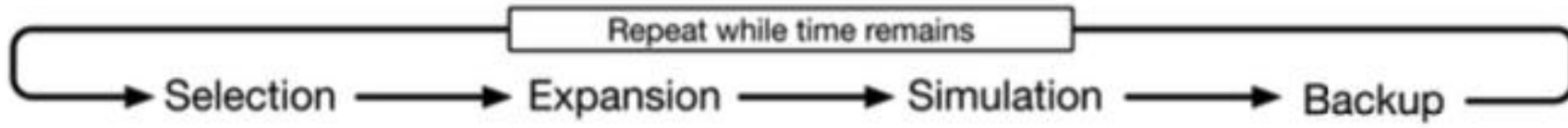
Consider set of random walks from initial to terminal state

Set average reward

No further steering, no further knowledge

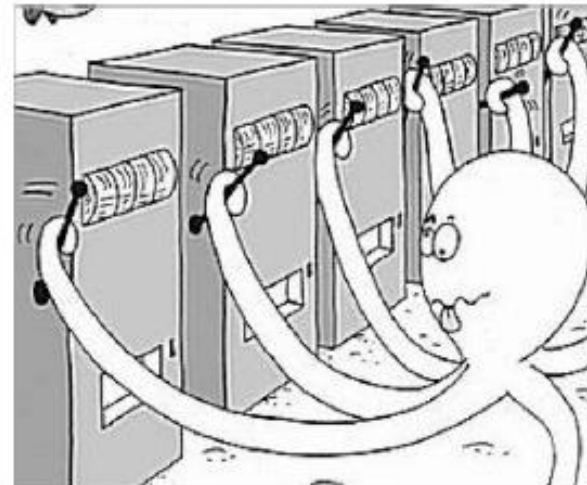
At end: Choose move with best evaluation





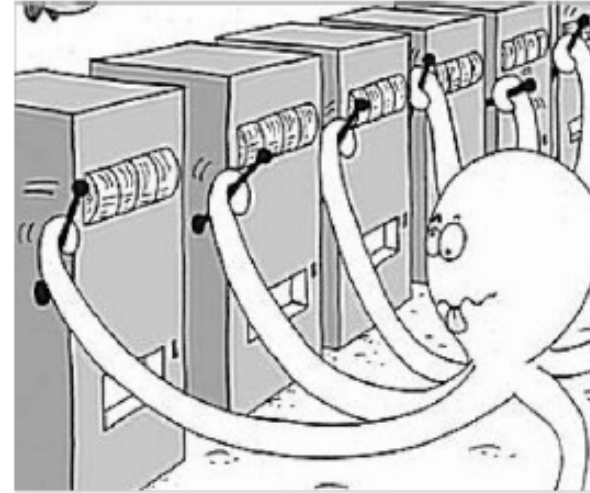
Multi-Arm Bandit Problem

- Statistical model of sequential experiments
 - Name comes from a traditional slot machine (one-armed bandit)
- Multiple actions a_1, a_2, \dots, a_n
 - Each a_i provides a reward from an unknown (but stationary) probability distribution p_i
 - Objective: maximize expected utility of a sequence of actions
- Exploitation vs exploration dilemma:
 - **Exploitation**: choose action that has given you high rewards in the past
 - **Exploration**: choose action that you don't know much about, in hopes that it might produce a higher reward



UCB (Upper Confidence Bound) Algorithm

- Assume all rewards are between 0 and 1
 - If they aren't, normalize them
- For each action a_i , let
 - r_i = average reward you've gotten from a_i
 - t_i = number of times you've tried a_i
 - $t = \sum_i t_i$



loop

if there are one or more actions that you haven't tried

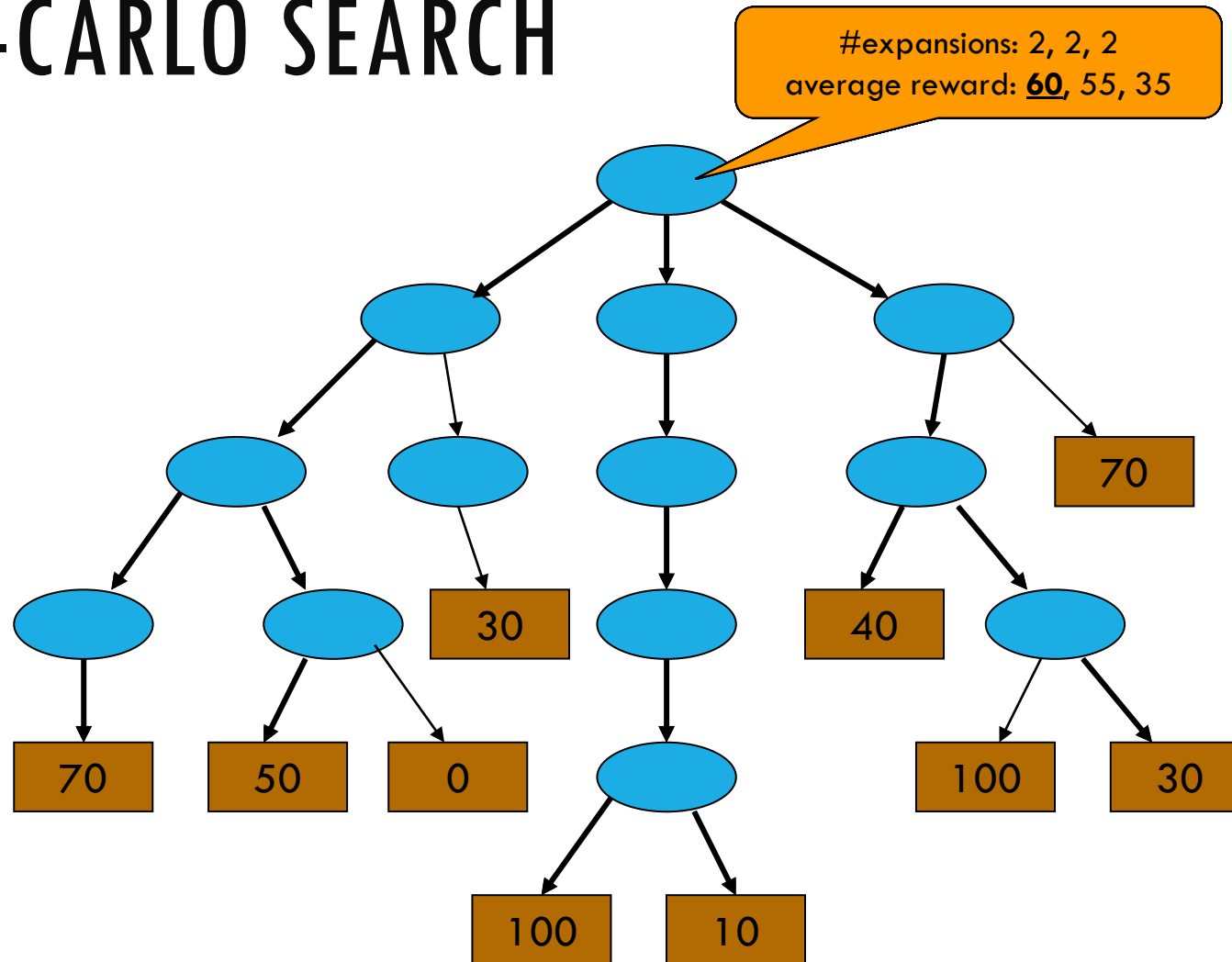
then choose an untried action a_i at random

else choose an action a_i that has the highest value of $r_i + \sqrt{2(\ln t)/t_i}$

perform a_i

update r_i , t_i , t

MONTE-CARLO SEARCH



MONTE-CARLO SEARCH

Advantage:

- Easy to implement
- Better than random
- Small memory requirement

Disadvantage:

- Repeated expansion of same states
 - Slowing down search
- No guidance
- Results potentially bad
- Information lost in further runs



MONTE-CARLO SEARCH WITH MEMORY

first + last problem resolved by using an explicit tree

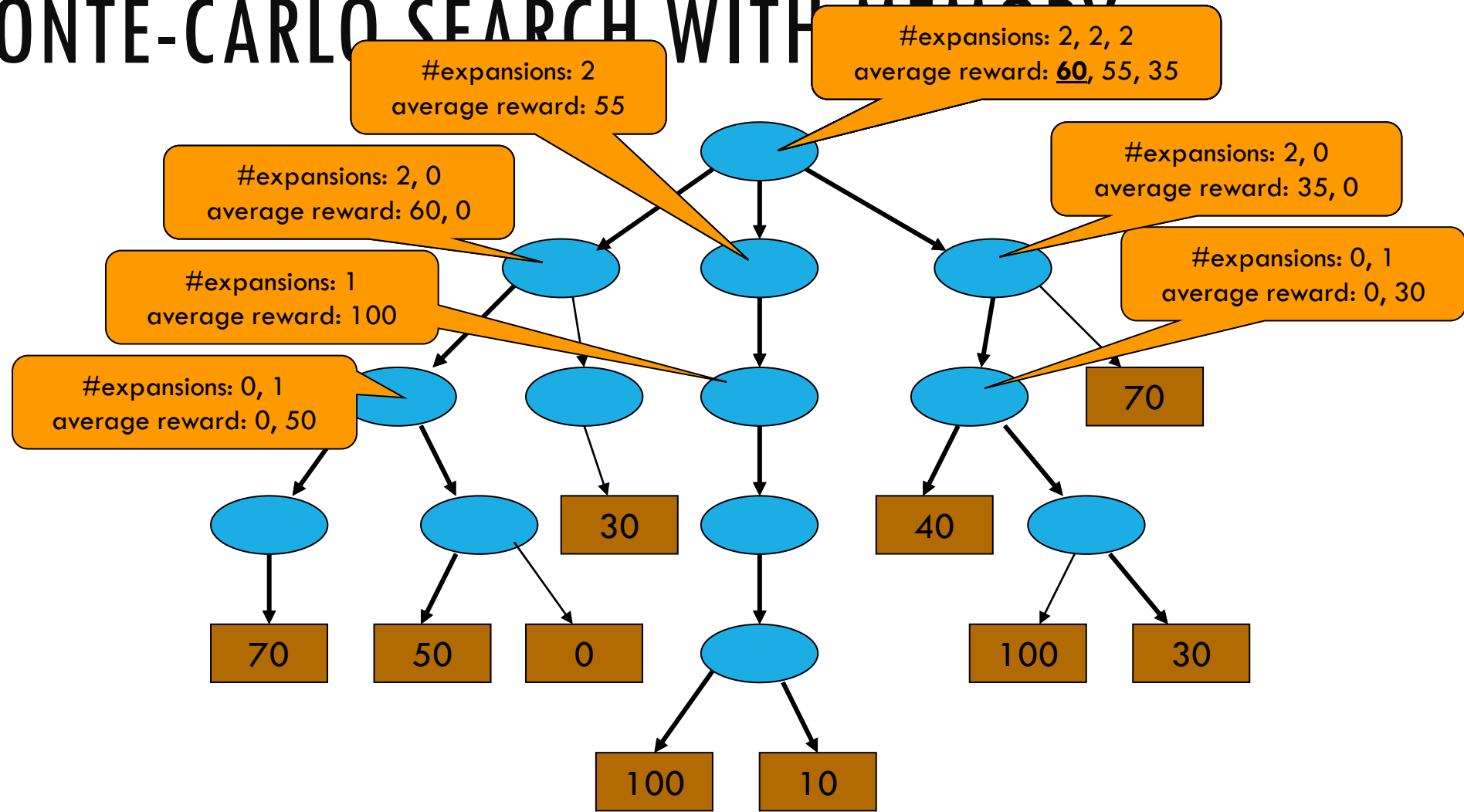
- Encapsulation of memory
- Instead forgetting about everything after Monte-Carlo run, insert node in tree
- Only on node at a time because of memory requirement

Search in the tree remains random

- Expansions faster if successor stored in tree sneller
- Update for all nodes in the tree
 - Stored node used later on in the search



MONTE-CARLO SEARCH WITH MEMORY



MONTE-CARLO SEARCH WITH MEMORY

Advantages:

- Storage → less expansions
- Information of successors can be used in upcoming moves

Disadvantages:

- More memory requirements
- Unguided search
- Results often not good



UCT [KOCSIS & SZEPESVÁRI, 2006]

“Upper Confidence Bounds applied to Trees”

Additional Information

- In the Tree
 - If ≥ 1 unexpanded move, chose a random one
 - Choose successor, maximizing UCT value
- $$Q(s, m) + C \sqrt{\frac{\ln(N(s))}{N(s, m)}}$$
- $Q(s, m)$ average reward of move m in state s
 - C : constant
 - $N(s)$: number visits of state s
 - $N(s, m)$: number visits of state s with chosen move m
 - If leaf found
 - Normal Monte-Carlo Suche d
 - Backpropagate terminal evaluation up the tree.



UCT Algorithm

- Recursive UCB computation to compute $Q(s,a)$
- Anytime algorithm: call repeatedly until time runs out
 - Then choose action $\operatorname{argmin}_a Q(s,a)$

UCT(s, h)

if $s \in S_g$ then return 0

if $h = 0$ then return $V_0(s)$

if $s \notin \text{Envelope}$ then do

 add s to *Envelope*

$n(s) \leftarrow 0$

 for all $a \in \text{Applicable}(s)$ do

$Q(s,a) \leftarrow 0$; $n(s,a) \leftarrow 0$

$\text{Untried} \leftarrow \{a \in \text{Applicable}(s) \mid n(s,a) = 0\}$

 if $\text{Untried} \neq \emptyset$ then $\tilde{a} \leftarrow \text{Choose}(\text{Untried})$

 else $\tilde{a} \leftarrow \operatorname{argmin}_{a \in \text{Applicable}(s)} \{Q(s,a) - C \times [\log(n(s))/n(s,a)]^{\frac{1}{2}}\}$

$s' \leftarrow \text{Sample}(\Sigma, s, \tilde{a})$

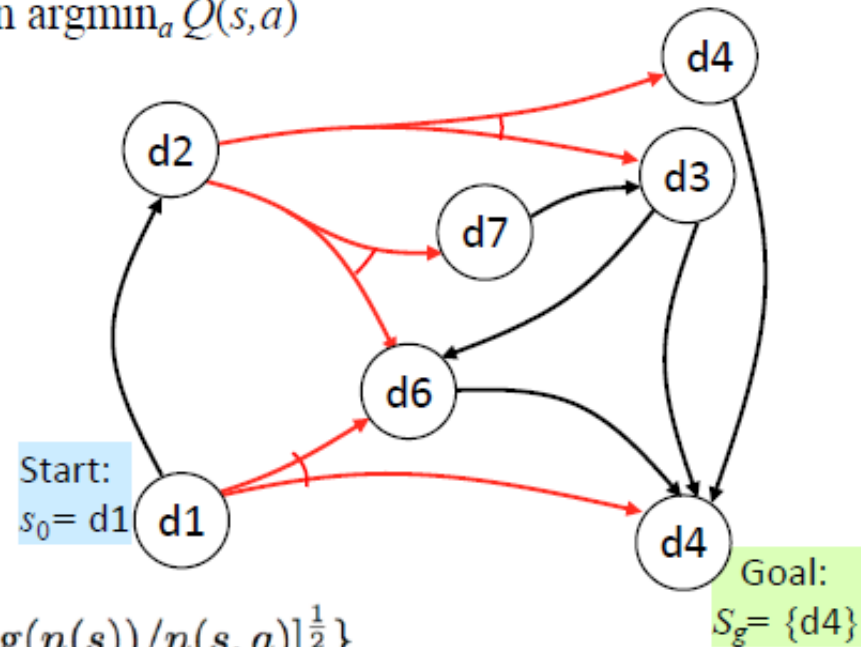
$\text{cost-rollout} \leftarrow \text{cost}(s, \tilde{a}) + \text{UCT}(s', h - 1)$

$Q(s, \tilde{a}) \leftarrow [n(s, \tilde{a}) \times Q(s, \tilde{a}) + \text{cost-rollout}] / (1 + n(s, \tilde{a}))$

$n(s) \leftarrow n(s) + 1$

$n(s, \tilde{a}) \leftarrow n(s, \tilde{a}) + 1$

 return cost-rollout

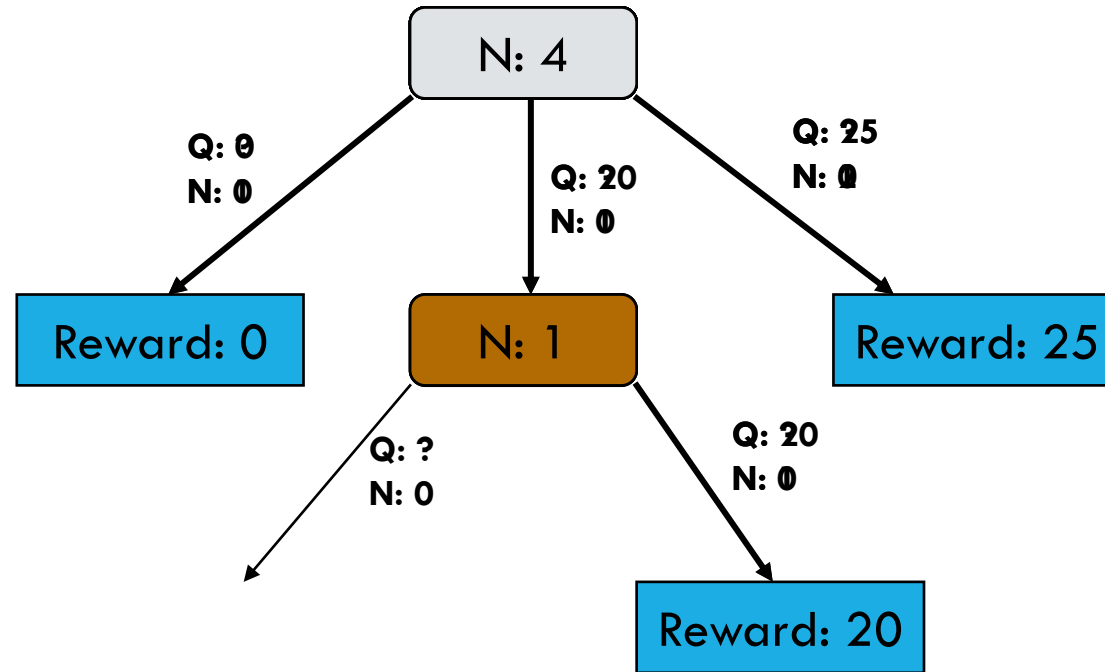


UCT in Two-Player Games

- Generate Monte Carlo rollouts using a modified version of UCT
- Main differences:
 - Instead of choosing actions that minimize accumulated cost, choose actions that maximize payoff at the end of the game
 - UCT for player 1 recursively calls UCT for player 2
 - Choose opponent's action
 - UCT for player 2 recursively calls UCT for player 1
- This produced the first computer programs to play go well
 - \approx 2008–2012
- Monte Carlo rollout techniques similar to UCT were used to train AlphaGo



UCT



$$Q(s, m) + C \sqrt{\frac{\ln(N(s))}{N(s, m)}}$$

for **C = 10**:
 Move 1: 10
 Move 2: 30
 Move 3: 33

for **C = 100**:
 Move 1: 108
 Move 2: 138
 Move 3: 136

for **C = 1,000**:
 Move 1: 1078
 Move 2: 1098
 Move 3: 8983



IMPROVEMENTS UCT

[Finnsson & Björnsson, 2010]

- Move-Average Sampling Technique (MAST)
- Tree-Only MAST (TO-MAST)
- Predicate-Average Sampling Technique (PAST)
- Features-to-Action Sampling Technique (FAST)
- Rapid Action Value Estimation (RAVE)

MAST

Idea:

- In every UCT-run, refine knowledge about all the moves
- Use knowledge to improve knowledge to steer search outside the tree

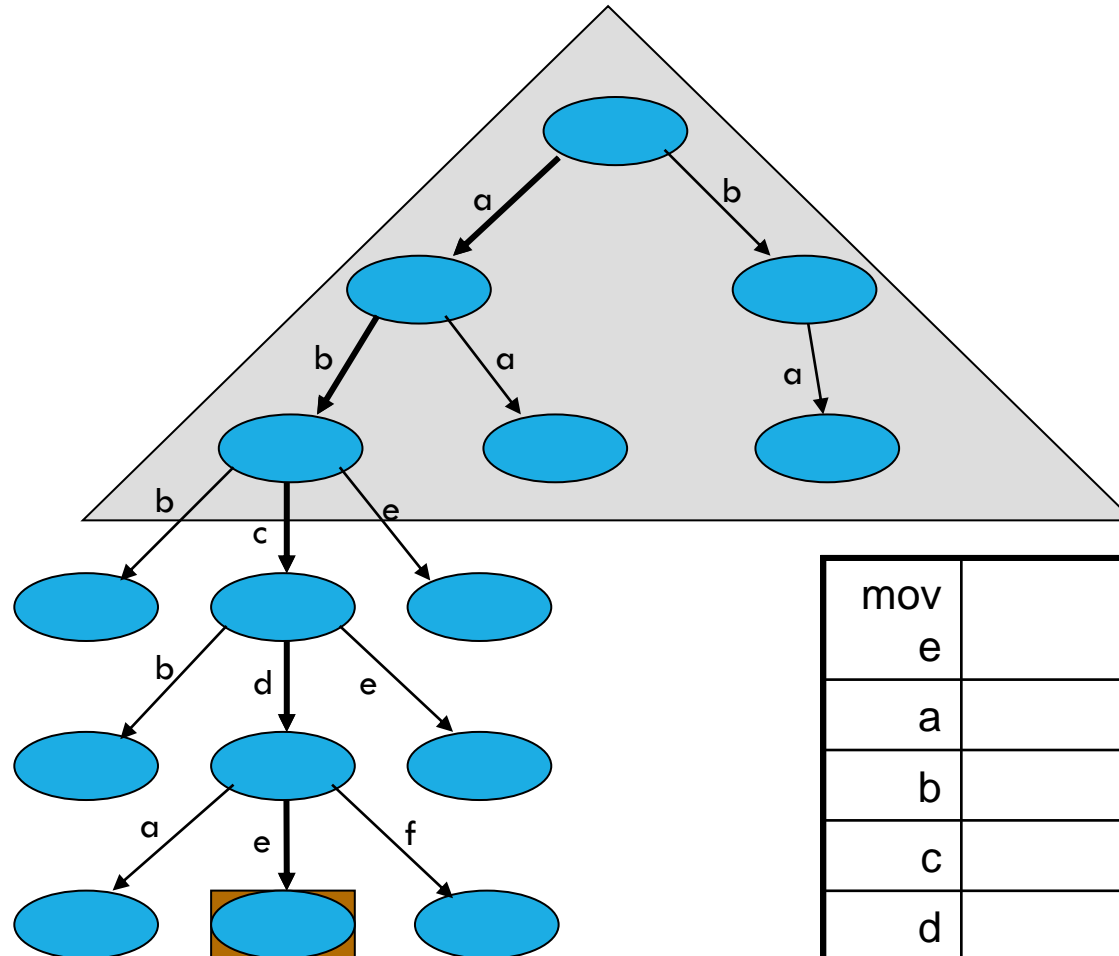
Maintain average reward for every move

- Independent of state

After UCT-run: update average reward of all the moves

- Moves, chosen frequently (independent of state) receive better score
- Hope: Moves are often good if available
 - e.g. placement of a stone in the corner of Reversi
 - e.g. Taking an opponent stone before the own base in Breakthrough

MOVE-AVERAGE SAMPLING TECHNIQUE



mov	average	#Visits
e		
a	25 41,67	2 3
b	60 63,75	3 4
c	15 35	2 3
d	85 83	4 5
e	10 20,83	5 6
f	70	1

MOVE-AVERAGE SAMPLING TECHNIQUE

Within Monte-Carlo run (outside UCT tree)

Choose move according to Gibbs sampling

$$P(m) = \frac{e^{Q_h(m)/\tau}}{\sum_{b=1}^n e^{Q_h(b)/\tau}}$$

with

m: chosen move

$Q_h(m)$: average reward of m

τ : constant to tune (large value: closer to uniform distribution)



TREE-ONLY MAST

(initial) Monte-Carlo runs random

may influence reward of actions negatively

Idea: use only results in UCT-tree

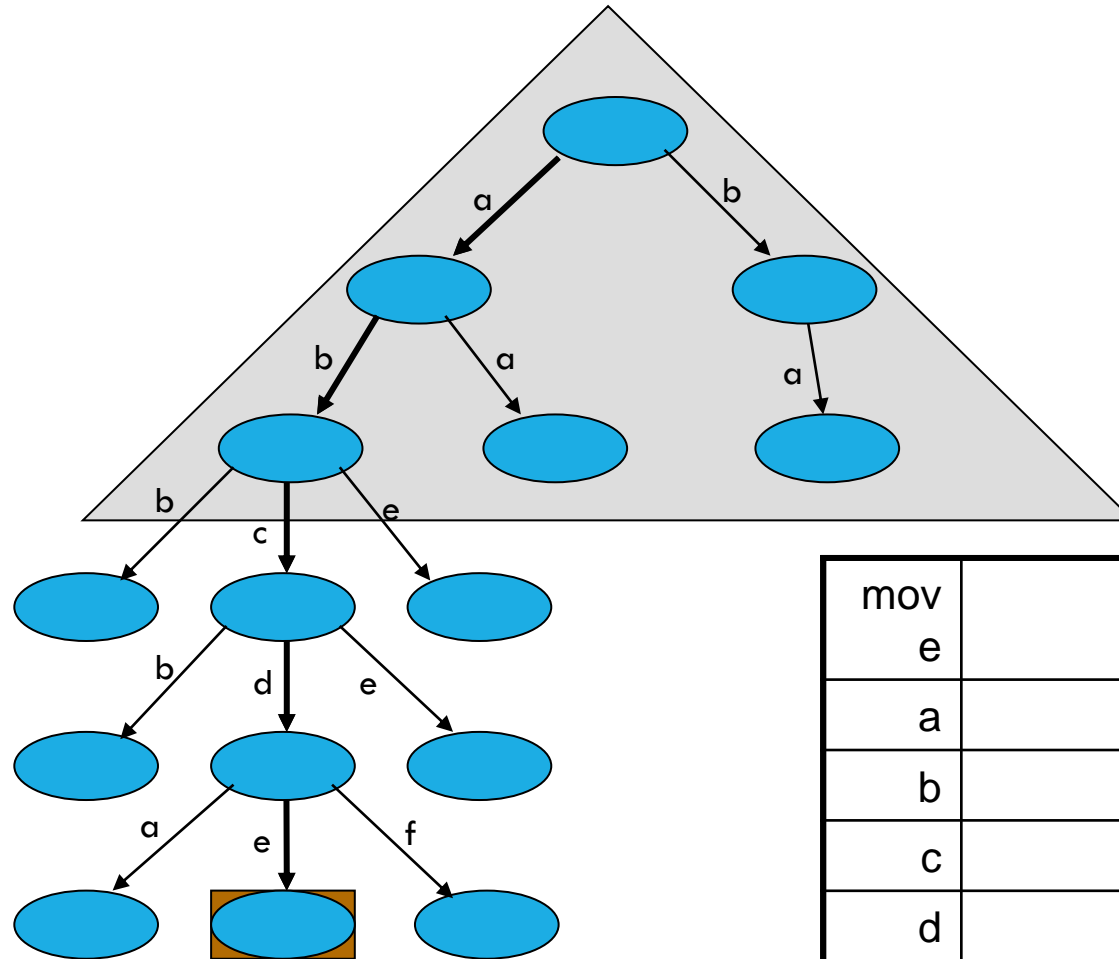
precise: apply UCT (with MAST) as before

but: update only moves, chosen within UCT tree

(ignore Monte-Carlo run)

In Monte-Carlo run: choose move according to MAST distribution

TO-MOVE-AVERAGE SAMPLING TECHNIQUE



mov	average	#visits
e		
a	25 41,67	2 3
b	60 63,75	3 4
c	15	2
d	85	4
e	10	5
f	70	1

RAPID ACTION VALUE ESTIMATION

First used in UCT Go AI [Gelly & Silver, 2007]

- known as all-moves-as-first geuristik

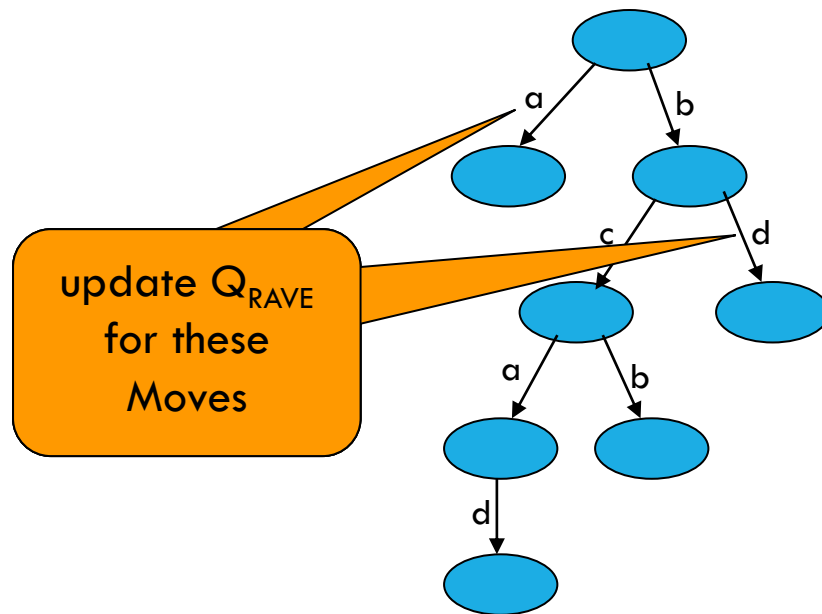
Accelerates learning withing UCT tree

Use later chosen moves, to derive more move samplings for the same but not chosen moves

RAPID ACTION VALUE ESTIMATION

Within UCT Tree

- update $Q(s, m)$, if move m chosen in state s (as before)
- update $Q_{RAVE}(s, m')$ for state s for non-chosen move m' , if m' chosen later on in UCT tree



RAPID ACTION VALUE ESTIMATION

Issue: drift of average

- Good if less samples are available
- Choose only in case of higher variance in $Q(s, m)$
- Later on, $Q(s, m)$ is more reliable
- Then ignore $Q_{\text{RAVE}}(s, m)$



- Store RAVE values on top
- Choose weighted combination: $\beta(s) * Q_{\text{RAVE}}(s, m) + (1 - \beta(s)) * Q(s, m)$ in UCT choice
- with

$$\beta(s) = \sqrt{\frac{k}{3N(s) + k}}$$

- k: parameter (deciding after how many samples weighted is equal)
- $N(s)$: number of visits of s



COMPARISON

[Finnsson & Björnsson, 2010] Comparisons UCT

Game	MAST win %	TO-MAST win %	PAST win %	RAVE win %	FAST win %
Breakthrough	90.00 (\pm 3.40)	85.33 (\pm 4.01)	85.00 (\pm 4.05)	63.33 (\pm 5.46)	81.67 (\pm 4.39)
Checkers	56.00 (\pm 5.37)	82.17 (\pm 4.15)	57.50 (\pm 5.36)	82.00 (\pm 4.08)	50.33 (\pm 5.36)
Othello	60.83 (\pm 5.46)	50.17 (\pm 5.56)	67.50 (\pm 5.24)	70.17 (\pm 5.11)	70.83 (\pm 5.10)
Skirmish	41.33 (\pm 5.18)	48.00 (\pm 5.29)	42.33 (\pm 5.16)	46.33 (\pm 5.30)	96.33 (\pm 1.86)

Comparison MAST

Game	TO-MAST win %	PAST win %	RAVE win %	FAST win %
Breakthrough	52.33 (\pm 5.66)	45.67 (\pm 5.65)	20.33 (\pm 4.56)	39.67 (\pm 5.55)
Checkers	82.00 (\pm 4.18)	55.83 (\pm 5.35)	78.17 (\pm 4.36)	46.17 (\pm 5.33)
Othello	40.67 (\pm 5.47)	49.17 (\pm 5.60)	58.17 (\pm 5.49)	56.83 (\pm 5.55)
Skirmish	56.00 (\pm 5.31)	43.33 (\pm 5.26)	59.83 (\pm 5.15)	97.00 (\pm 1.70)

RESULTS

[Finnsson & Björnsson, 2010]

UCT vs RAVE+MAST (RM) vs RAVE+FAST (RF)

Game	RM win %	RF win %
Breakthrough	89.00 (\pm 4.35)	76.50 (\pm 5.89)
Checkers	84.50 (\pm 4.78)	77.00 (\pm 5.37)
Othello	79.75 (\pm 5.52)	81.00 (\pm 5.32)
Skirmish	45.00 (\pm 6.55)	96.00 (\pm 2.34)

MAST vs RAVE+MAST (RM) vs. RAVE+FAST (RF)

Game	RM win %	RF win %
Breakthrough	50.50 (\pm 6.95)	38.50 (\pm 6.76)
Checkers	83.50 (\pm 4.87)	74.00 (\pm 5.81)
Othello	73.75 (\pm 6.01)	66.00 (\pm 6.43)
Skirmish	53.00 (\pm 6.47)	97.00 (\pm 2.04)

NESTED ROLLOUT POLICY ADAPTATION

Rosin 2011 IJCAI – Best Paper, Morpion Solitaire with new Record

MCS tree based on Complete Rollout and Recursive Search

Not really MCTS, No Search Tree.

In Each Level a Policy is Maintained, Updated and Refreshed

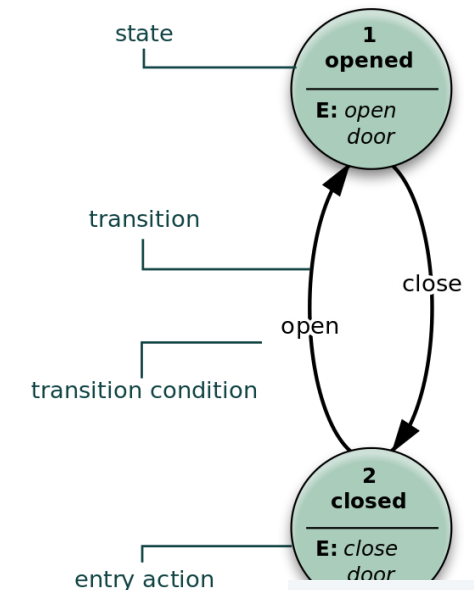
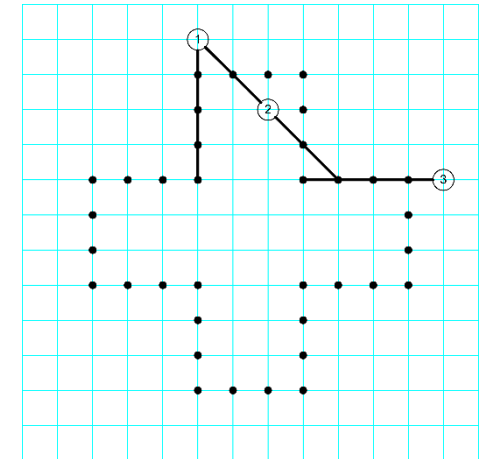
Updating Policy based on better Solutions Coming in from below

Policy in Turn Influences the Rollouts

Parameters: Level of Recursion, and Iteration Width

Effective for TSPTW and many other Approaches

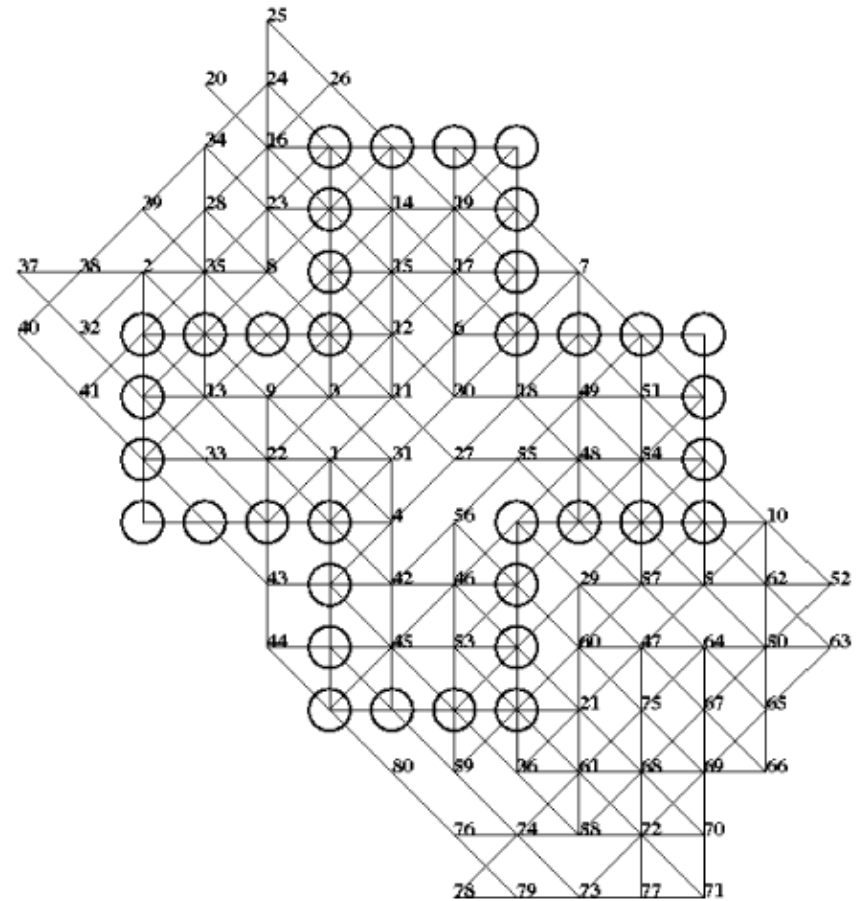
Refinements: Beam / Diversity / Generalization



NEW RECORD

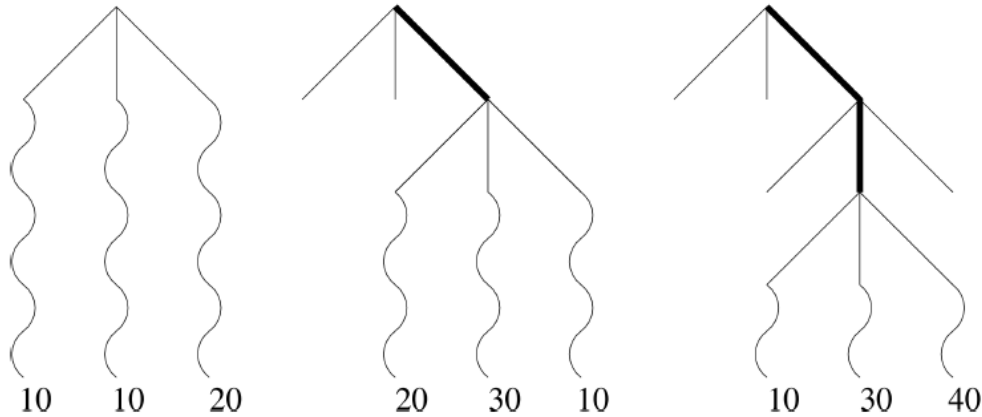
Morpion Solitaire

- 80 moves :



NESTED MCS

Nested Monte-Carlo Search



Nested Monte-Carlo Search

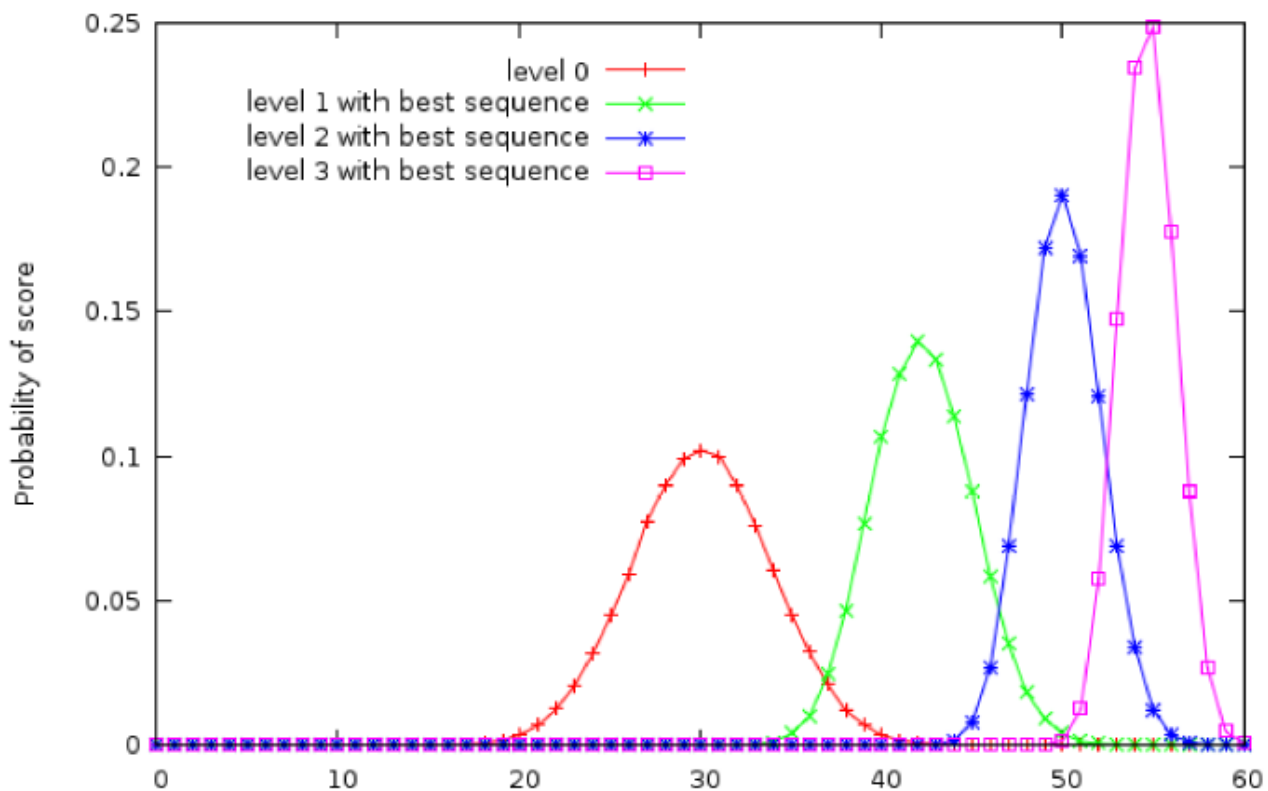
- Play random games at level 0
- For each move at level $n+1$, play the move then play a game at level n
- Choose to play the move with the greatest associated score
- Important : memorize and follow the best sequence found at each level

NMCS PSEUDO-CODE

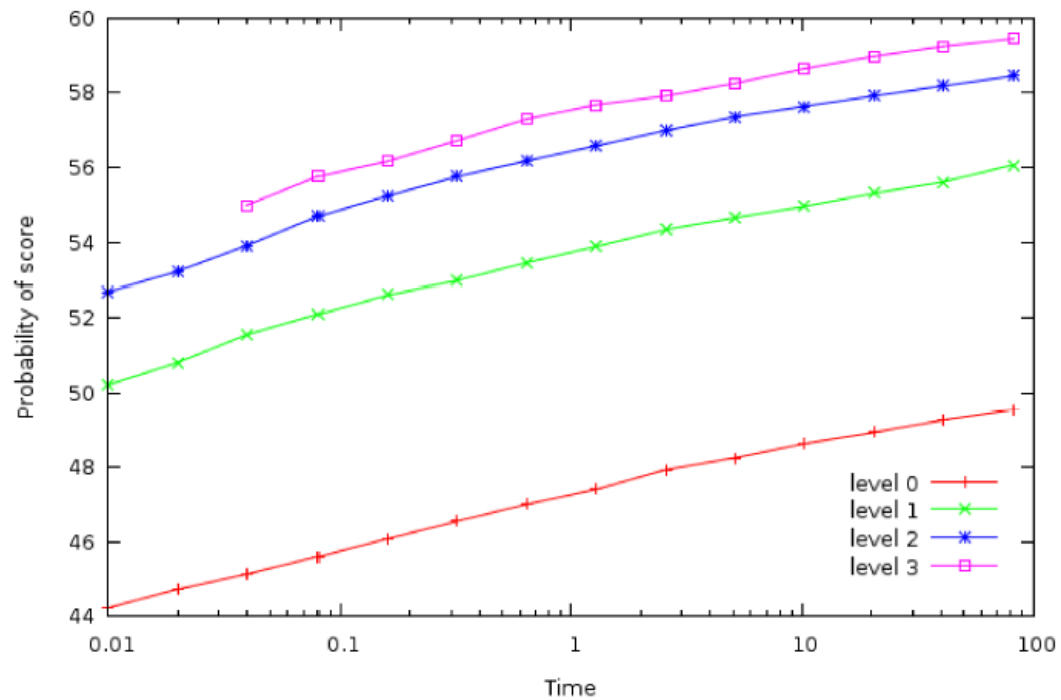
Algorithm 4 The NMCS algorithm.

```
NMCS (state, level)
  if level == 0 then
    return playout (state, uniform)
  end if
  BestSequenceOfLevel  $\leftarrow \emptyset$ 
  while state is not terminal do
    for m in possible moves for state do
      s  $\leftarrow$  play (state, m)
      NMCS (s, level - 1)
      update BestSequenceOfLevel
    end for
    bestMove  $\leftarrow$  move of the BestSequenceOfLevel
    state  $\leftarrow$  play (state, bestMove)
  end while
```

TYPICAL NMCS



- Mean score in real time



TSP

Given a map, compute a *minimum-cost* round trip visiting certain cities

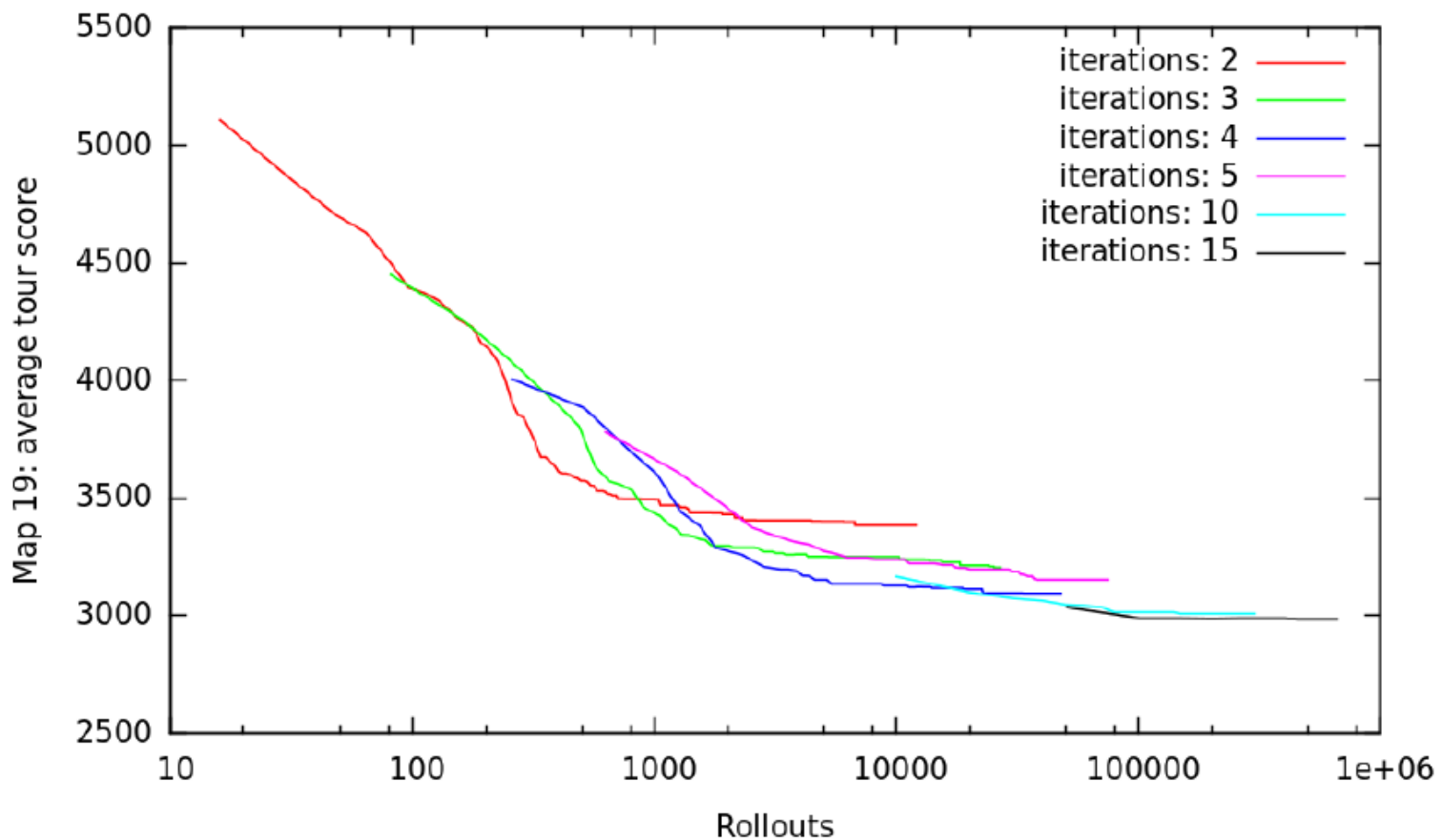


Graph reduction: Precompute *all-pairs-shortest-paths*

Given a distance matrix, compute a *minimum-cost* trip

- Model problem as an IP and call solver (CPLEX, IPSolve, ...)
- Neighborhood search (xOPT: SA; GA; AA; PSO; LNS, ...)
- (Depth-First) Branch and Bound with
 - DFBnB₀ No Heuristic – incremental $O(1)$ time
 - DFBnB₁ Column/Row Minima – incremental $O(1)$ time
 - DFBnB₂ Assignment Problem – incremental $O(n^2)$ time
- ...

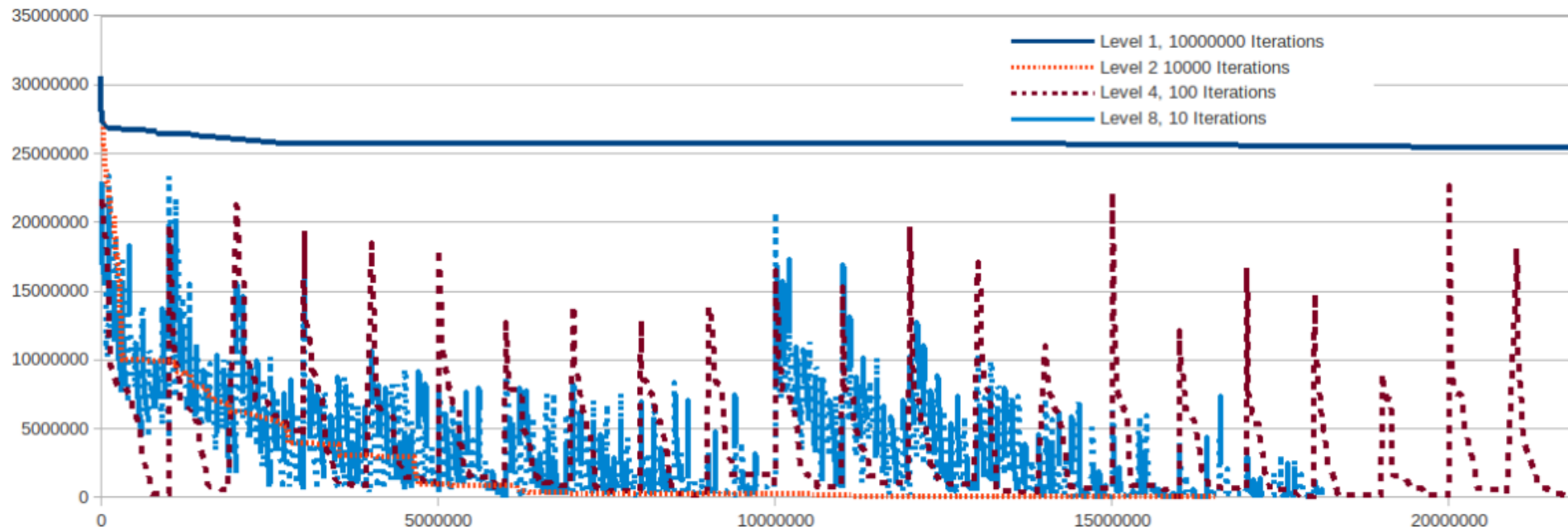
LEARNING CURVE



NESTED ROLLOUT POLICY ADAPTATION

Input: Iteration width (exploitation), nestedness (exploration)

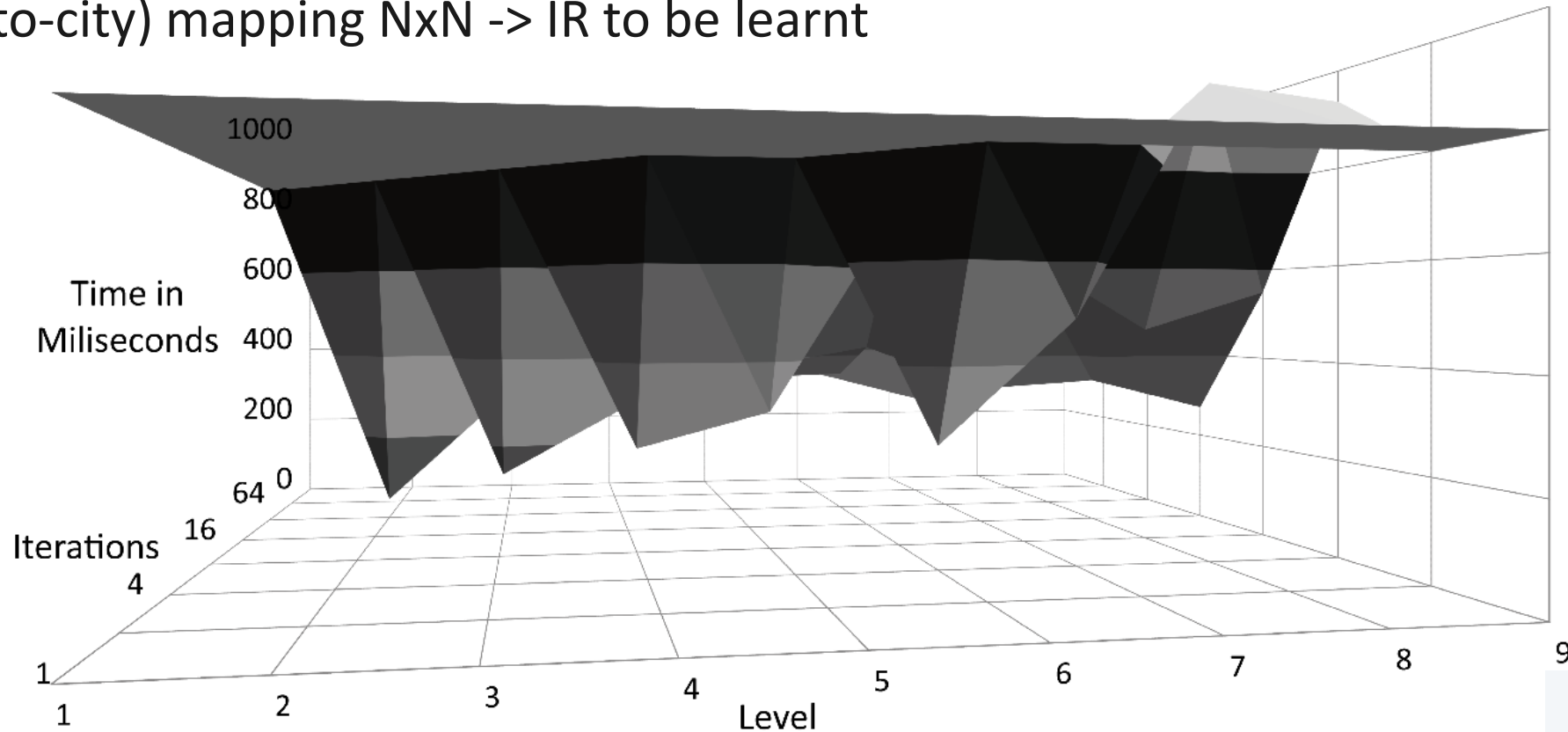
Policy: (city-to-city) mapping $N \times N \rightarrow IR$ to be learnt



NESTED ROLLOUT POLICY ADAPTATION

Input: Iteration width (exploitation), nestedness (exploration)

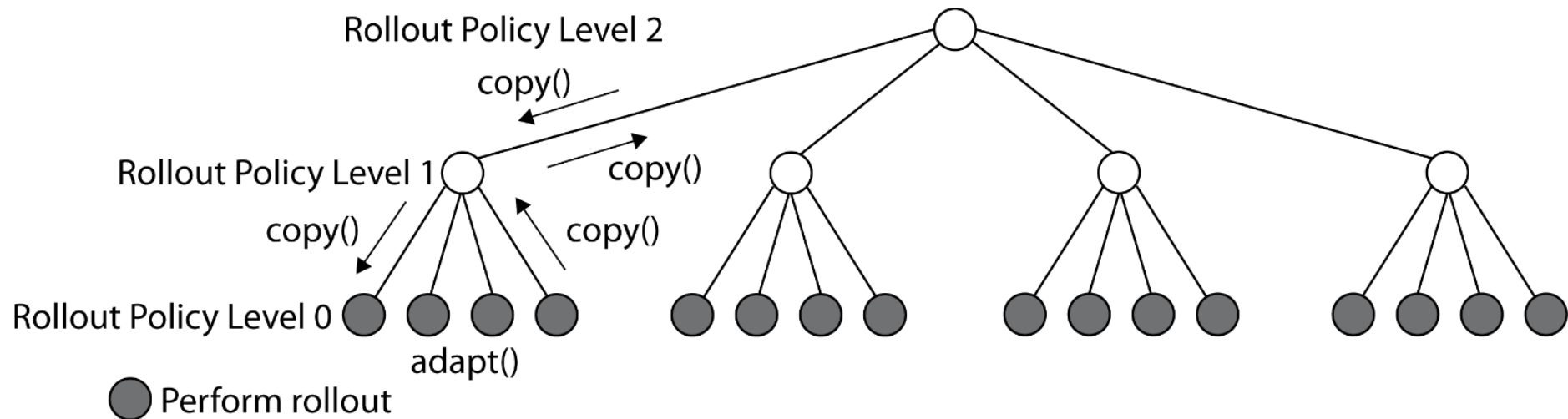
Policy: (city-to-city) mapping $N \times N \rightarrow IR$ to be learnt



NESTED ROLLOUT POLICY ADAPTATION

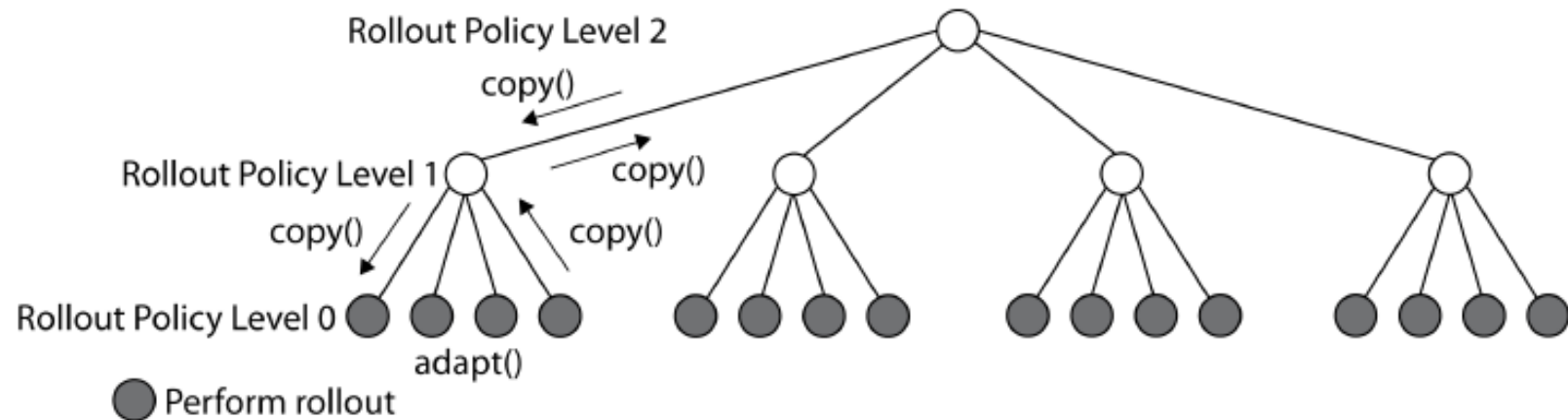
Input: Iteration width (exploitation), nestedness (exploration)

Policy: (city-to-city) mapping $N \times N \rightarrow IR$ to be learnt



NRPA

```
NRPA(level)
  best.score = MAX;
  if (level == 0) best = ROLLOUT();
  else
    backup[level] = global;
    for (i=0; i<iterations; i++)
      current = NRPA(level-1);
      if (current.score < best.score)
        best = current; ADAPT(best);
    global = backup[level];
  return best;
```



PLAYOUT

Algorithm 1 The playout algorithm

```
1: playout (state, policy)
2:   sequence  $\leftarrow$  []
3:   while true do
4:     if state is terminal then
5:       return (score (state), sequence)
6:     end if
7:      $z \leftarrow 0.0$ 
8:     for m in possible moves for state do
9:        $z \leftarrow z + \exp(\text{policy}[\text{code}(m)])$ 
10:    end for
11:    choose a move with probability  $\frac{\exp(\text{policy}[\text{code}(move)])}{z}$ 
12:    state  $\leftarrow$  play (state, move)
13:    sequence  $\leftarrow$  sequence + move
14:  end while
```

ADAPT

Algorithm 2 The Adapt algorithm

```
1: Adapt (policy, sequence)
2:   polp  $\leftarrow$  policy
3:   state  $\leftarrow$  root
4:   for move in sequence do
5:     polp [code(move)]  $\leftarrow$  polp [code(move)] +  $\alpha$ 
6:     z  $\leftarrow$  0.0
7:     for m in possible moves for state do
8:       z  $\leftarrow$  z + exp (policy [code(m)])
9:     end for
10:    for m in possible moves for state do
11:      polp [code(m)]  $\leftarrow$  polp [code(m)] -  $\alpha * \frac{\text{exp}(\text{policy}[\text{code}(\textit{m})])}{z}$ 
12:    end for
13:    state  $\leftarrow$  play (state, move)
14:  end for
15:  policy  $\leftarrow$  polp
```

SEARCH

Algorithm 3 The NRPA algorithm.

```
1: NRPA (level, policy)
2:   if level == 0 then
3:     return playout (root, policy)
4:   else
5:     bestScore  $\leftarrow -\infty$ 
6:     for N iterations do
7:       (result, new)  $\leftarrow$  NRPA(level - 1, policy)
8:       if result  $\geq$  bestScore then
9:         bestScore  $\leftarrow$  result
10:        seq  $\leftarrow$  new
11:       end if
12:       policy  $\leftarrow$  Adapt (policy, seq)
13:     end for
14:     return (bestScore, seq)
15:   end if
```

THEORY...

The probability p_{ik} of choosing the move m_{ik} in a playout is the softmax function:

$$p_{ik} = \frac{e^{w_{ik}}}{\sum_j e^{w_{ij}}}$$

The cross-entropy loss for learning to play move m_{ib} is $C_i = -\log(p_{ib})$. In order to apply the gradient we calculate the partial derivative of the loss: $\frac{\delta C_i}{\delta p_{ib}} = -\frac{1}{p_{ib}}$. We then calculate the partial derivative of the softmax with respect to the weights:

$$\frac{\delta p_{ib}}{\delta w_{ij}} = p_{ib}(\delta_{bj} - p_{ij})$$

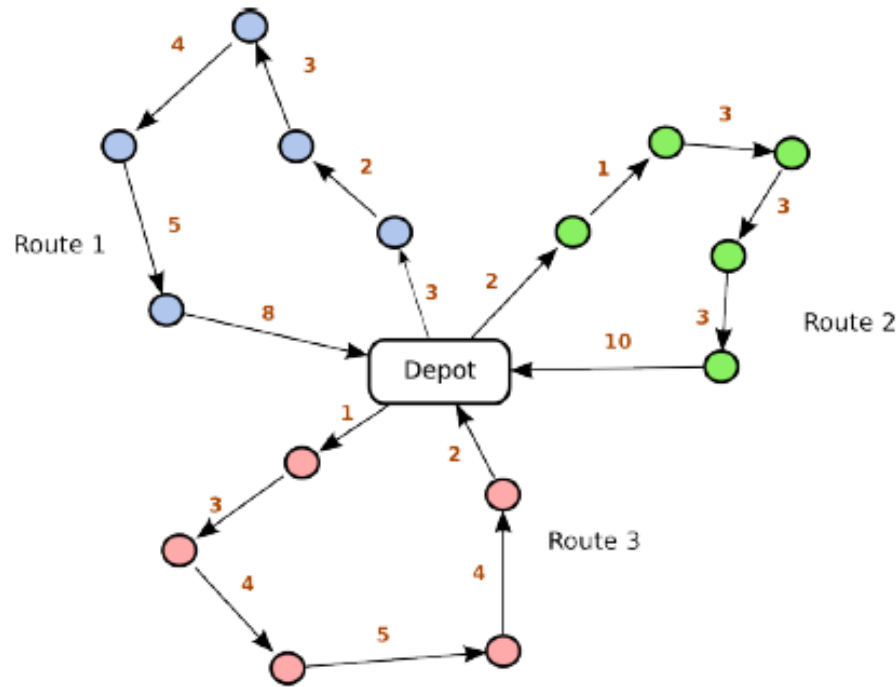
Where $\delta_{bj} = 1$ if $b = j$ and 0 otherwise. Thus the gradient is:

$$\nabla w_{ij} = \frac{\delta C_i}{\delta p_{ib}} \frac{\delta p_{ib}}{\delta w_{ij}} = -\frac{1}{p_{ib}} p_{ib}(\delta_{bj} - p_{ij}) = p_{ij} - \delta_{bj}$$

If we use α as a learning rate we update the weights with:

$$w_{ij} = w_{ij} - \alpha(p_{ij} - \delta_{bj})$$

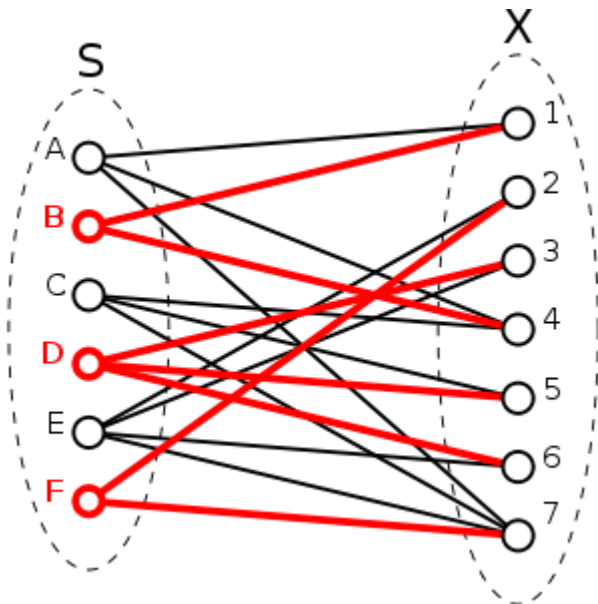
MONTE-CARLO SEARCH FOR VRP



- No additional knowledge needed: rollout function fits in paper.
- Solving Solomon's 100-city problem R101 optimally in 10min (score: 1650.7823, rollouts: 3240000), new highscore in R102
- Very particular benchmark: Even number precision matters.
- Solving Homberger's r2_2_4 (200 packages, 4 vehicles), r2_2_8 (200 packages, 8 vehicles), and r2_4_4 (400 packages, 4 vehicles) in 2h

MONTE-CARLO SEARCH 4 HITTING SET

Given a bipartite graph $G = (V, E)$ with $V = V_1 \cup V_2$, $V_1 \cap V_2 = \emptyset$, and $E \subseteq (V_1 \times V_2)$, find a subset V' of V_1 of minimal cardinality, so that all nodes in V_2 are covered, i.e., for every $v_2 \in V_2$ there is a $v_1 \in V_1$ such that $(v_1, v_2) \in E$.



```
class Game {
public:
    int length, size;
    Move rollout[MaxPlyoutLength];
    Game () {
        for (int j=0; j<SET; j++) visited[j] = 0;
        for (int j=0; j<HITTING; j++) chosen[j] = 0;
        length = size = 0;
    }
    int code (Move m) { return m; }
    bool terminal() { return size >= SET*ALPHA; }
    double score() { return length - size*100; }
    void play(Move m) {
        rollout[length++] = m;
        chosen[m] = 1;
        for (int j = 0; j<SET; j++)
            if (visited[j] == 0 && adjacent[m][j]) {
                visited[j] = 1;
                size++;
            }
    }
    int legalMoves(Move moves[MaxLegalMoves]) {
        int succs = 0;
        for (int m = 0; m < HITTING; m++)
            if (chosen[m] == 0) moves[succs++] = m;
        return successors;
    }
};
```

PRAXIS...

<https://nms.kcl.ac.uk/stefan.edelkamp/lectures/pi1/programs/VRP.java>

The screenshot displays a Java IDE window titled "Blue: Konsole - Einführung CV". The console window shows the following output:

```
Optionen
Level: 2,13, score: 2502165.1351650227, runs: 883320
Level: 2,18, score: 2502164.690862609, runs: 883470
Level: 2,0, score: 2502195.3907905878, runs: 883830
Level: 2,29, score: 2502189.536582865, runs: 884700
Level: 2,0, score: 2502186.811441709, runs: 884730
Level: 2,13, score: 2502177.7208089014, runs: 885120
Level: 2,0, score: 2502187.4106493737, runs: 885630
Level: 2,1, score: 2502186.3647505976, runs: 885660
Level: 2,0, score: 2502178.3200165667, runs: 886530
Level: 2,7, score: 2502177.7208089014, runs: 886740
Level: 2,0, score: 2502177.7208089014, runs: 887430
Level: 2,3, score: 2502165.1351650227, runs: 887520
Level: 2,12, score: 2502160.826497041, runs: 887790
Level: 2,0, score: 2502178.3200165667, runs: 888330
Level: 2,3, score: 2502165.1351650227, runs: 888420
Level: 2,0, score: 2502165.1351650227, runs: 889230
Can only enter input while your programming is running
```

Below the console, there is a graphical user interface (GUI) for the `Mill.Pair` class. The GUI is divided into two main sections:

- Mill.Pair:** This section contains a `double score` field with the value `2502119.8597814734`, an `int[] assignment` field, and buttons for `Inspiziere`, `Hole`, `Zeige statische Variablen`, and `Schließen`.
- assignment : int[]:** This section displays the `int length` as `101` and a list of 11 input fields for the array elements, indexed from `[0]` to `[10]`. The values in the fields are: `[0]: 0`, `[1]: 1`, `[2]: 2`, `[3]: 3`, `[4]: 0`, `[5]: 0`, `[6]: 0`, `[7]: 3`, `[8]: 3`, `[9]: 2`, and `[10]: 0`. Buttons for `Inspiziere` and `Hole` are also present.

The Windows taskbar at the bottom shows the system time as 17:13 on 28.04.2020.

SUMMARY

Monte Carlo Search is a simple algorithm that gives state of the art results for multiple problems:

- Games
- Puzzles
- Snake in the box
- Pancake
- Logistics
- Multiple Sequence Alignment
- RNA Inverse Folding