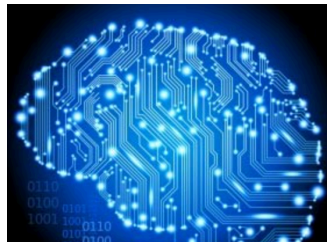




# Lecture 5: Search 4

[http://cs.nju.edu.cn/yuy/course\\_ai16.ashx](http://cs.nju.edu.cn/yuy/course_ai16.ashx)



# Previously...



## Path-based search

### Uninformed search

Depth-first, breadth first, uniform-cost search

### Informed search

Best-first, **A\* search**

## Adversarial search

**Alpha-Beta search**



## Beyond classical search

Bandit search

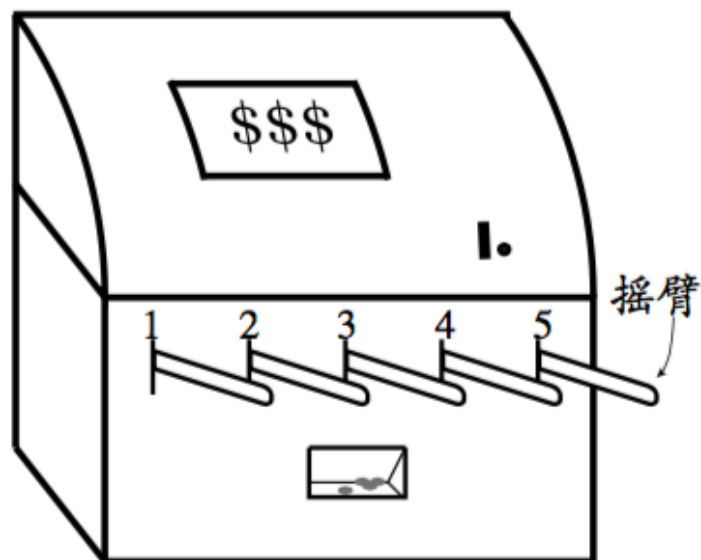
Tree search: Monte-Carlo Tree Search

General search:

Gradient decent

Metaheuristic search

# Bandit



Multiple arms

Each arm has an expected reward,  
but unknown, with an unknown distribution

Maximize your award in fixed trials

# Simplest strategies



Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem?

# Simplest strategies



Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem? waste on suboptimal arms

# Simplest strategies



## Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem? waste on suboptimal arms

Exploitation-only:

# Simplest strategies



## Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem? waste on suboptimal arms

Exploitation-only:

1. try each arm once
2. try the observed best arm  $T-K$  times





# Simplest strategies

## Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem? waste on suboptimal arms

Exploitation-only:

1. try each arm once
2. try the observed best arm  $T-K$  times

problem?



# Simplest strategies

## Two simplest strategies

Exploration-only:

for  $T$  trials and  $K$  arms, try each arm  $T/K$  times

problem? waste on suboptimal arms

Exploitation-only:

1. try each arm once

2. try the observed best arm  $T-K$  times

problem? risk of wrong best arm

# $\epsilon$ -greedy



Balance the exploration and exploitation:

with  $\epsilon$  probability, try a random arm  
with  $1-\epsilon$  probability, try the best arm

$\epsilon$  controls the balance

---

输入: 摇臂数  $K$ ;

奖赏函数  $R$ ;

尝试次数  $T$ ;

探索概率  $\epsilon$ .

过程:

1:  $r = 0$ ;

2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, \text{count}(i) = 0$ ;

3: **for**  $t = 1, 2, \dots, T$  **do**

4:   **if**  $\text{rand}() < \epsilon$  **then**

5:      $k =$  从  $1, 2, \dots, K$  中以均匀分布随机选取

6:   **else**

7:      $k = \arg \max_i Q(i)$

8:   **end if**

9:    $v = R(k)$ ;

10:  $r = r + v$ ;

11:  $Q(k) = \frac{Q(k) \times \text{count}(k) + v}{\text{count}(k) + 1}$ ;

12:  $\text{count}(k) = \text{count}(k) + 1$ ;

13: **end for**

输出: 累积奖赏  $r$

---

# Softmax



Balance the exploration and exploitation:  
Choose arm with probability

$$P(k) = \frac{e^{\frac{Q(k)}{\tau}}}{\sum_{i=1}^K e^{\frac{Q(i)}{\tau}}}, \quad (16.4)$$

$\tau$  controls the balance

---

输入: 摇臂数  $K$ ;  
      奖赏函数  $R$ ;  
      尝试次数  $T$ ;  
      温度参数  $\tau$ .

过程:

- 1:  $r = 0$ ;
- 2:  $\forall i = 1, 2, \dots, K : Q(i) = 0, \text{count}(i) = 0$ ;
- 3: **for**  $t = 1, 2, \dots, T$  **do**
- 4:    $k =$  从  $1, 2, \dots, K$  中根据式(16.4)随机选取
- 5:    $v = R(k)$ ;
- 6:    $r = r + v$ ;
- 7:    $Q(k) = \frac{Q(k) \times \text{count}(k) + v}{\text{count}(k) + 1}$ ;
- 8:    $\text{count}(k) = \text{count}(k) + 1$ ;
- 9: **end for**

输出: 累积奖赏  $r$

---

# Upper-confidence bound

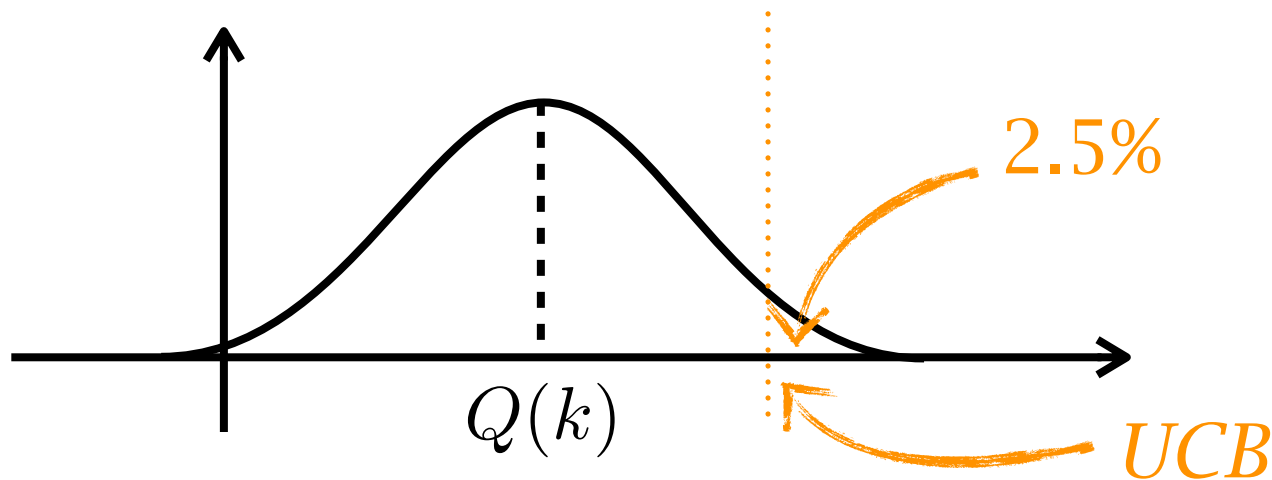


Balance the exploration and exploitation:

Choose arm with the largest value of

average reward + upper confidence bound

$$Q(k) + \sqrt{\frac{2 \ln n}{n_k}},$$



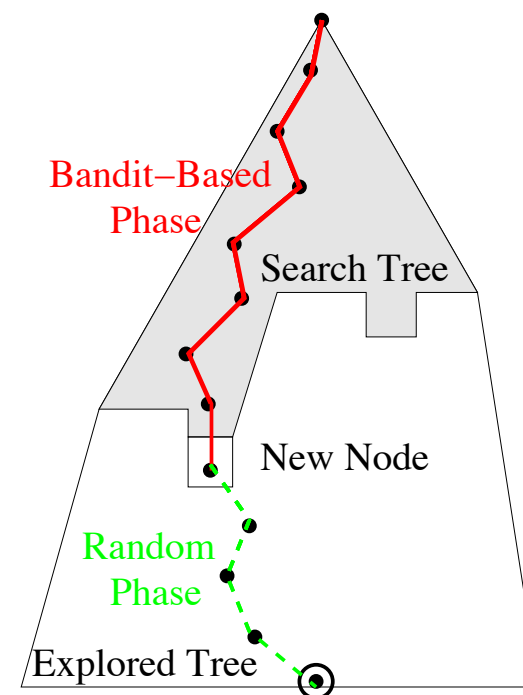
# Monte-Carlo Tree Search



Kocsis Szepesvári, 06

Gradually grow the search tree:

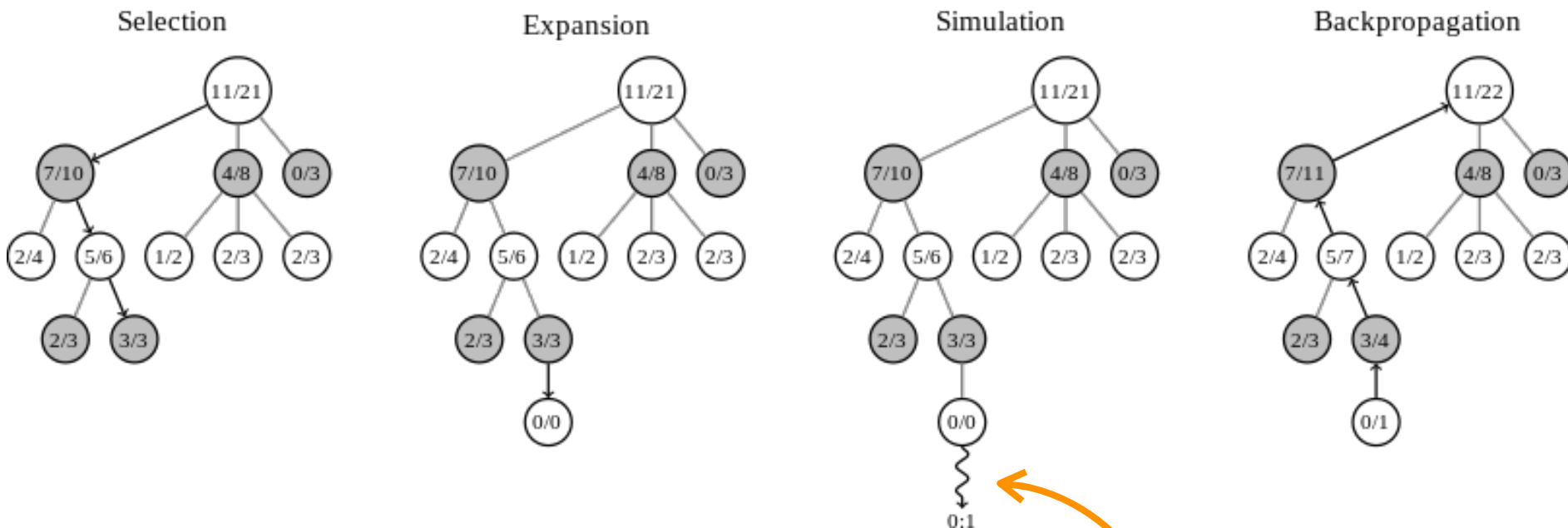
- ▶ Iterate Tree-Walk
  - ▶ Building Blocks
    - ▶ Select next action **Bandit phase**
    - ▶ Add a node **Grow a leaf of the search tree**
    - ▶ Select next action bis **Random phase, roll-out**
    - ▶ Compute instant reward **Evaluate**
    - ▶ Update information in visited nodes **Propagate**
  - ▶ Returned solution:
    - ▶ Path visited most often



# Monte-Carlo Tree Search



Example:



Pic from [https://en.wikipedia.org/wiki/Monte\\_Carlo\\_tree\\_search#cite\\_note-Kocsis-Szepesvari-5](https://en.wikipedia.org/wiki/Monte_Carlo_tree_search#cite_note-Kocsis-Szepesvari-5)

rollout

How to select the leave? As bandit

# Monte-Carlo Tree Search



```
public class TreeNode {
    static Random r = new Random();
    static int nActions = 5;
    static double epsilon = 1e-6;

    TreeNode[] children;
    double nVisits, totValue;

    public void selectAction() {
        List<TreeNode> visited = new LinkedList<TreeNode>();
        TreeNode cur = this;
        visited.add(this);
        while (!cur.isLeaf()) {
            cur = cur.select();
            visited.add(cur);
        }
        cur.expand();
        TreeNode newNode = cur.select();
        visited.add(newNode);
        double value = rollOut(newNode);
        for (TreeNode node : visited) {
            // would need extra logic for n-player game
            node.updateStats(value);
        }
    }
}
```

```
public void expand() {
    children = new TreeNode[nActions];
    for (int i=0; i<nActions; i++) {
        children[i] = new TreeNode();
    }
}
```

```
public void updateStats(double value) {
    nVisits++;
    totValue += value;
}
```



# Monte-Carlo Tree Search



```
private TreeNode select() {
    public
    st
    st
    st
    Tr
    do
    pu
    TreeNode selected = null;
    double bestValue = Double.MIN_VALUE;
    for (TreeNode c : children) {
        double uctValue = c.totValue / (c.nVisits + epsilon) +
            Math.sqrt(Math.log(nVisits+1) / (c.nVisits + epsilon)) +
                r.nextDouble() * epsilon;
        // small random number to break ties randomly in unexpanded nodes
        if (uctValue > bestValue) {
            selected = c;
            bestValue = uctValue;
        }
    }
    return selected;
}
```

```
    cur = cur.select();
    visited.add(cur);
```

```
    }
    cur.expand();
    TreeNode newNode = cur.select();
    visited.add(newNode);
    double value = rollOut(newNode);
    for (TreeNode node : visited) {
        // would need extra logic for n-player game
        node.updateStats(value);
    }
}
```

```
    totValue += value;
}
```

# Monte-Carlo Tree Search



optimal? Yes, after infinite tries

compare with alpha-beta pruning  
no need of heuristic function

# Monte-Carlo Tree Search



## Improving random rollout

### Monte-Carlo-based

Brügman 93

1. Until the goban is filled, add a stone (black or white in turn) at a uniformly selected empty position
2. Compute  $r = \text{Win}(\text{black})$
3. The outcome of the tree-walk is  $r$



### Improvements ?

- ▶ Put stones randomly in the neighborhood of a previous stone
- ▶ Put stones matching patterns
- ▶ Put stones optimizing a value function

Silver et al. 07



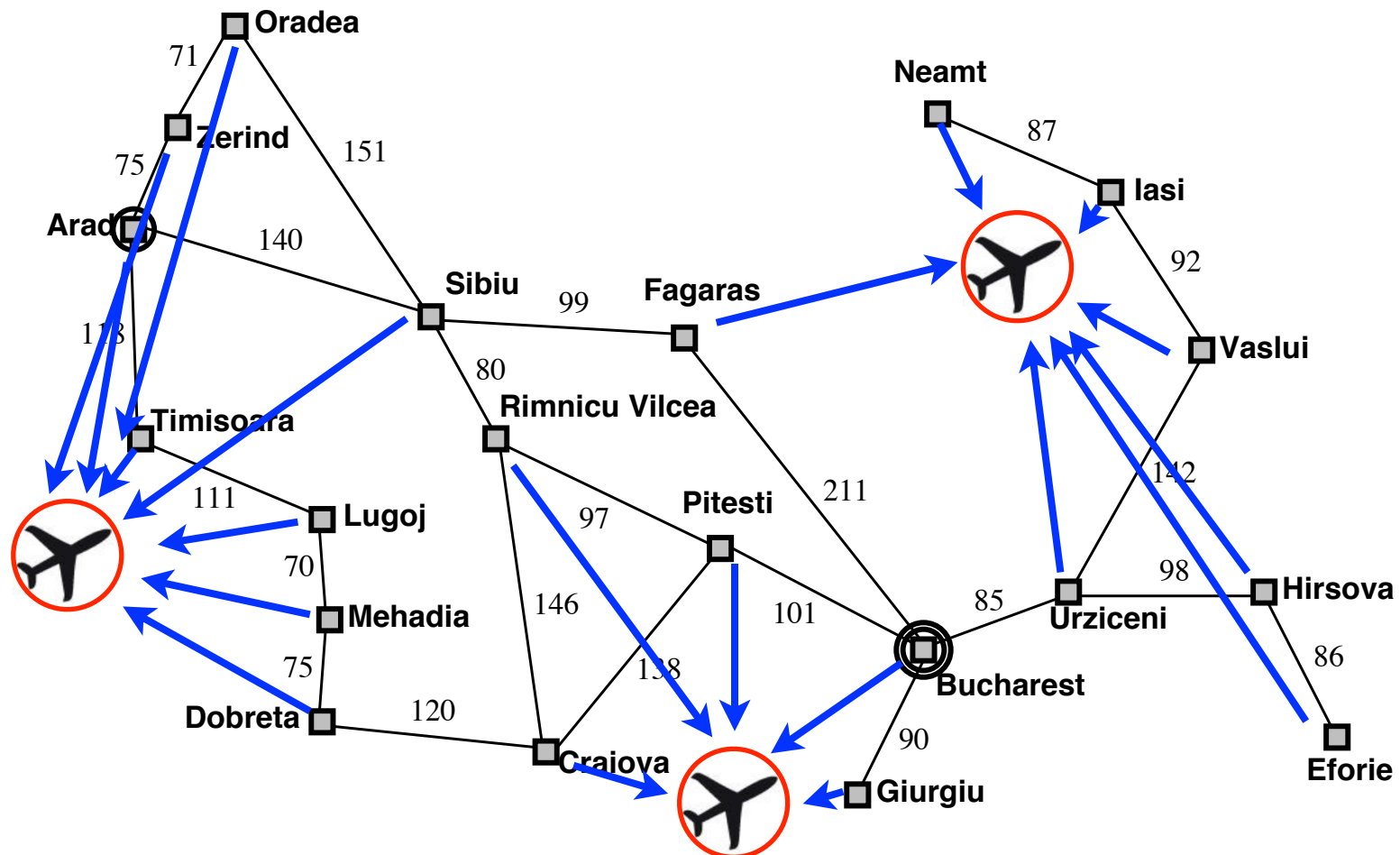
# General search

# Greedy idea in continuous space



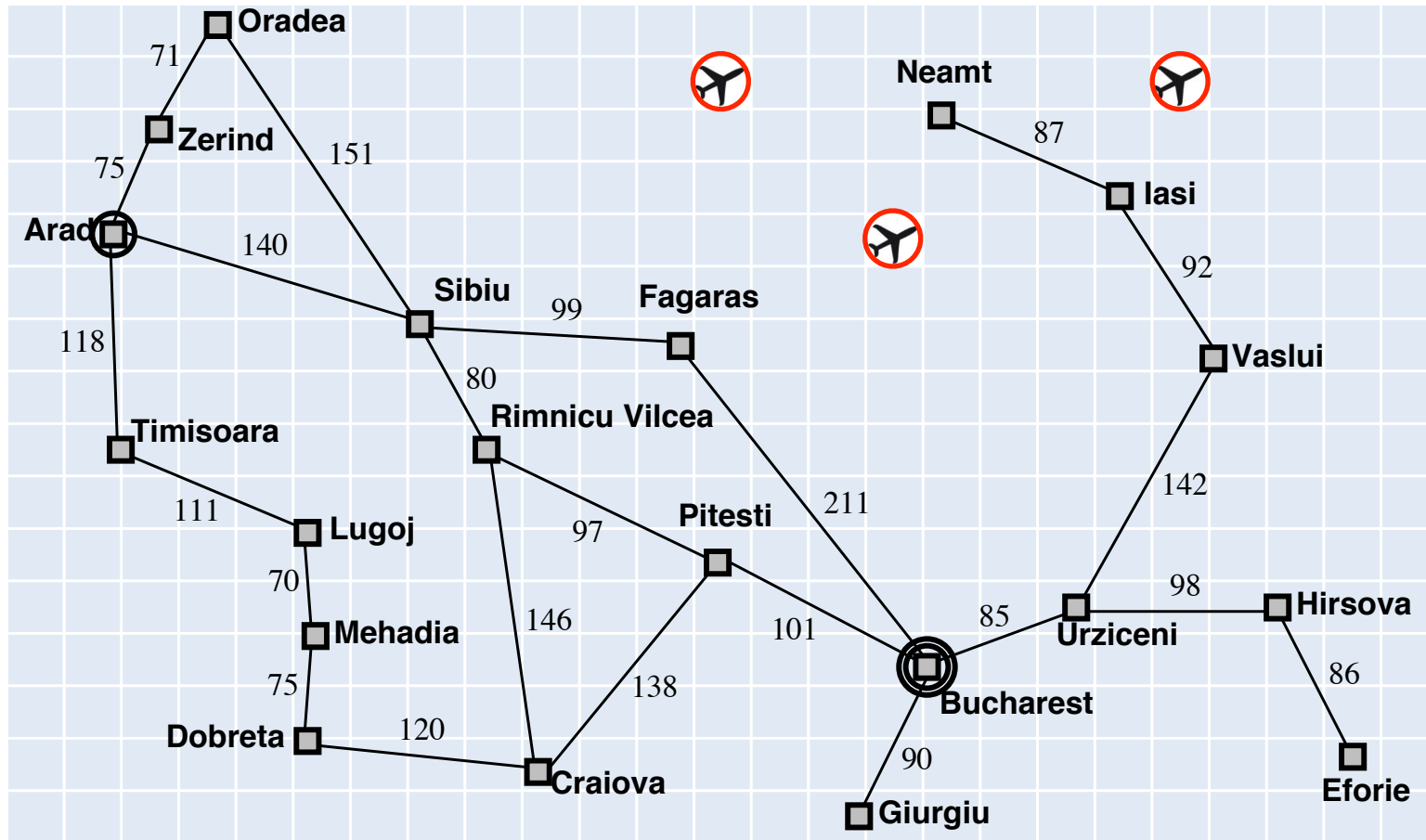
Suppose we want to site three airports in Romania:

- 6-D state space defined by  $(x_1, y_1), (x_2, y_2), (x_3, y_3)$
- objective function  $f(x_1, y_1, x_2, y_2, x_3, y_3) =$   
sum of squared distances from each city to nearest airport



# Greedy idea in continuous space

discretize and use hill climbing





# Greedy idea in continuous space

## gradient decent

- 6-D state space defined by  $(x_1, y_1), (x_2, y_2), (x_3, y_3)$
- objective function  $f(x_1, y_1, x_2, y_2, x_3, y_3) =$   
sum of squared distances from each city to nearest airport

Gradient methods compute

$$\nabla f = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial y_1}, \frac{\partial f}{\partial x_2}, \frac{\partial f}{\partial y_2}, \frac{\partial f}{\partial x_3}, \frac{\partial f}{\partial y_3} \right)$$

to increase/reduce  $f$ , e.g., by  $\mathbf{x} \leftarrow \mathbf{x} + \alpha \nabla f(\mathbf{x})$

*1-order method*



# Greedy idea in continuous space

## gradient decent

- 6-D state space defined by  $(x_1, y_2), (x_2, y_2), (x_3, y_3)$
- objective function  $f(x_1, y_2, x_2, y_2, x_3, y_3) =$   
sum of squared distances from each city to nearest airport

Sometimes can solve for  $\nabla f(\mathbf{x}) = 0$  exactly (e.g., with one city).  
Newton–Raphson (1664, 1690) iterates  $\mathbf{x} \leftarrow \mathbf{x} - \mathbf{H}_f^{-1}(\mathbf{x})\nabla f(\mathbf{x})$   
to solve  $\nabla f(\mathbf{x}) = 0$ , where  $\mathbf{H}_{ij} = \partial^2 f / \partial x_i \partial x_j$

## 2-order method

Taylor's series:

$$f(x) = f(a) + (x - a)f'(a) + \frac{(x - a)^2}{2}f''(a) + \dots = \sum_{i=0}^{\infty} \frac{(x - a)^i}{i!} f^{(i)}(a).$$

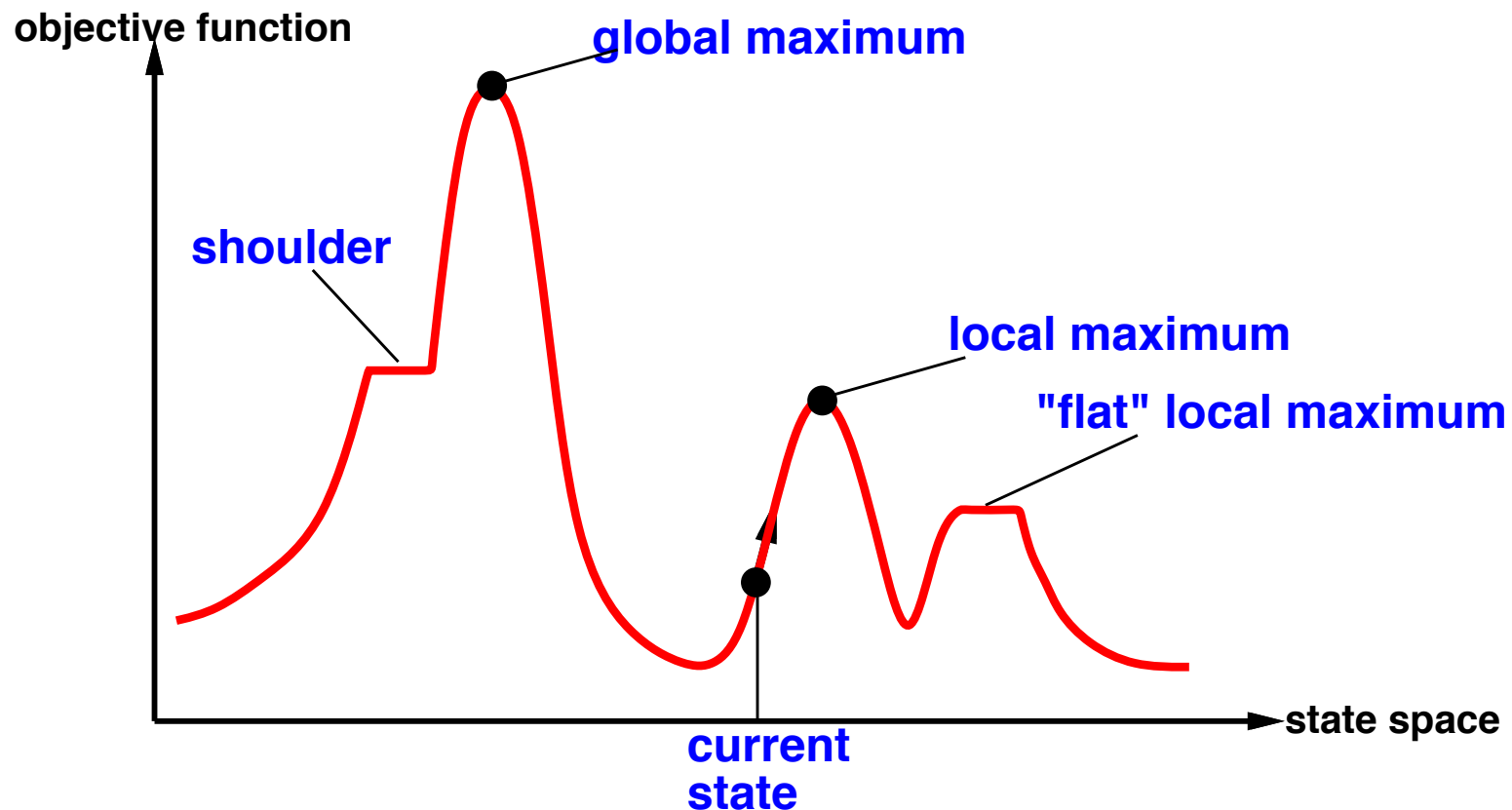


# Greedy idea



1st and 2nd order methods may not find global optimal solutions

they work for convex functions



# Meta-heuristics



“problem independent

“black-box

“zeroth-order method

...

and usually inspired from nature phenomenon

# Simulated annealing



temperature from high to low

when high temperature, form the shape  
when low temperature, polish the detail

# Simulated annealing



Idea: escape local maxima by allowing some “bad” moves  
but gradually decrease their size and frequency

**function** SIMULATED-ANNEALING(*problem*, *schedule*) **returns** a solution state

**inputs:** *problem*, a problem

*schedule*, a mapping from time to “temperature”

**local variables:** *current*, a node

*next*, a node

*T*, a “temperature” controlling prob. of downward steps

*current* ← MAKE-NODE(INITIAL-STATE[*problem*])

**for**  $t \leftarrow 1$  to  $\infty$  **do**

*T* ← *schedule*[*t*]

**if**  $T = 0$  **then return** *current*

*next* ← a randomly selected successor of *current* *the neighborhood range*

$\Delta E \leftarrow \text{VALUE}[\textit{next}] - \text{VALUE}[\textit{current}]$  *shrinks with T*

**if**  $\Delta E > 0$  **then** *current* ← *next*

**else** *current* ← *next* only with probability  $e^{\Delta E/T}$  *the probability of accepting a bad solution decreases with T*

# Simulated annealing

a demo



# Local beam search



**Idea:** keep  $k$  states instead of 1; choose top  $k$  of all their successors

Not the same as  $k$  searches run in parallel!

Searches that find good states recruit other searches to join them

**Problem:** quite often, all  $k$  states end up on same local hill

**Idea:** choose  $k$  successors randomly, biased towards good ones

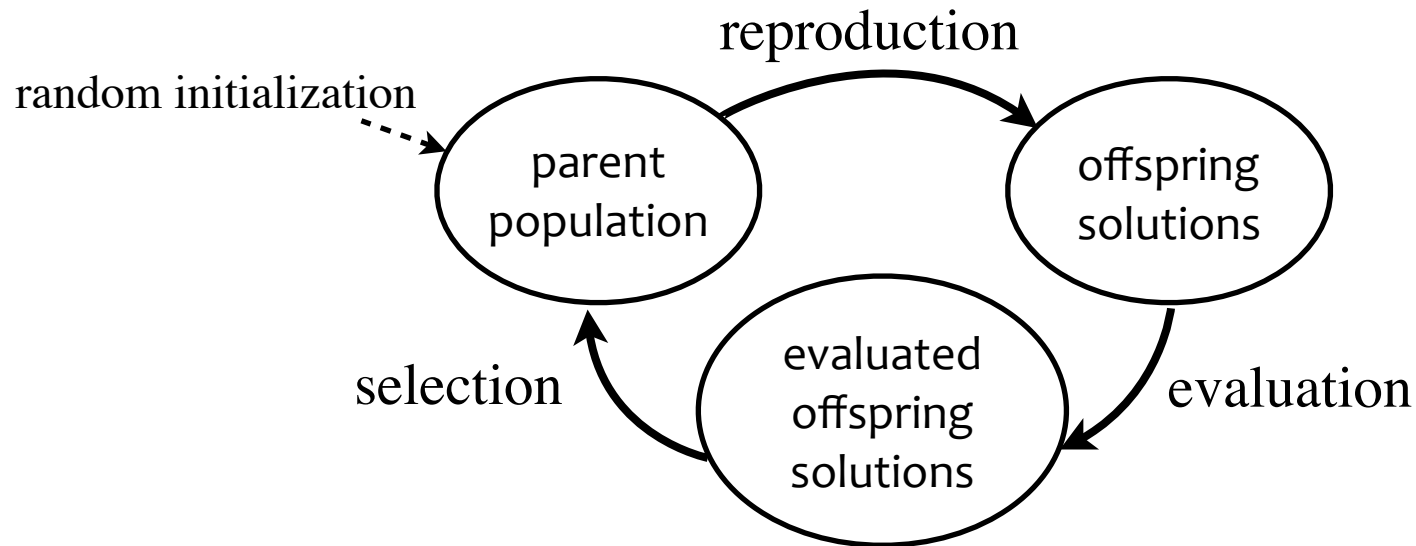
Observe the close analogy to natural selection!

# Genetic algorithm



a simulation of Darwin's evolutionary theory

over-reproduction with diversity  
nature selection



# Genetic algorithm

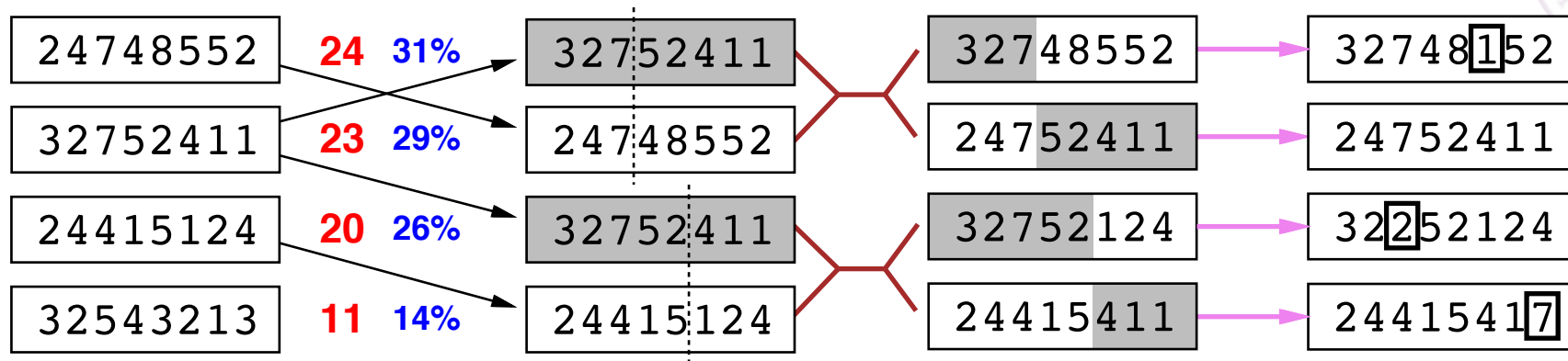


Encode a solution as a vector,

- 1:  $Pop \leftarrow n$  randomly drawn solutions from  $\mathcal{X}$
- 2: **for**  $t=1,2,\dots$  **do**
- 3:      $Pop^m \leftarrow \{mutate(s) \mid \forall s \in Pop\}$ , the mutated solutions
- 4:      $Pop^c \leftarrow \{crossover(s_1, s_2) \mid \exists s_1, s_2 \in Pop^m\}$ , the recombined solutions
- 5:     evaluate every solution in  $Pop^c$  by  $f(s)(\forall s \in Pop^c)$
- 6:      $Pop^s \leftarrow$  selected solutions from  $Pop$  and  $Pop^c$
- 7:      $Pop \leftarrow Pop^s$
- 8:     **terminate** if meets a stopping criterion
- 9: **end for**



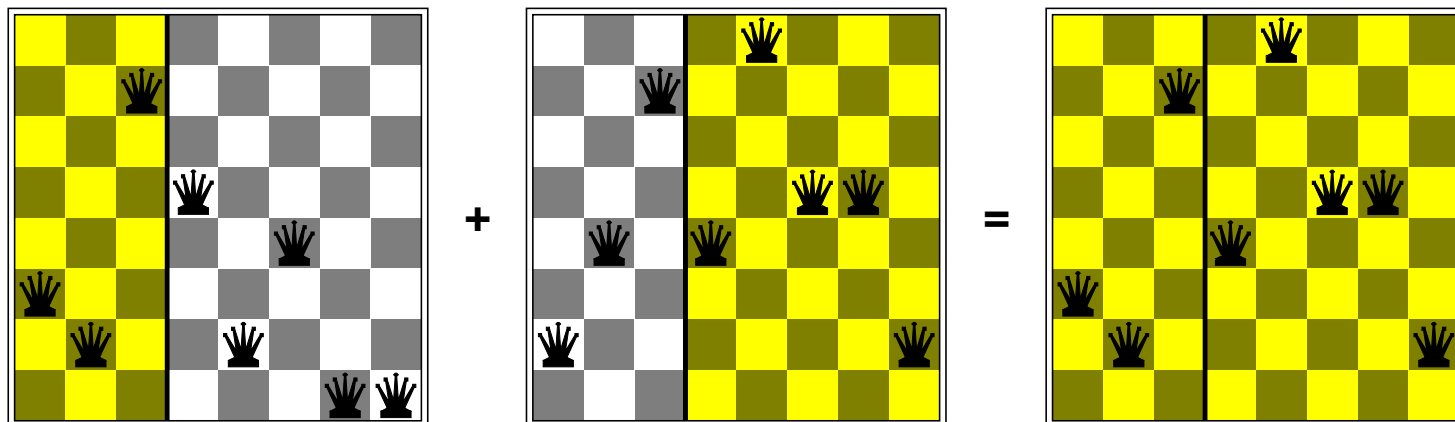
# Genetic algorithm



**Fitness**   **Selection**   **Pairs**   **Cross-Over**   **Mutation**

GAs require states encoded as strings (GPs use programs)

Crossover helps **iff substrings are meaningful components**





# Example

Encode a solution as a vector with length  $n$

each element of the vector can be chosen from  $\{1, \dots, V\}$

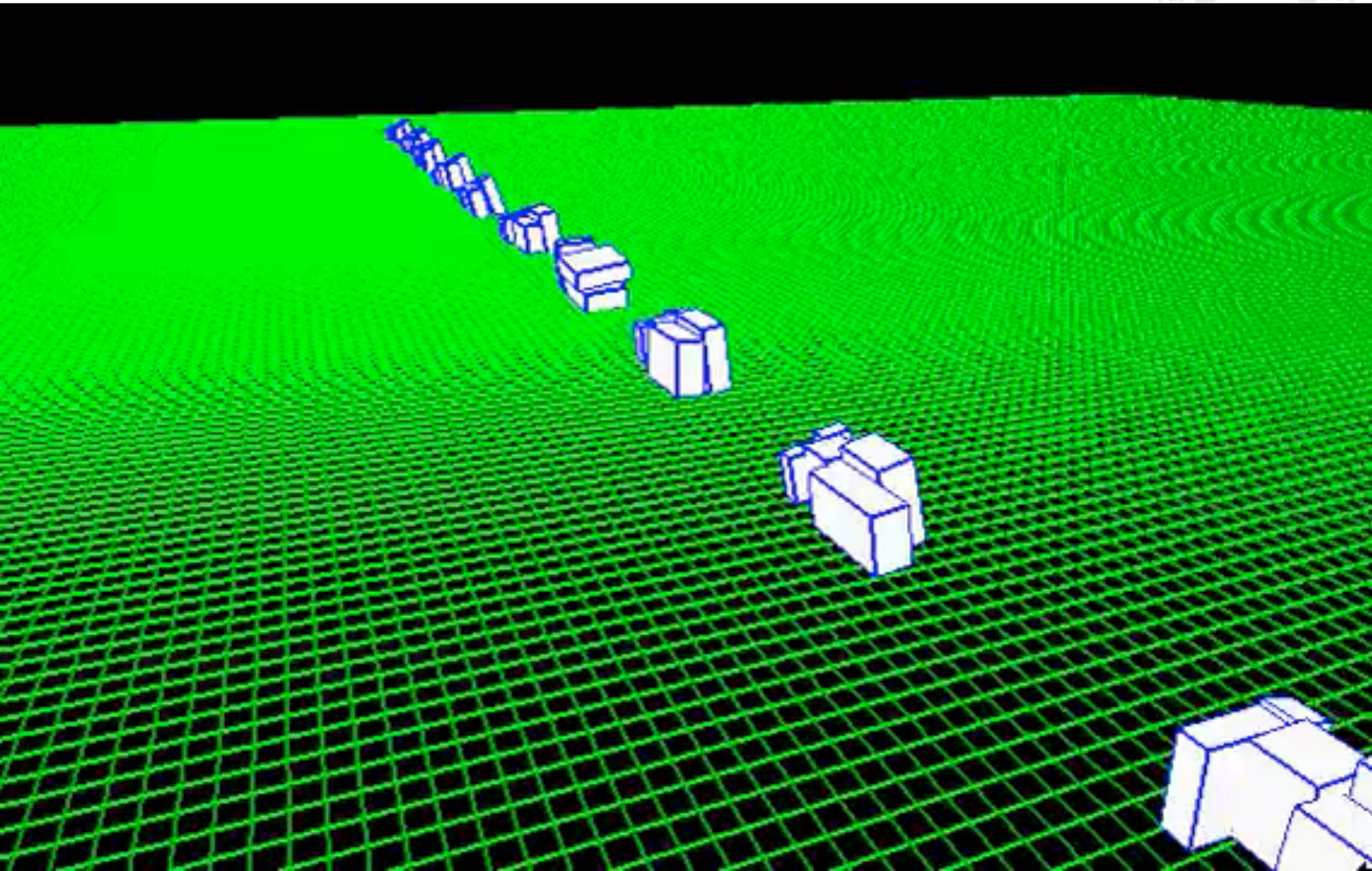
parameters: mutation probability  $p_m$ , crossover probability  $p_c$

- 1:  $Pop$  = randomly generate  $n$  solutions from  $\{1, \dots, V\}^n$
- 2: for  $t=1, 2, \dots$  do
- 3:      $Pop^m = \text{emptyset}, Pop^c = \text{emptyset}$
- 4:     for  $i = 1$  to  $n$
- 5:         let  $x$  be the  $i$ -th solution in  $Pop$
- 6:         for  $j = 1$  to  $n$ : with probability  $p_m$ , change  $x_j$  by a random value from  $\{1, \dots, V\}$
- 7:         add  $x$  into  $Pop^m$
- 8:     end for
- 9:     for  $i = 1$  to  $n$
- 10:         let  $x$  be the  $i$ -th solution in  $Pop^m$
- 11:         let  $x'$  be a randomly selected solution from  $Pop^m$
- 12:         with probability  $p_c$ , exchange a random part of  $x$  with  $x'$
- 13:         add  $x$  into  $Pop^c$
- 14:     end for
- 15:     evaluate solutions in  $Pop^c$ , select the best  $n$  solutions from  $Pop$  and  $Pop^c$  to  $Pop$
- 16:     terminal if a good solution is found
- 17: end for

# An evolutionary of virtual life



# An evolutionary of virtual life



# Properties of meta-heuristics



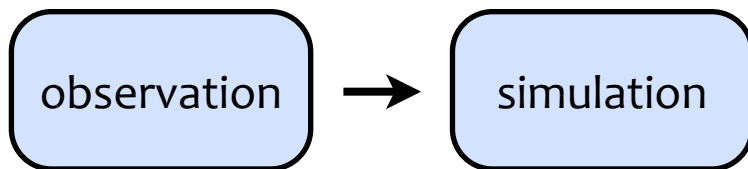
zeroth order

do not need differentiable functions

convergence

will find an optimal solution if  $P(x^* | x) > 0$   
or  $P(x \rightarrow x_1 \rightarrow \dots \rightarrow x_k \rightarrow x^*) > 0$

a missing link



# Properties of meta-heuristics



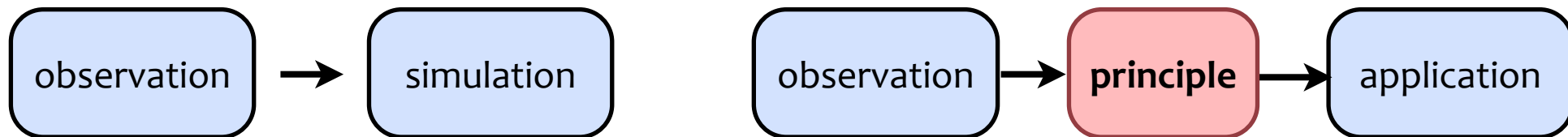
zeroth order

do not need differentiable functions

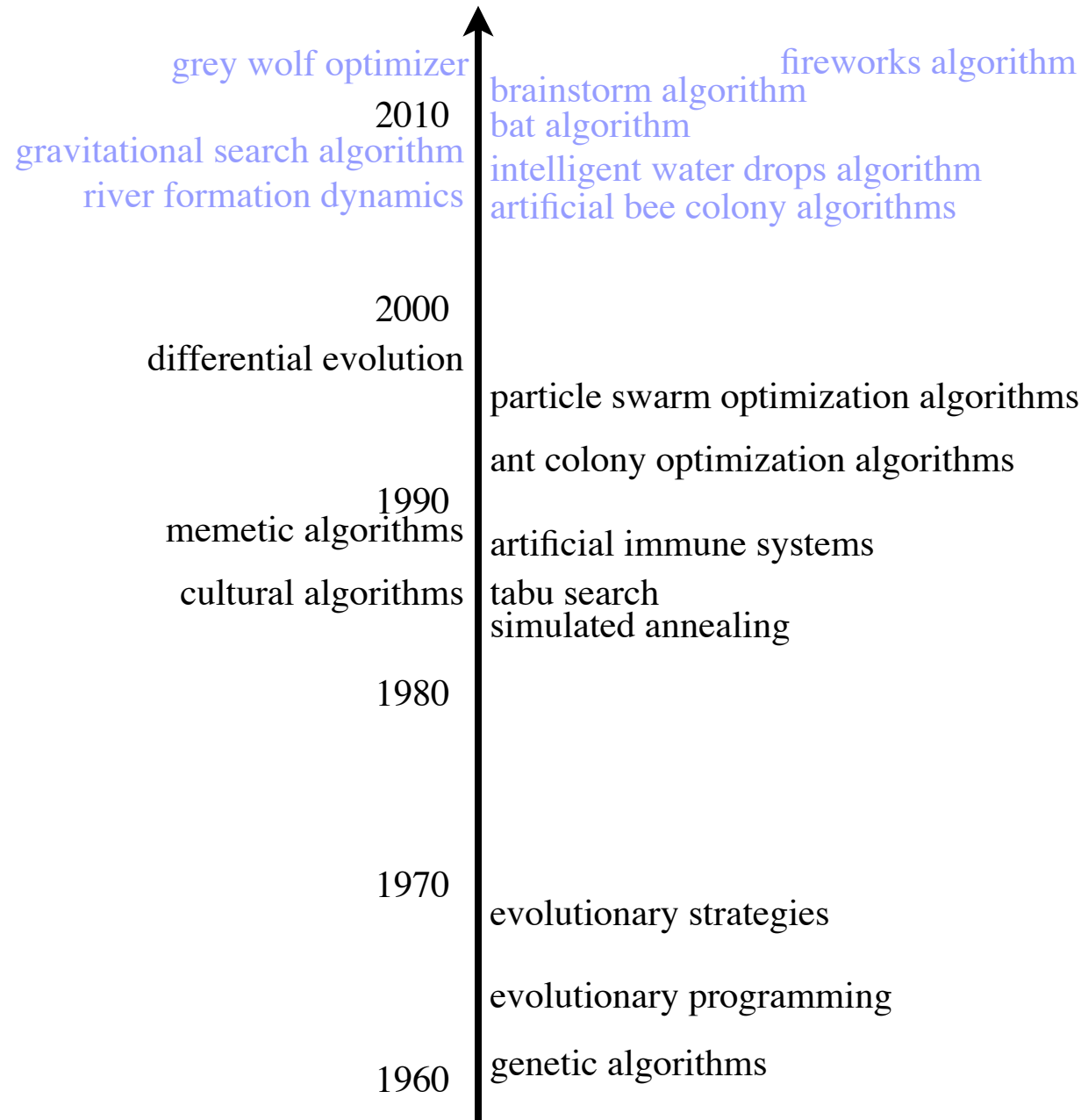
convergence

will find an optimal solution if  $P(x^* | x) > 0$   
or  $P(x \rightarrow x_1 \rightarrow \dots \rightarrow x_k \rightarrow x^*) > 0$

a missing link



# Properties of meta-heuristics



# Example

*hard to apply traditional optimization methods  
but easy to test a given solution*



## Representation:



parameterize  
→

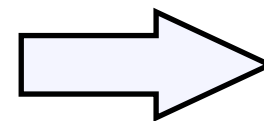
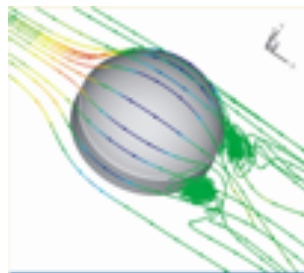
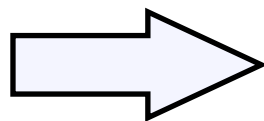


represented as a vector of parameters

## Fitness:



$x_i$



$f(x_i)$

test by simulation/experiment



# Example



Series 700



Series N700

## Technological overview of the next generation Shinkansen high-speed train Series N700

M. Ueno<sup>1</sup>, S. Usui<sup>1</sup>, H. Tanaka<sup>1</sup>, A. Watanabe<sup>2</sup>

<sup>1</sup>Central Japan Railway Company, Tokyo, Japan, <sup>2</sup>West Japan Railway Company, Osaka, Japan

### Abstract

In March 2005, Central Japan Railway Company (JR Central) has completed prototype trainset of the Series N700, the next generation Shinkansen high-speed rolling stock developed

aiming at the aerodynamic system, and subjected to the problem of aerodynamic pressure waves and other issues related to environmental compatibility such as external noise. To combat this, an aero double-wing-type has been adopted for nose shape (Fig. 3). This nose shape, which boasts the most appropriate aerodynamic performance, has been newly developed for railway rolling stock using the latest analytical technique (i.e. genetic algorithms) used to develop the main wings of airplanes. The shape resembles a bird in flight, suggesting a feeling of boldness and speed.

suppress aerodynamic noise. [1]

On the Tokaido Shinkansen line, Series N700 cars save 19% energy than Series 700 cars, despite a 30% increase in the output of their traction equipment for higher-speed operation (Fig. 4).

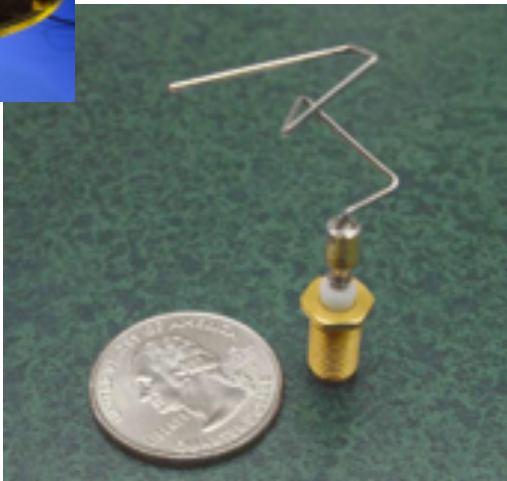
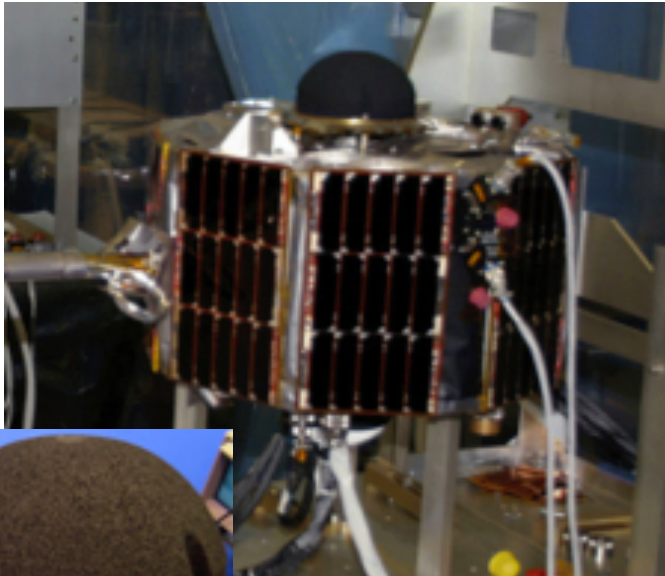
This is a result of adopting the aerodynamically excellent nose shape, reduced running resistance thanks to the drastically smoothed car body and under-floor equipment, effective

this nose ... has been newly developed ... using the latest analytical technique (i.e. **genetic algorithms**)

N700 cars save **19%** energy ... **30%** increase in the output... This is a result of adopting the ... nose shape

# Example

NASA ST5 satellite



*hard to apply traditional optimization methods  
but easy to test a given solution*

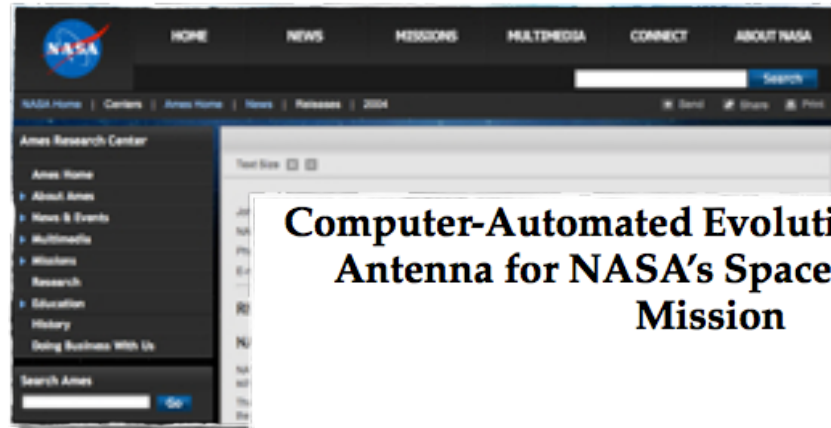
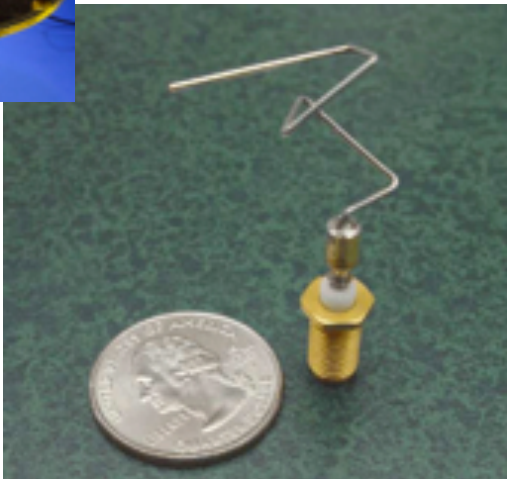
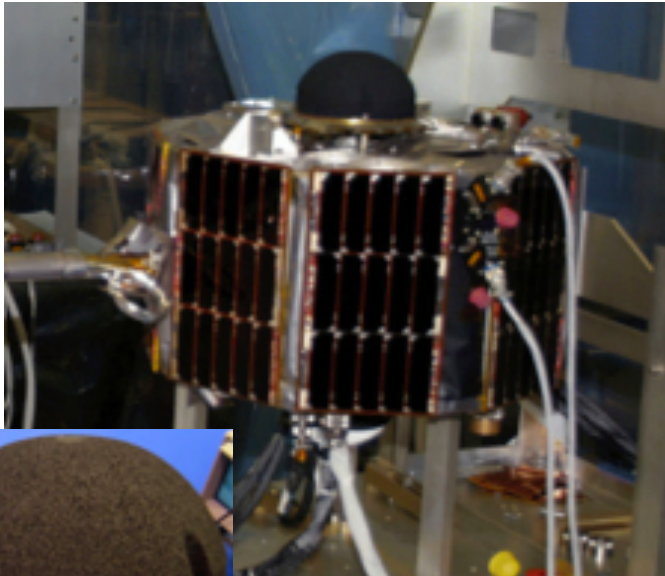


# Example



## NASA ST5 satellite

*hard to apply traditional optimization methods  
but easy to test a given solution*



### Computer-Automated Evolution of an X-Band Antenna for NASA's Space Technology 5 Mission

**Gregory S. Hornby**

Gregory.S.Hornby@nasa.gov  
Mail Stop 269-3, University Affiliated Research Center, UC Santa Cruz, Moffett Field, CA, 94035, USA

**Jason D. Lohn**

Jason.Lohn@west.cmu.edu  
Carnegie Mellon University, Mail Stop 23-11, Moffett Field, CA 94035, USA

**Derek S. Linden**

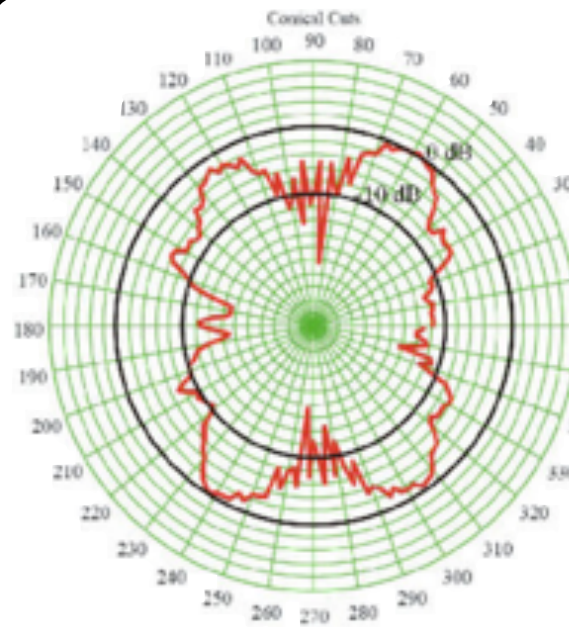
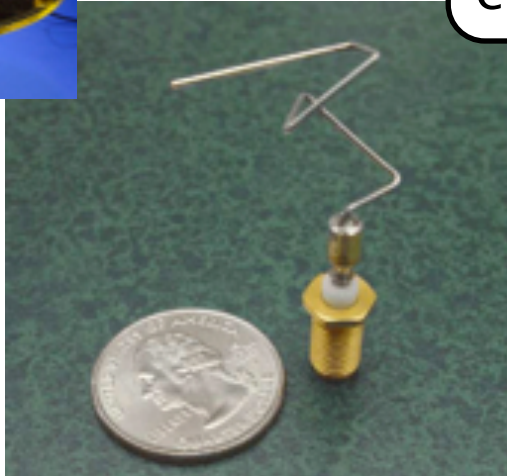
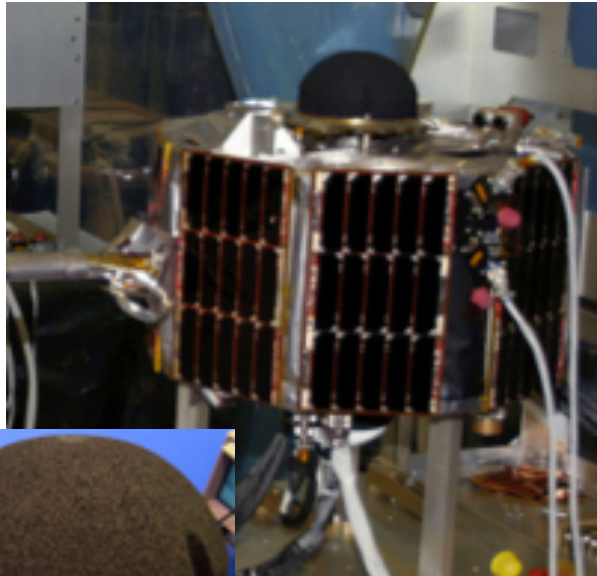
dlinden@jemengineering.com  
JEM Engineering, 8683 Cherry Lane, Laurel, MD 20707, USA Moffett Field, CA 94035, USA

Since there are two antennas on each spacecraft, and not just one, it is important to measure the overall gain pattern with two antennas mounted on the spacecraft. For this, different combinations of the two evolved antennas and the QHA were tried on the the ST5 mock-up and measured in an anechoic chamber. **With two QHAs 38% efficiency was achieved, using a QHA with an evolved antenna resulted in 80% efficiency, and using two evolved antennas resulted in 93% efficiency.** Here "efficiency" means how much power is being radiated versus how much power is being eaten up in resistance, with greater efficiency resulting in a stronger signal and greater range. Figure 11

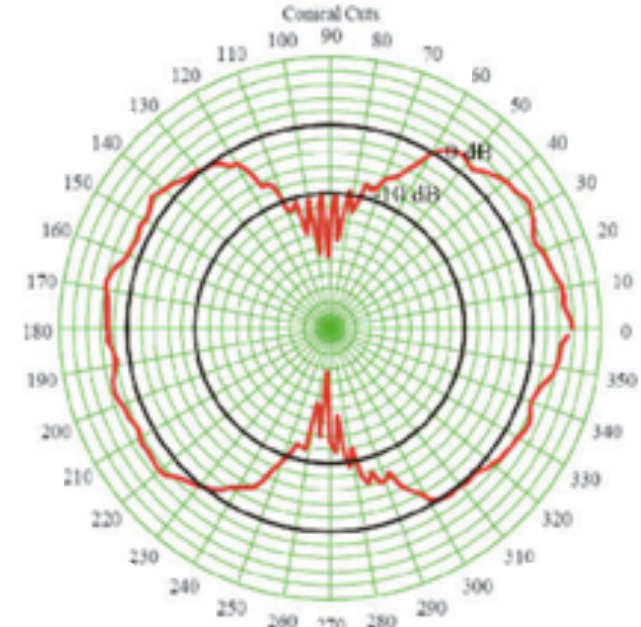
# Example



NASA ST5 satellite



QHAs ( ) 38% efficiency



evolved antennas resulted in 93% efficiency

**Jason D. Lohn** [JDLohn@west.cmu.edu](mailto:JDLohn@west.cmu.edu)  
Carnegie Mellon University, Mail Stop 23-11, Moffett Field, CA 94035, USA

**Derek S. Linden** [dlinden@jemengineering.com](mailto:dlinden@jemengineering.com)  
JEM Engineering, 8683 Cherry Lane, Laurel, MD 20707, USA Moffett Field, CA 94035, USA

Since there are two antennas on each spacecraft, and not just one, it is important to measure the overall gain pattern with two antennas mounted on the spacecraft. For this, different combinations of the two evolved antennas and the QHA were tried on the the ST5 mock-up and measured in an anechoic chamber. **With two QHAs 38% efficiency was achieved, using a QHA with an evolved antenna resulted in 80% efficiency, and using two evolved antennas resulted in 93% efficiency.** Here "efficiency" means how much power is being radiated versus how much power is being eaten up in resistance, with greater efficiency resulting in a stronger signal and greater range. Figure 11



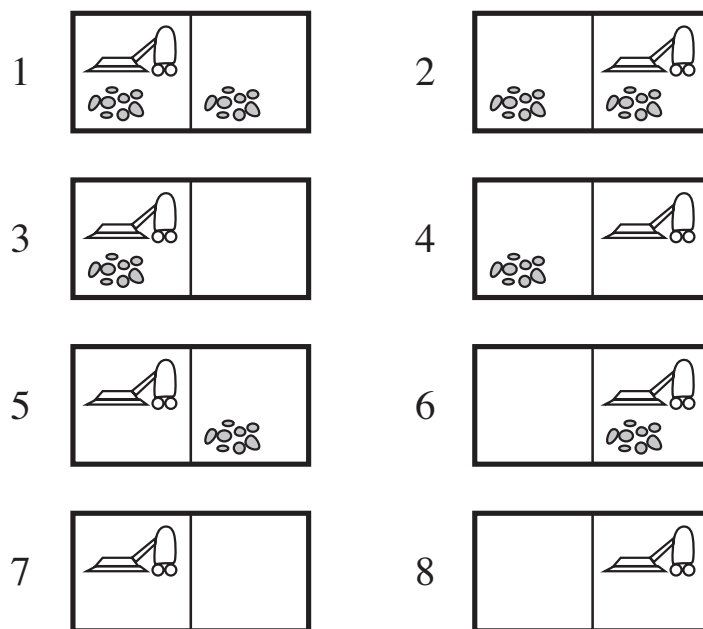
# Different Environment Properties

# Nondeterministic actions



In the **erratic vacuum world**, the *Suck* action works as follows:

- When applied to a dirty square the action cleans the square and sometimes cleans up dirt in an adjacent square, too.
- When applied to a clean square the action sometimes deposits dirt on the carpet.<sup>9</sup>

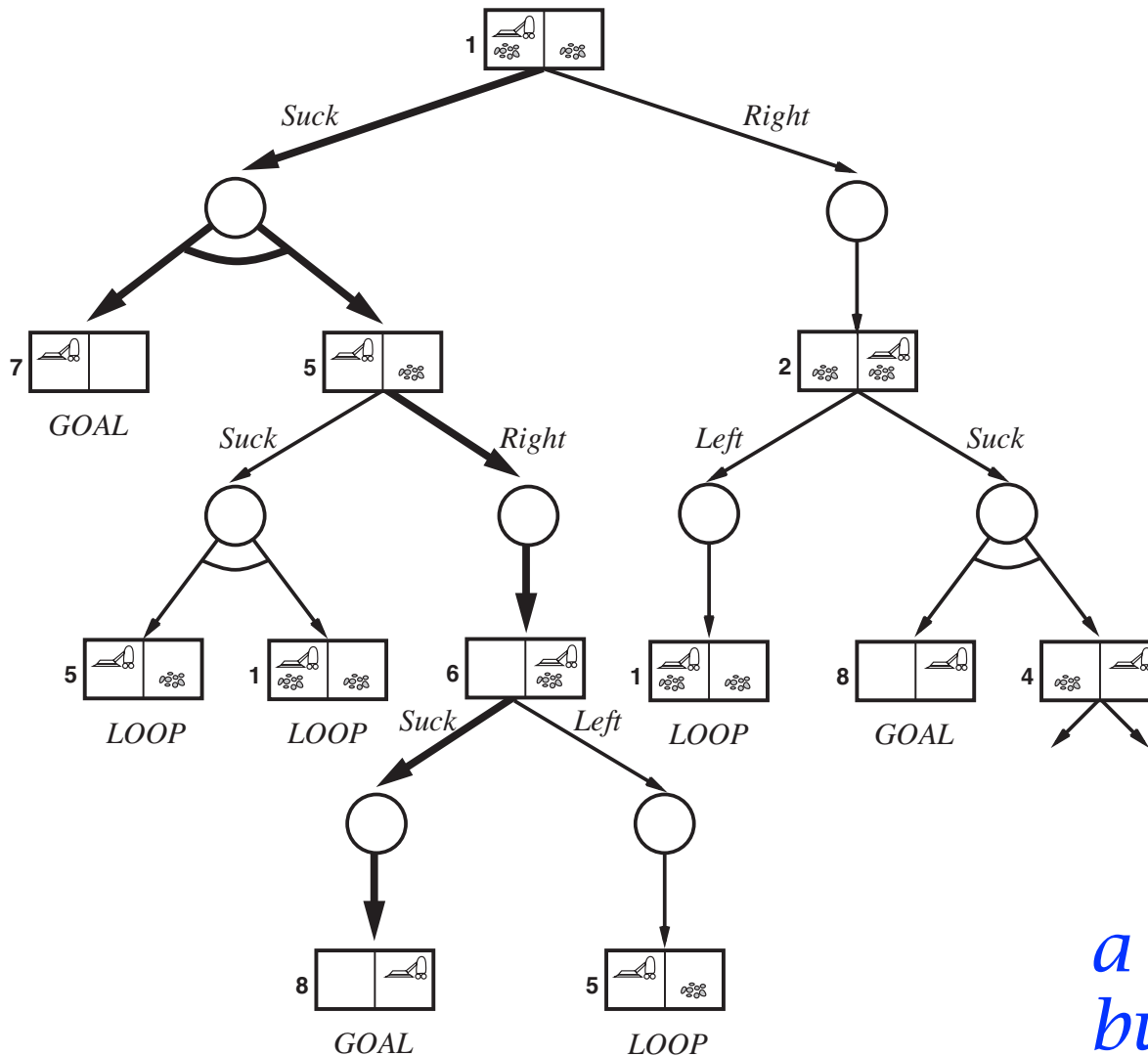


*almost all real-world problems are nondeterministic*  
*how do you solve this problem?*



# AND-OR tree search

OR node: different actions (as usual)  
AND node: different transitions



*a solution is not a path  
but a tree*

# Depth-first AND-OR tree search



**function** AND-OR-GRAPH-SEARCH(*problem*) **returns** a conditional plan, or failure  
OR-SEARCH(*problem*.INITIAL-STATE, *problem*, [])

---

**function** OR-SEARCH(*state*, *problem*, *path*) **returns** a conditional plan, or failure  
**if** *problem*.GOAL-TEST(*state*) **then return** the empty plan  
**if** *state* is on *path* **then return** failure  
**for each** *action* **in** *problem*.ACTIONS(*state*) **do**  
    *plan* ← AND-SEARCH(RESULTS(*state*, *action*), *problem*, [*state* | *path*])  
    **if** *plan* ≠ failure **then return** [*action* | *plan*]  
**return** failure

---

**function** AND-SEARCH(*states*, *problem*, *path*) **returns** a conditional plan, or failure  
**for each**  $s_i$  **in** *states* **do**  
    *plan*<sub>*i*</sub> ← OR-SEARCH( $s_i$ , *problem*, *path*)  
    **if** *plan*<sub>*i*</sub> = failure **then return** failure  
**return** [**if**  $s_1$  **then** *plan*<sub>1</sub> **else if**  $s_2$  **then** *plan*<sub>2</sub> **else** ... **if**  $s_{n-1}$  **then** *plan* <sub>$n-1$</sub>  **else** *plan* <sub>$n$</sub> ]



# Search with no observations



search in **belief** (in agent's mind)

