

Computational Intelligence

Master in Artificial Intelligence

2023-24

Lluís A. Belanche

Introduction to Genetic Algorithms



Soft Computing Research Group



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

A fundamental difference ...

- **Genetic Algorithms (bitstrings)**
- Evolutionary Programming (finite-state automata)
- Evolution Strategies (real-valued vectors)
- Genetic Programming (computer programs)
- Neuro Evolution (neural networks)
- Learning Classifier Systems (if-then rules)

Genetic Algorithms

Canonical Genetic Algorithms (CGAs) are a class of Evolutionary Algorithms for which individuals are bitstrings of length n and:

- Strategy is usually (μ, λ) with $\mu = \lambda$
- Selection is usually proportional to fitness
- Mutation is seen as “transcription error”
(secondary discovery force)
- Recombination is called “crossover”
(primary discovery force)

A **coding function** may be needed to represent potential solutions as bitstrings

Genetic Algorithms

When solving a given optimization problem with a CGA, arguably the first task would be **devising an adequate description/coding of a solution** in terms of a **bitstring**:

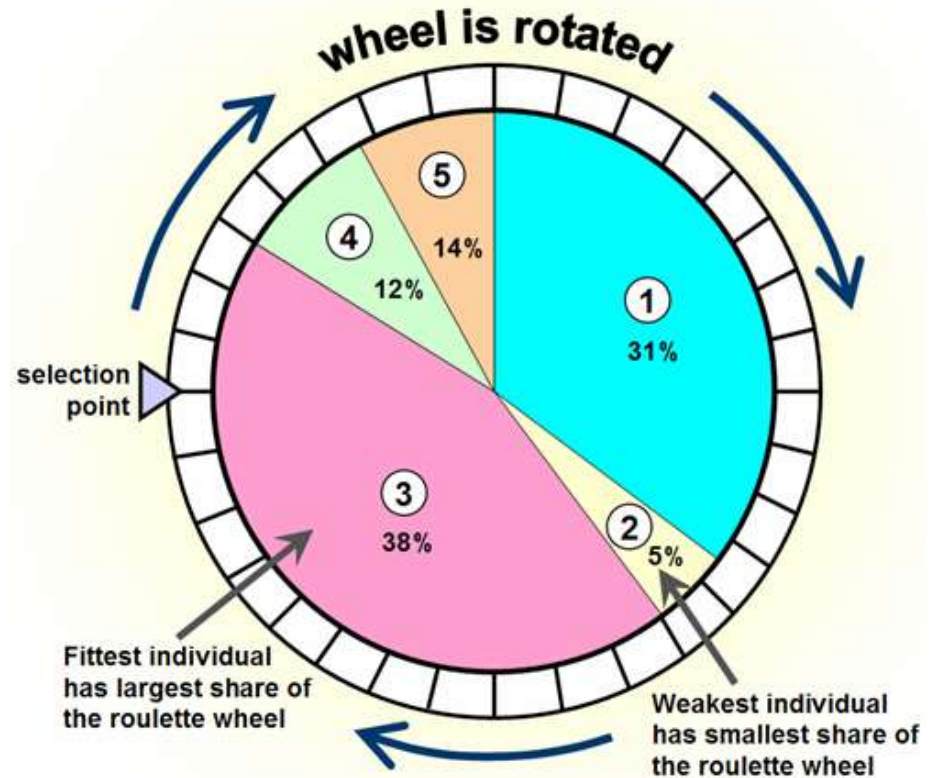
- How many loci? (chromosome length n) Is it scalable?
- What exactly are the genes? Do they come in groups?
- Do we favour low epistasis?
- Is it “continuous”?
- Do we allow polyploidy? (dominant/recessive traits)
- Is it bijective (or, at least, injective)?
- Does it allow/facilitate feasible crossover & mutations?

Selection – first ideas

- Selection completely at **random**
(no benefit for being better than others)
→ little or no exploitation
- Always select (only) the **best**
(very high selective pressure)
→ little or no exploration
- **Stochastic** selection (better fitness → higher chance of reproduction)
(leads to fitness-proportional diversity)
→ more balanced exploitation-exploration searches

Selection: the Roulette Wheel

- Each solution gets a region on a roulette wheel according to its fitness
- Spin the wheel, select solution marked by roulette pointer
- Stochastic selection**
(better fitness = higher chance of reproduction)



Selection: alternatives to Roulette Wheel

Standard Roulette Wheel (SRW) may be very sensitive to the “details” of the fitness function

- Spin wheel once with as many equally-spaced pointers as individuals (stochastic universal sampling), or
- Give a sure number of copies for the above-average individuals and perform standard SRW with remainder (remainder stochastic sampling), or
- Rank selection (probability of selection is proportional to *rank*), which can be used to cope with:
 - High selective pressure in the first few generations
 - Low selective pressure due to a decrease in fitness variance as the search proceeds

Selection/Replacement: Elitism

The best A individuals from the last B generations are passed on *unchanged* for the next generation

Example:

- keep the best individual of past population ($A=1, B=1$)
- “unrealistic”, but ensures best fitness of a generation never decreases
- entails a small decrease in diversity (more exploitation)
- they can be subject to specific local improvement
 - Hill-climbing
 - Specific mutation steps

Selection: Tournament

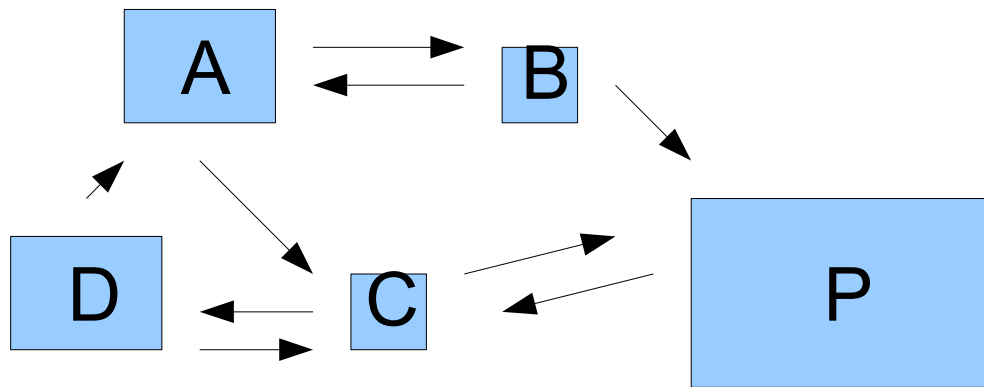
- Randomly select k individuals (with replacement)
- Select the individual with best fitness among these k individuals
- **Example:**
 - randomly pick two individuals and keep the best of the two; do this $\mu-1$ times (these parents will then undergo recombination & mutation)
 - use **elitism** (once) to restore population size μ

Selection: Niching

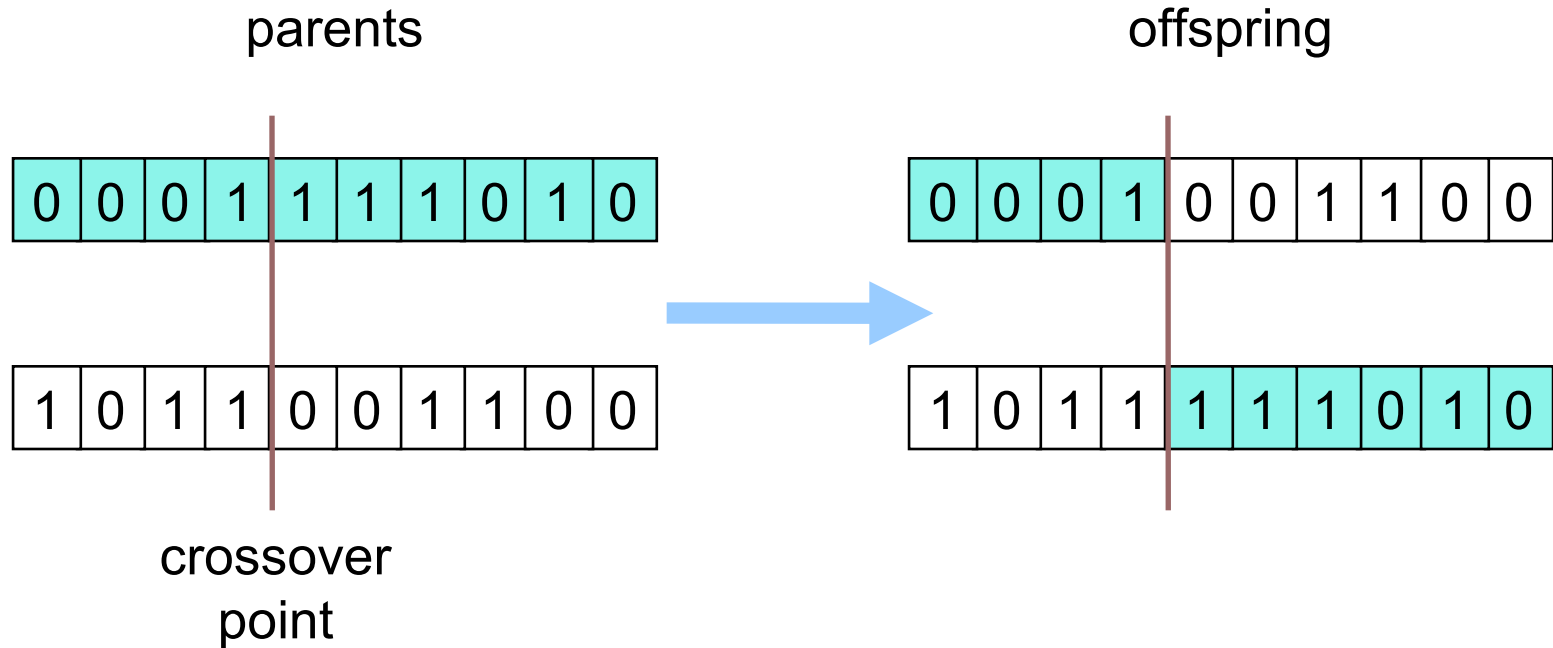
Start with a single population P

Split it into parts (niches) around promising parts of the search space (eg. A,B,C,D)

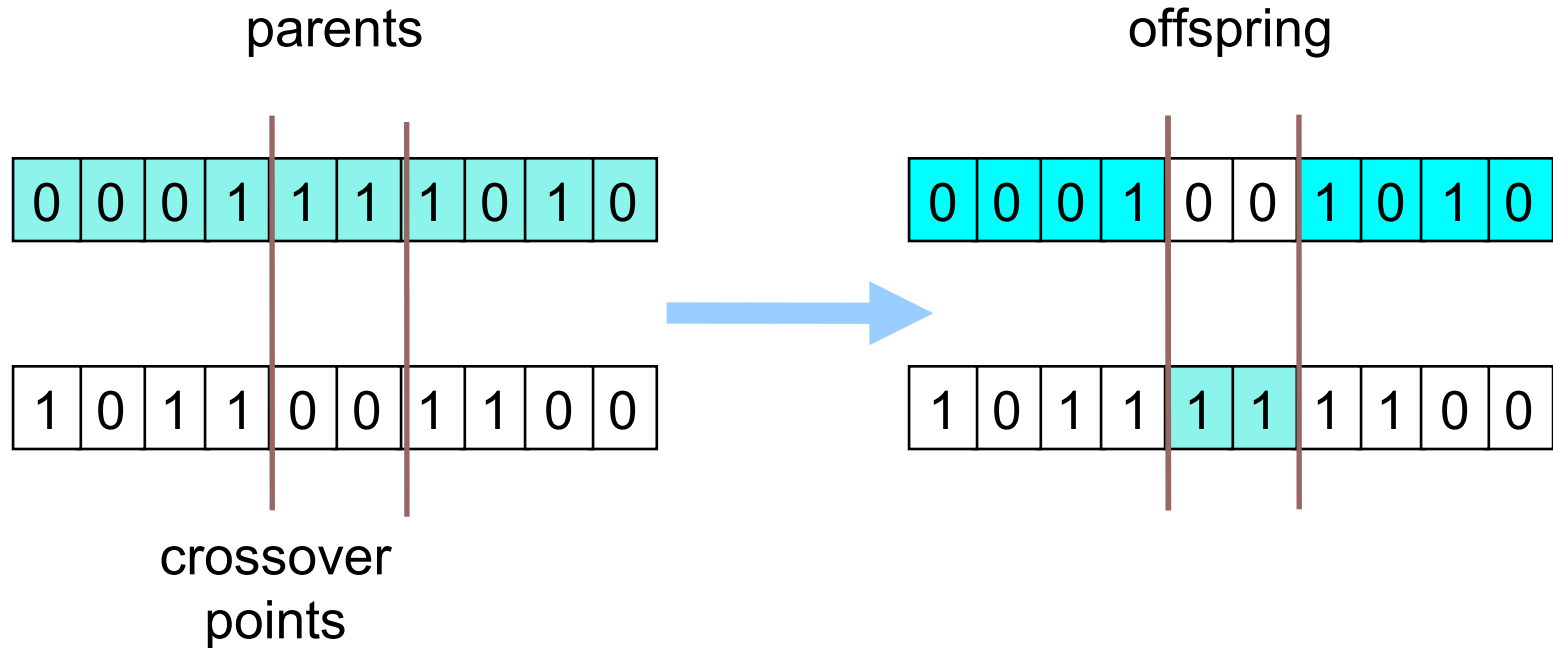
Reproduction takes place within each niche, but letting them interchange information and merge back:



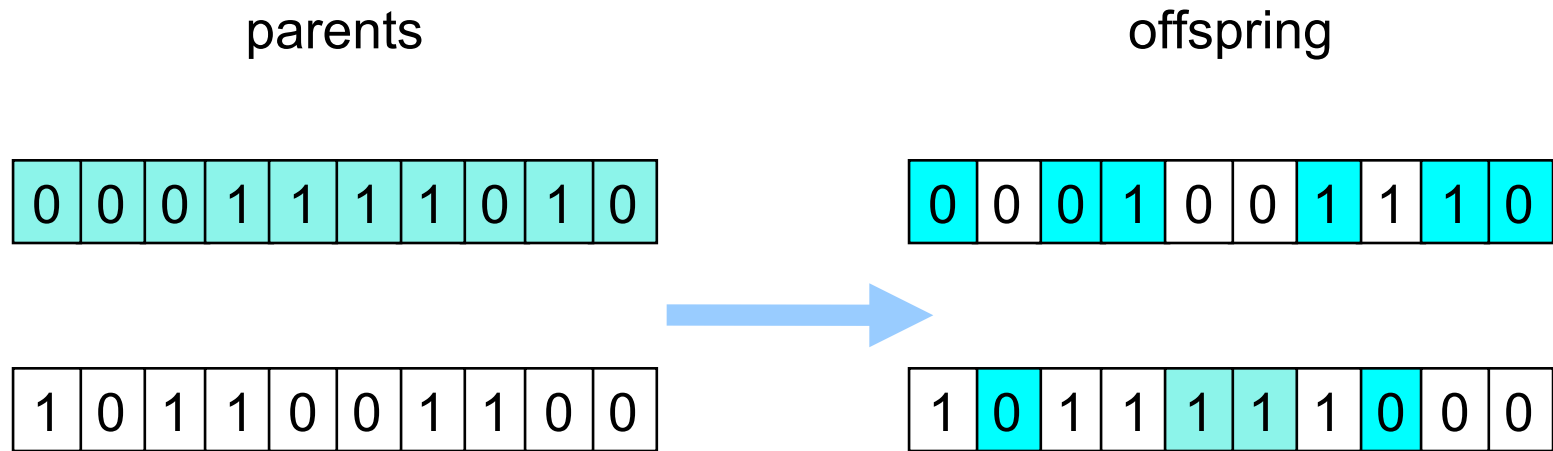
Crossover: One-Point crossover (1P)



Crossover: Two-Point Crossover (2P)



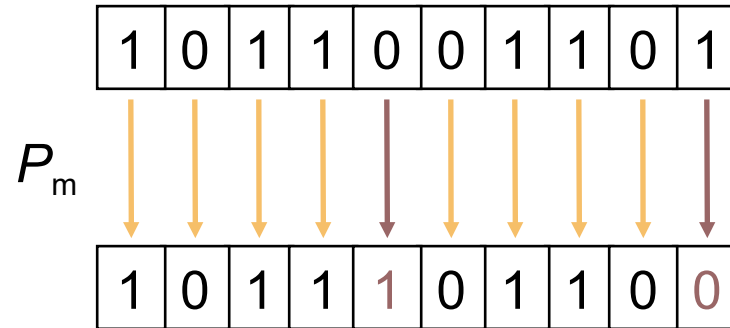
Crossover: Uniform Crossover (UX)



UX evaluates each bit in the parent strings for exchange with a probability of 0.5

HUX (Half-Uniform Crossover) is as UX, but exactly half of the non-matching bits are swapped

Mutation



Binomial = independent Bernoulli “transcription errors”
Default choice is Bin (n , $1/n$)

Mutations occur with some small probability every time a chromosome is duplicated

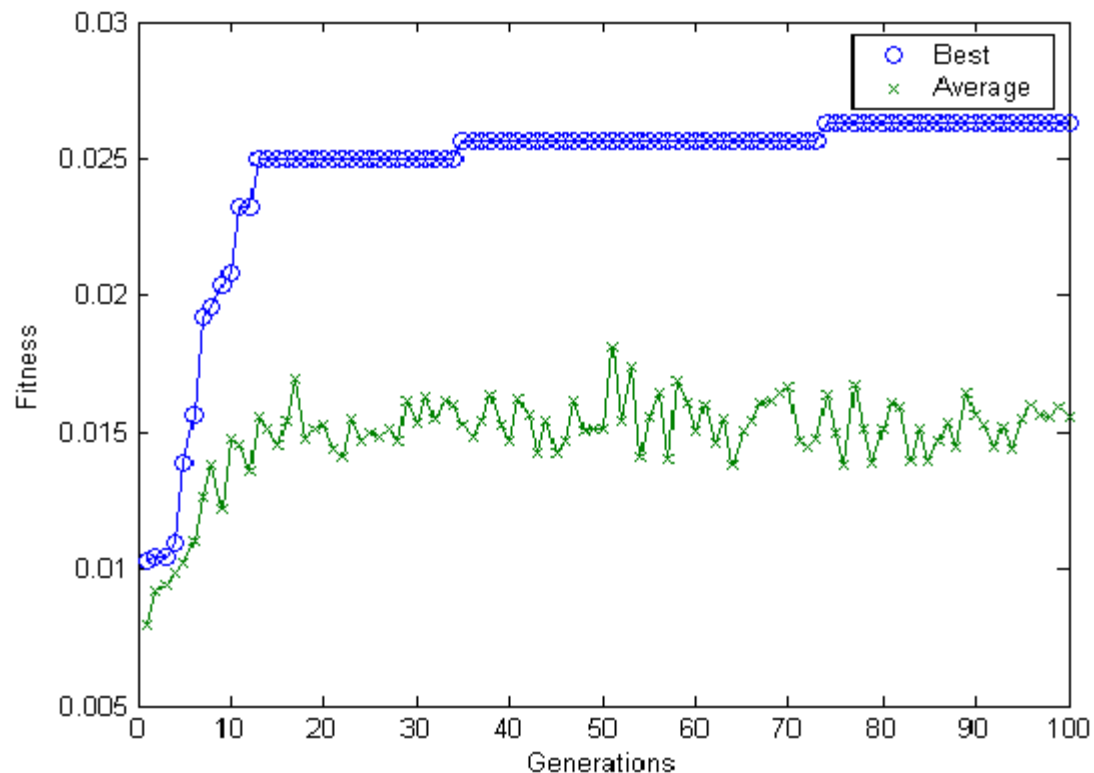
In humans, this probability (the chance that a given nucleotide is mutated in a generation) is around 10^{-8}

Knowing when to stop

- Limit the **number of generations** or fitness evaluations
- Detect **convergence**:
 - No relative improvement in mean/max fitness
 - No change in best (watch out elitism!)
 - Low diversity of fitness evaluations or individuals
 - Fitness evaluation beyond a predefined value
 - Closeness to optimum (if the optimal fitness is known)

Premature convergence

- Population almost converged to (very) suboptimal solution, and
- Not being able to generate better offspring anymore.



Prevention is best option. Detect and counteract (eg. injecting new genetic material, modifying population size, mutation rate, ...)

Replacement strategies (1)

- **Generational** genetic algorithms (GGA):
 - replace all parents with their offspring, (μ, λ) with $\mu = \lambda$
 - no overlap between populations of different generations
- **Steady-state** genetic algorithms (SSGA):
 - Offspring is used to replace some parents in the old generation
 - some overlap exists between populations of different generations

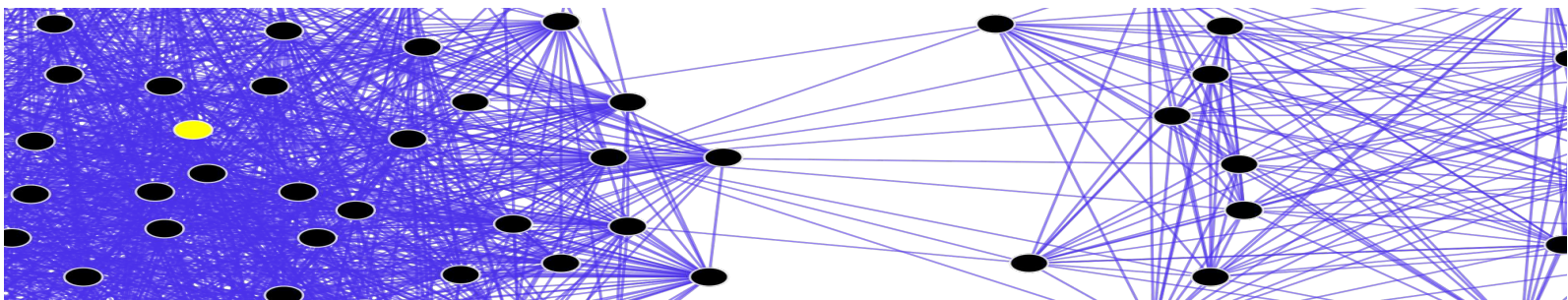
Replacement strategies (2)

- Replacement strategies for SSGA:
 - the offspring replace the **worst** individuals of the current population (elitism)
 - the offspring replace **random** individuals of the current population
 - the individuals to be replaced are selected using «reverse» **tournament selection** (now for the worst one)
 - the offspring replace the **oldest** individuals of the current population

Schools of “thought”

In GAs, originally there were two different views of what a population represents:

- The **Michigan** view: individuals represent parts of the solution; the whole population represents a full solution
- The **Pittsburgh** view: each individual represents a full solution (this view has prevailed and is the common approach)
- **Example:** evolving neural networks



Binary representations

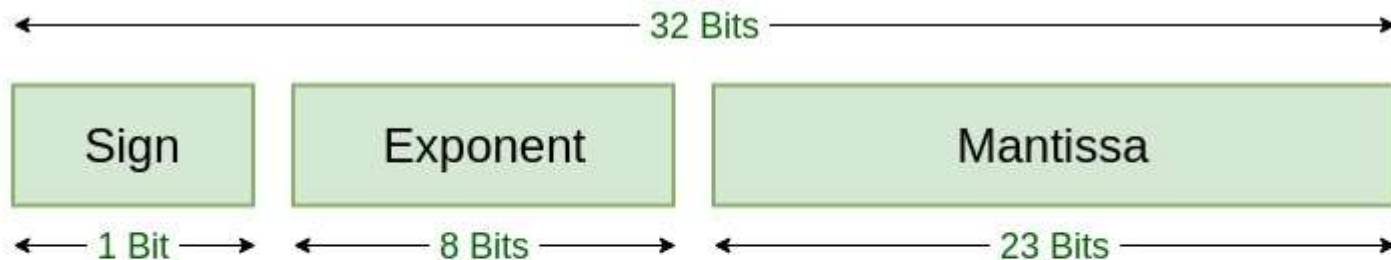
1010101110011110100000101000100110101001010 ?

The (variables in the) optimization problem may:

- admit a natural binary representation
- be given by nominal-valued variables (aka polychotomous)
- be given by a discrete structure (tree, graph, ...)
- be given by integer or rational-valued variables
- be given by continuous information (real numbers)

Binary representations for reals

Is it really a good idea?



Single Precision
IEEE 754 Floating-Point Standard

Binary representations for reals

Example. $85.125 = 1010101.001 = 1.010101001 \times 2^6$

biased exponent $127+6 = 133 = 10000101$

Normalised mantisa = 010101001, sign = 0

add 0's to complete the 23 bits

The IEEE 754 Single precision is:

0 10000101 01010100100000000000000

(hexadecimal **42AA4000**)

Binary representations for reals

Standard code for real numbers

We want to represent real numbers $r \in [a, b]$ with an attainable precision ε as $\text{enc}(r)$, s.t. $|r - \text{enc}(r)| < \varepsilon$:

- ① We create 2^k segments with $k = \lceil \log_2 \frac{b-a}{2\varepsilon} \rceil$, which are mapped to the numbers $0, \dots, 2^k - 1$
- ② $\text{enc}(r) := \lfloor \frac{r-a}{b-a}(2^k - 1) + \frac{1}{2} \rfloor$

Decoding is performed by $a + \text{enc}(r) \frac{b-a}{2^k-1}$.

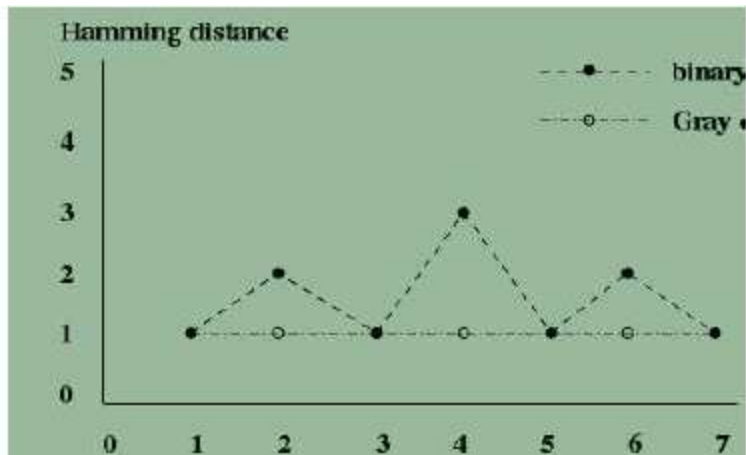
Example: $r = 0.637197$ in $[-1, 2]$ with precision $\varepsilon = 10^{-6}$

Solution: **1000101110110101000110**

Binary representations

Problem using this binary representation:

- 1) maximum attainable precision may be less than needed to represent the optimum
- 2) introduction of «Hamming cliffs»



Hamming cliff is formed when two numerically adjacent values have bit representations that are far apart

$$7_{10} = 0111_2 \text{ vs } 8_{10} = 1000_2$$

Large change in solution is needed for small change in fitness

Gray Coding

In a **Gray coding**, the Hamming distance between the representation of consecutive numerical values is 1

Example for $n_b=3$:

	Binary	Gray
0	000	000
1	001	001
2	010	011
3	011	010
4	100	110
5	101	111
6	110	101
7	111	100

Converting binary strings to Gray bit strings:

$$g_1 = b_1$$

$$g_l = b_{l-1}\bar{b}_l + \bar{b}_{l-1}b_l$$

where b_l is bit l of the binary number

$$b_1b_2\cdots b_{n_b}$$

How a GA *might* work: schemata theory

A **schema** S is a template describing a subset of individuals:

*	*	*	1	0	*	1	*	*	*
---	---	---	---	---	---	---	---	---	---

wildcard symbol: *

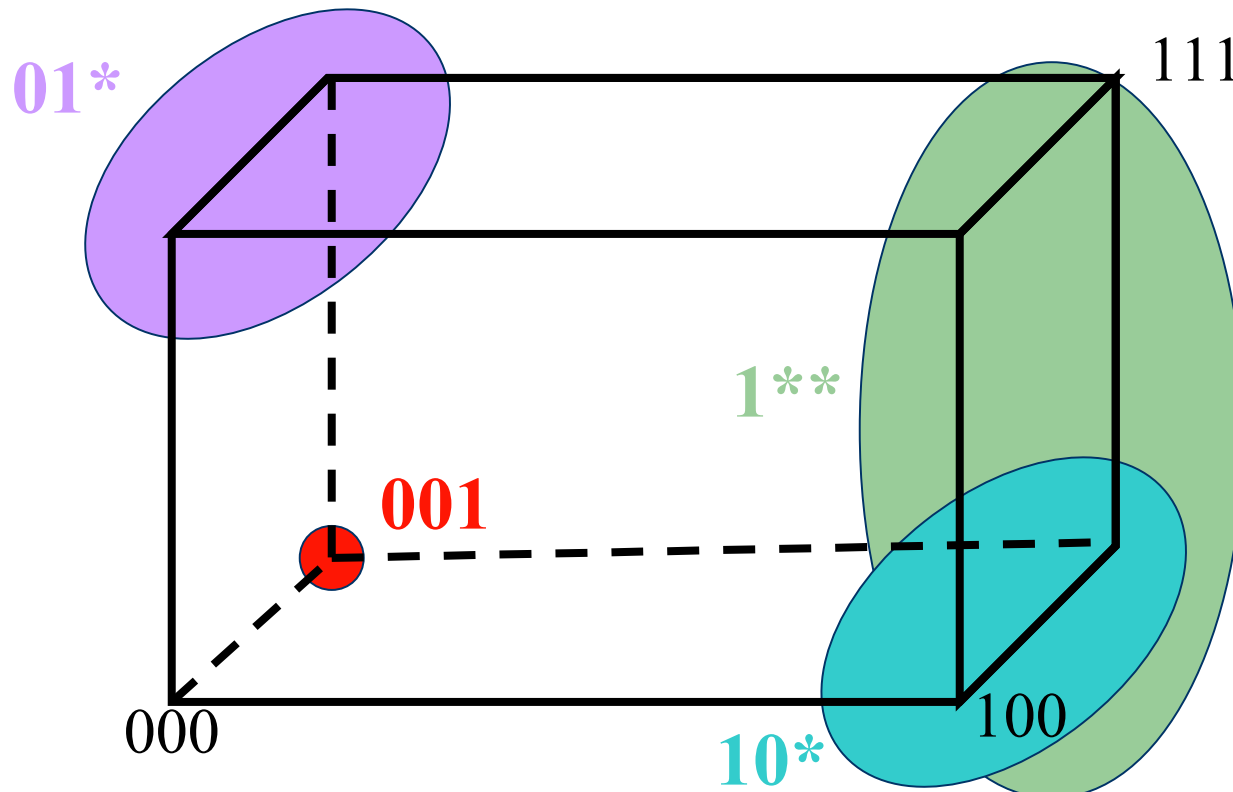
order $\omega(S) = \# \text{ fixed positions}$

defining length $\delta(S) = \text{distance between first and last fixed positions}$

- there are 3^n possible schemata
- there are $\text{Choose}(\delta, \omega) \cdot 2^\omega$ schemata of order ω and defining length δ
- a schema S includes (represents) $2^{n - \omega(S)}$ individuals
- an individual of length n belongs to (i.e. *represents*) 2^n schemata
(this result does not depend on the alphabet being binary)

Schemata as hyperplanes

A **schema** (*pl.* schemata) is a string in the ternary alphabet $\{0,1,*\}$ representing a hyperplane in the search space



Implicit Parallelism

In a population of μ individuals of length n :

$$2^n \leq \text{number of processed schemata} \leq \mu \cdot 2^n$$

→ far more schemata than individuals are *implicitly* processed (not disrupted by crossover and mutation) [Holland 1989]:

Lower bound of the order of $\mu^3 / \sqrt{\log_2 \mu}$

--see Bertoni & Dorigo (1993) "Implicit Parallelism in Genetic Algorithms"
Artificial Intelligence **61**(2), p. 307–314

Schemata interpretations

- **Geometric** view: hyperplanes in $\{0,1\}^n$
- **Evolutionary** view: what survives is not the individuals, but their genetic stuff (the schemata)

Thus the GA is a **schema processor**

The Schema theorem

Theorem. Short, low-order, above-average schemata receive exponentially increasing trials (individuals represented by the schema) in subsequent generations.

Schema theory is (generally) accepted to be (essentially) correct, providing an intuitive explanation for the good (and bad) performance observed in practical GA applications

Theoretical work on GAs includes:

- Markov chains as an alternative formalism
- Synthetic search problems (e.g. Walsh functions)
- Effects of genetic operators on exploration/exploitation
- Studies on (kinds and effect of) deceptive problems

The Building Blocks Hypothesis

Closely related to the Schema theorem is the Building Block hypothesis (BBH):

“Genetic Algorithms work by discovering and exploiting **building blocks** --groups of closely interacting genes-- and then combining these blocks (via crossover) to produce successively larger blocks until the problem is solved.”

D. Goldberg. *Genetic Algorithms in Search, Optimization and Machine Learning*.
Addison-Wesley Boston, MA, USA 1989.

--see *An Overview of Schema Theory*, D. White.
arxiv.org/abs/1401.2651 for a modern analysis.

Deception

A problem is said to be **deceptive** if the building blocks identified actually lead the GA away from the global optima.

Example. fully deceptive order-3 problem: all fitness evaluations of order-1 and order-2 schemata lead the search away from the global optimum (111), as given by evaluation of order-3:

$$f(0**) > f(1**)$$

$$f(*0*) > f(*1*)$$

$$f(**0) > f(**1)$$

$$f(00*) > f(11*), f(01*), f(10*)$$

$$f(0*0) > f(1*1), f(0*1), f(1*0)$$

$$f(*00) > f(*11), f(*01), f(*10)$$

$$f(000) = 28$$

$$f(010) = 22$$

$$f(110) = 0$$

$$f(101) = 0$$

$$f(001) = 26$$

$$f(100) = 14$$

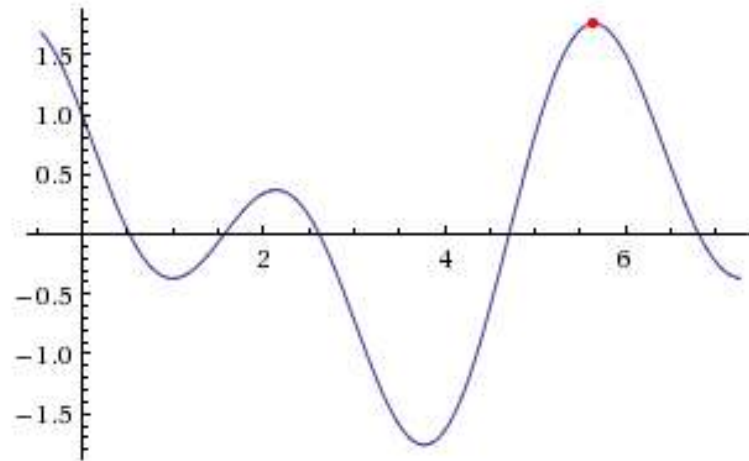
$$f(011) = 0$$

$$f(111) = 30$$

Examples in function optimization

Let $f(x) = \cos(x) - \sin(2x)$

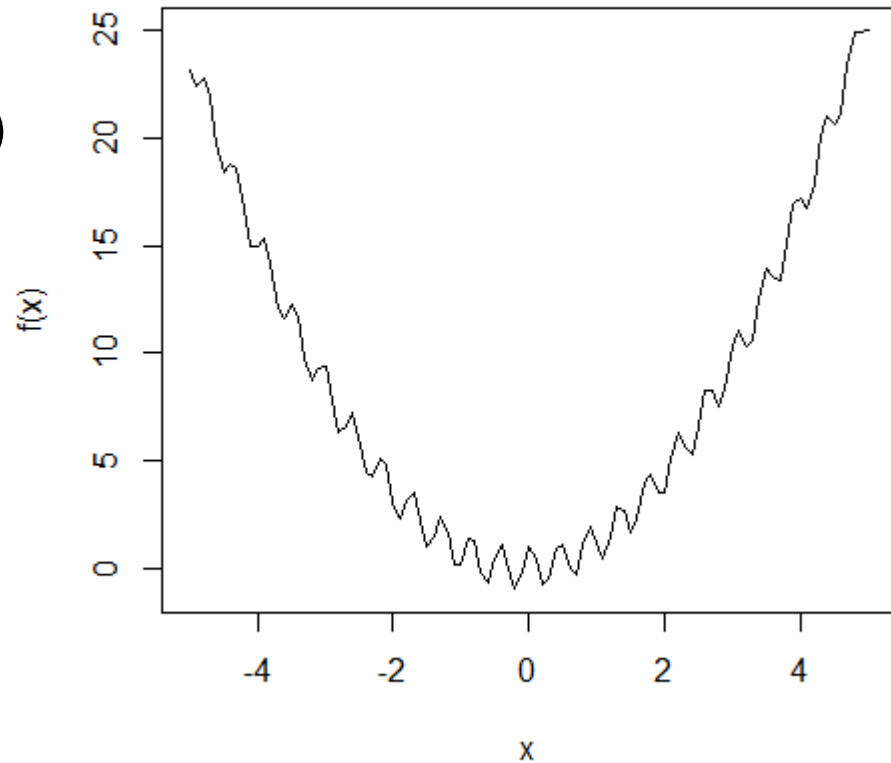
(to be maximized in $[0, 2\pi]$)



Examples in function optimization

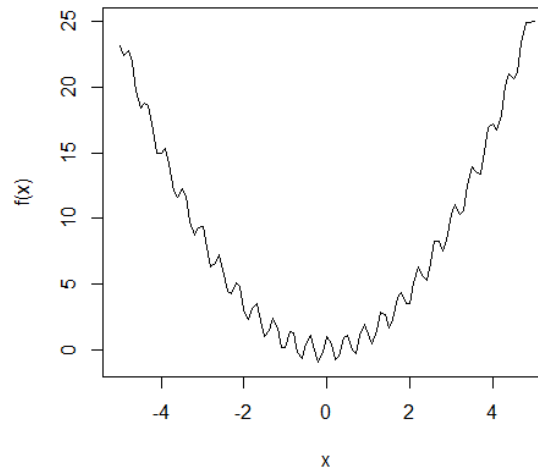
Let $f(x) = ?$

(to be minimized in $[-4, 4]$)



Examples in function optimization

Let $f(x) = (x+0.2)*x + \cos(14.5*x-0.3)$
(to be minimized in $[-4, 4]$)

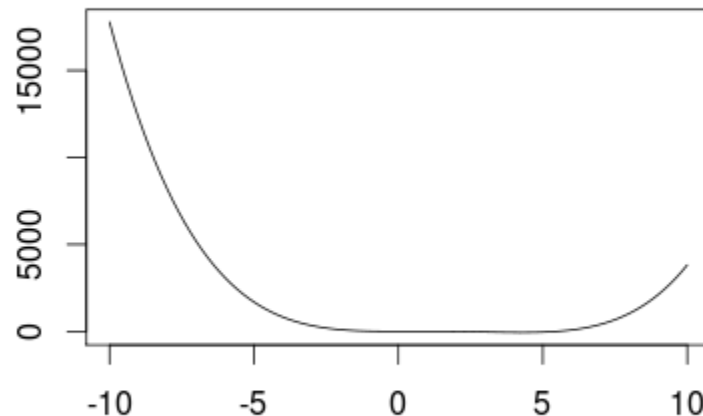


Examples in function optimization

Finding the roots of a polynomial $p(x)$ in an interval

Fitness function: maximize $F(x) = -p(x) \cdot p(x)$

Example: $p(x) = x^4 - 7x^3 + 8x^2 + 2x - 1$ in $[0, 10]$

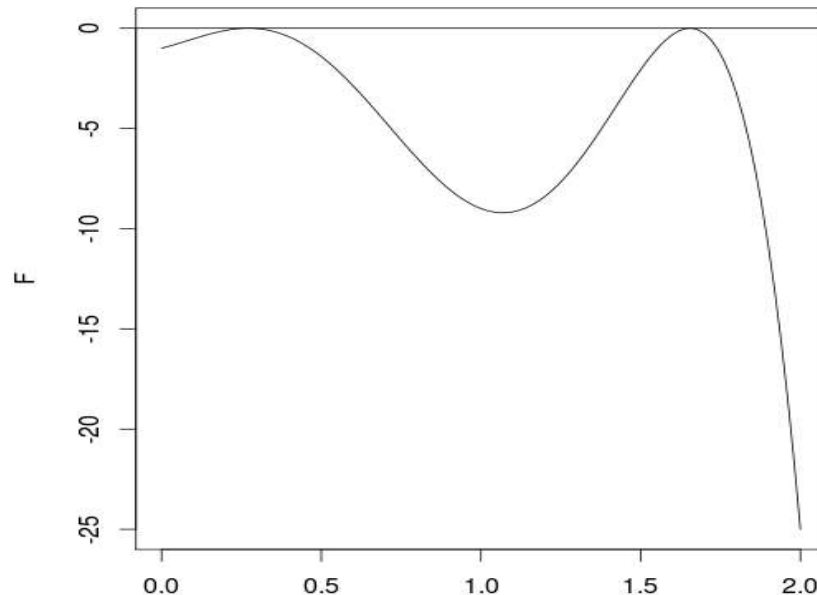


Examples in function optimization

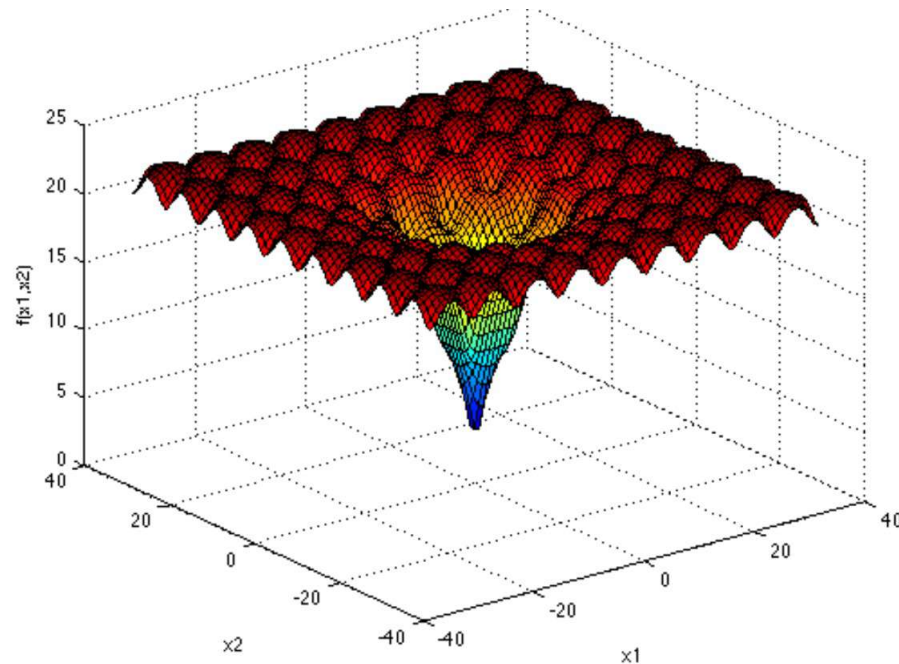
Most methods need a **range** for the possible values

Example: $p(x) = x^4 - 7x^3 + 8x^2 + 2x - 1$ in $[0, 10]$ can be constrained to within $[1/9, 9]$:

- $x_1 \approx 0.27190$
- $x_2 \approx 1.65434$



Examples in function optimization



$$f(\mathbf{x}) = -a \exp \left(-b \sqrt{\frac{1}{d} \sum_{i=1}^d x_i^2} \right) - \exp \left(\frac{1}{d} \sum_{i=1}^d \cos(cx_i) \right) + a + \exp(1)$$

https://en.wikipedia.org/wiki/Test_functions_for_optimization

<https://www.sfu.ca/~ssurjano/optimization.html>

Example: the (iterated) Prisoner's Dilemma

- The 'Prisoner's dilemma game' was invented by Merrill Flood & Melvin Dresher in the 50s
- Albert Tucker formalized the game with prison sentence rewards and gave it the name "prisoner's dilemma"
- Studied in game theory, economics, political science, cold war
- The story:
 - Alice and Bob get **arrested**, no communication between them
 - They are offered a **deal**:
 - If only **one** person accuses the other then he/she gets suspended sentence while the other gets 5 years in prison
 - If **both** confess & testify against the other, they both get 4 years
 - If **none** of them confesses then they both get 2 years

What is the best strategy for maximising one's own payoff?

Example: the (iterated) Prisoner's Dilemma

- Humans play against each other with the goal of maximising one's payoff in the long run
- A winner strategy is **TIT-FOR-TAT***:
 - Cooperates as long as the other player does not defect
 - Defects on defection until the other player begins to cooperate again
 - Can a GA discover this strategy?
 - Can a GA evolve a better strategy?

(*) actions done intentionally to punish other people because they have done something unpleasant to you
(<http://dictionary.cambridge.org>)

Example: the (iterated) Prisoner's Dilemma

Define a representation of cases:

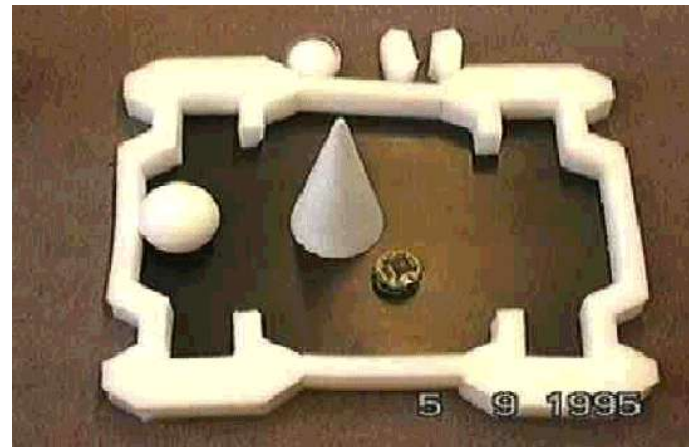
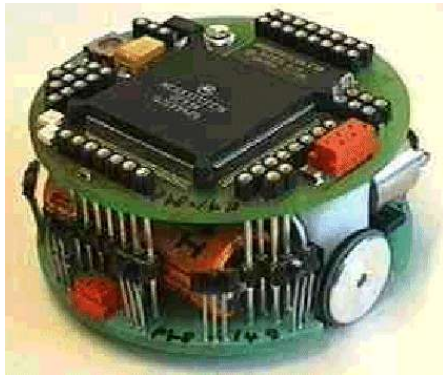
Case 1: CC	C	
Case 2: CD	D	← encoding of TIT FOR TAT
Case 3: DC	C	
Case 4: DD	D	

● Results

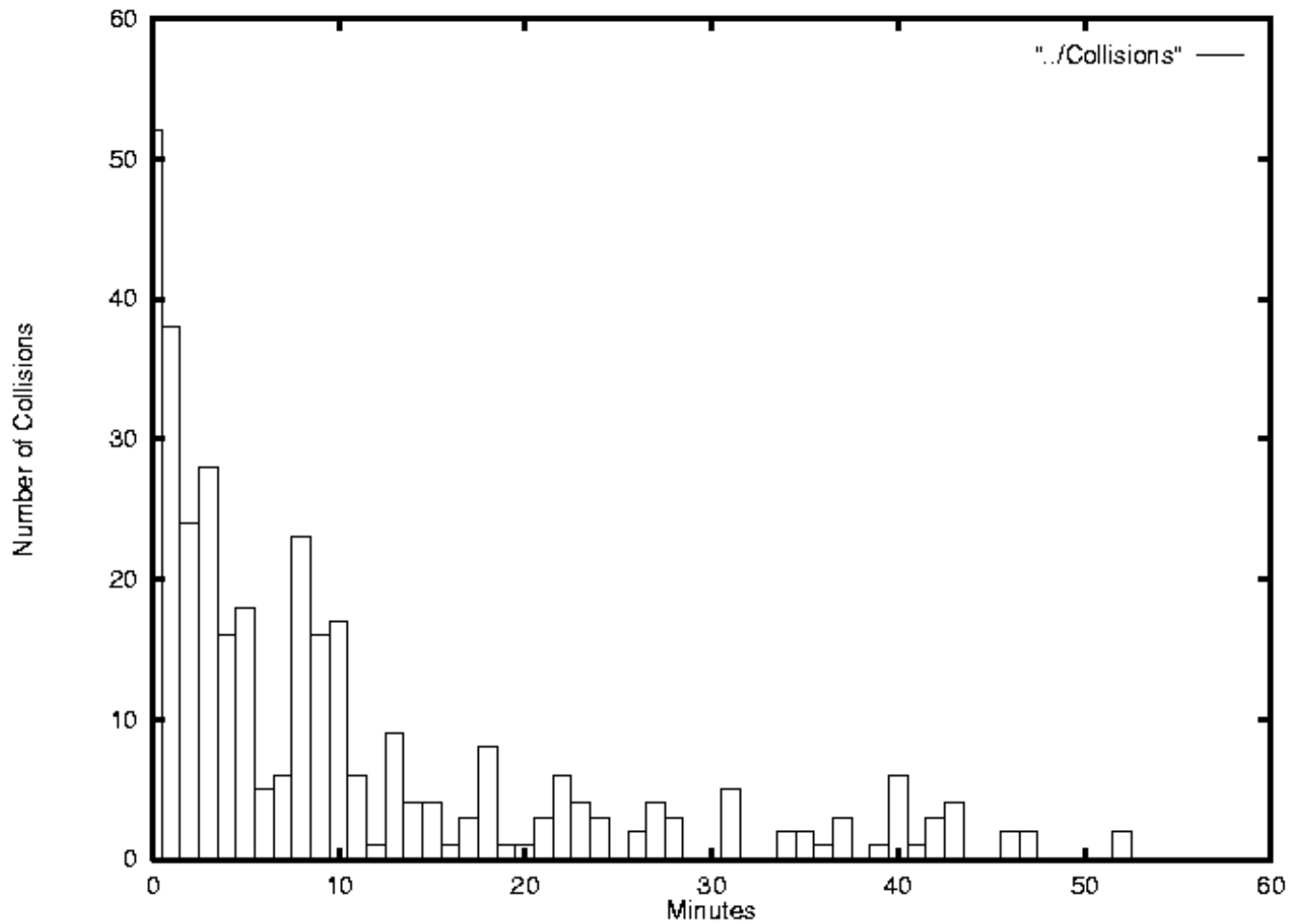
- For a environment formed by 8 human-designed strategies, a GA found a better strategy (i.e., highly adapted to *that* environment) [memory of 3 previous games, 40 runs, 50 generations each, population of 20, fitness: average score over all games played]
- Experiments in *changing* environments: the evolving strategies played against each other: found strategies similar to the winner human-designed strategies
- Idealized model of evolution & co-evolution

Example: Evolving programs for a small mobile robot

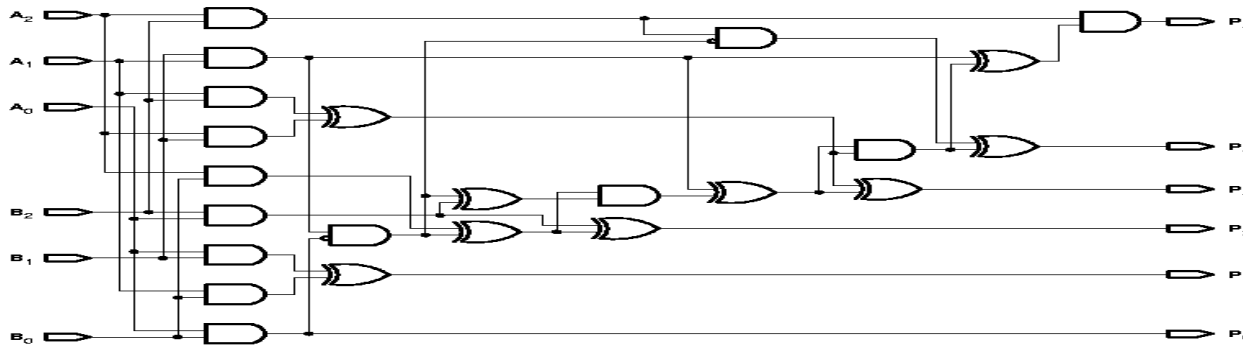
- **Goal:** obstacle avoidance
- Inputs from eight sensors on robot $\{s_1, \dots, s_8\}$ with values in $\{0, \dots, 1023\}$ (higher values mean closer obstacle)
- Output to two motors (speeds) with values in $\{0, \dots, 15\}$



Results

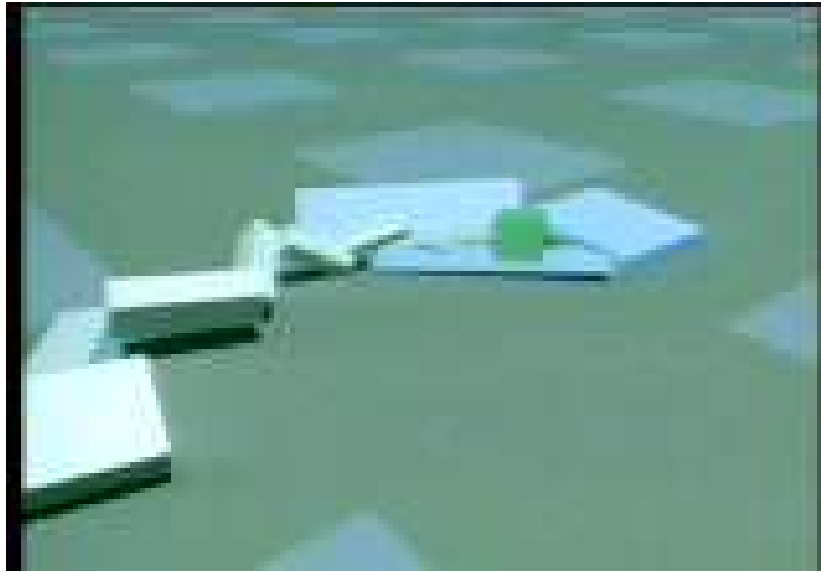


Example: Electronic circuit design



- Individuals are programs that transform an initial circuit to a final circuit by adding/subtracting components and connections
- Fitness: computed by simulating the circuit
- Population of 640,000 individuals run on a parallel processor
- After only 137 generations, some discovered circuits exhibit performance competitive with best human designs

Example: Karl Sims' virtual creatures



<http://www.karlsims.com/evolved-virtual-creatures.html>

Key ideas (1)

- The encoding should respect the schemata: it must allow discovery of small building blocks from which larger, more complete solutions can be formed
- The encoding should reflect functional interactions, as proximity on the genome (*linkage bias*)
- You should devise appropriate genetic operators:
 - all individuals correspond to feasible solutions and vice versa
 - genetic operators preserve feasibility

Key ideas (2)

- High flexibility and adaptability because of many options:
 - Problem representation
 - Genetic operators with parameters
 - Mechanism of selection
 - Size of the population
 - Fitness function
- These decisions are highly problem-dependent
- Parameters are not independent: you cannot optimize them one by one
- Parameters can be adaptable, e.g. high in the beginning (more exploration), going down (more exploitation), **or even be subject to evolution themselves!**