
Using Evolutionary Algorithms To Solve Hard Problems (ralph.morelli@trincoll.edu)

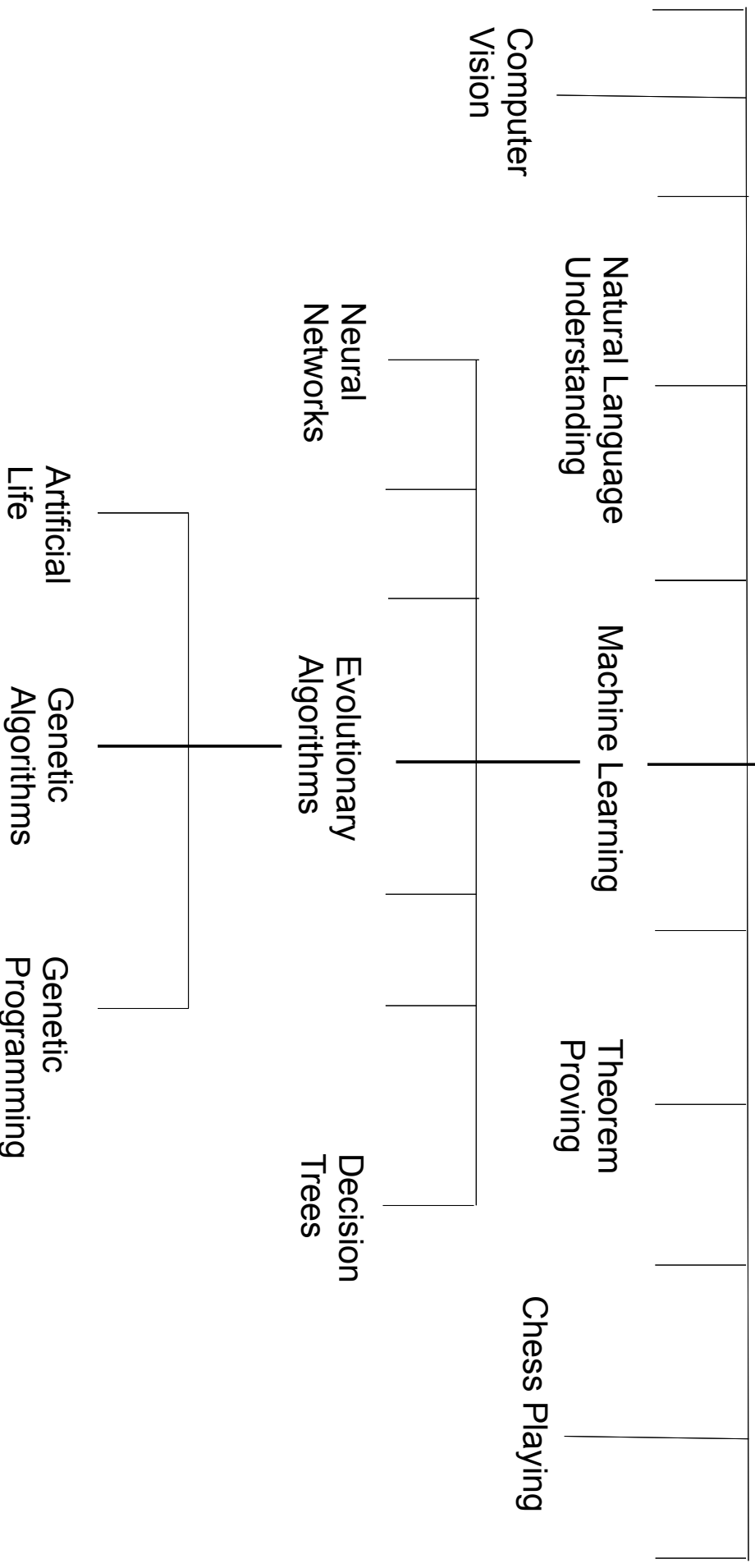


Outline

- What are Genetic Algorithms.
- Some examples from different domains.
- Themes
 - GA: inspired by natural selection.
 - Simple model: Random variation plus selection.
 - Complexity: emergent property.
- Implications for secularism? (discussion)
 - Matches without watch makers.

What are GAs?

Artificial Intelligence



GAs Competitive with Humans

- AI researchers keep score (Turing Test)
- 1997: Deep blue beat Kasparov.
- <http://www.genetic-programming.com/humancompetitive.html>
- Some examples from Genetic Programming (out of 36)
 - Creation of quantum algorithm for 'early promise' problem.
 - Creation of quantum algorithm for Grover's DB search.
 - Rediscovery of Cauer elliptical topology for filters.
 - Synthesis of a NAND circuit.
 - Rediscovery of negative feedback.

Solving Short Cryptograms

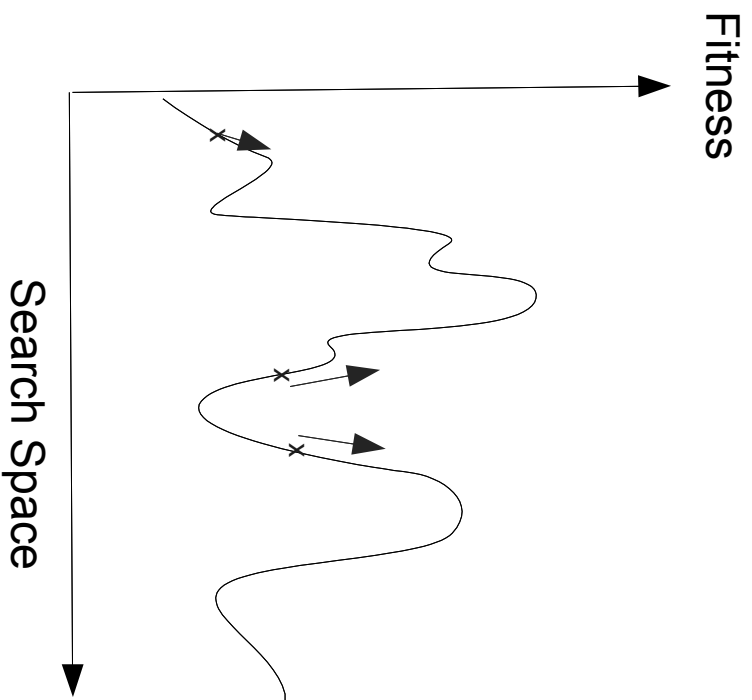
tmlk sceq adqeq e dj blekb sl qclv ylt clv d atkrc ln
apebcs yltkb nlhgq wew nekw d rcdjmelk; d jdk vesc alvq
dkw bephq ln ceq lvk.

- Short cryptograms are 'hard'.
 - Simple substitution cipher: a/e, b/g, c/t, ..., x/c, y/m, z/j
 - Key: A permutation of the alphabet, a..z.
ZWUSQOMKIGECABDFHJLNPRTVXY
 - Possible keys: $26! = 4.0329 * 10^{26}$
 - Unsolvable if too short
 - Theoretical minimum length ~ 28 characters
 - To we or not to we that is the question
 - Xabcx = River = Raver = Rover = Saves = David ...

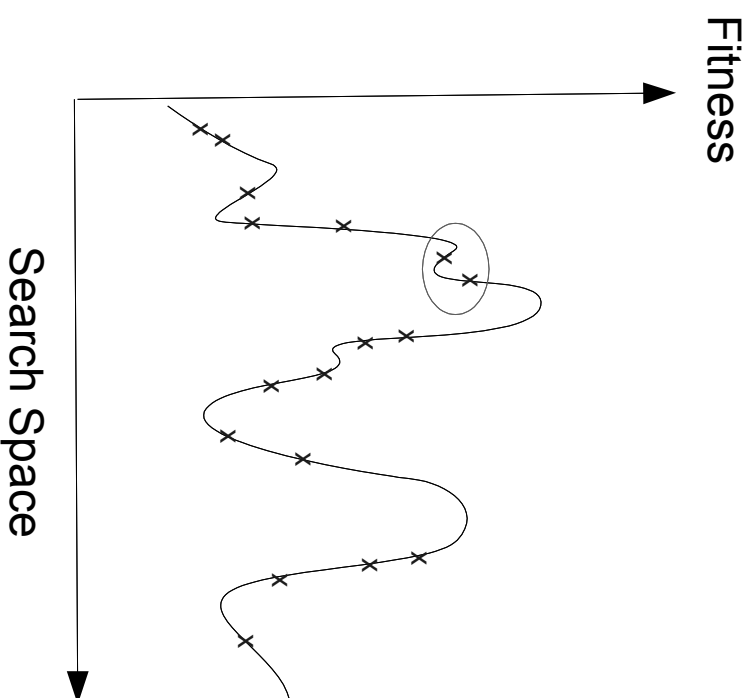
Searching a Problem Space

- *Algorithmic search*
 - Try every key (brute force).
 - Guaranteed to work (theoretically).
 - Intractable: $10^{26}/10^9 = 10^{17}$ seconds $\approx 10^9$ years
- *Heuristic search*
 - Use a rule of thumb to reduce search space.
 - Not guaranteed to succeed.
 - E.G. Traveling salesman nearest neighbor heuristic.
 - E.G. Hill climbing optimization.
 - E.G. Genetic algorithm.

Hill Climbing vs GA

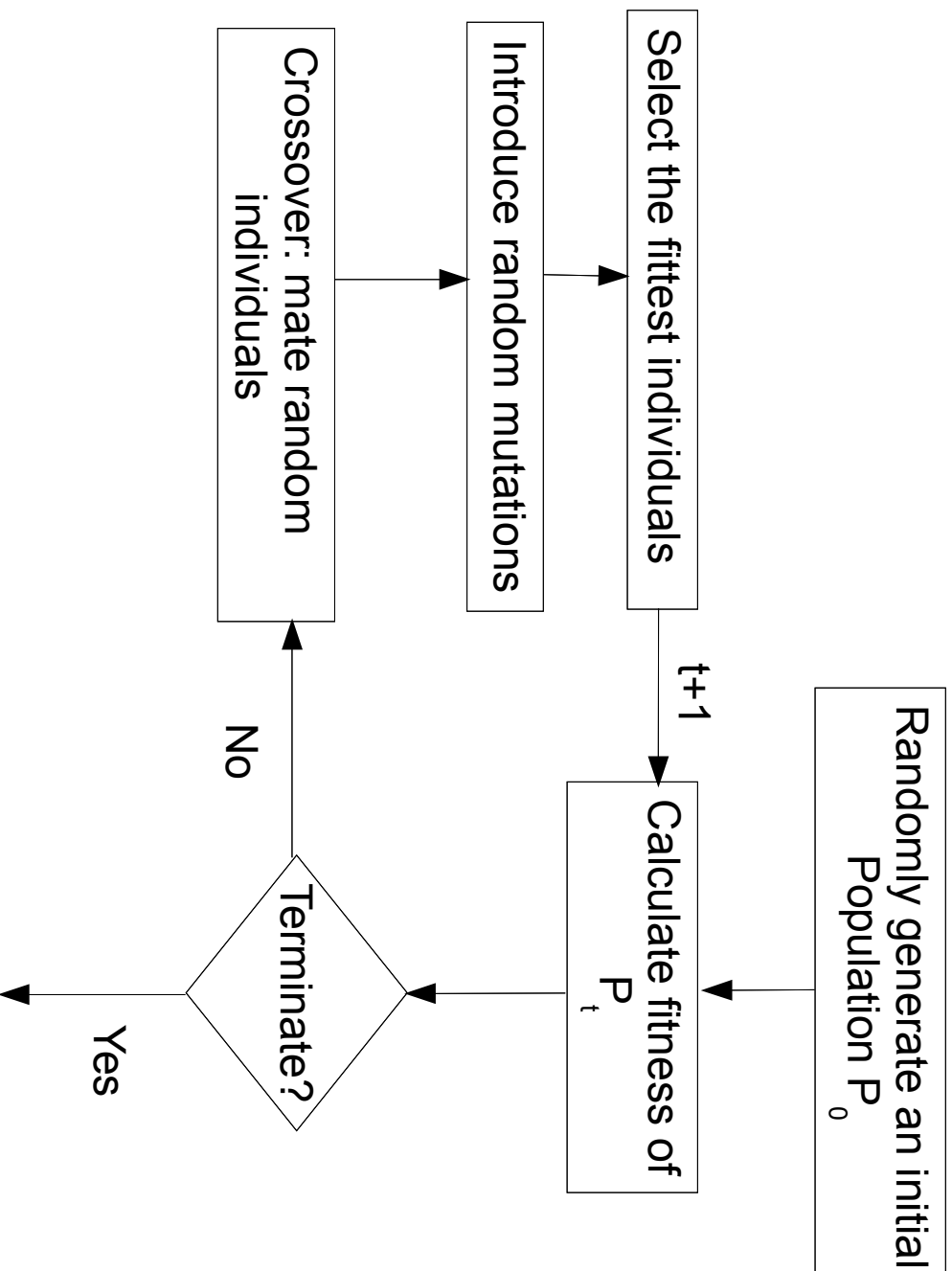


Hill Climber



Genetic Algorithm

The Genetic Algorithm



GA Basics

Population of Chromosomes

```
01101011010010101101000110
11011010010001110000101010
11001010010100010001000010
01010010001000101010100010
10100000100001100001010100
01010101010001000001011000
```

0101001000100 | 01010101000010

Crossover

0110101101001 | 0101101000110

0101001000100 | 0101101000110

Mutation

0101001000100 | 0101100100110

$$f(0101001000100 | 0101100100110) = N$$

Fitness:

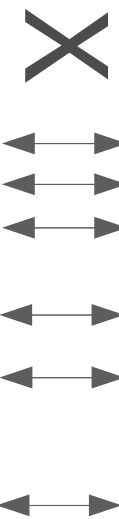
$$f(0101001000100 | 0101100100110) = N$$

Cryptogram: Crossover and Mutation

ABCDEFGHIJKLMNOPQRSTUVWXYZ

SUQHCLADFGJKNOP~~TVWX~~ZBEMRIY

Parent1 Chromosome (words: and)



AZBYXCWDE~~VUFG~~TSHRIQJPKLOMN

Parent2 Chromosome (words: big)

SZQHCLWDEGJKNOP~~TVAX~~UYFMRI~~B~~

Child1 Chromosome (words: and big)

SZBHXCWDE~~VUFG~~OAYRIQJPKLTMN

Child2 Chromosome (words: and big)

f(SZBHXCWDE~~VUFG~~OAYRIQJPKLTMN) = " ... gx ... and ny cd ... big ... "

Results

Parameter	Setting
Population size	512
Seed dictionary size	50
Fitness dictionary size	3500
Crossover rate	0.9
Mutation rate	0.2

msg	nToks	nWds	%Found	nGens	nTrials
1	24	19	97.4	21.1	50
2	33	22	99.5	16.2	50
3	23	14	61.5	46.9	50
4	23	12	58.6	52.4	50
5	26	17	95.7	16.1	50
6	24	22	93.1	57.8	50
7	25	16	74.0	46.2	50

Table 3: Performance results for seven cryptograms.

Multiple Sequence Alignment

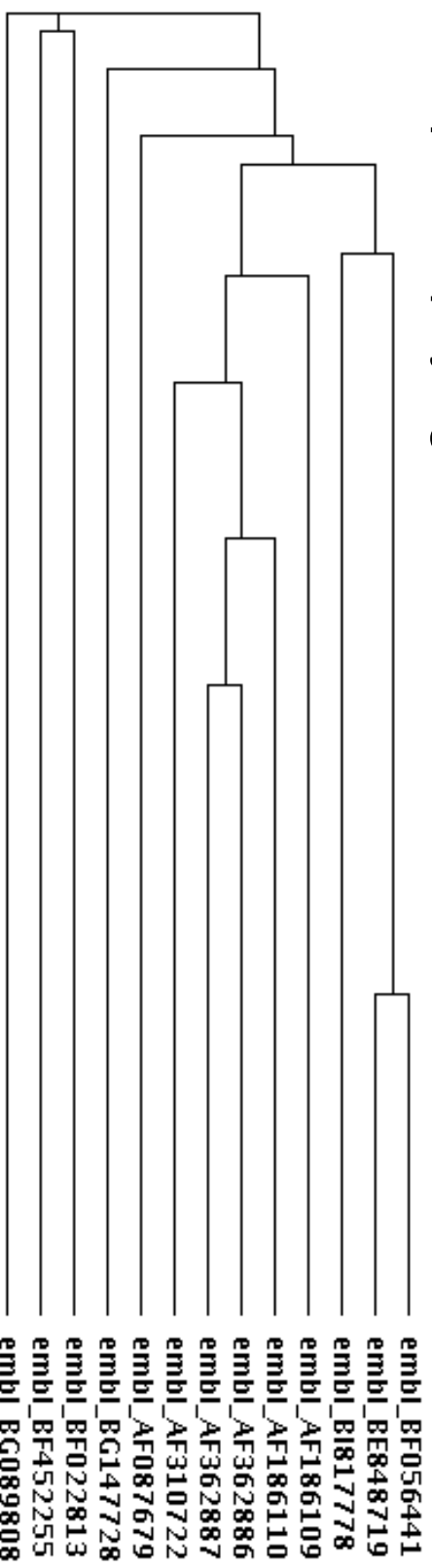
- Fundamental molecular biology problem:
 - Identify common structure in a string of nucleotides (DNA, RNA) or amino acids (in proteins).
- For “deciphering”
 - Evolutionary history
 - Phylogenetic relationships among organisms.
- Very hard ~ on the order of m^n operations for n sequences of length m .

Clustal-W Example

- From a set of DNA sequences such as these:

```
emb1_BF022813  ATGGCCCGCCCTCAACTCACTGGAGGCAGTGAAGCCGCAAGATCCAGGCCCTGCAGCAGCAG  86
emb1_BF452255  ATGGCCCGCCCTCAACTCACTGGAGGCAGTGAAGCCGCAAGATCCAGGCCCTGCAGCAGCAG  83
emb1_BG089808  ATGGCCCGCCCTCAACTCACTGGAGGCAGTGAAGCCGCAAGATCCAGGCCCTGCAGCAGCAG  94
emb1_BG147728  -----CAACTCACTGGAGGCAGTGAAGCCGCAAGATCCAGGCCCTGCAGCAGCAG  49
emb1_AF087679  ATGGCCCGCCCTCAACTCCCTGGAAGCGGTGAAAACGCAAGATCCAGGCCCTGCAGCAGCAG  60
emb1_AF362886  -----
emb1_AF362887  -----CGAGAAGTTGAGGGAGAAAGCGCGGCC  27
emb1_AF186110  ATGAAGGATGAGGAGAAGATGGAGATTCAGGAGATGCAGCTCAAAAGAGGCCAAGCACATTT  285
emb1_AF310722  ATGAAGGATGAGGAGAAGATGGAGATTCAGGAGATGCAGCTCAAAAGAGGCCAAGCACATTT  360
emb1_AF186109  ATGAAGGATGAGGAGAAGATGGAGATTCAGGAGATGCAGCTCAAAAGAGGCCAAGCACATTT  357
emb1_BI817778  ATGTCGGGTGGCAGTTCATCGATGCGGTGAAGAAAGAGATCCAGAGCCCTTCAGCAGGTTG  137
emb1_BF056441  GATTCATTAATTTGCTTGACATTTCCACGCAAGCCGAAGATGGCCATAACCAAAAGGAACTT  271
emb1_BE848719  -TGTTACCAATCTGCTTGGCATTTTCCTGCAAGGTGAAACC-TGGTAATAAGCGGAACCTT  58
```

- Compute a phylogram such as this:



The Algorithm

- The chromosomes were candidate alignments, represented in a matrix with n rows.
- Create a random population of candidate alignments.
- For each candidate, apply *variation* operators to derive a child candidate.
- Apply a selection operator (fitness test) to generate the next generation.
- Repeat until 200 generations, no change in 100 generations, or number of gaps fell below a certain threshold.

Evolutionary Programming

- Chellapilla and Fogel's EP is comparable to Clustal-W.

Table 1: Data sets used for testing the proposed evolutionary programming procedure for multiple sequence alignment. Information regarding the data sets is provided in the Appendix.

Data Set	Number of Sequences	Mean Sequence Length in nucleotides (min,max)	Percentage of matched columns based on the best alignment using ClustalW	Number of Matched Columns in ClustalW solution	Number of Matched Columns in EP solution	Number of Generations	EP Score
1. S1 (Zhang, 1997)	10	211.9 (211, 212)	93.39	198	198	200	4082
2. 16S rRNA	8	457.0 (457, 457)	98.25	449	449	400	7233
3. Histone H3	21	122.0 (122, 122)	89.34	109	109	160	4766
4. Histone H3-H4 Intergenic Region	21	333.4 (322, 346)	27.41	91	102	180	7033

- EP Advantages
 - Flexibility in fitness functions.
 - Tackles harder and longer sequences.
 - Outperforms Clustal-W for low-similarity sequences.

Thank You!

