

brief introduction to bioinformatics and computational biology



from systems science to systems biology

luis m. rocha

Indiana university

school of informatics and cognitive science program

901 East Tenth Street, Bloomington IN 47408

and

Instituto Gulbenkian de Ciencia

Computational and Mathematical Biology

Oeiras, Portugal

rocha@indiana.edu
<http://informatics.indiana.edu/rocha>



**INDIANA
UNIVERSITY**



informatics
luis rocha 2006



- Information Processes in Biology
- Systems Biology, Computational Biology, Bioinformatics
- Synthetic, Multi- Disciplinary Approach to Biology
- Grand Challenges of Systems Biology
- Components of Bioinformatics & Computational Biology
- Some traditional components of Bioinformatics
- Literature Discussion and Useful Resources



from systems science to post-genome informatics

The word “system” is almost never used by itself, it is generally accompanied by an adjective or other modifier: physical system; biological system; social system [...] The adjective describes what is specific and particular, i.e., it refers to the specific “thinghood” of the system; the “system” describes those properties which are independent of this specific “thinghood.” [Rosen, 1986]



- Systems Science is the methodology used to study *systemhood* not *thinghood* properties in Nature.
 - ▶ General Principles of Life (and other systems)
 - ▶ Modeling and Simulation of systems measured from and validated in real things.
 - ▶ It accumulates knowledge via Mathematical and Computational analysis of classes of systems, models, and problems.
 - Dynamical Systems, Automata Theory, Pattern Recognition, etc.
- Interdisciplinary Meta-Methodology
 - ▶ Comparative, Integrative, Non-reductionist
- Historically Related to Cybernetics
 - ▶ Complex Systems, Artificial Life

dealing with general principles of complex systems

- Weaver [1948] identified 3 types of problems in Science
 - ▶ Organized Simplicity: systems with small number of components
 - Classical mathematical tools: calculus and differential equations
 - ▶ Disorganized Complexity: systems with large number of erratic components
 - Stochastic, Statistical Methods
 - ▶ Organized Complexity: systems with a fair number of components with some functional identity
 - When the behavior of components depends on the organization and function of the whole
 - Techniques depend on Computer Science and Informatics. Require massive combinatorial searches, simulations, and knowledge integration.
 - The realm of Systems Science
 - ▶ Complex Systems are systems of many components which cannot be completely understood by the behavior of their components.
 - Complementary models, Hierarchical Organization, Functional decomposition [See Klir, 1991]



And its Involvement with Systems Science

- People

- ▶ Von Bertalanffy [1952, 1968], Mesarovic [1968], Rosen [1972, 1978, 1979, 1991], Pattee [1962, 1979, 1982, 1991, 2001], Maturana and Varela [1980], Kauffman [1991], Conrad [1983], Matsuno [1981], Cariani [1987].
- ▶ Leading Journal: *Biosystems*

- Biology is the most Fundamental Inspiration for Systems Science

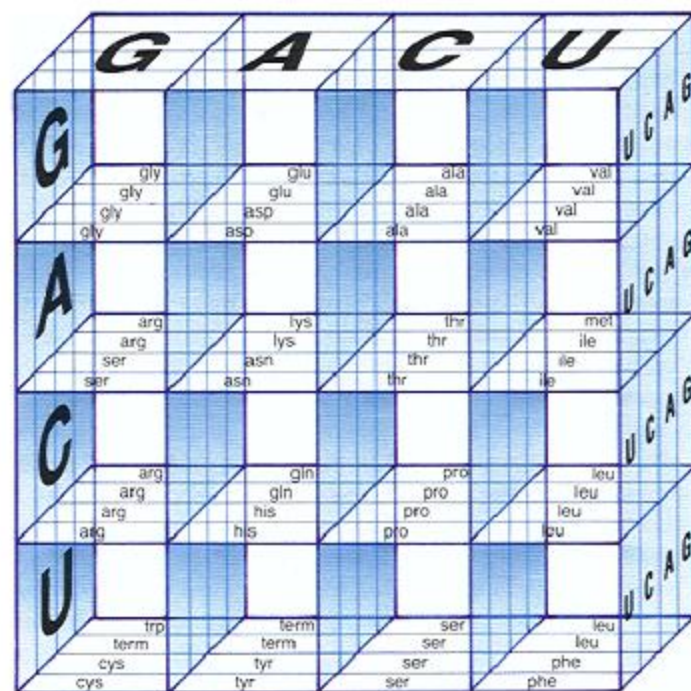
- ▶ Cybernetics and Control Theory derive Feedback Control from the physiological concept of Homeostasis
- ▶ Automata Theory, Artificial Intelligence, Artificial Life derived from attempts (by Turing, McCulloch and Pitts) to study the behavior of the Brain and Evolution (Von Neumann)
- ▶ Self-Organizing, Autopoiesis, Complex Adaptive Systems, Artificial Life, Embodied Cognition from developmental and evolutionary biology.

- But Systems Science has had a Small impact in the practice of Biology

- ▶ Due to a large gap between theoretical and experimental biologists.
 - Systems-based theoretical Biology versus a reductionist view
 - Theoretical biology has had more impact on other areas (AI, Alife, Complexity, Systems Science) than Biology itself.



general principles and metaphors from life



Life and Information

how to identify it?

■ List of properties

- Growth
- Metabolism
- Reproduction
- Adaptability
- Self-maintenance (autonomy)
- Self-repair
- Reaction
- Evolution
- Choice

■ Threshold of complexity

- Categorization and Control
- Function (self-reference)
- Open-ended evolution
- Information

Is life
Fuzzy?



viruses, candle
flames, the
Earth, certain
robots?



Is there a synthetic
criteria? How
general can it be?

in the living organization

- organisms act according to information they perceive in an environment
- organisms reproduce and develop from genetic information
 - genetic information is **transmitted** “vertically” (inherited) in phylogeny and cell reproduction, and **expressed** “horizontally” within a cell in ontogeny and plain functioning
- Self-reference
 - Information relevant to organism: **function**
 - Only in **reference** to an organism does a piece of DNA **function** as a gene
 - Biology is contextual, physics is universal



“Life is a dynamic state of matter organized by information”. Manfred Eigen [1992]



“Biology and physics have nothing to do with each other because biological evolution is essentially historical, and physical laws must be independent of history”. Ernst Mayer

- impossibility of epistemological reduction of the properties of a system to its components
 - Wave-particle duality
 - Information and function are contextual and historical
 - “Clockness”: many possible material implementations
 - Several biological designs for similar function (e.g. flying)
 - The function of DNA does not lie in its dynamic (bio-chemical) characteristics

“First, nothing in biology contradicts the laws of physics and chemistry; any adequate biology must be consonant with the ‘basic’ sciences. Second, the principles of physics and chemistry are not sufficient to explain complex biological objects because new properties emerge as a result of organization and interaction. These properties can only be understood by the direct study of the whole, living systems in their normal state. Third, the insufficiency of physics and chemistry to encompass life records no mystical addition, no contradiction to the basic sciences, but only reflects the hierarchy of natural objects and the principle of emergent properties at higher levels of organization”. Stephen Jay Gould



How much is specific bio-chemistry?

- **Can there be several implementations of life?**
 - To study life do we need to find and synthesize the necessary threshold of complexity?
 - Hard and wet Artificial Life
 - Or is it enough to simulate the behavior of life?
 - Soft Artificial Life
- **Important to study the living organization**
 - What can be abstracted and implemented in a different medium?
 - Understanding organization and design principles
 - Scientific advancement of the essential principles of life
 - Systems Biology, Artificial Life
 - Solving engineering and design problems
 - Bio-inspired computing

■ Genetic System

- Construction (expression, development, and maintenance) of cells ontogenetically: horizontal transmission
- Heredity (reproduction) of cells and phenotypes: vertical transmission

■ Immune System

- Internal response based on accumulated experience (information)

■ Nervous and Neurological system

- Response to external cues based on memory

■ Language, Social, Ecological, Eco-social, etc.



“Life is a complex system for information storage and processing”.
Minoru Kanehisa
[2000]

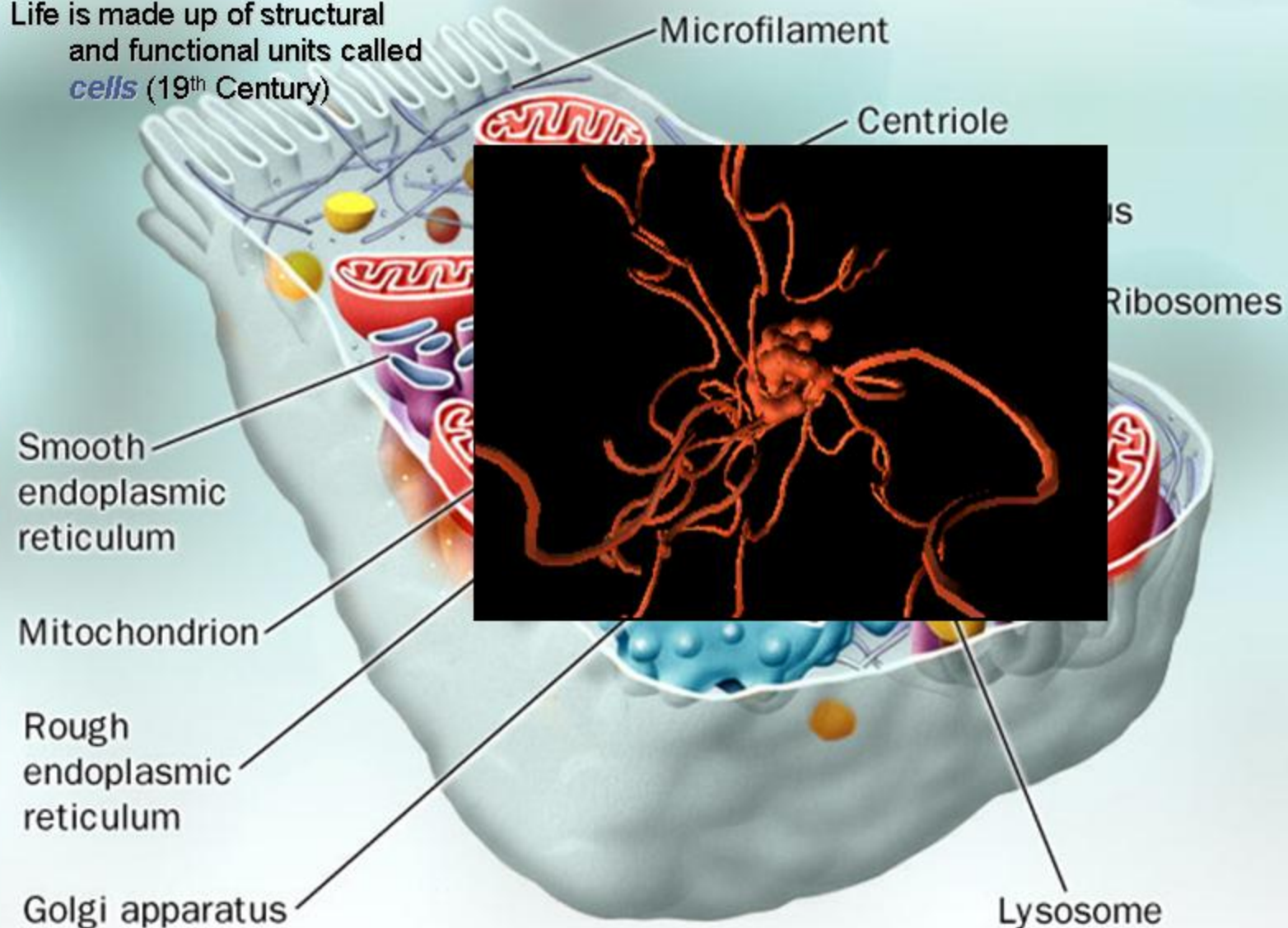
memory, structural and functional

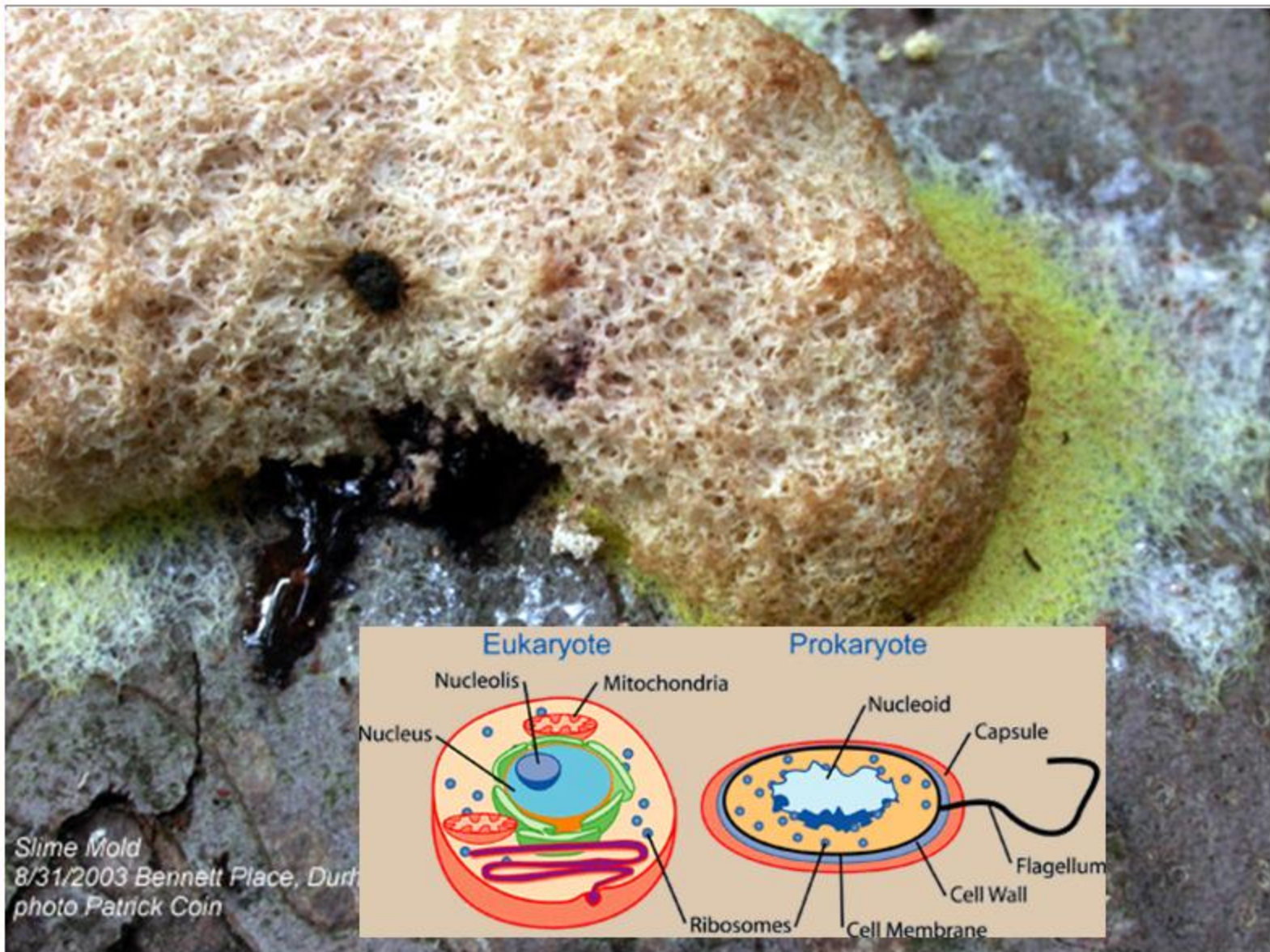
- **Mendelian Gene**
 - Hereditary unit responsible for a particular characteristic or trait
- **Molecular Biology Gene**
 - Unit of information expression via *transcription* and *translation*
 - Horizontal information expression (semantics, active)
- **Genome**
 - Unit of information transmission via *DNA replication*
 - Vertical information transmission (syntactic, passive)
 - Set of genes in the chromosome of a species
- **Genotype**
 - Instance of the genome for an individual
- **Phenotype**
 - Expressed and developed genotype
 - Genes have different alleles
- **Transcriptome**
 - Set of expressed genes (mRNA transcripts) in a given context
- **Proteome**
 - Set of proteins that are encoded and expressed by a genome

“Biology is the science of life that aims at understanding both functional and structural aspects of living organisms”. Minoru Kanehisa [2000]

Structure and organelles

Life is made up of structural and functional units called *cells* (19th Century)



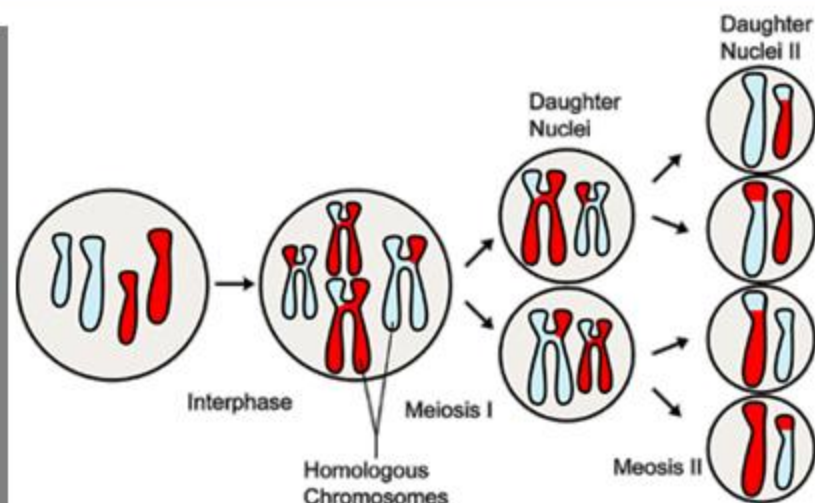


■ Meiosis

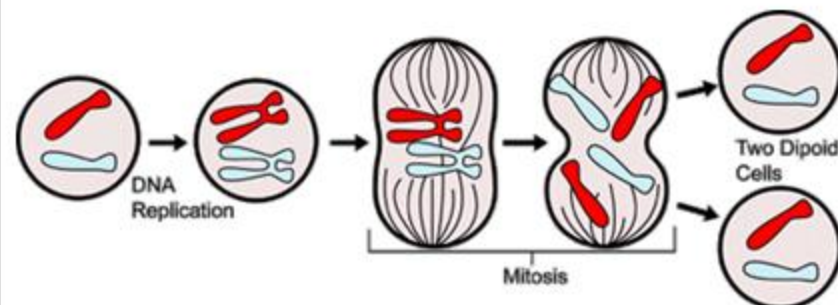
- Diploid cell's genome is replicated once and split twice
- produces four haploid (*germ*) cells each with half the chromosomes
- Sexual reproduction combines germ cells from two individuals to produce diploid (*zygocyte*) cells
- Vertical genetic information transmission
 - Offspring with a new genotype

■ Mitosis

- Eukaryotic cell separates its duplicated genotype into two identical halves
 - somatic cells in multicellular organisms
- Horizontal genetic information expression
 - development

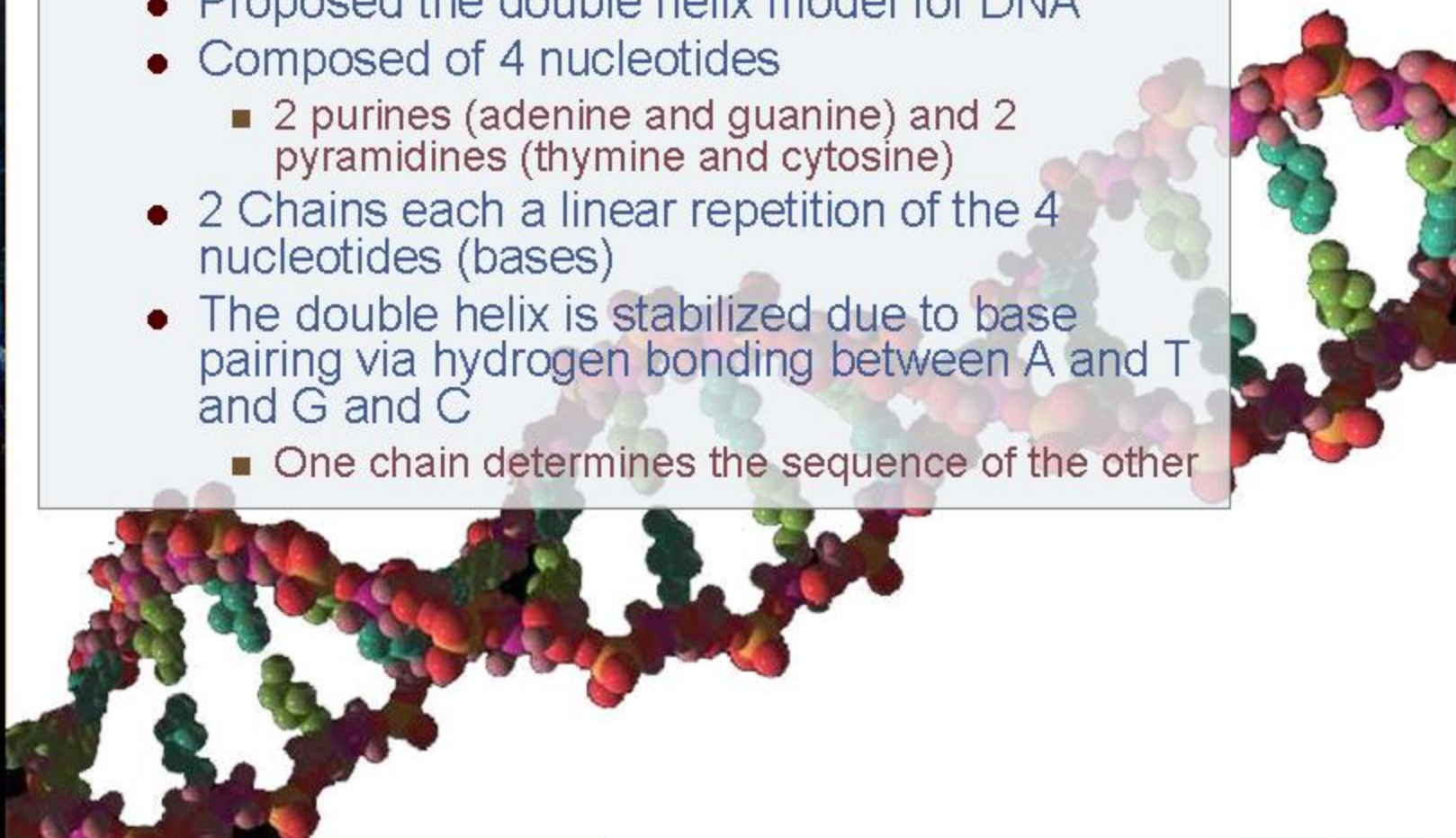


Crossovers may occur in meiosis

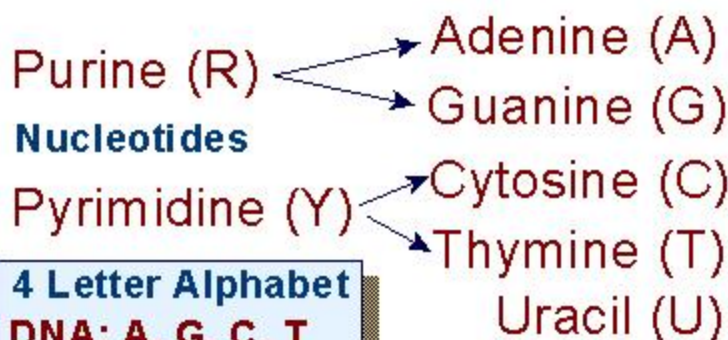


deoxyribonucleic acid

- The chromatin contains DNA and protein
- James Watson and Francis Crick (1953)
 - Proposed the double helix model for DNA
 - Composed of 4 nucleotides
 - 2 purines (adenine and guanine) and 2 pyrimidines (thymine and cytosine)
 - 2 Chains each a linear repetition of the 4 nucleotides (bases)
 - The double helix is stabilized due to base pairing via hydrogen bonding between A and T and G and C
 - One chain determines the sequence of the other



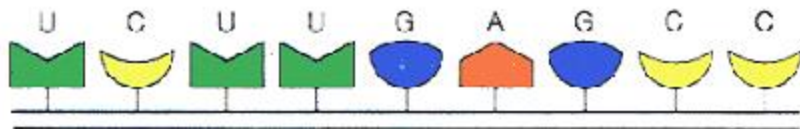
a molecular language system



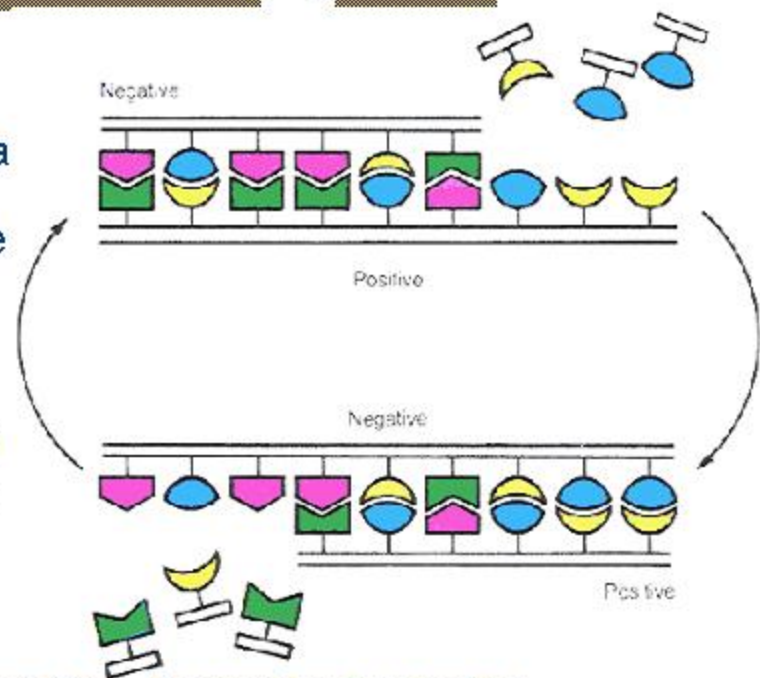
4 Letter Alphabet
DNA: A, G, C, T
RNA: A, G, C, U

Form sequences that can store information

Linear molecules with a phosphate-sugar backbone (deoxyribose and ribose)

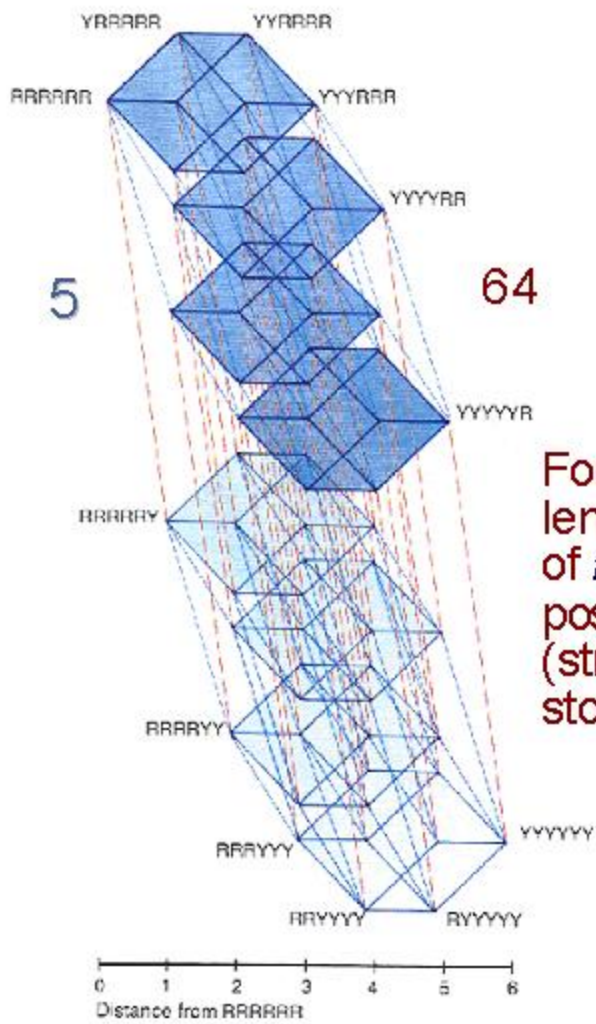
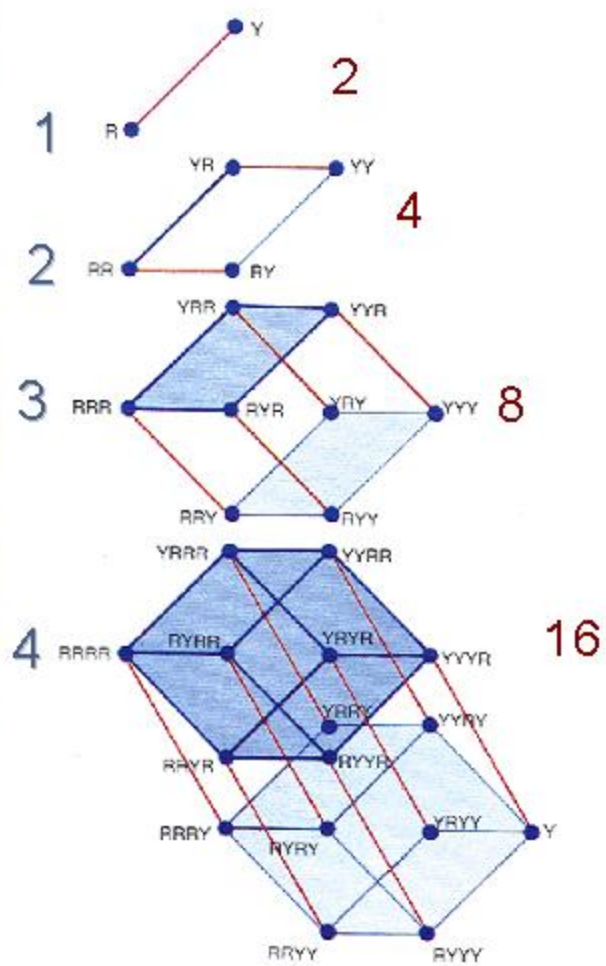


Complementary base pairing
(Hydrogen-bonding between purines and pyrimidines)



Requirements for structural information

Possibility of repeated copying



For a sequence of length n , composed of m -ary symbols, m^n possible values (structures) can be stored

biologically
Inspired
computing

functional products

Polypeptide chains of aminoacids

Primary Structure

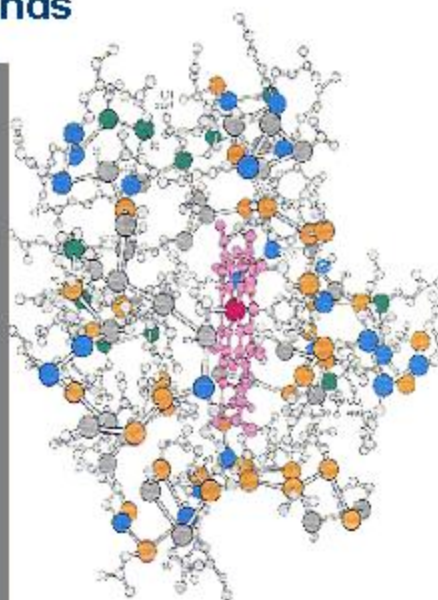


Folding

3-dimensional structure

Secondary and tertiary bonds

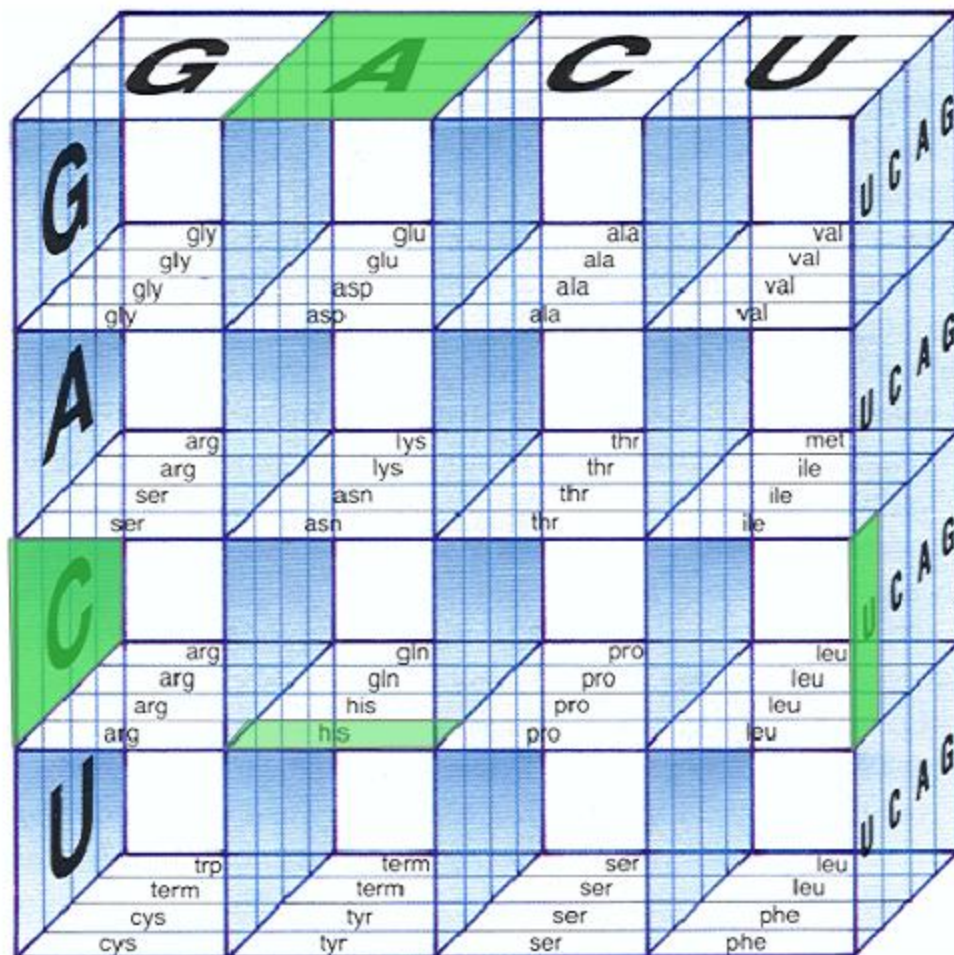
- In proteins, it is the 3-dimensional structure that dictates function
 - ▶ The specificity of enzymes to recognize and react on substrates
- The functioning of the cell is mostly performed by proteins
 - ▶ Though there are also ribozymes

**Table 1.4.** Amino acid codes

Ala	A	Alanine
Arg	R	Arginine
Asn	N	Asparagine
Asp	D	Aspartic acid
Cys	C	Cysteine
Gln	Q	Glutamine
Glu	E	Glutamic acid
Gly	G	Glycine
His	H	Histidine
Ile	I	Isoleucine
Leu	L	Leucine
Lys	K	Lysine
Met	M	Methionine
Phe	F	Phenylalanine
Pro	P	Proline
Ser	S	Serine
Thr	T	Threonine
Trp	W	Tryptophan
Tyr	Y	Tyrosine
Val	V	Valine
Asx	B	Asn or Asp
Glx	Z	Gln or Glu
Sec	U	Selenocysteine
Unk	X	Unknown

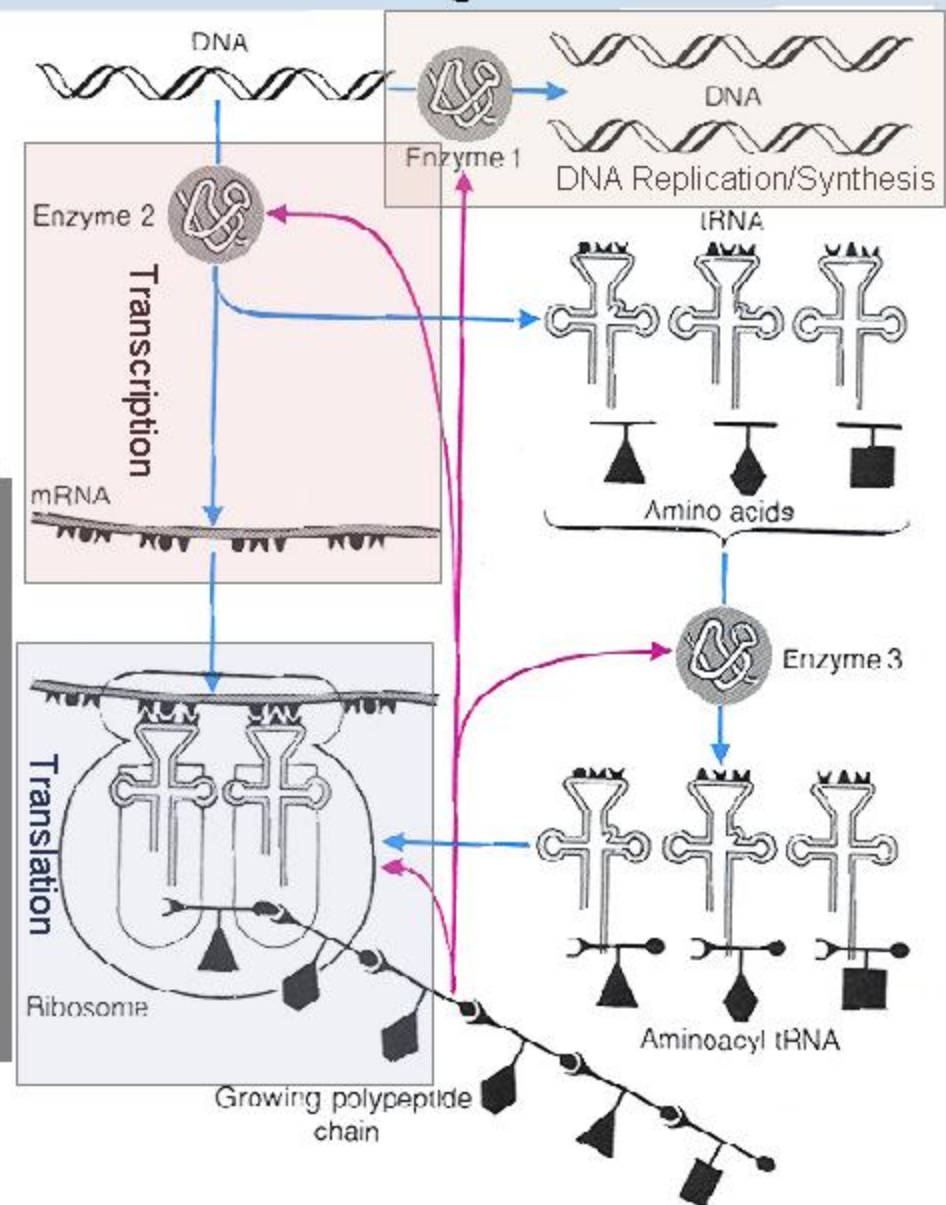
3.5.14
2005biologically
Inspired
computing

- The genetic code maps information stored in the genome into functional proteins
 - ▶ Triplet combinations of nucleotides into amino acids

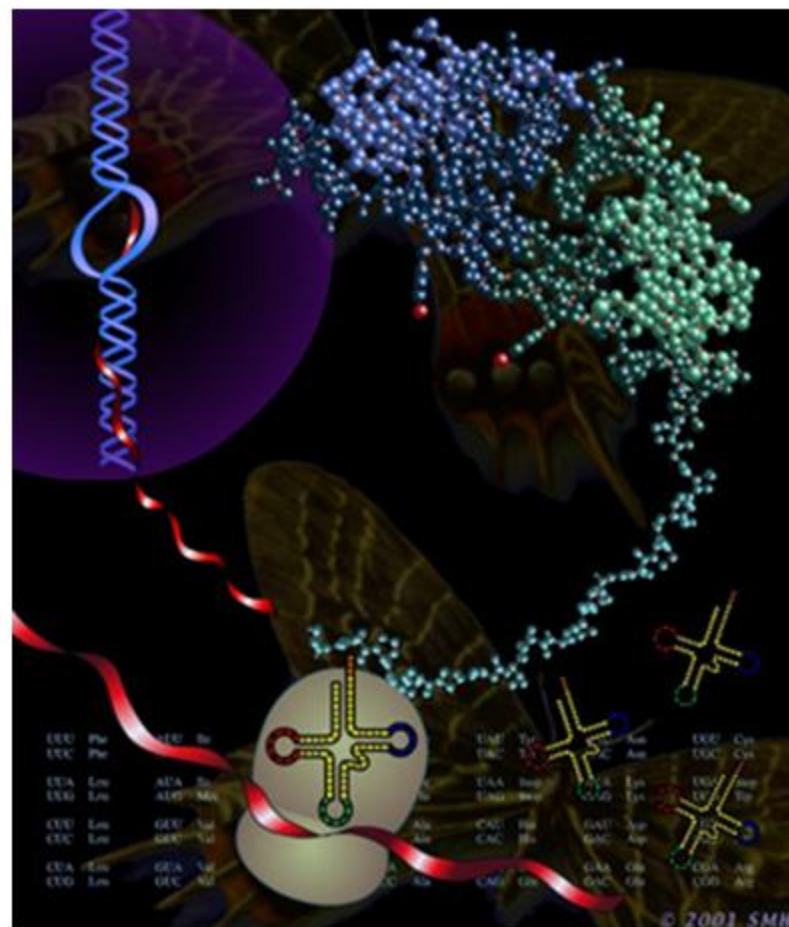
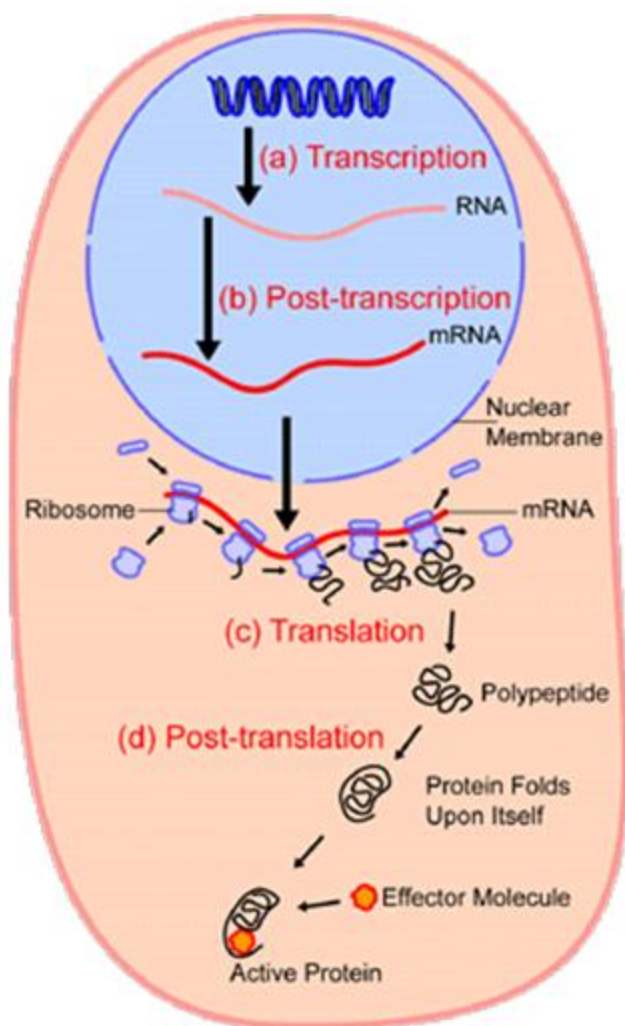


Triplets of 4 Nucleotides can define 64 possible codons, but only 20 amino acids are used (redundancy)

- **Reproduction**
 - ▶ DNA Polymerase
- **Transcription**
 - ▶ RNA Polymerase
- **Translation**
 - ▶ Ribosome
- **Coupling of AA's to adaptors**
 - ▶ Aminoacyl Synthetase

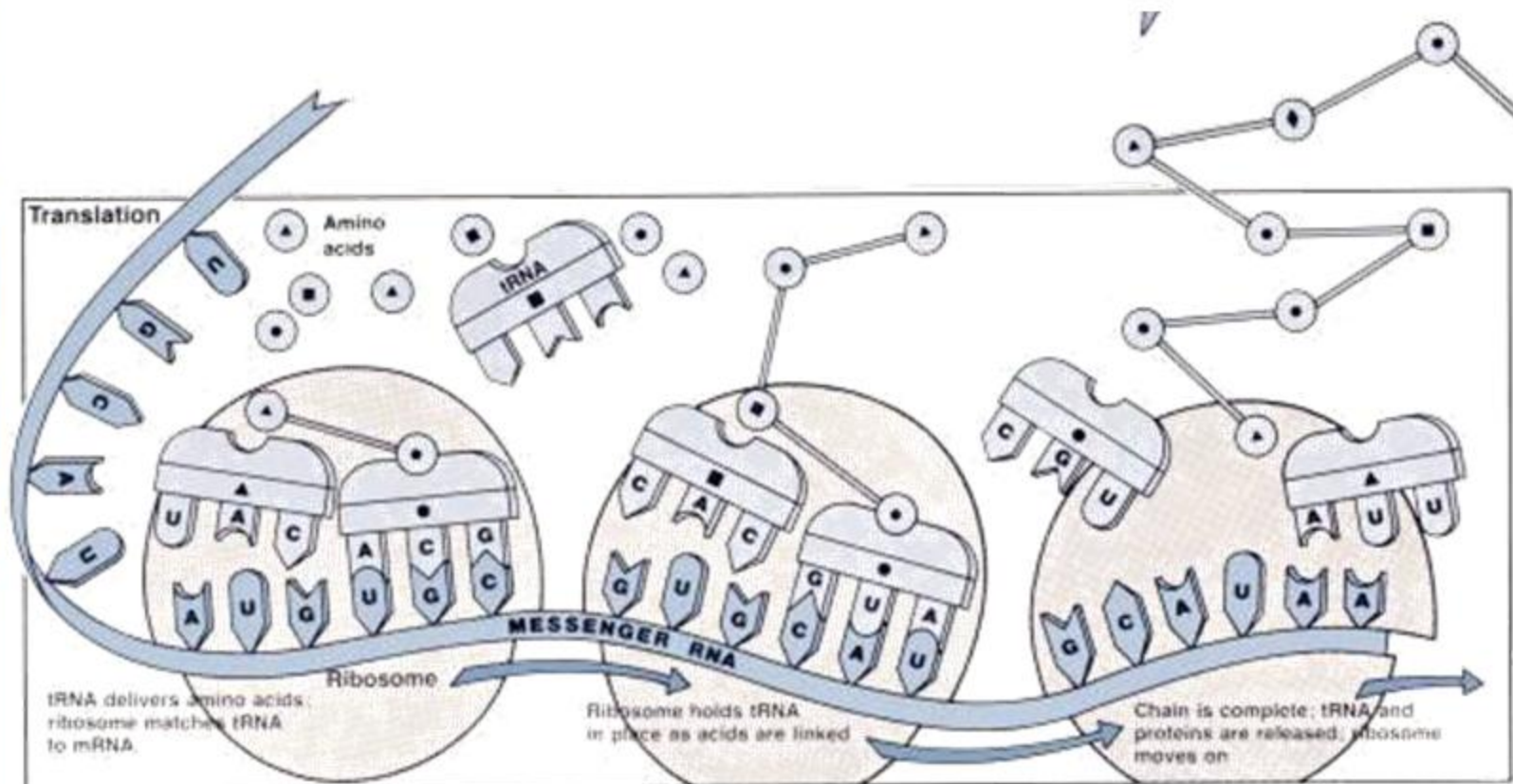


transcription and translation

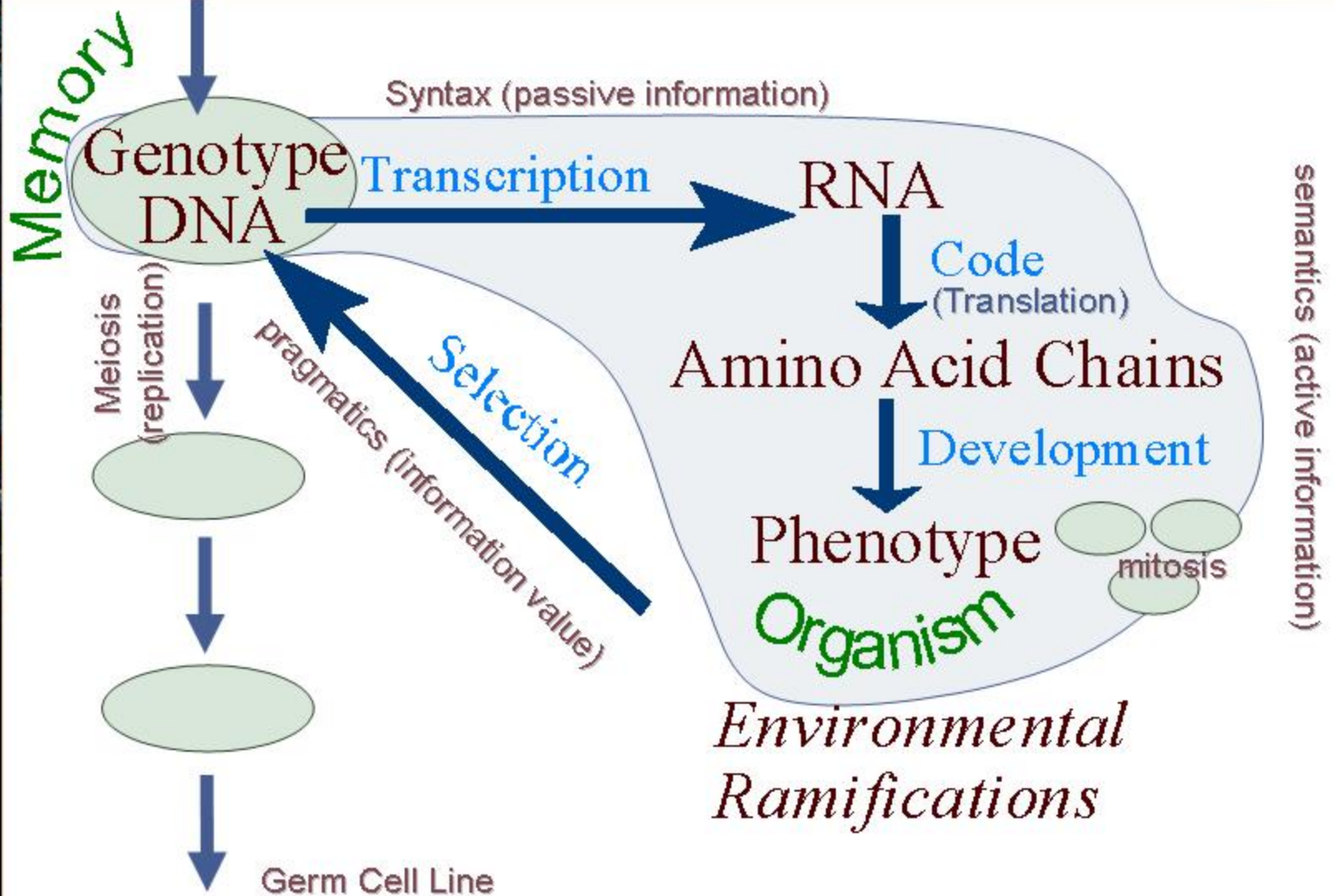


biologically
Inspired
computing

constructing (decoding) the message



biologically
Inspired
computing



life-as-it-could-be



■ Chris Langton

- Artificial Life can contribute to theoretical biology by locating *life-as-we-know-it* within the larger picture of *life-as-it-could-be*
- life as a property of the *organization* of matter, rather than a property of the matter which is so organized
 - The way information is processed
- Whereas biology has largely concerned itself with the material basis of life, Artificial Life is concerned with the formal basis of life.
 - views an organism as a large population of *simple* machines
 - *Synthetic approach or emergent behavior*

■ Analytical

- Reduction to (non-living) components
 - Reductionism
- Life is complicated chemistry
- Tied to specific materiality
- Does not allow emergence
 - Function, control, measurement, categorization, information are unnecessary "illusions"

■ Synthetic

- Construction from components
 - Holist
- Life is Organization
 - Networks of components
- Universal or implementation independent
- Emergence
 - "bottom-up" approach

■ Hard Alife

- Logical mechanisms of life
- Discover and synthesize the design principles of life
 - Threshold of complexity
 - Lists of characteristics

■ Soft or weak Alife

- To simulate life
- Compare design principles of life with simulations
- Extract design principles to solve problems
 - Bio-inspired computing

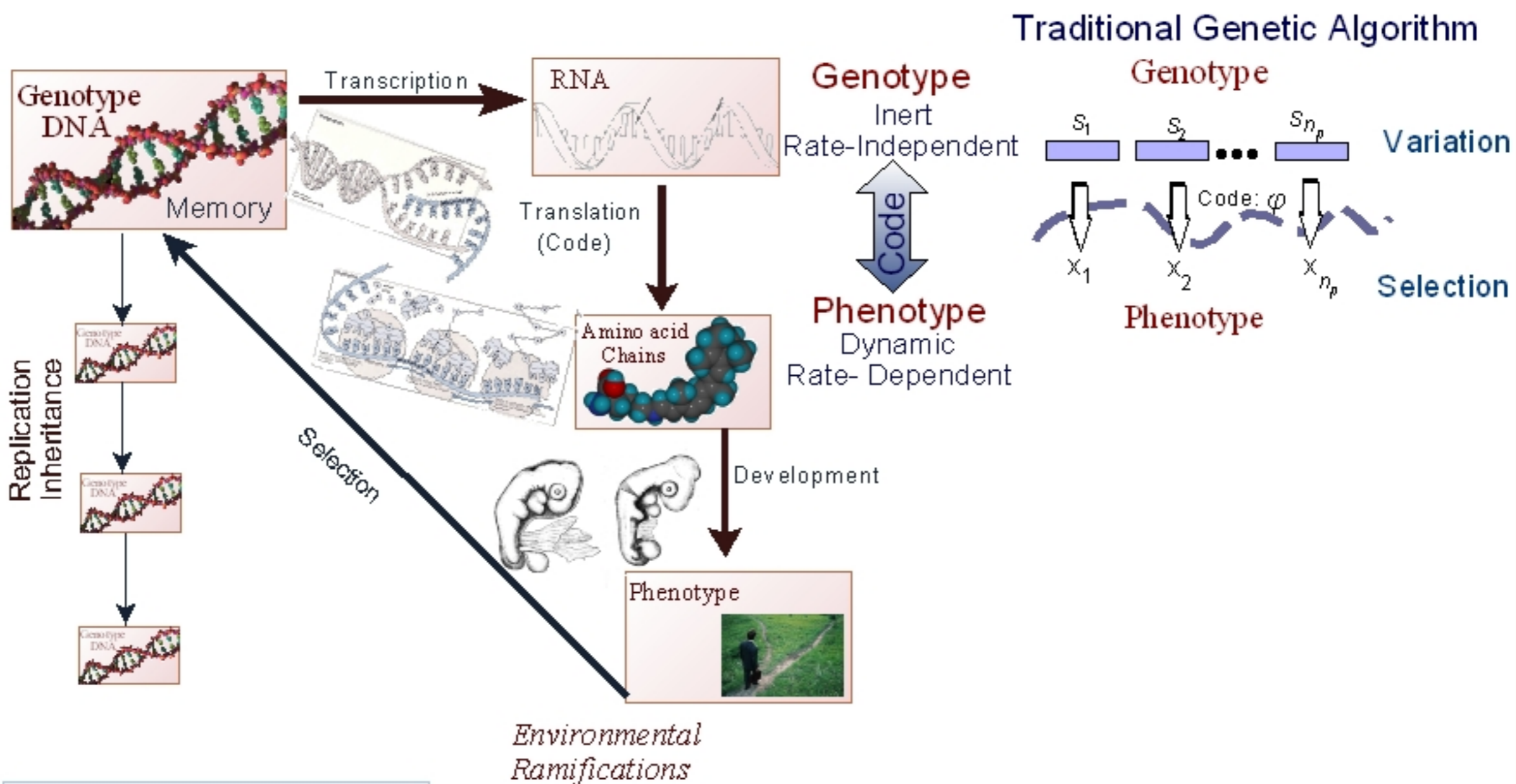
■ Bottom-up methodology

Systemhood

- A system possesses *systemhood* and *thinghood* properties
 - Thinghood refers to the specific material that makes up the system
 - Systemhood are the abstracted properties
 - E.g. a clock can be made of different things, but there are implementation-independent properties of “clockness”
 - Systems science deals with the implementation-independent aspects of systems
 - Robert Rosen, George Klir...

genotype/phenotype mapping in artificial life

informatics
luis rocha 2006



ncRNA: a regulatory hidden layer in Eukaryotes

- Evidence for non-protein coding RNA (ncRNA) in complex organisms (higher eukaryotes)
 - ▶ “ncRNA dominates the genomic output of the higher organisms and has been shown to control chromosome architecture, mRNA turnover and the developmental timing of protein expression, and may also regulate transcription and alternative splicing.”
 - Mattick, J. S. (2003). *BioEssays*. 25: 930-939
 - ▶ *A Hidden Layer* of Non-protein-coding RNAs in Complex Organisms.

- Two types of genetic information

- ▶ **mRNA** for proteins
- ▶ **ncRNA** for RNA products

- Three types of genes in eukaryotes

- ▶ Encoding only proteins
- ▶ Encoding only ncRNA
- ▶ Encoding both

- Many types of ncRNA

- ▶ tRNA, rRNA, SnoRNA, miRNA, siRNA, eRNA, etc.

**Single output
Simple system**

Prokaryotic gene



**Multiplex output
Parallel processing**

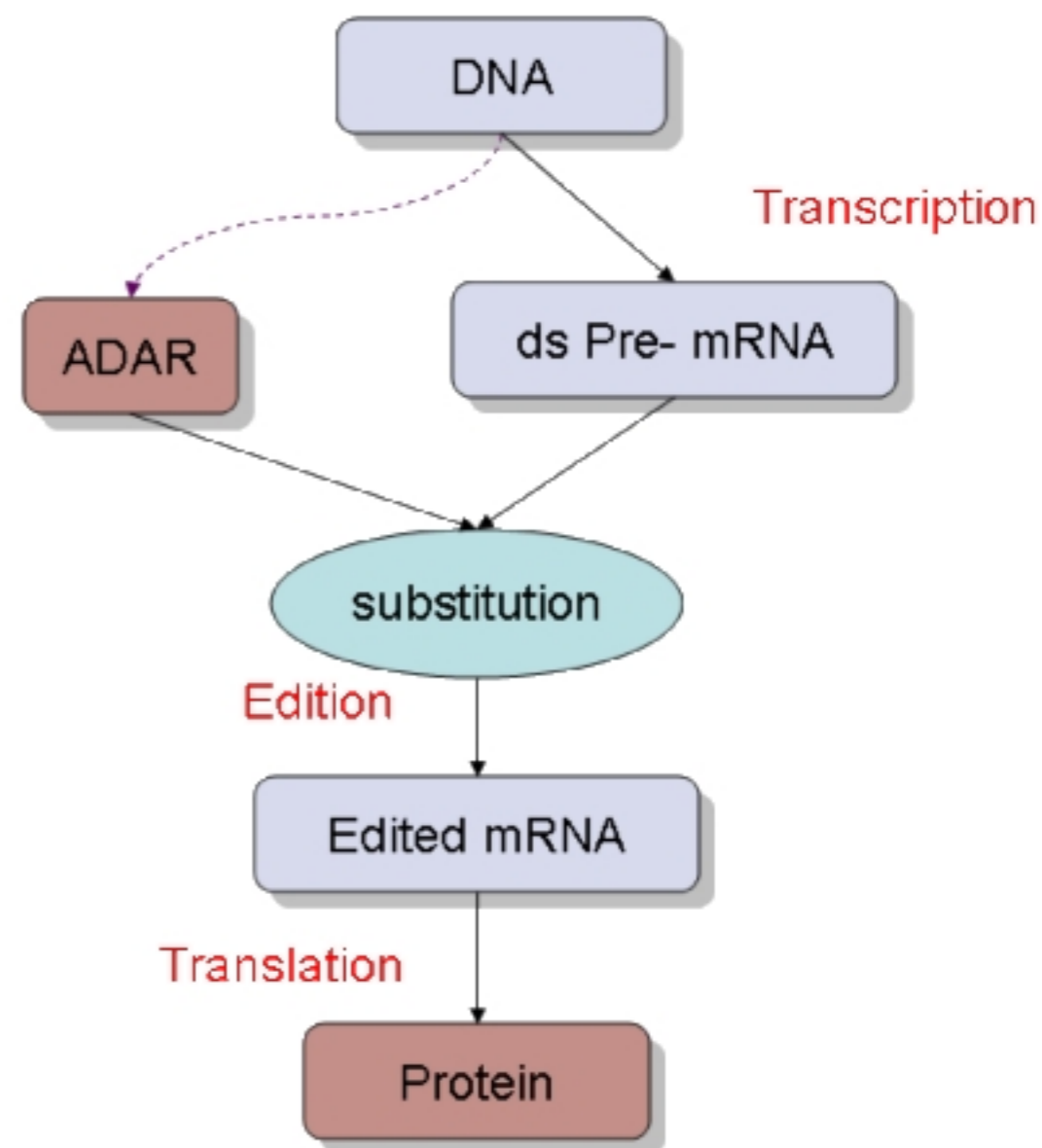
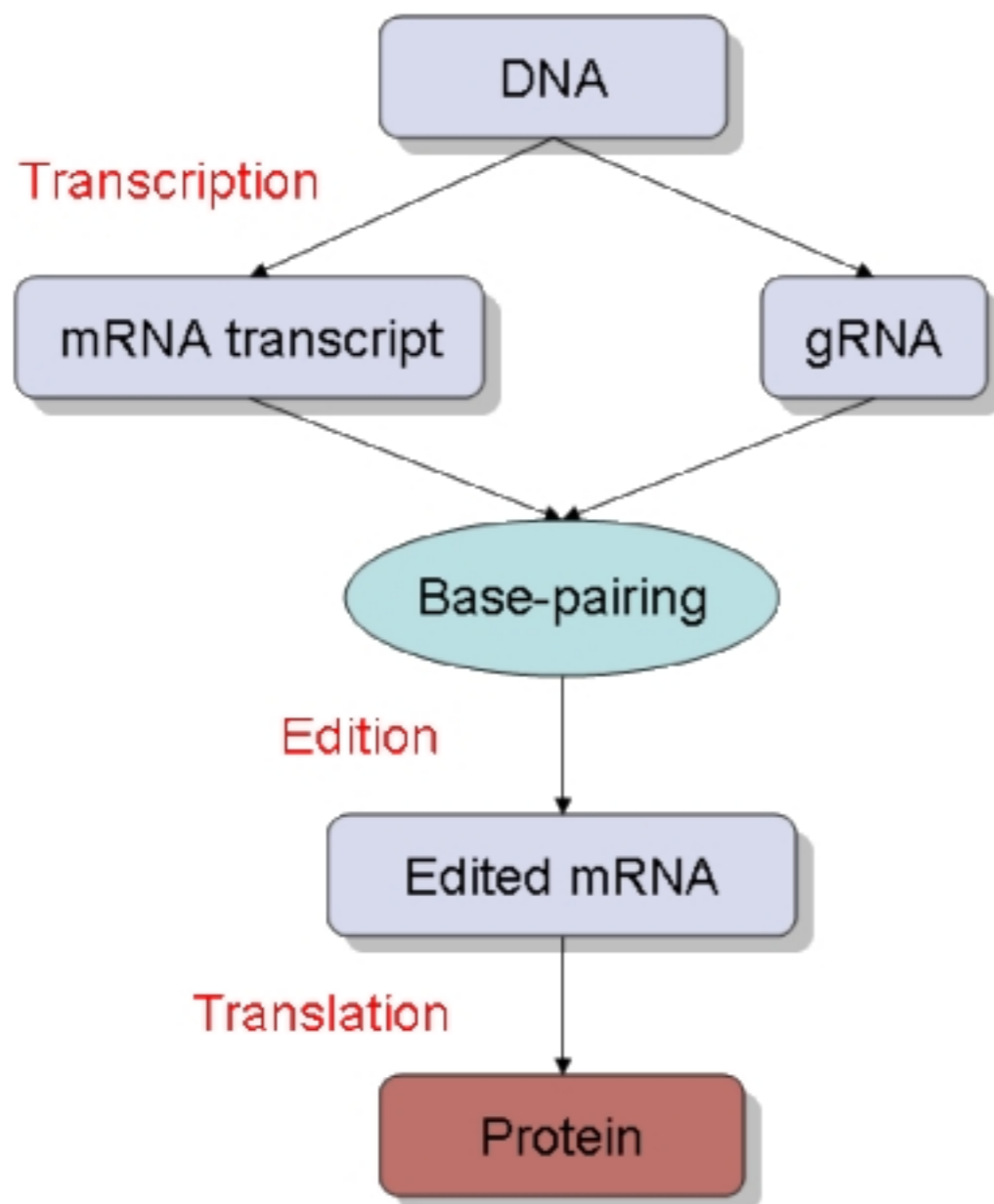
Eukaryotic gene



Mattick, J. S. [2001]. *EMBO Reports* 2, 11, 986–991

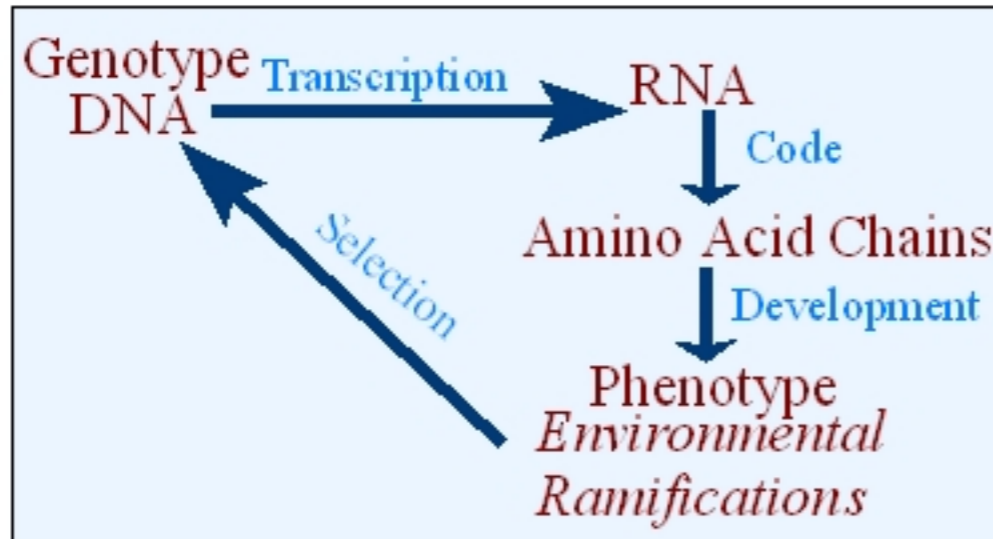
Mattick, J. S. And V. Makunin [2005]. *Human Molecular Genetics* 14, 11, R121-R132

U-insertion and A-to-I substitution



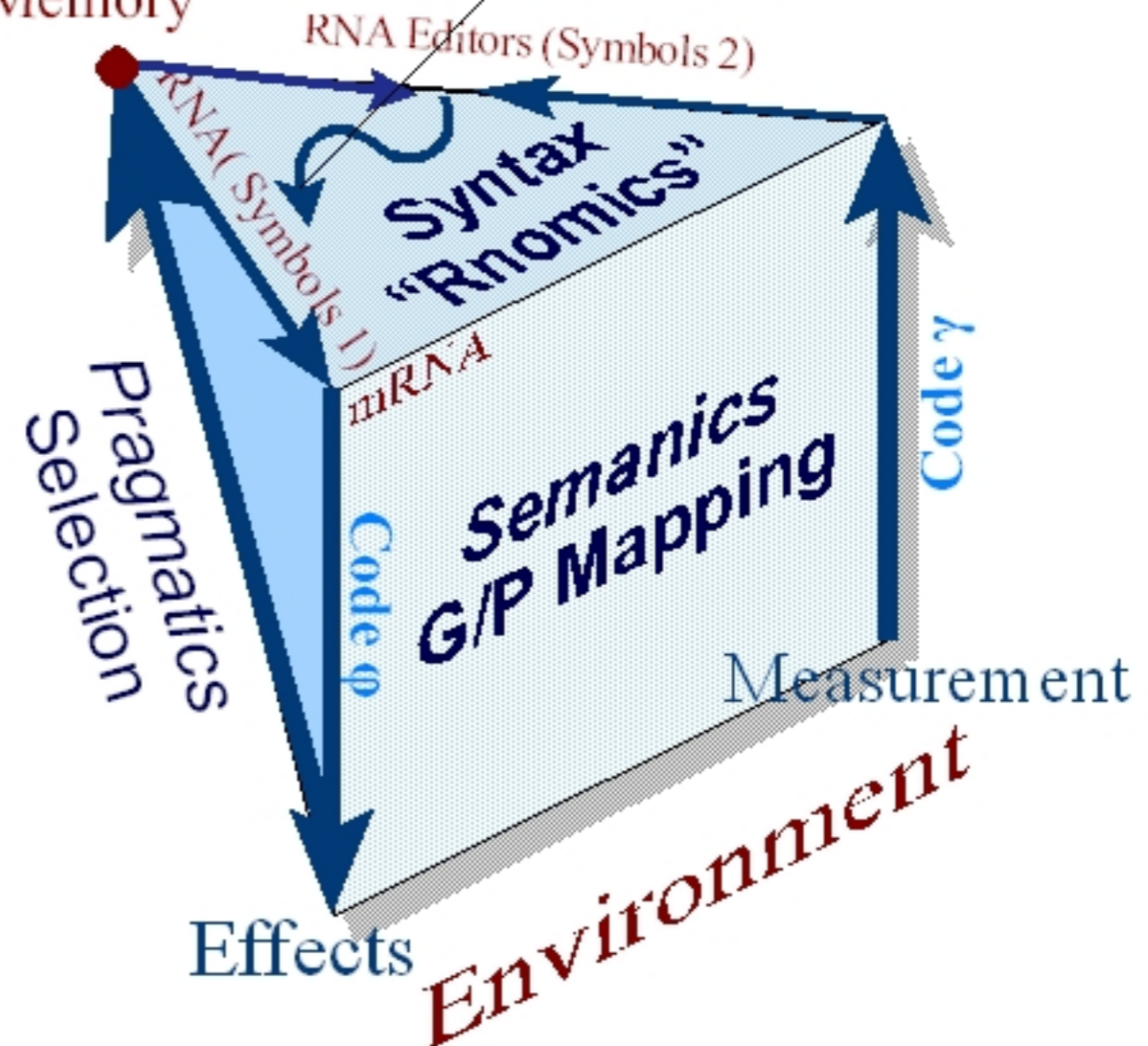
RNA editing modulates gene expression

genes may encode different proteins depending on environment



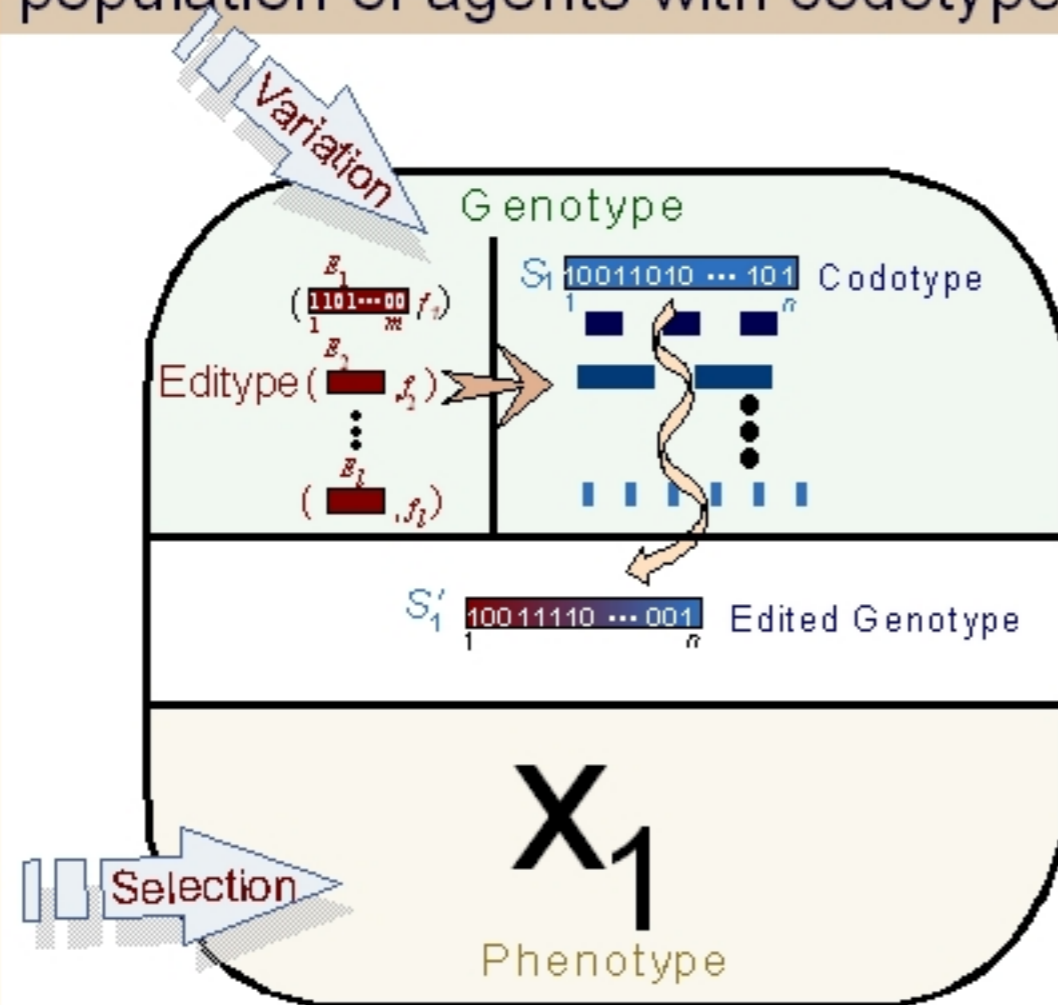
DNA
(Symbolic)
Memory

A Richer informational
process---re-programmable
genotype??

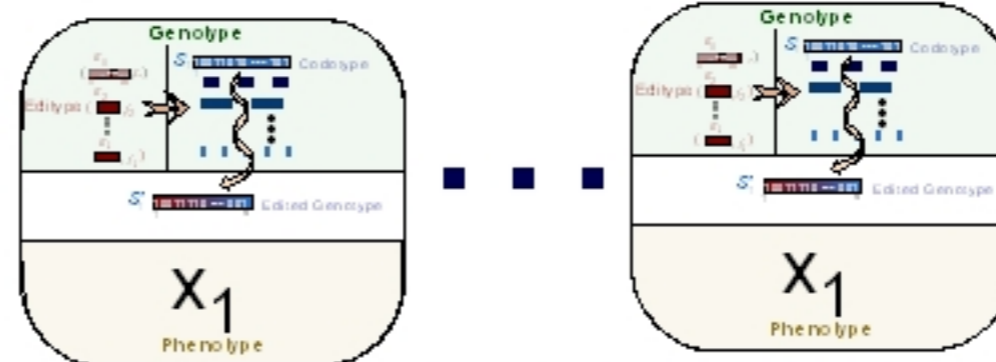


- Only mutations that occur during DNA replication can become permanent and heritable
- RNA Editing may produce different mRNA's (and thus proteins), but editions are not inherited.
 - ▶ What is inheritable, and subjected to variation, is the genetic material (both coding and non-coding) which is ultimately selected and transmitted to the offspring of the organism

population of agents with codotype and editype

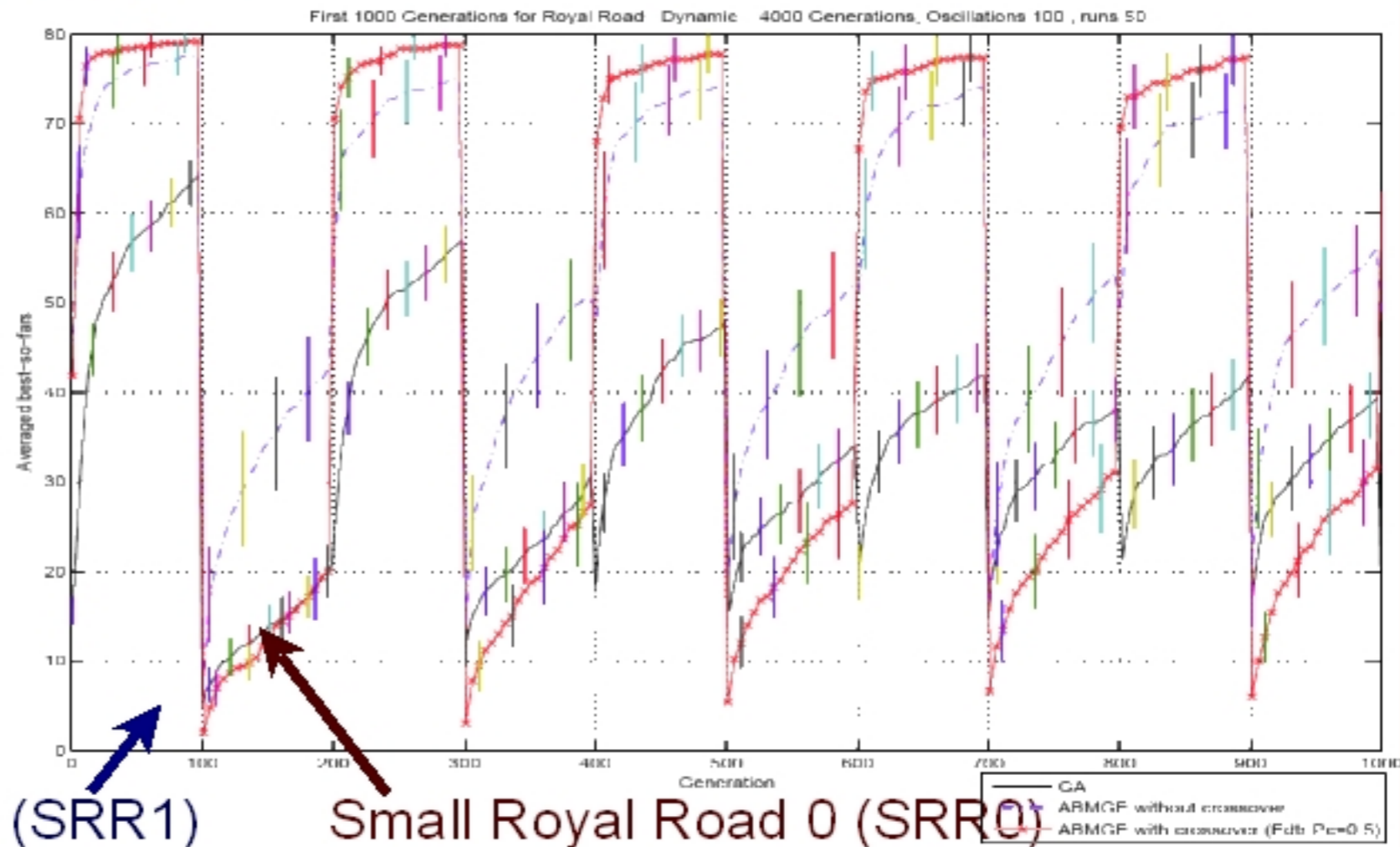


- **Genome contains both coding and non-coding portions**
 - ▶ Codome and Editome (Editosome)
- **For each agent**
 - ▶ Codotype edited by editype before "translation"
- **Modeling pre-translation information (syntactic) processes**
 - ▶ no RNA/DNA distinction
 - ▶ a process of *non-inheritable alteration of genotypes via edition*, not any specific type of RNA Editing.
 - ▶ Not mutation
- **co-evolution of editype and codotype**
 - ▶ Not in the EC sense of independent populations
 - ▶ Independent variation



dramatic environmental changes

- Oscillation period
 - ▶ **100** (50, 200) generations
- First 1000 generations
 - ▶ Same parameters as in static case



Small Royal Road 1 (SRR1)

```

s1 = 11111***** c1 = 10
s2 = *****11111***** c2 = 10
s3 = *****11111***** c3 = 10
s4 = *****11111***** c4 = 10
s5 = *****11111***** c5 = 10
s6 = *****11111***** c6 = 10
s7 = *****11111***** c7 = 10
s8 = *****11111***** c8 = 10
    
```

Small Royal Road 0 (SRR0)

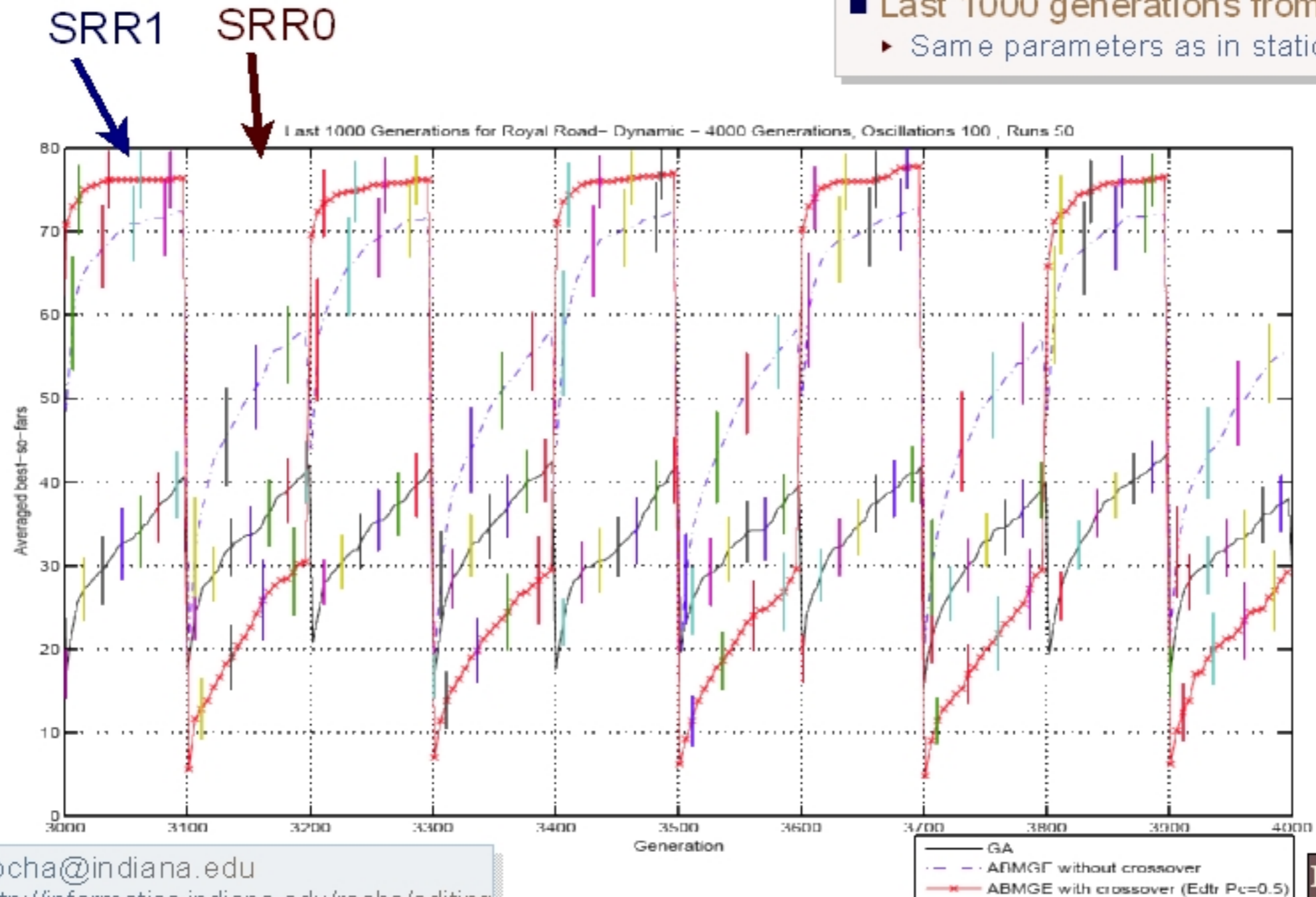
```

s1 = 00000***** c1 = 10
s2 = *****00000***** c2 = 10
s3 = *****00000***** c3 = 10
s4 = *****00000***** c4 = 10
s5 = *****00000***** c5 = 10
s6 = *****00000***** c6 = 10
s7 = *****00000***** c7 = 10
s8 = *****00000***** c8 = 10
    
```

oscillatory royal road function

dramatic environmental changes

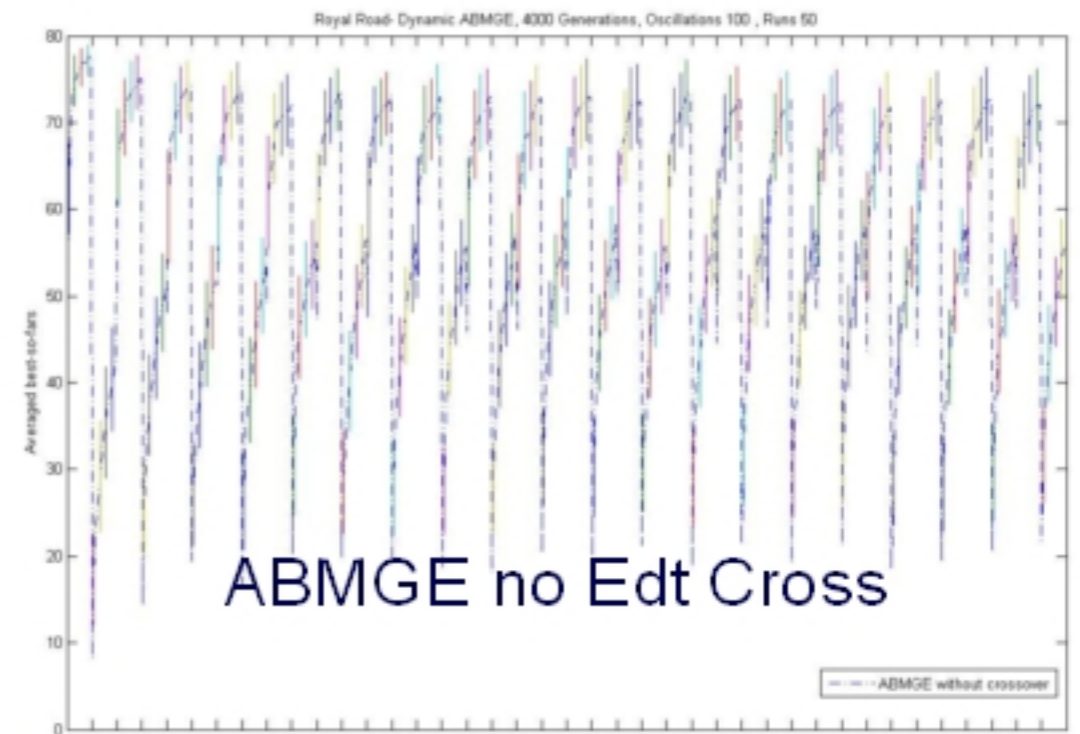
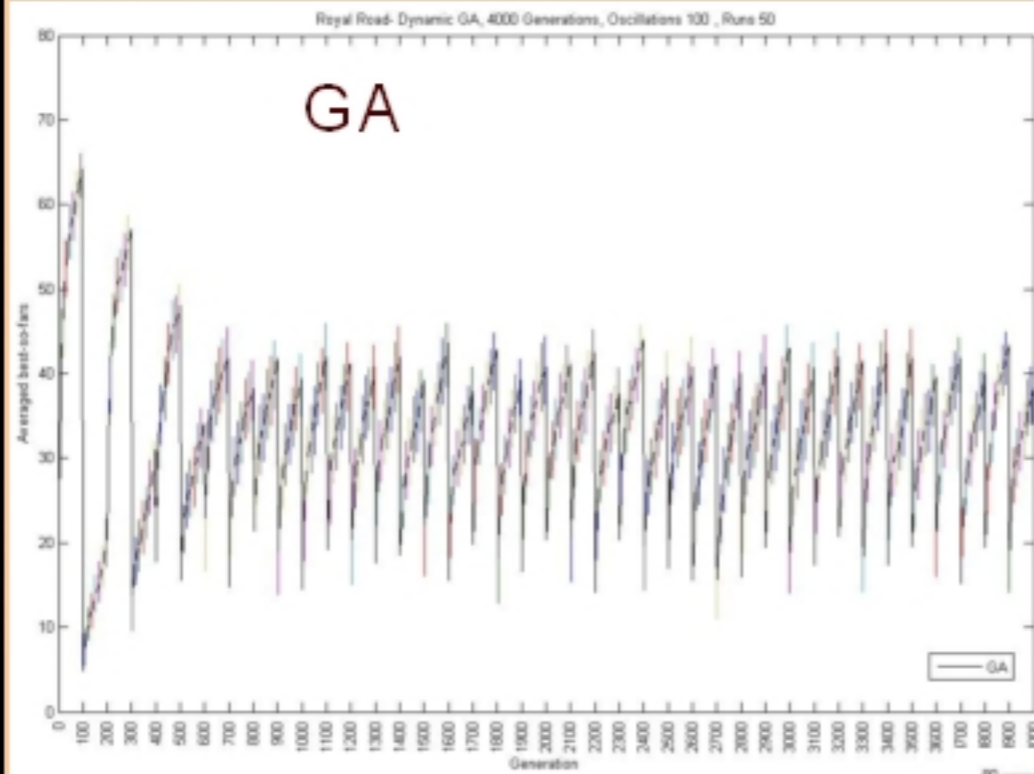
- Oscillation period
 - ▶ 100 generations
- Last 1000 generations from 4000
 - ▶ Same parameters as in static case



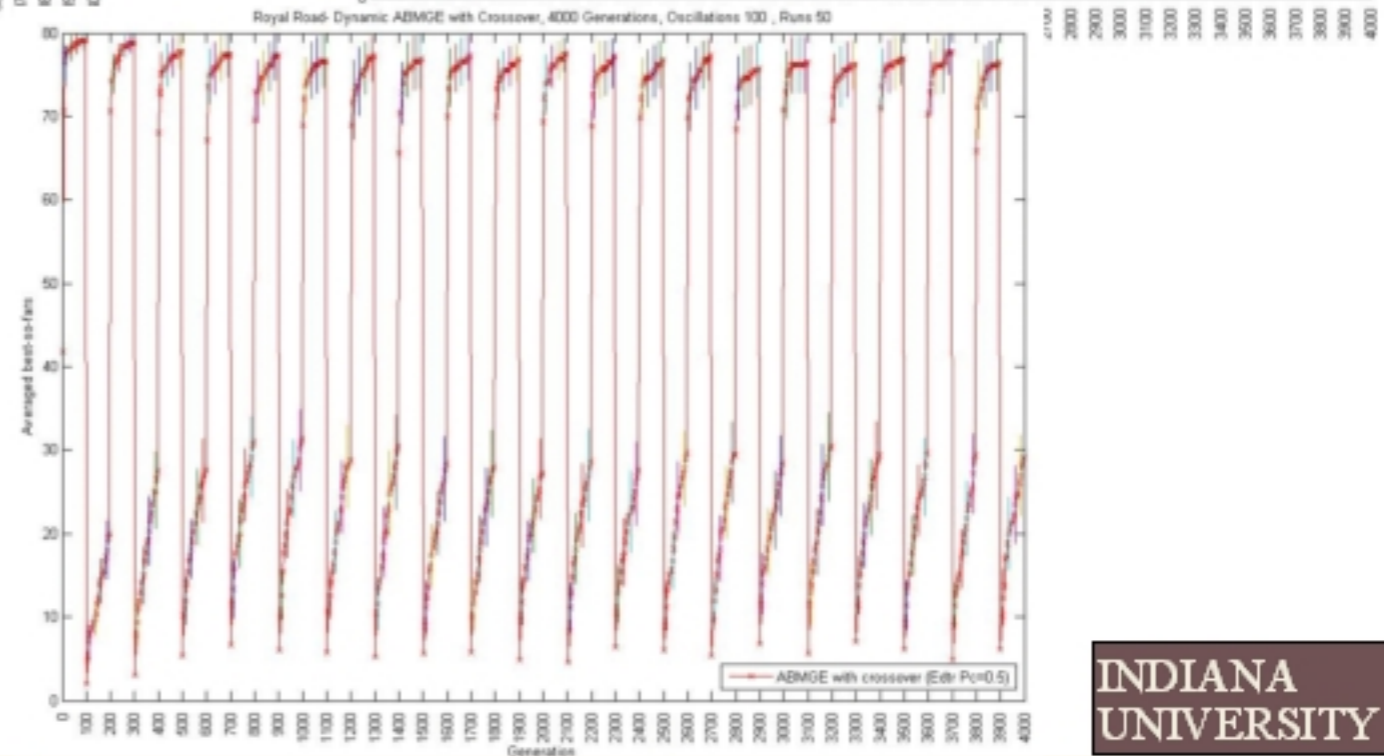
rocha@indiana.edu
<http://informatics.indiana.edu/rocha/editing>

INDIANA
UNIVERSITY

behavior of 3 algorithms

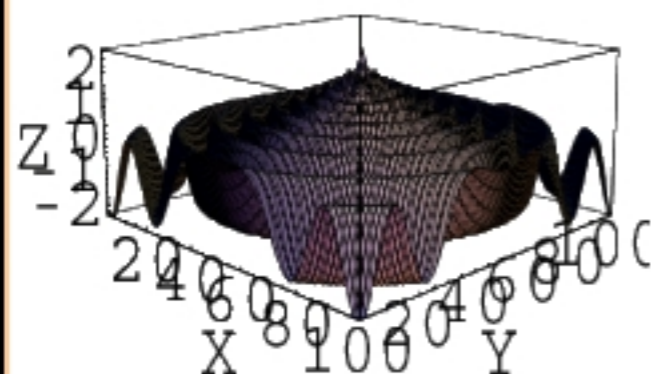


ABMGE

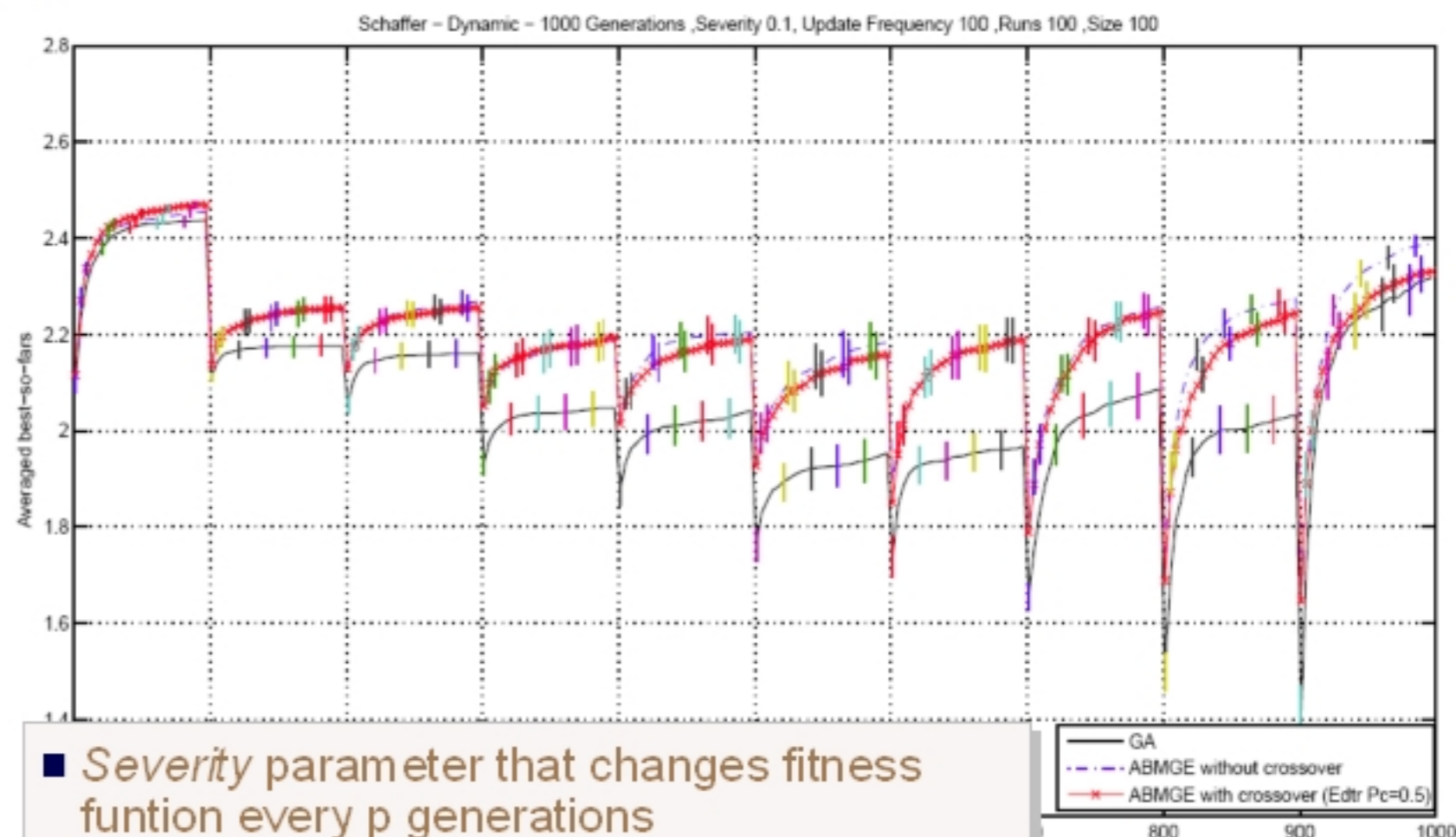
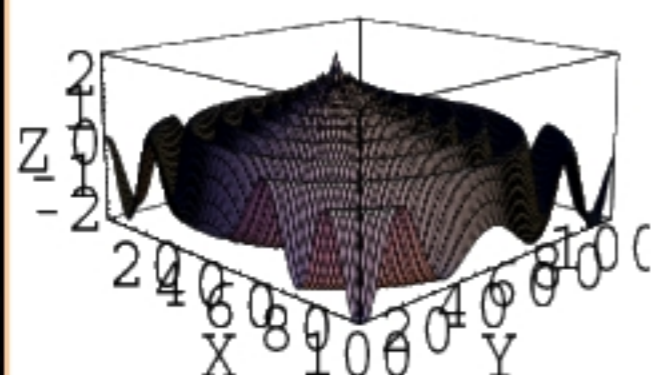


dynamic Shaffer function

l: delta=0



delta=.5



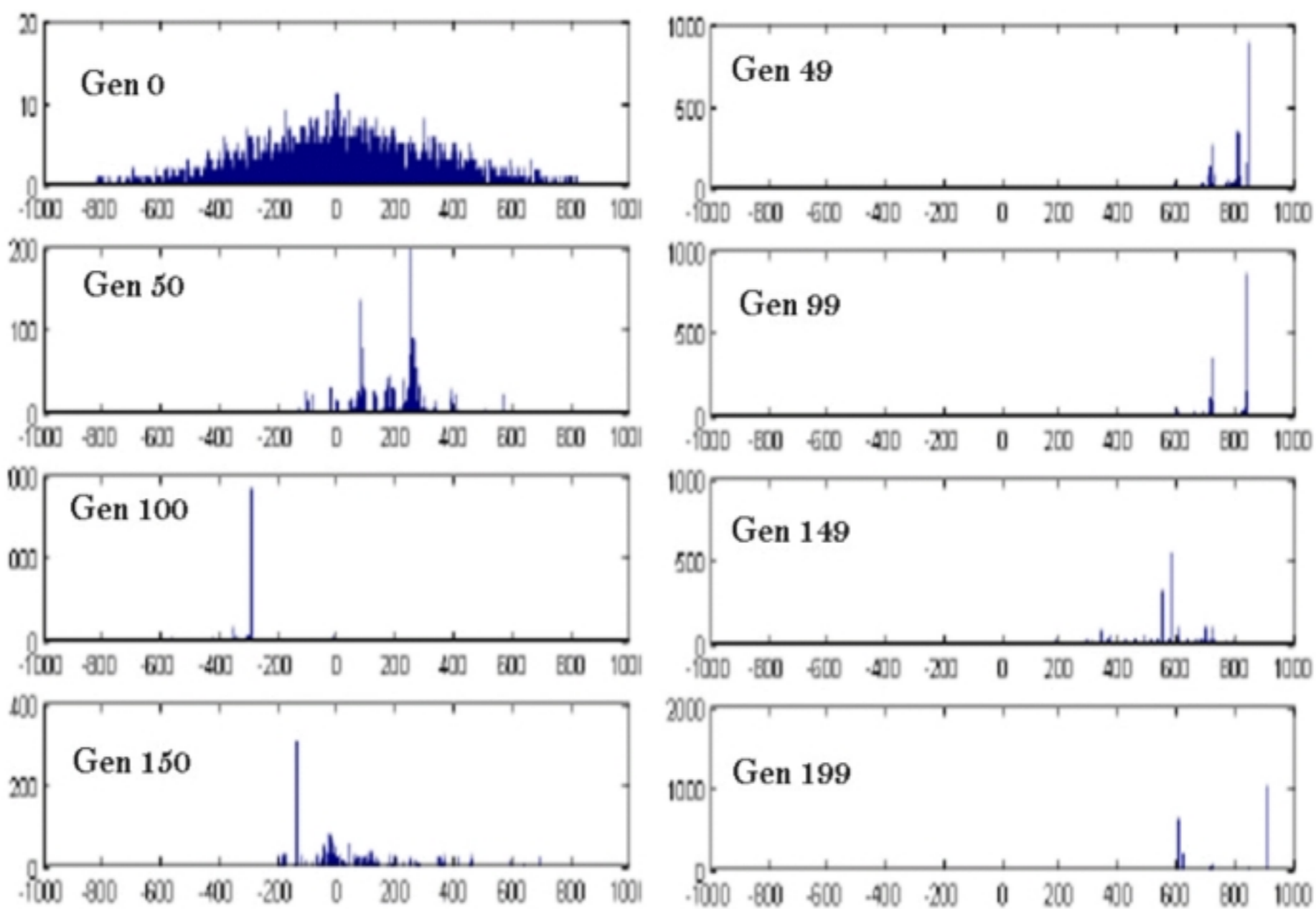
■ *Severity* parameter that changes fitness function every p generations

- ▶ **100** (50, 200)
- ▶ Linear and **jumping** dynamics

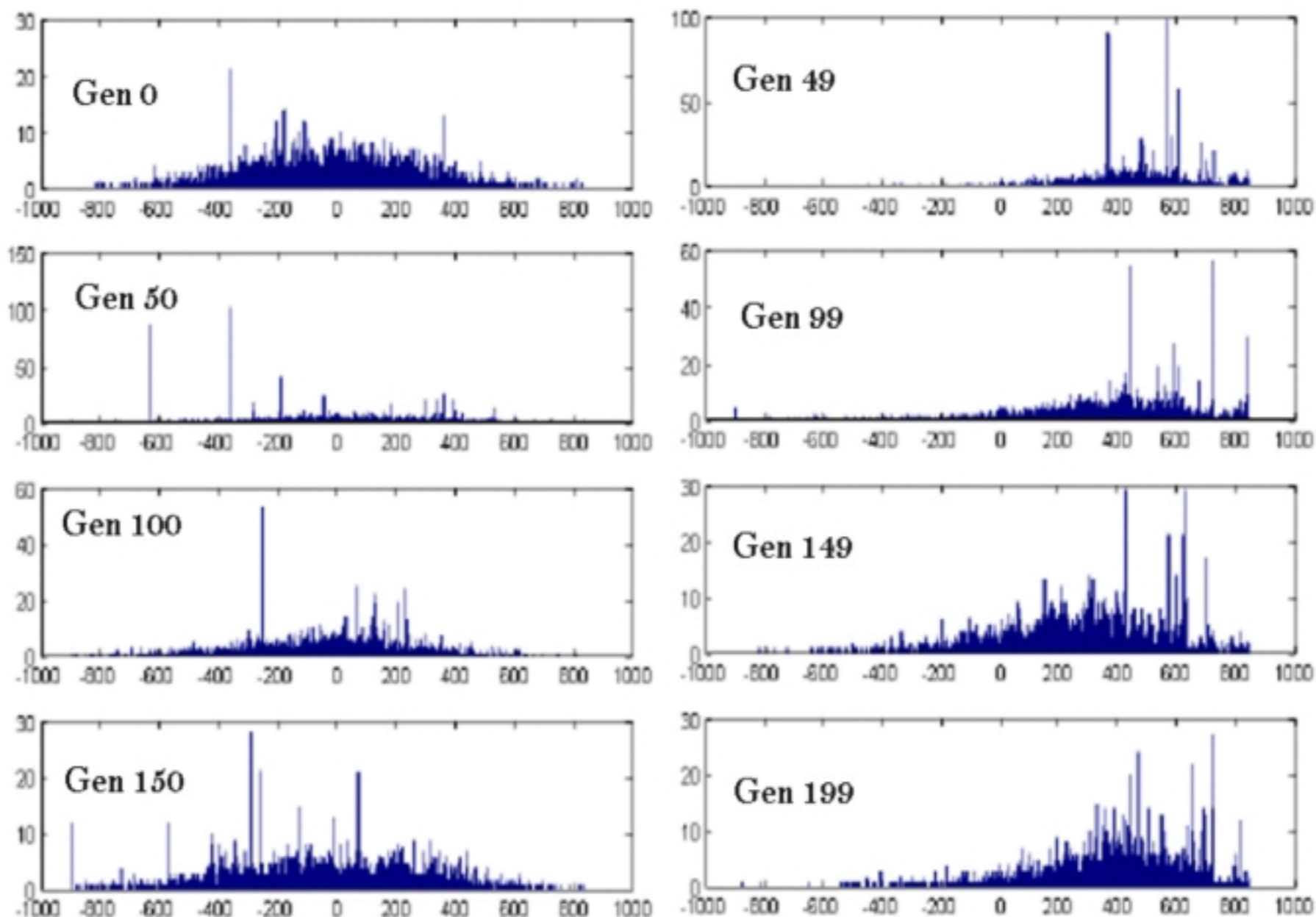
■ **50 agents**

- ▶ 1000 generations
- ▶ 100 runs

dynamic Schwafel function example (GA)

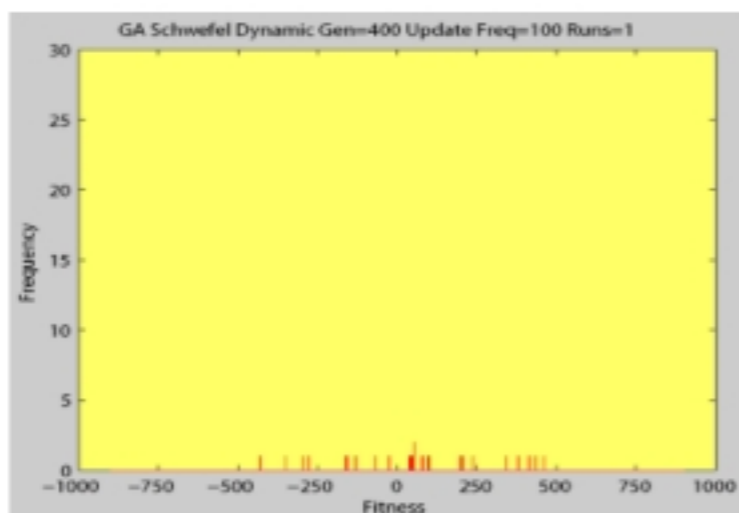


dynamic Schwafel function example (ABMGE)

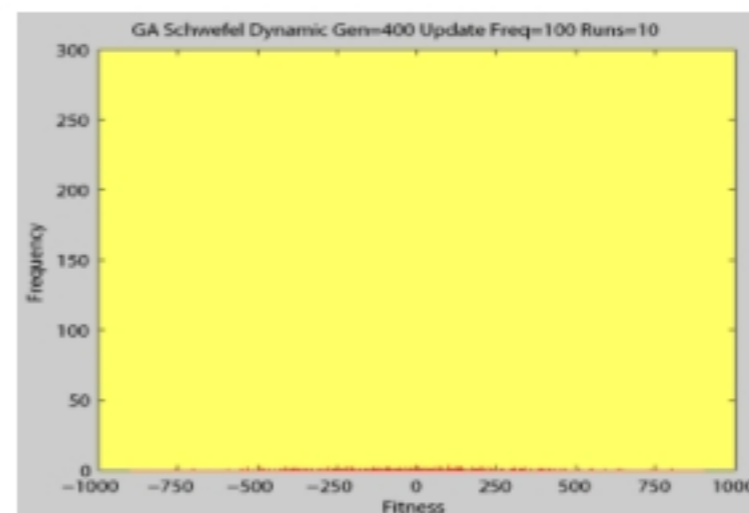


exploration and exploitation with genotype editing

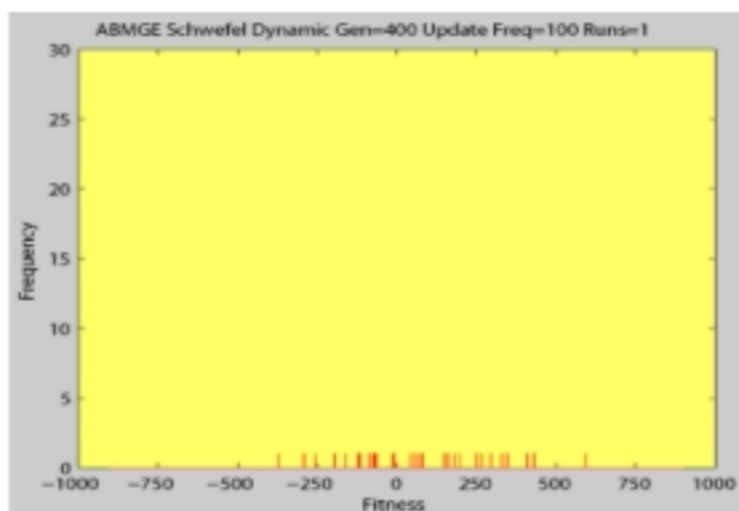
dynamic Schwafel function fitness distribution videos



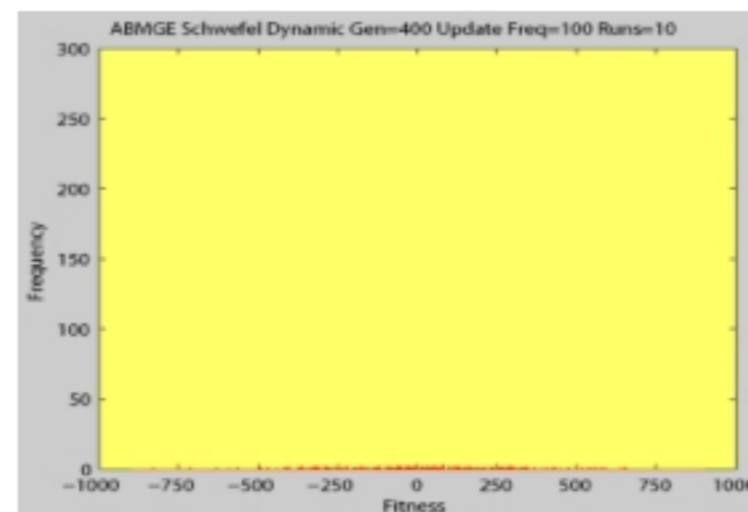
1 run



10 run



ABMGA



the gap between experimental reductionism vs. Systems view

The only consensus found among biologists about their subject is that biological systems are complicated, by any criterion of complexity that one may care to specify. [Rosen, 1972]

- Biology must simplify organisms to study them – some type of abstraction or modeling is needed.
 - ▶ External (Functional) description (favored by Systems Thinking)
 - *Blackbox*, input-output behavior of observables
 - Tells us what the system does
 - Function depends on repercussions in an environment
 - ▶ Internal (structural) description (favored by Experimentalists)
 - State description, trajectory behavior
 - Tells us how the system does what it does
 - Structural information can be measured for any component
 - ▶ Ideally, we would like to move between the two descriptions
 - But in Biology, the structural states we can measure, are not obviously related to the observed functional activities (and vice versa).
 - Thus, Systems Biology has mostly been relegated to deal with evolutionary problems, and Experimental Biology to increase our knowledge of the molecular components of organisms



Why Structural Reductionism is Not Sufficient

Destruction of Dynamical Properties

■ Naive Structural Decomposition

- ▶ Breaks an organism into simpler components, gathers information about those, and attempts to assemble information about the organism from the components
- ▶ But some properties of the original system cannot be reconstructed from components
 - E.g. the crucial stability properties of 3-body system cannot be reconstructed from knowledge of 2-body or 1-body constituents
 - the dynamics is destroyed.
 - Think what this means for the methodologies of molecular biology!

<http://faculty.ifmo.ru/butikov/Projects/Collection2.html>

<http://www.freewebz.com/vitaliy/triApplet/triGrav.html>

<http://www.dynamical-systems.org/threebody>

Coupling Structural Data with Functional Decomposition

- Biological Systems require “function-preserving” and “dynamics-preserving” Decompositions
 - ▶ In biology, the same physical structure typically is simultaneously involved in several functional activities
 - E.g. unlike airplanes, birds use the same structure (wing) as both propeller and airfoil
 - ▶ We must allow the simplifying decompositions to be dictated by system dynamics
 - Iterative Design of Experiments from Knowledge of Dynamics
 - Data accumulated from experiments based on naive structural decompositions are simply the first iteration!
 - ▶ Search for Global Patterns and Juxtaposed Functional Modes
 - E.g. studying global patterns of antigens rather than specific molecular interactions [Coutinho et al]
 - Spectral, PCA-like, Fourier Analysis approaches
 - ▶ Build Integrative Technology to Disseminate and Utilize Structural Data – for a diverse group of scientists

Integrative Link for bridging Experimental and Systems Biology

- Genome Informatics initially as enabling technology for the genome projects
 - ▶ Support for experimental projects
 - ▶ Genome projects as the ultimate reductionism: search and characterization of the function of information building blocks (genes)
- Post-genome informatics [Kanehisa 2000] aims at the synthesis of biological knowledge from genomic information
 - ▶ Towards an understanding of basic principles of life (while developing biomedical applications) via the search and characterization of networks of building blocks (genes and molecules)
 - The genome contains information about building blocks but, given the knowledge of Systems Biology, it is naive to assume that it also contains the information on how the building blocks relate, develop, and evolve.
 - ▶ Interdisciplinary: biology, computer science, mathematics, and physics



"Life is a complex system for information storage and processing".
Minoru Kanehisa
[2000]

Enabling a Systems Approach to Biology

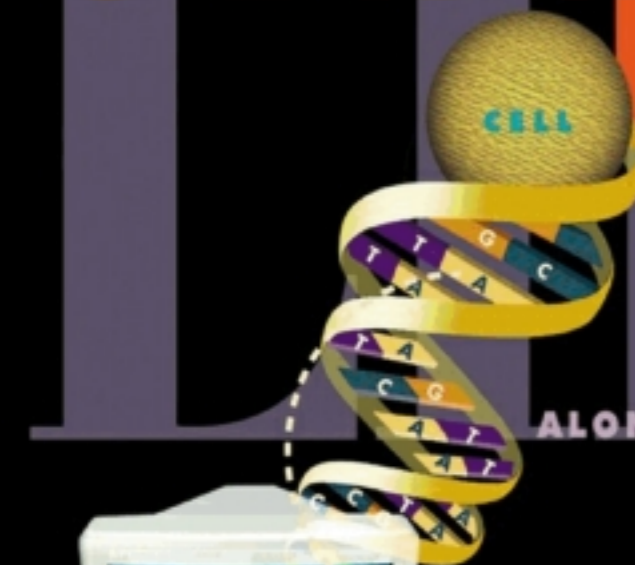
- Not just support technology but involvement in the systematic design and analysis of experiments
 - ▶ *Functional genomics*
 - ▶ Where, when, how, and why of gene expression
 - ▶ *Post-genome informatics* aims to understand biology at the molecular network level using all sources of data: sequence, expression, diversity, etc.
 - ▶ Cybernetics, Systems Theory, Artificial Life, Complex Systems approach to Theoretical Biology
 - ▶ *Synthetic Biology*: to engineer novel life forms and bio-technology
- **Grand Challenge of Computational Biology**
 - ▶ Given a complete genome sequence, reconstruct in a computer the functioning of a biological organism
 - Regards Genome more as set of initial conditions for a dynamic system, not as complete blueprint (Pattee, Rosen, Atlan). The genome can be contextual and dynamically accessed and even modified by the complete network of reactions in the cell (e.g. editing).
 - Uses additional knowledge for integration comparative analysis: Comparative Biology
- **Grand Challenge of Synthetic Biology**
 - ▶ **If we understand it, we can build it!**
 - intentional design of real biological systems
 - reversal of aging and innovative medical treatments such as beneficial bacterial infections programmed to augment immunity

- *Systems biology* is a unique approach to the study of genes and proteins which has only recently been made possible by rapid advances in computer technology. Unlike traditional science which examines single genes or proteins, systems biology studies the complex interaction of all levels of biological information: genomic DNA, mRNA, proteins, functional proteins, informational pathways and informational networks to understand how they work together. Systems biology embraces the view that most interesting human organism traits such as immunity, development and even diseases such as cancer arise from the operation of complex biological systems or networks.
 - ▶ Institute for Systems Biology: <http://www.systemsbiology.org>
 - ▶ Kitano Symbiotic Systems Project: <http://www.symbio.jst.go.jp/>
- The “New” Systems Biology is not novel per se, it is rather a result of new enabling technology for doing “Old” Systems Biology
 - ▶ But it is finally allowing experimentalists to work with theorists.

GENOMES to LIFE

BIOLOGICAL SOLUTIONS FOR ENERGY CHALLENGES

INNOVATIVE APPROACHES ALONG UNCONVENTIONAL PATHS

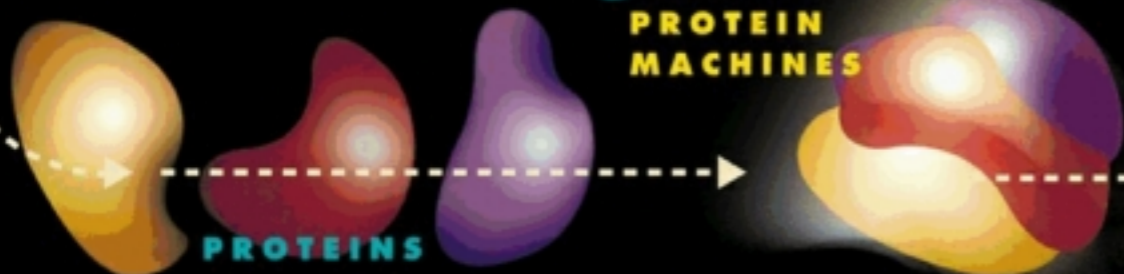


DNA SEQUENCE DATA FROM GENOME PROJECTS



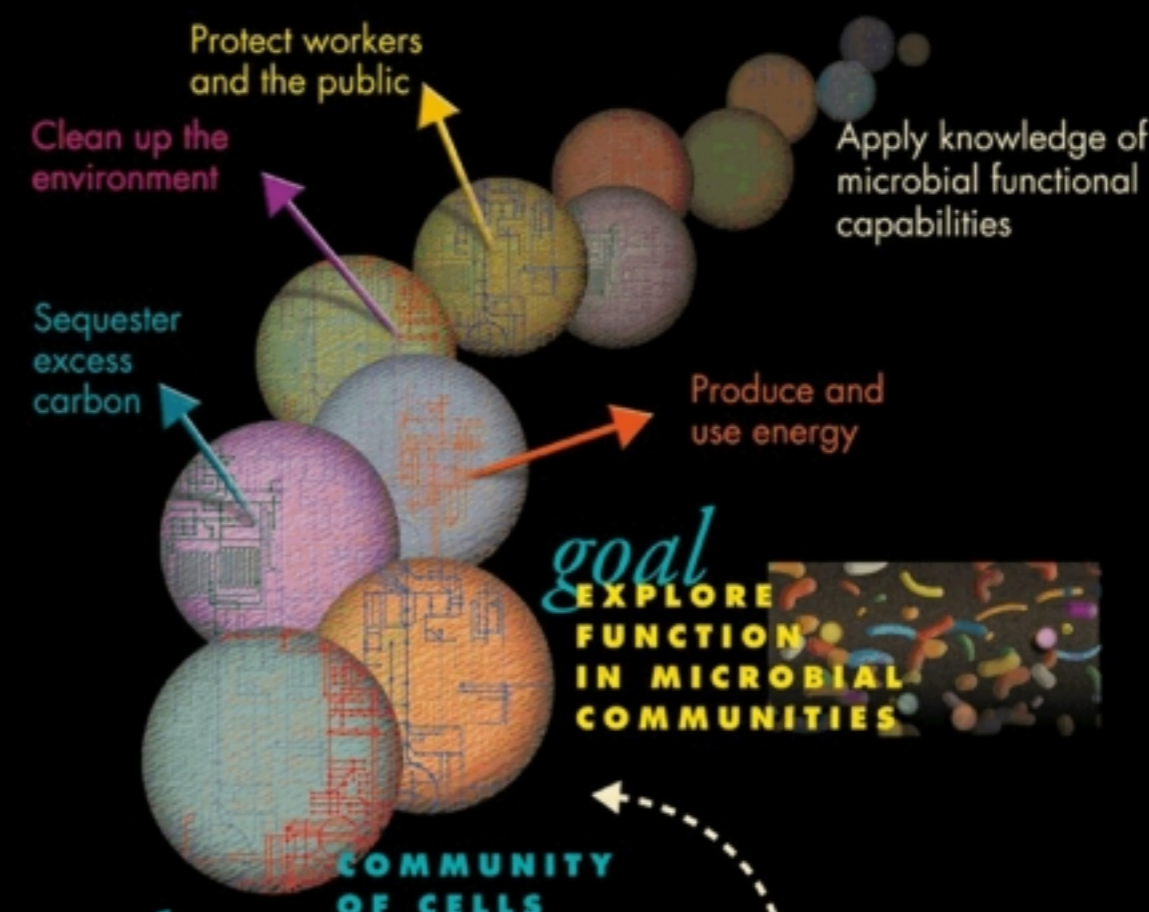
Genes and other DNA sequences contain instructions on how and when to build proteins

goal
IDENTIFY PROTEIN MACHINES



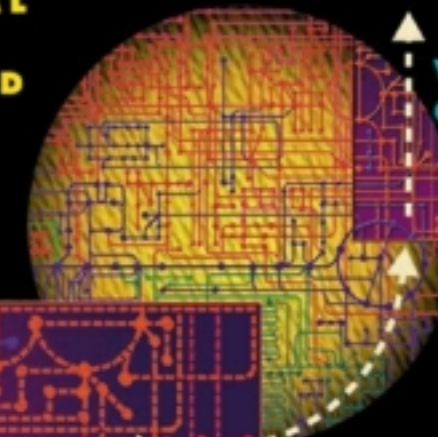
PROTEINS

Proteins perform many of life's most essential functions. To carry out their specific roles, they often work together in the cell as protein machines.



COMMUNITY OF CELLS

goal
DEVELOP COMPUTATIONAL CAPABILITIES TO UNDERSTAND COMPLEX BIOLOGICAL SYSTEMS



WORKING CELL

Many protein machines interact through complex, interconnected pathways. Analyzing these dynamic processes will lead to models of life processes.

goal
CHARACTERIZE GENE REGULATORY NETWORKS

URL DOEGenomesToLife.org

■ Experimental Side

- ▶ Improving cellular measurement methods
 - High-throughput identification of the components of protein complexes; Parallel, comparative, high-throughput identification of DNA fragments among microbial communities and for community characterization; Whole-cell imaging including in vivo measurements; Better Separation techniques.
- ▶ Measurements Based on Functional Decompositions
 - Functional assays? Flexible, fast, novel experimental design based on informatics results.

■ Computational Side

- ▶ Integrative Technology
 - Standardized formats, databases, and visualization methods
 - Automated collection, integration and analysis of biological data
 - Algorithms for genome assembly and annotation and measurement of protein expression and interactions;
- ▶ Simulation Technology
 - Improved methods for distributed simulation, analysis, and visualization of complex biological pathways;
 - Prediction of emergent functional capabilities of microbial communities



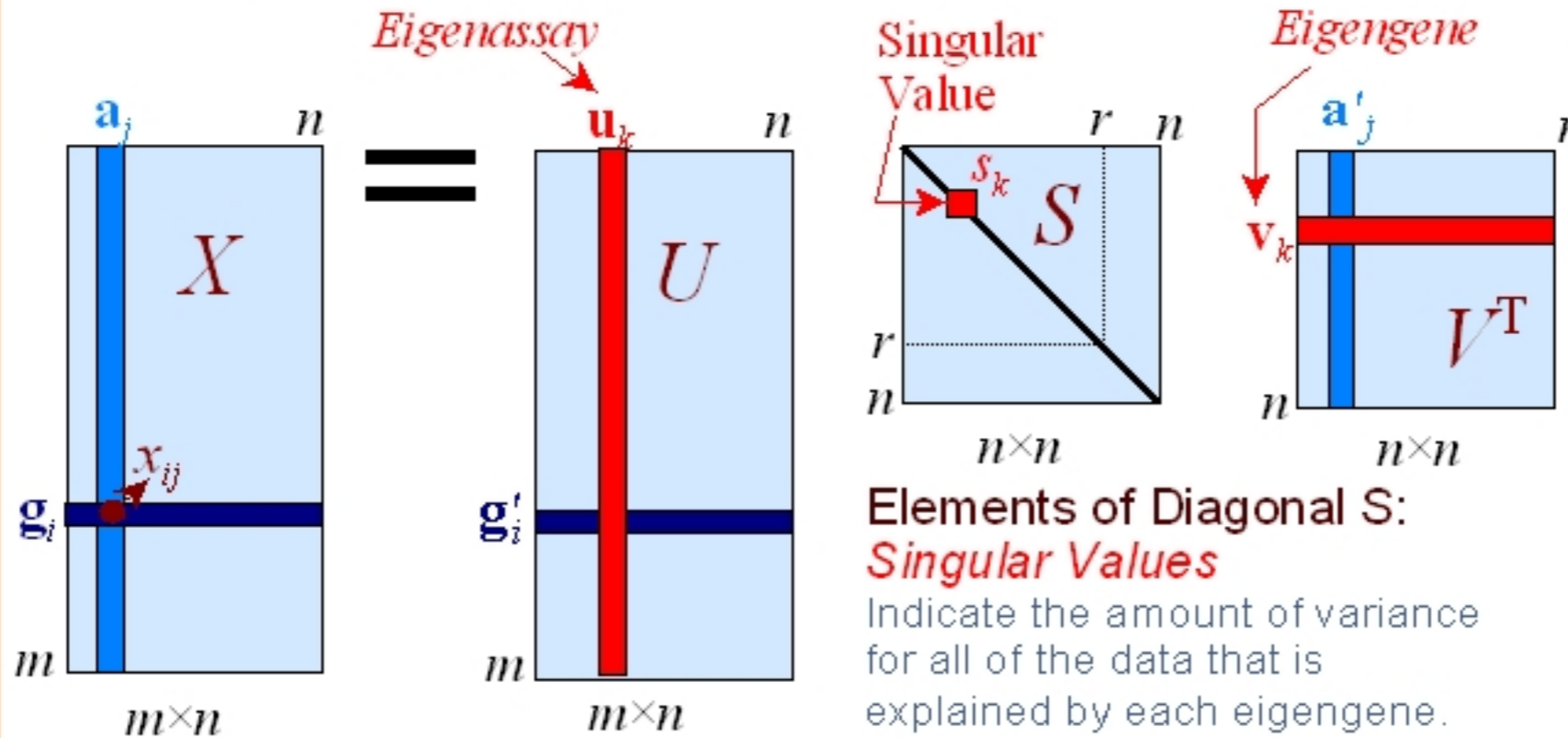
Continuation

■ Modeling Side

- ▶ Algorithms for Discovery of Global Patterns and Juxtaposed Functional Modes
 - Pattern Recognition, data-mining, “Spectral” methods.
- ▶ Network Models and Analysis
 - Predictive Models based on biochemical pathways of observed networks
 - Simplification Strategies for Network Modeling
 - Reduction of possible cell-behaviors from steady-state models of metabolic network models
 - High-Performance Algorithms to allow whole-system Kinetic models



for microarray analysis



Rows of V^T : *eigengenes* (columns are time steps). Each gene's expression pattern is a linear combination of the eigengene patterns.

$$\mathbf{g}_i = \sum_{k=1}^r u_{ik} s_k \mathbf{v}_k, \quad i:1, \dots, m$$

$$X = USV^T$$

Gene Expression Matrix: Columns are assays (time steps) and rows are genes

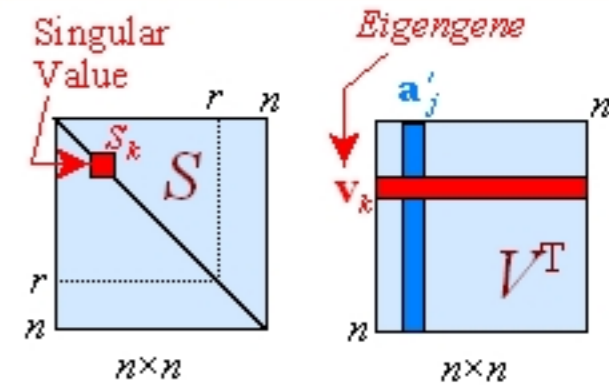
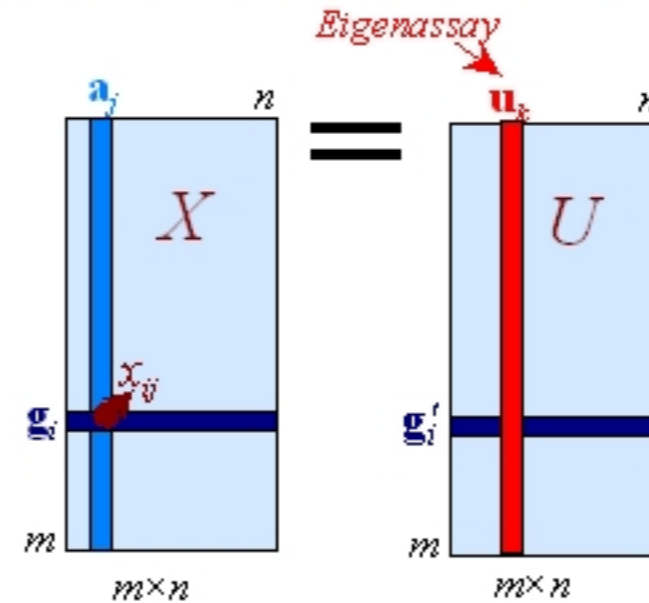
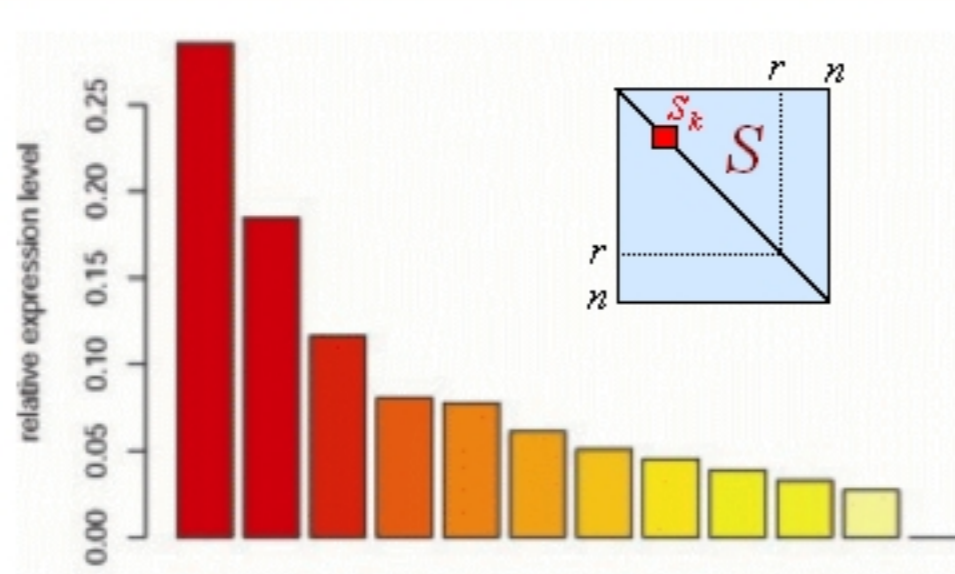
Columns of U : *eigenassays* (rows are genes) describe how each component contributes to a single gene's expression pattern

$$\mathbf{a}_j = \sum_{k=1}^r \mathbf{v}_{jk} s_k \mathbf{u}_k, \quad j:1, \dots, n$$

Wall, Rechtsteiner and Rocha [2002]. "Singular value decomposition and principal component analysis". In *Understanding and Using Microarray Analysis Techniques: A Practical Guide*. D.P. Berrar, W. Dubitzky, M. Granzow, eds.

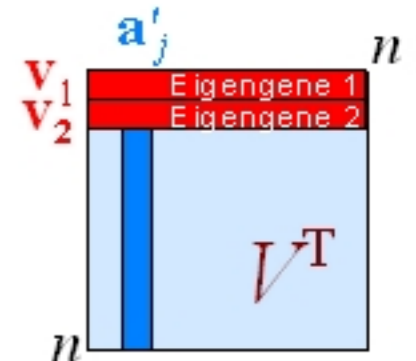
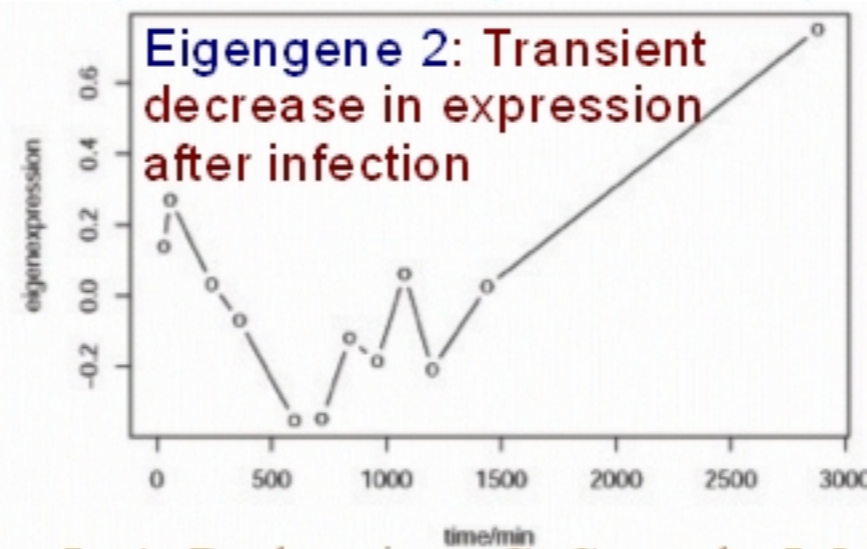
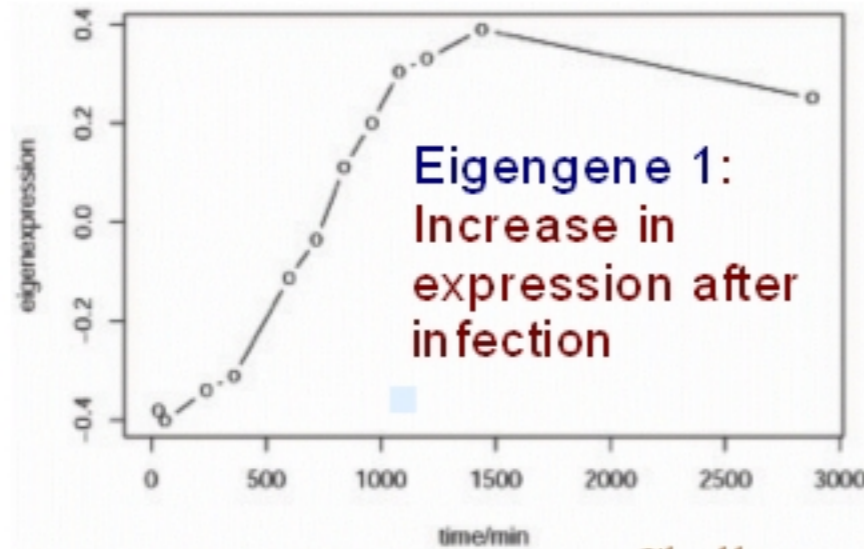
singular value decomposition

gene expression (13000 genes) after infection with herpes virus



$$X = USV^T$$

12 point time series (30min - 48hrs)



Challacombe, J., A. Rechtsteiner, G. Gottardo, L.M. Rocha, E.P. Brown, T. Shenk, M. Altherr, T. Brettin [2004]. "Evaluation of the host transcriptional response to human cytomegalovirus infection". *Physiol. Genomics*. 10.1152

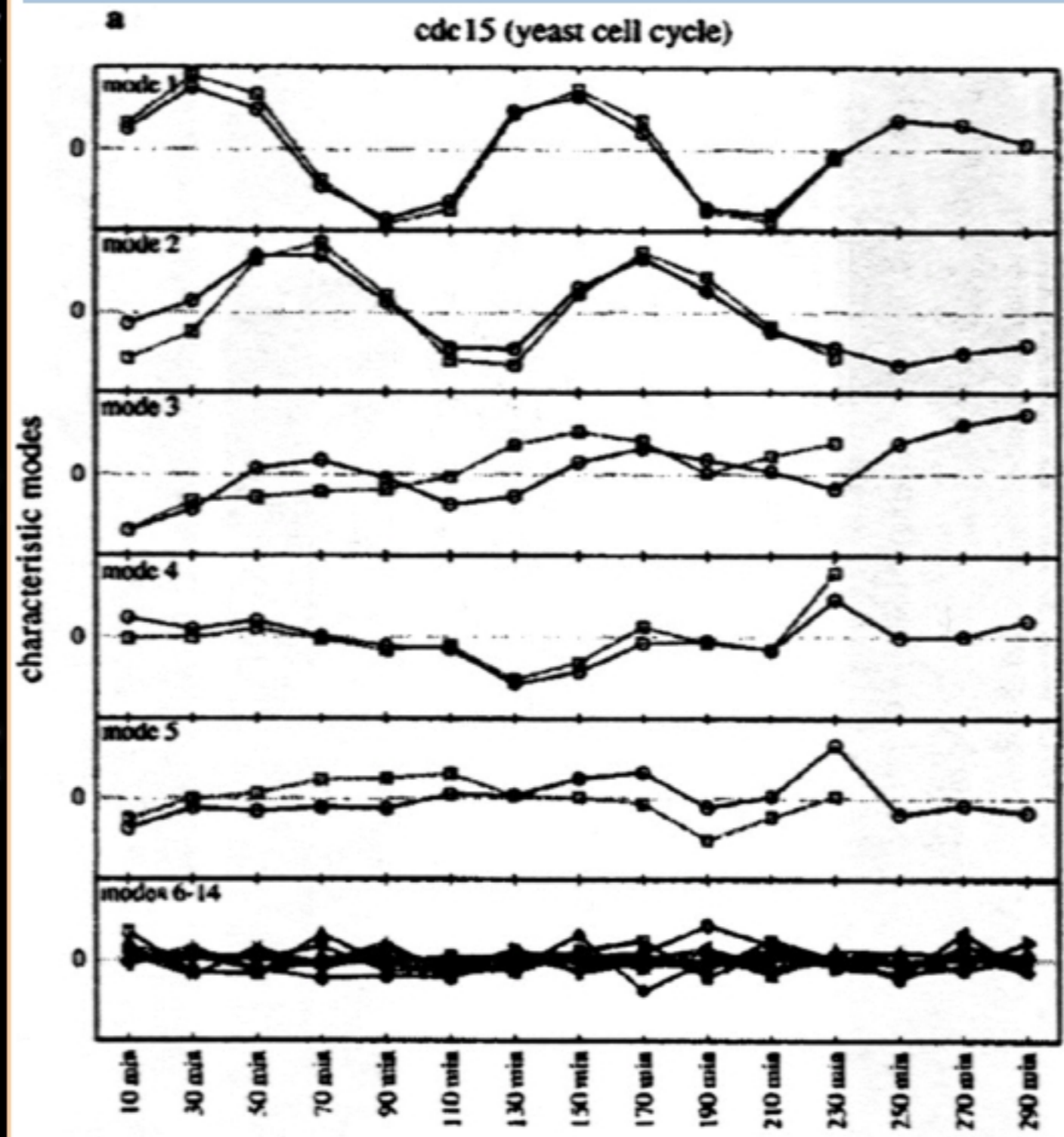
Discovery of Juxtaposed Functional Modes

■ Gene Expression Modes

- ▶ Cluster analysis provides little insight into inter-relationships among groups of co-regulated genes. Tends to demand separated groupings.
- ▶ Component ("spectral") analysis yields a description of superposed behavior of gene expression networks, rather than a partition.
 - PCA, SVD, etc.
 - Holter et al [2000] compares the superposed components to the characteristic vibration modes of a violin string which entirely specify the tone produced
- ▶ Holter et al [2000] compared SVD analysis of yeast cdc15 cell-cycle [Spellman et al 1998] and sporulation [Chu et al, 1998] data sets, as well as the data set from serum-treated human fibroblasts [Iyer et al, 1999].
 - Essential temporal behavior is captured by first 2 modes (sine and cosine)
 - Large group of genes with same sinusoidal period but dephased

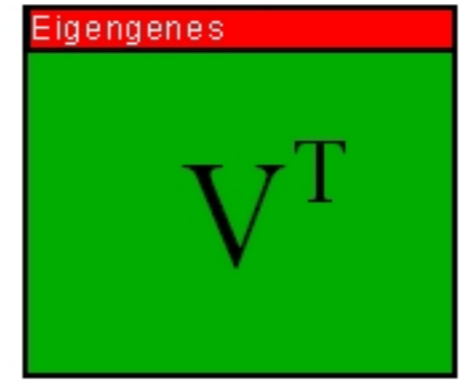


Holter et al SVD Analysis

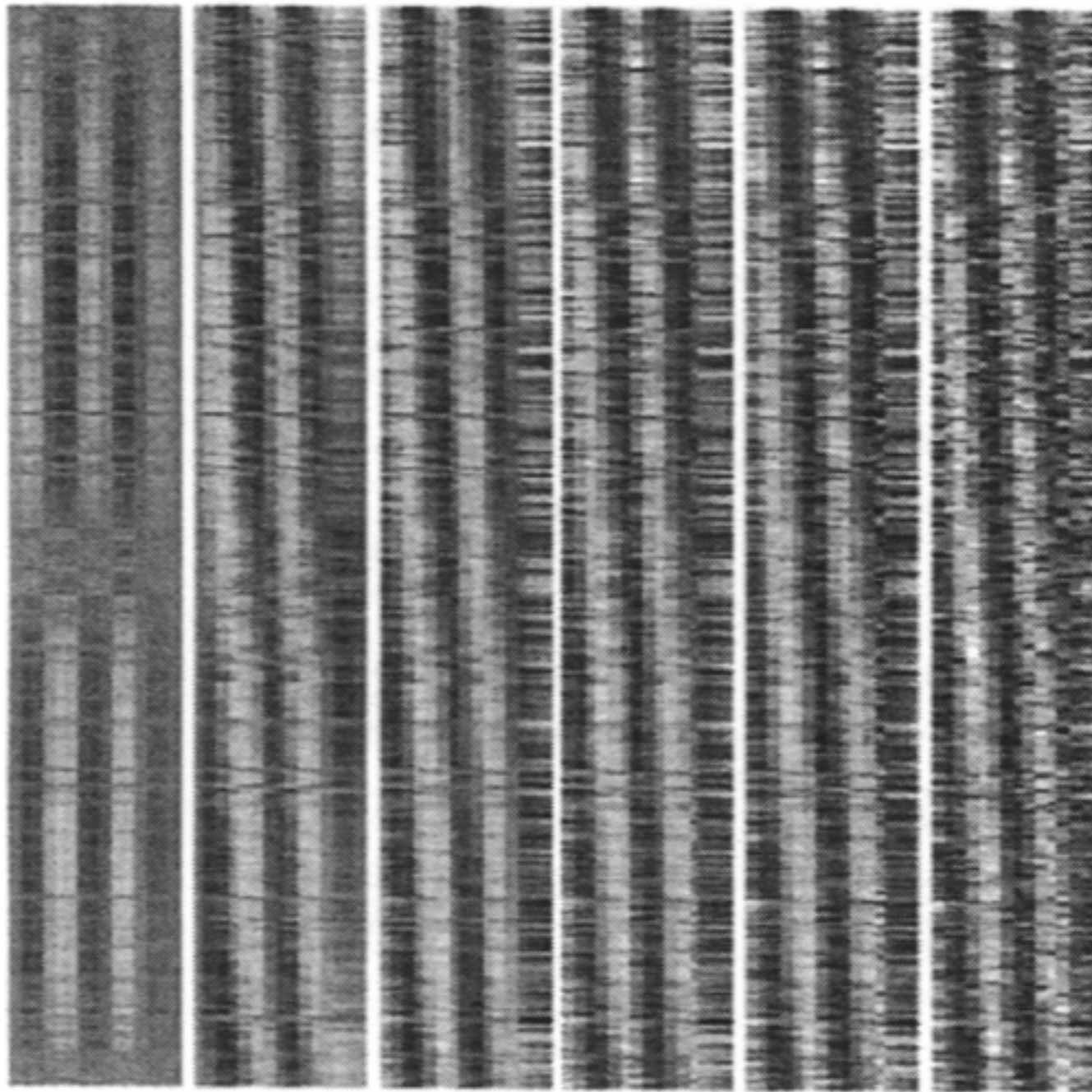


- 800 genes by 15 (12) time measurements
- 2 dominant modes
 - ▶ Approximately sinusoidal and out of phase
 - ▶ Less synchronized as cell enters 3rd cycle
 - ▶ If only 12 points are used, third SV loses relevance, but 2 first components remain largely unchanged

Eigengene: rows of V^T (each column is a time instance)



cdc15 Reconstruction with k-highest modes



1

2

3

4

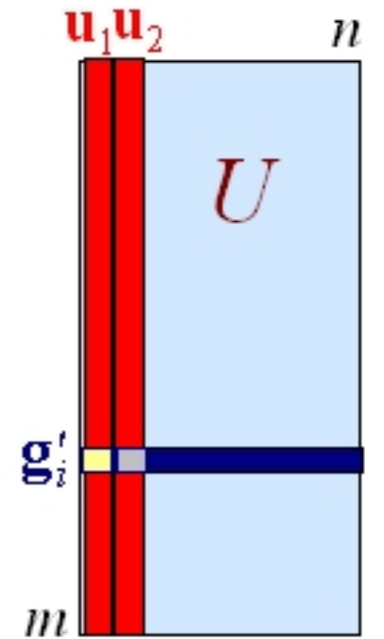
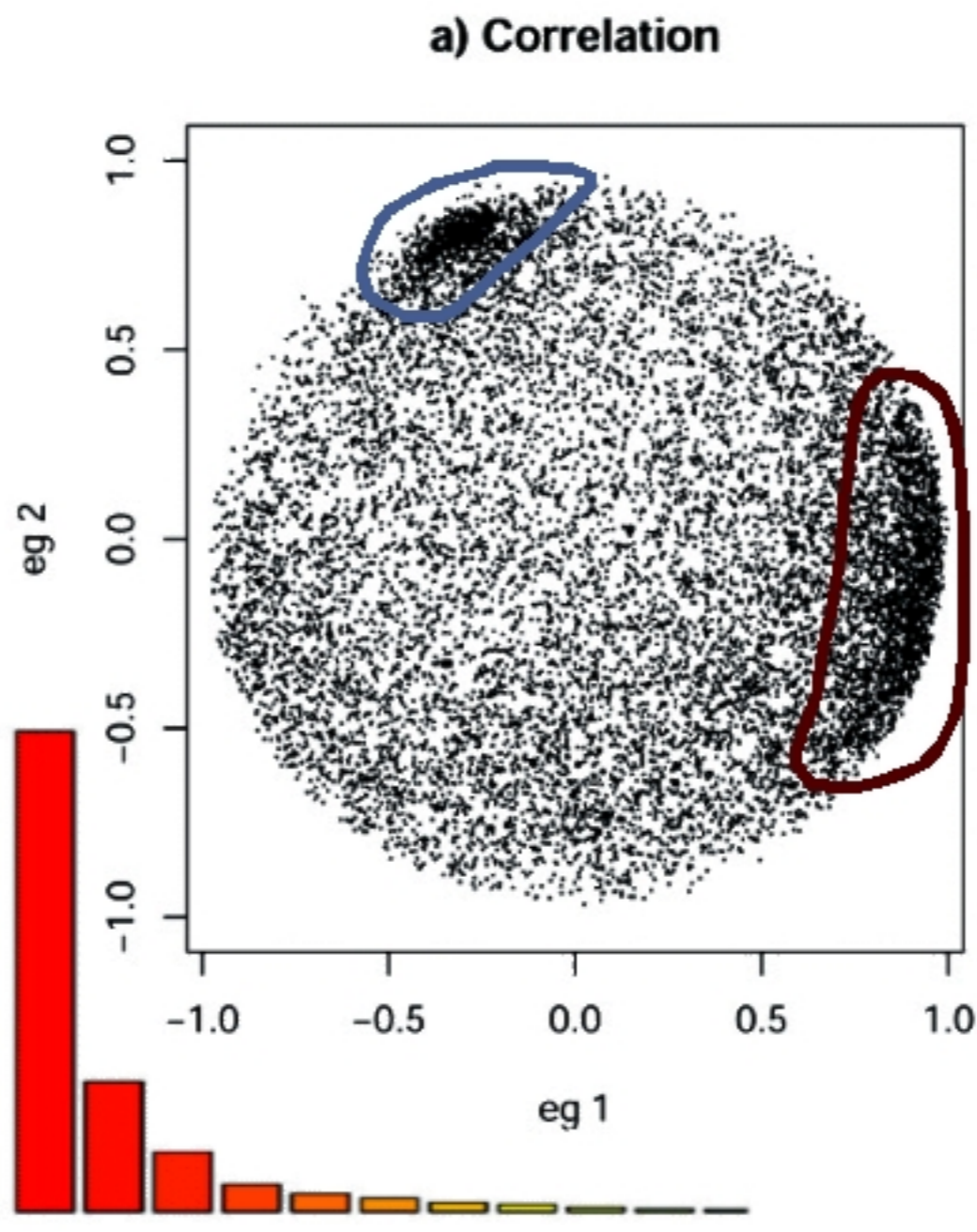
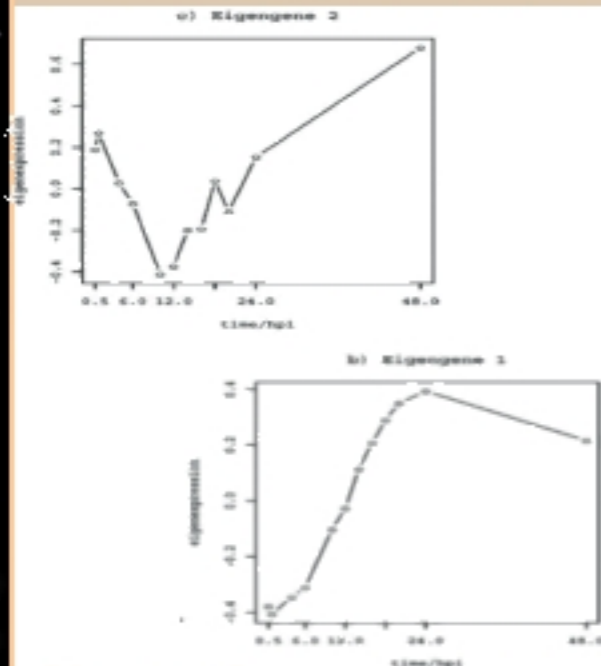
5

14

Rows are genes
Columns are time
points

It implies an
undelying simplicity in
genetic response

eigenassay coefficient plot: human cytomegalovirus infection

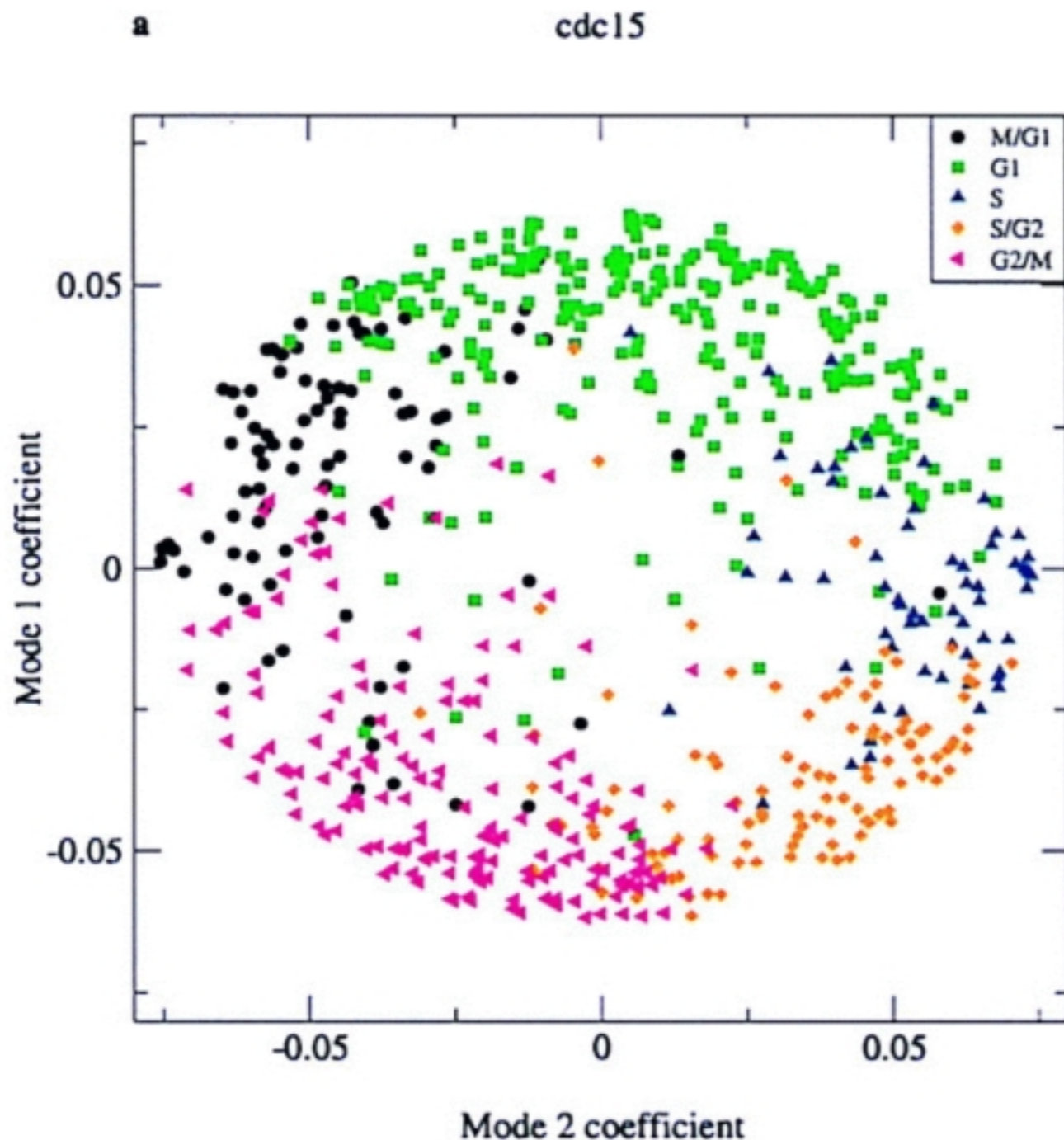


Cluster 2:
 Genes involved in immune system regulation, signal transduction and cell adhesion. Also mainly in cluster 2, genes targeted by HCMV's immune evasion strategies.

Cluster 1:
 Genes involved in *transcriptional regulation*, oncogenesis and cell cycle regulation. Also mainly in cluster 1, genes involved in the host response to HCMV infection.

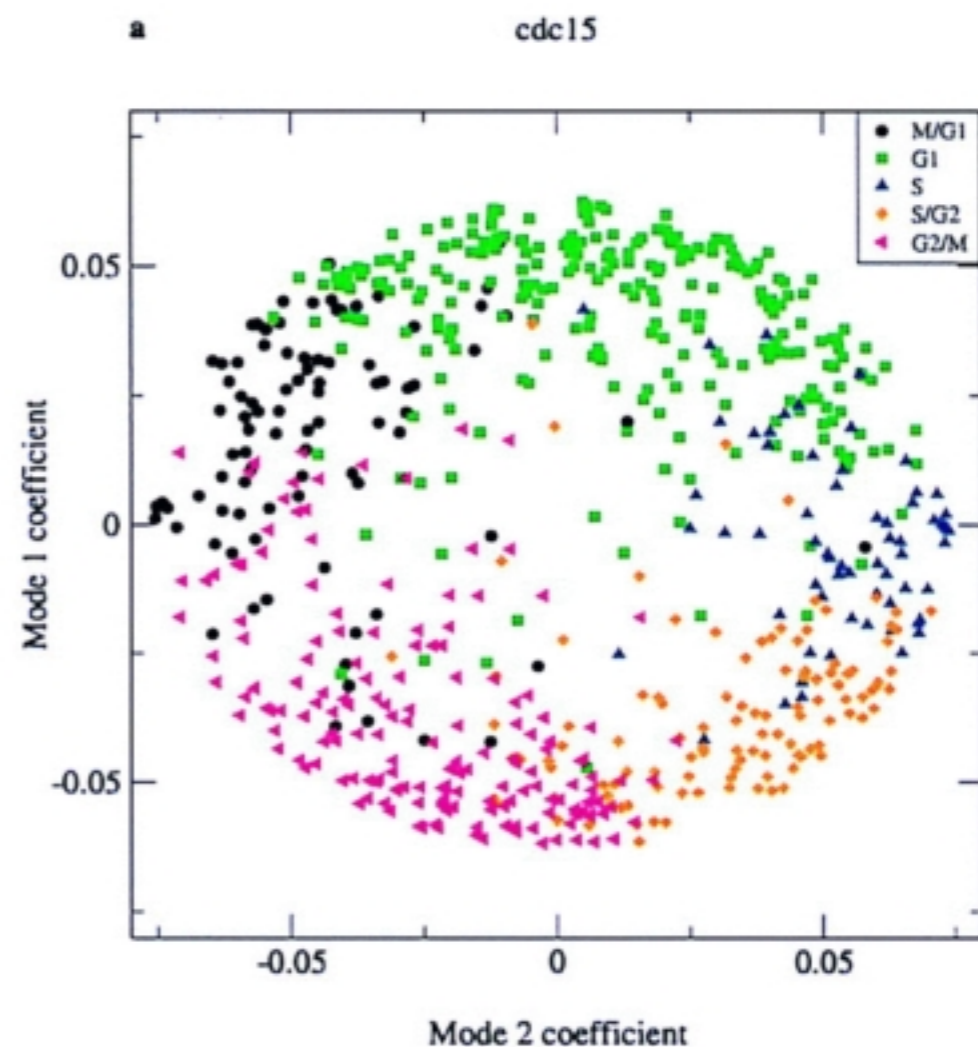
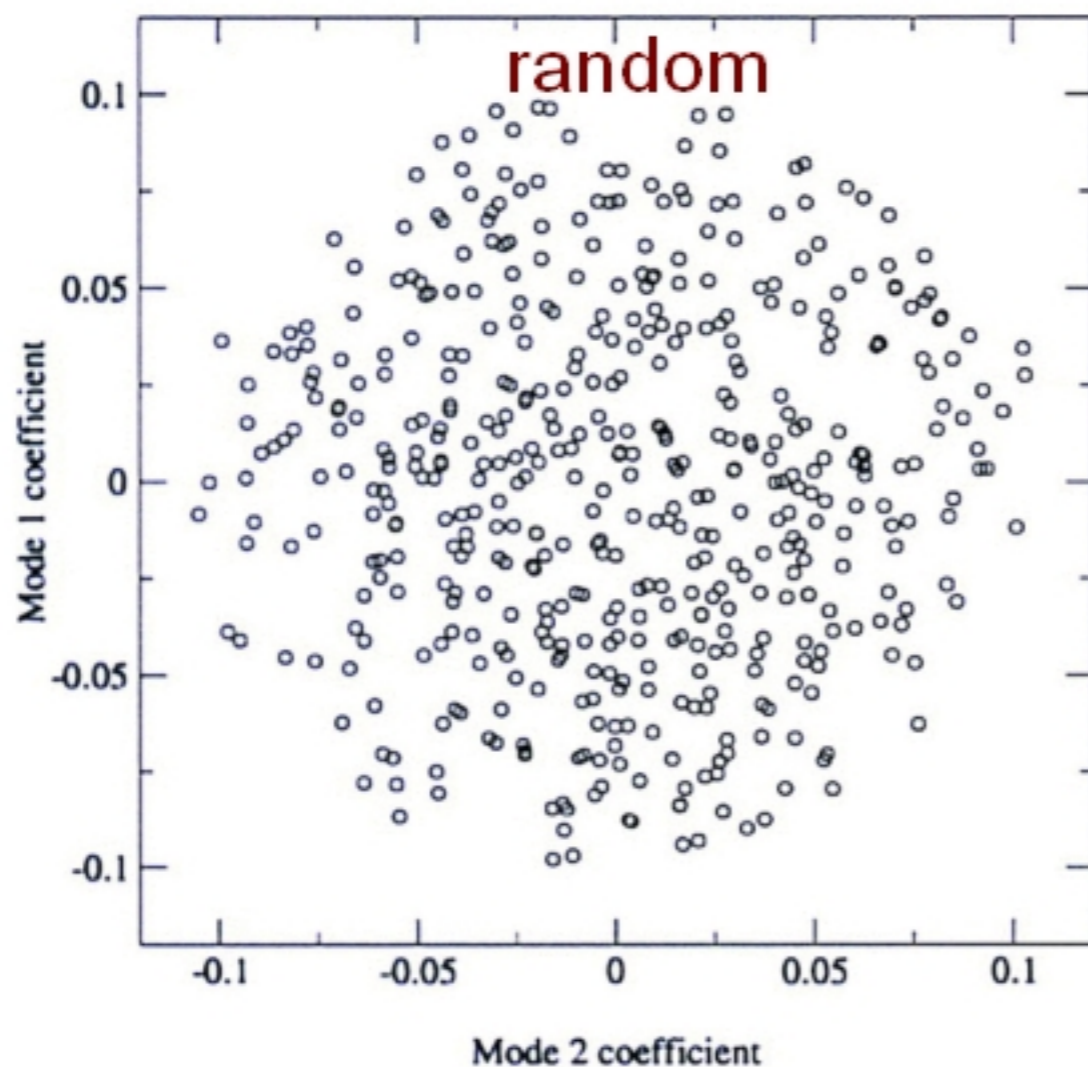


Plot of the coefficients of the first 2 modes for all genes



- Clusters of genes by other methods cluster in these plots, but the temporal progression in the cell cycle and in the course of sporulation is more evident in the SVD analysis
- Holter et al conclude that genes are not activated in discrete groups or blocks, as historically implied by the division of the cell cycle into phases or the sporulation response into temporal groups. There is a continuity in expression change

Random data



Fill most of the plot because genes are not very correlated with components. A circle implies equal contribution from each component (rather than an ellipse)

Most common areas

■ **Journals**

- ▶ Bioinformatics
- ▶ BMC Bioinformatics
- ▶ Journal of Theoretical Biology
- ▶ PNAS
- ▶ Biosystems
- ▶ Genome Research
- ▶ IEEE Transactions on Computational Biology and Bioinformatics

■ **Conferences**

- ▶ Intelligent Systems for Molecular Biology (ISMB)
- ▶ Research in Computational and Molecular Biology (RECOMB)
- ▶ Pacific Symposium on Biocomputing (PSB)

■ **Areas**

- ▶ Genome Analysis
- ▶ Sequence Analysis
- ▶ Systems Biology
- ▶ Data and Text Mining
- ▶ Structural Bioinformatics
- ▶ Gene Expression
- ▶ Genetics and Population Analysis



■ Direction

- ▶ Marie-France Sagot (Program Director)
- ▶ Jorge Carneiro (Program Deputy-Director)
- ▶ Luis Rocha (Collaboratorium Director)

■ Background

- ▶ Knowledge of empirical sciences (Physics, Chemistry, Biology) and quantitative technical disciplines (programming, applied mathematics, statistics)
 - Introduction Module to catch up on biology, modeling, and CS

■ Regular Syllabus

- ▶ Training in Biology
 - Molecular evolution and sequences
 - Theory of evolution and population genetics, sequence alignment, from pairwise to multiple, from genes to genomes, molecular phylogeny
 - Structures (DNA, RNA and proteins)
 - Introduction to biomolecular structures, determination and visualisation, biomolecular structure mechanics, dynamics, prediction and design



■ Regular Syllabus (Cont)

▶ Training in Biology

- Genome structure
- Genome evolution and genome dynamics
- Function classification
- Transcriptomics and proteomics
- Networks
 - Generic aspects, Protein interaction networks, Metabolic networks, Genetic networks
- Systems Biology
 - Population biology, epidemiology and immune system, Cytoskeleton and cell morphogenesis, motion and chemotaxis modelling, Development and whole organism modelling, Evolutionary development, Computational Neurobiology

▶ Training in Computer Science

- Algorithms in computational biology
- Statistical data mining and machine learning
- Database management systems, knowledge systems and integration
- Introduction to dynamical systems

<http://bc.igc.gulbenkian.pt/pdbc/syllabus.htm>



■ Monday, June 19, 2006

▶ 10:00 – 13:00 **Introduction: *From Bioinformatics to Systems Biology***

– by Luis M. Rocha, Indiana University and Instituto Gulbenkian de Ciencia

▶ 14:30 – 15:30 ***Microarray Data Analysis with Data Mining and Machine Learning Methods***

– by Miguel Rocha and Isabel Rocha, Universidade do Minho

▶ 15:30 – 16:30 ***Modeling and Optimization of Metabolic and Regulatory Networks in Systems Biology*** by Miguel Rocha and Isabel Rocha, Universidade do Minho

■ Tuesday, June 20, 2006

▶ 10:00 – 11:00 ***GENE-CBR: a Case-Based Reasoning Tool for Cancer Diagnosis using Microarray Datasets***

– by Florentino Fernández Riverola, Universidad de Vigo.

▶ 11:00 – 12:00 ***Bibliome Informatics***

– by Luis M. Rocha, Indiana University and Instituto Gulbenkian de Ciencia

▶ 12:00 – 13:00 ***Machine Learning Methods for Computational Proteomics and Beyond***

– by Pierre Baldi, University of California, Irvine