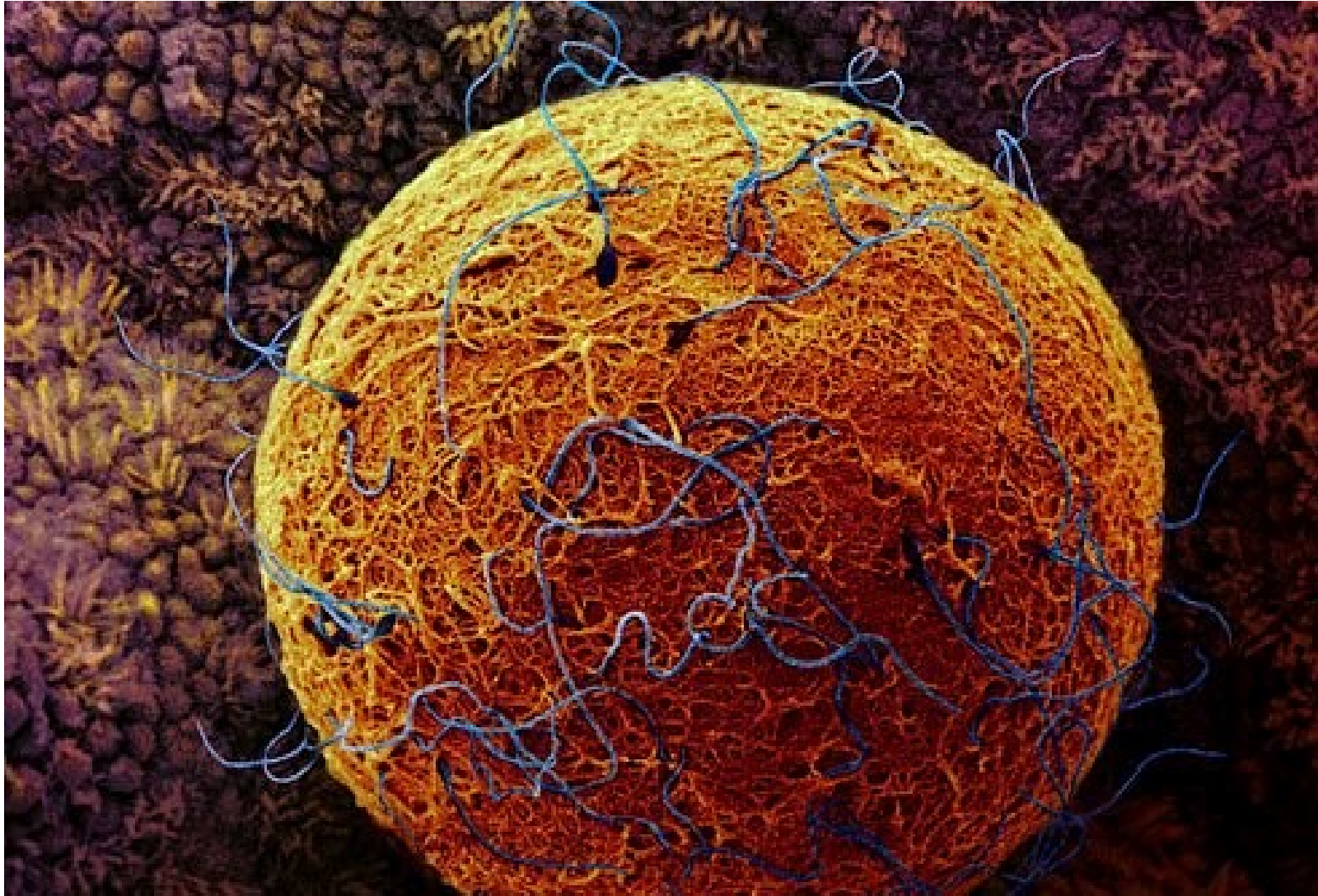


Computational biology

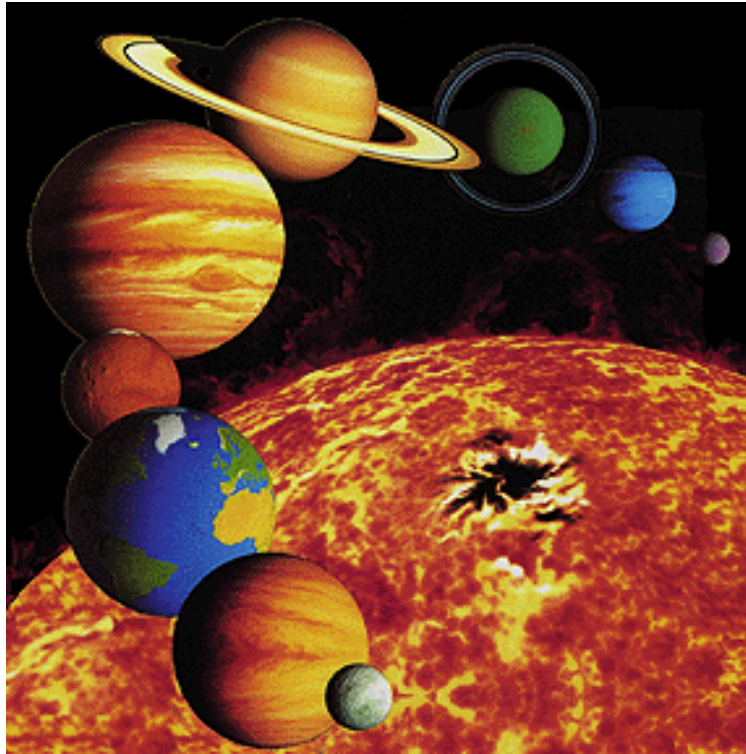


Il materiale di seguito riportato è destinato solo ed unicamente all'attività didattica nell'ambito del CdLM in Biotecnologie della Riproduzione, Università di Teramo

fertility

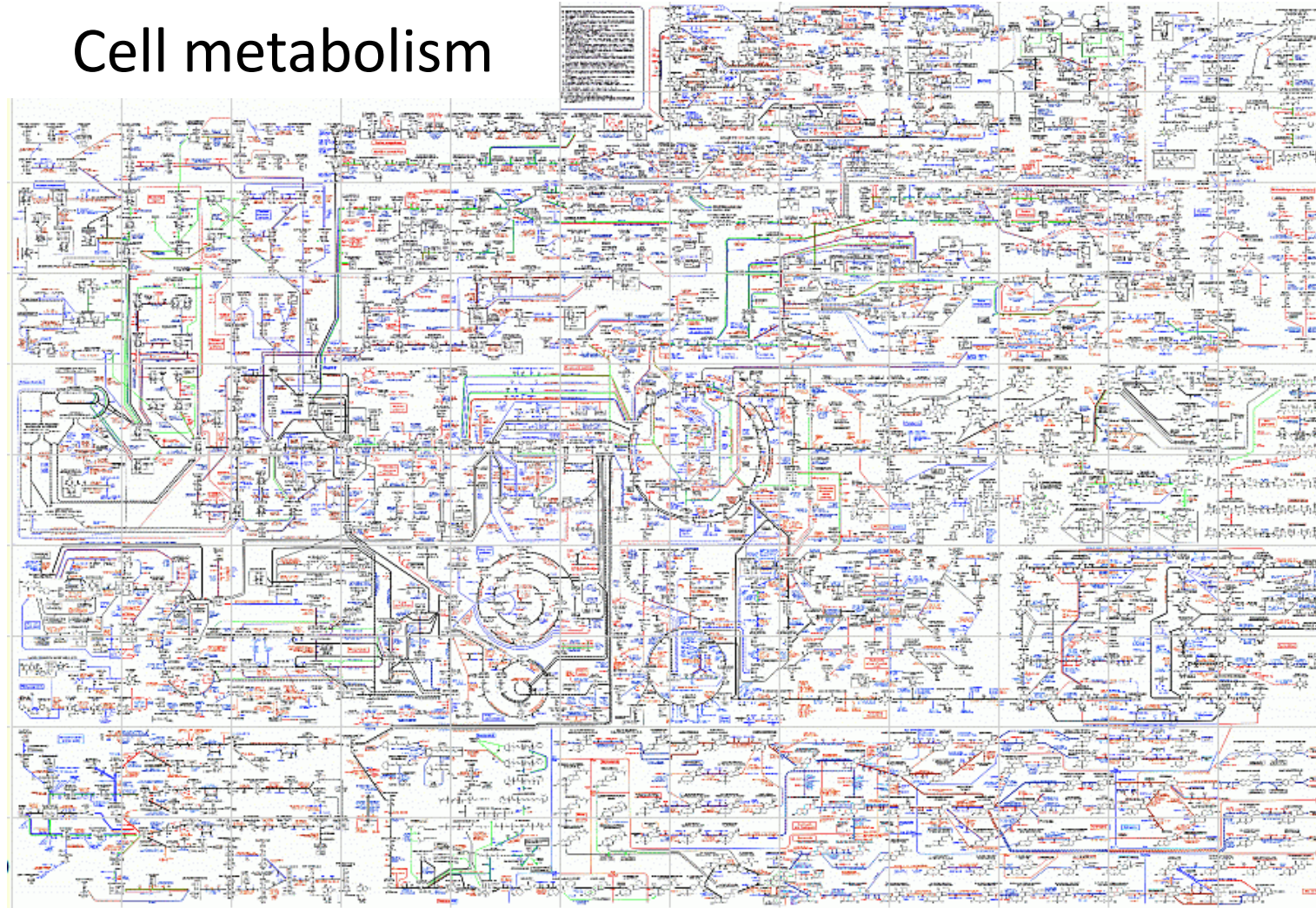


Complicated vs. complex



Biological complexity

Cell metabolism

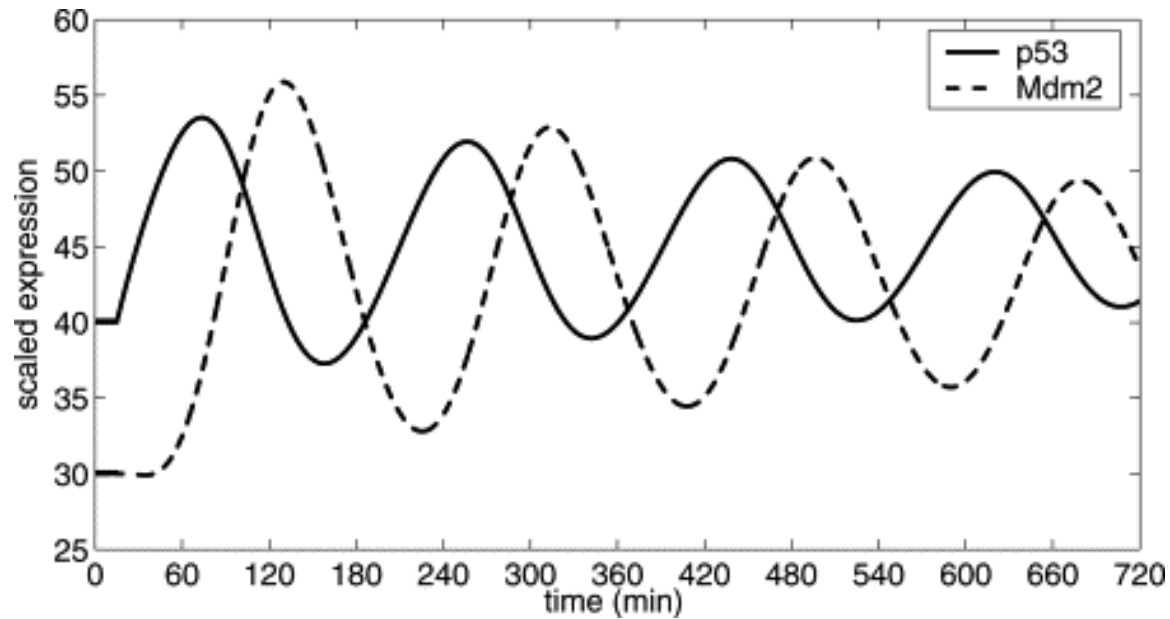
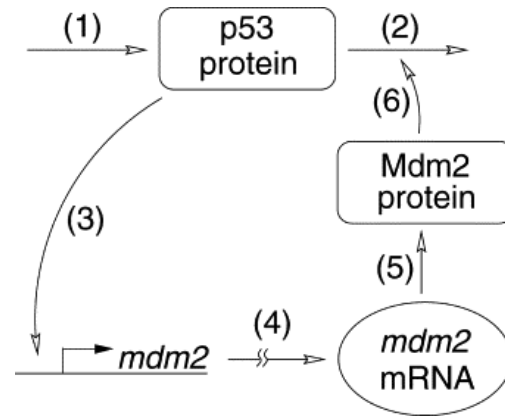


Complexity:

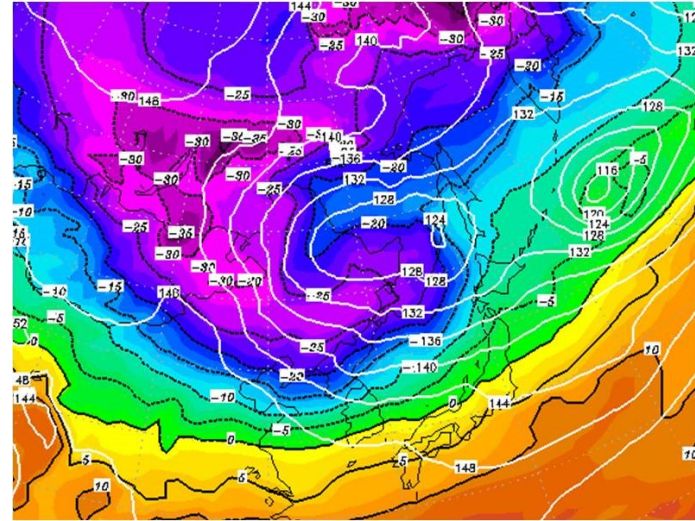
non-linearity of interactions:



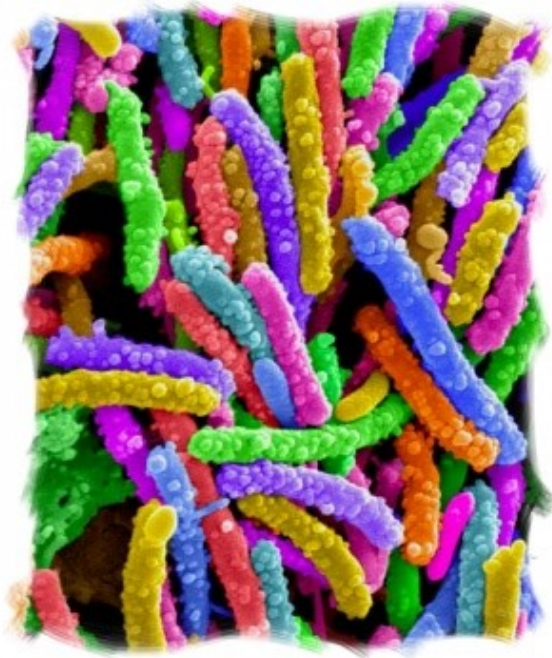
An example...



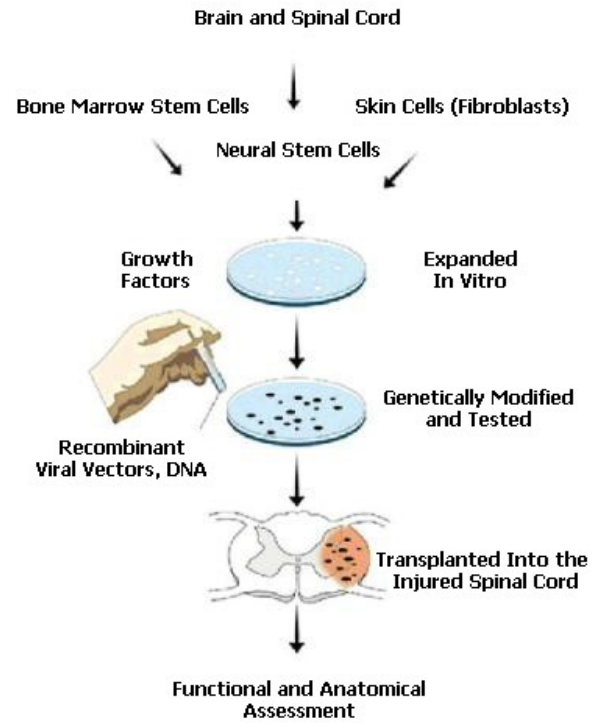
unpredictability



unpredictability



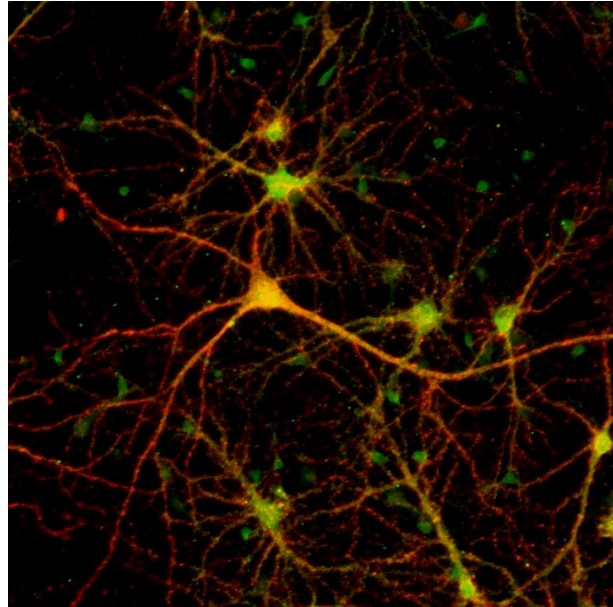
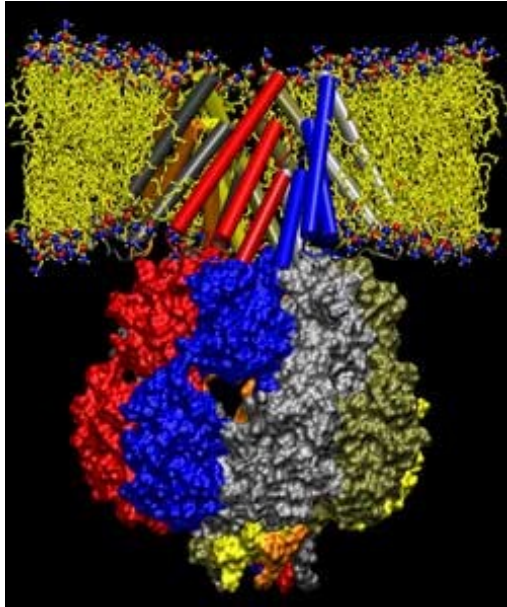
Strategies of spinal cord transplantation and gene therapy



Butterfly effect



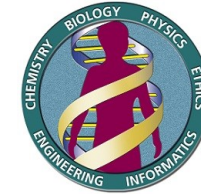
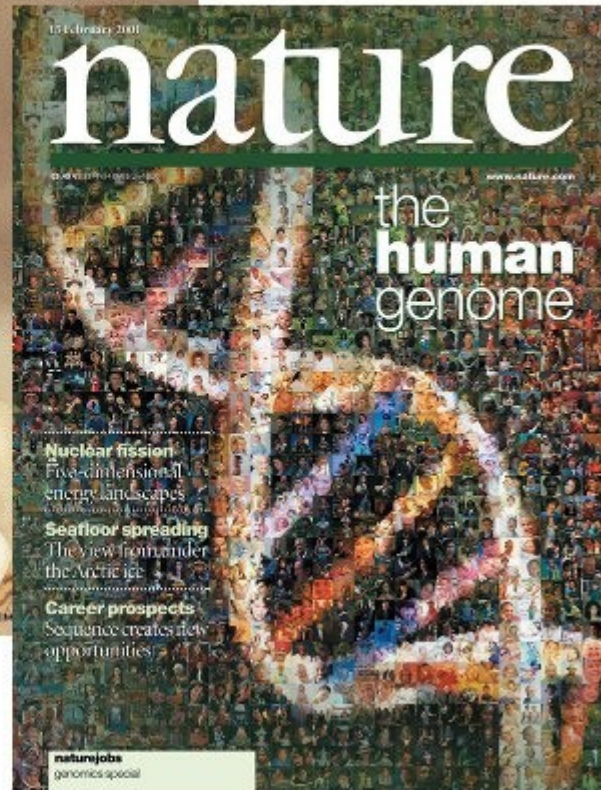
Emergence of properties



The whole is more (different) than the sum of the individual components



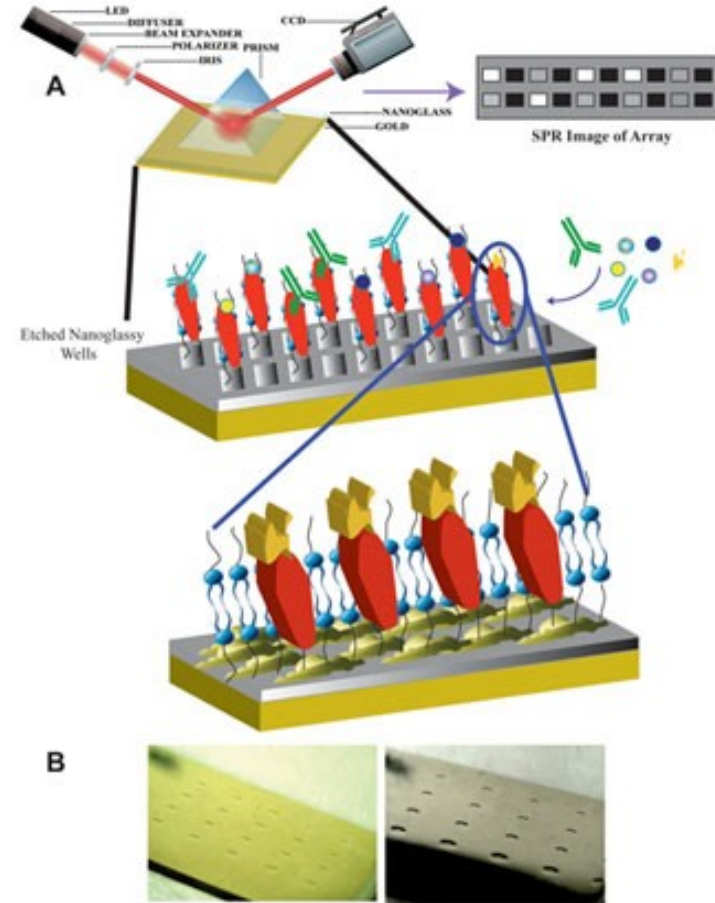
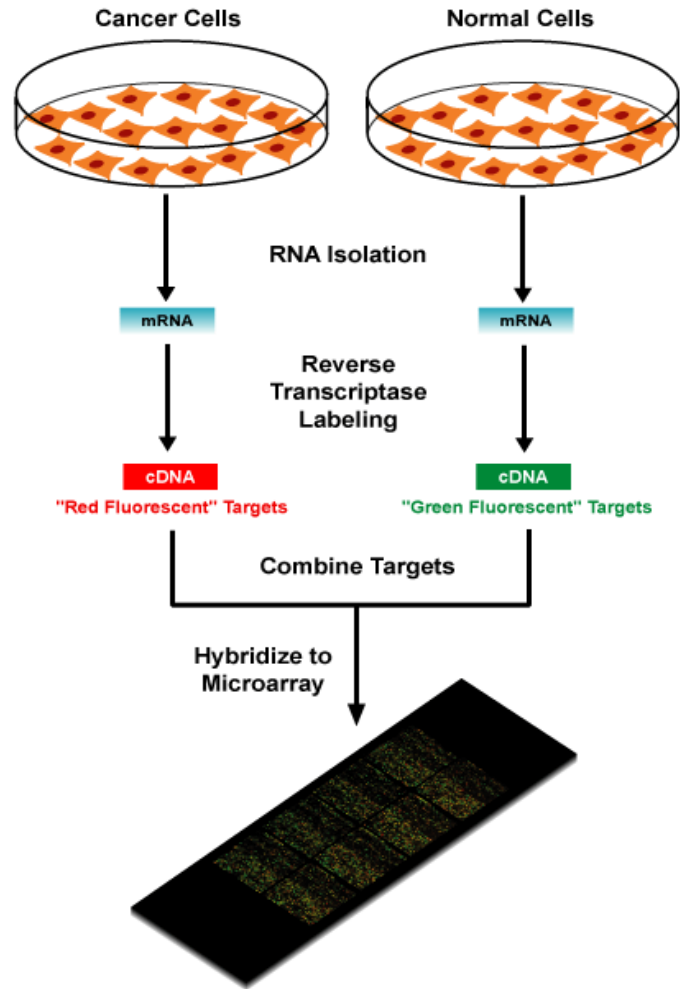
Human Genome Project



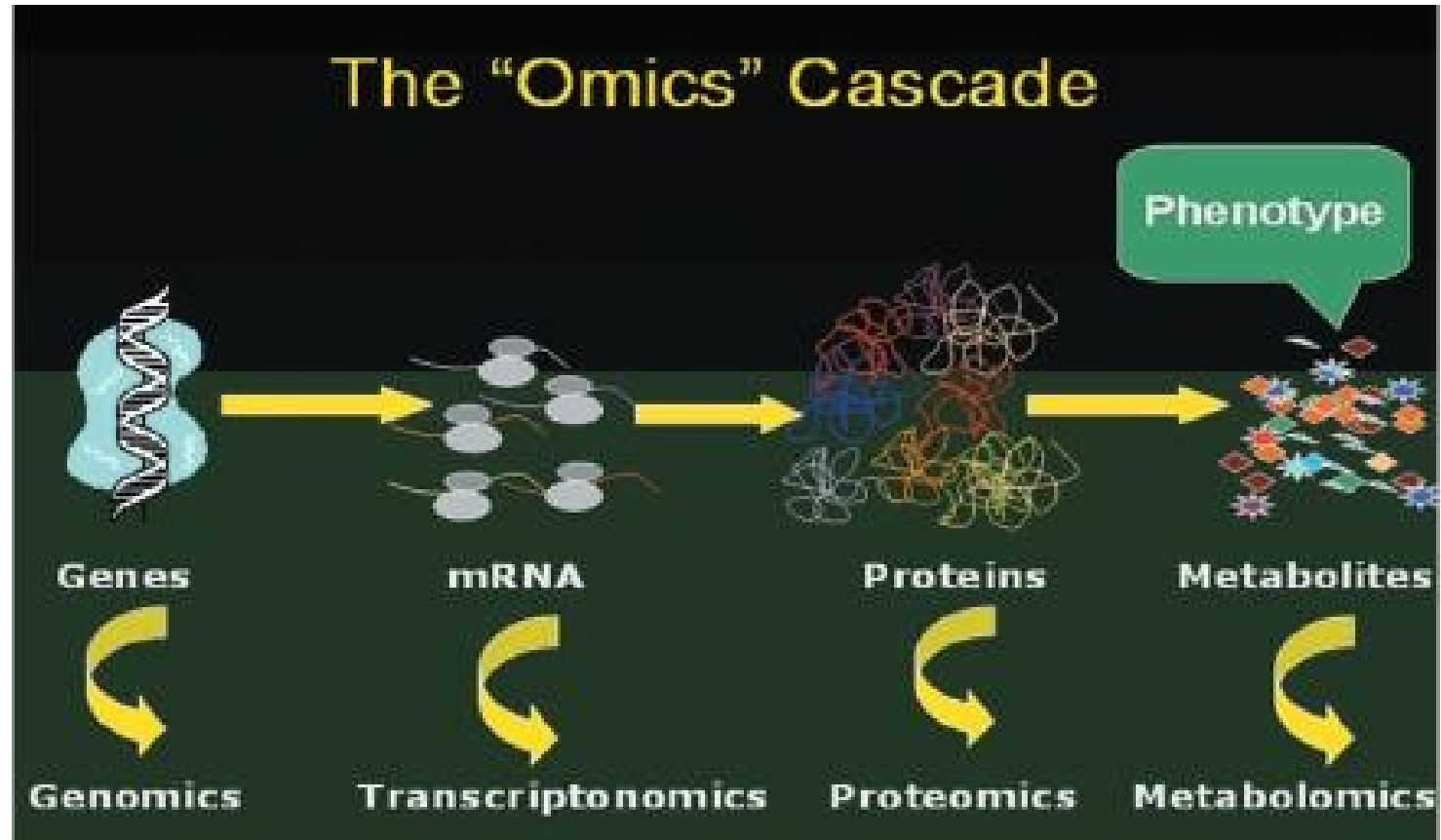
Begun formally in 1990, the U.S. Human Genome Project was a 13-year effort coordinated by the U.S. Department of Energy and the National Institutes of Health. The project originally was planned to last 15 years, but rapid technological advances accelerated the completion date to 2003. Project goals

- *identify* all the approximately 20,000-25,000 genes in human DNA,
- *determine* the sequences of the 3 billion chemical base pairs that make up human DNA,
- *store* this information in databases,
- *improve* tools for data analysis,
- *transfer* related technologies to the private sector, and
- *address* the ethical, legal, and social issues (ELSI) that may arise from the project.

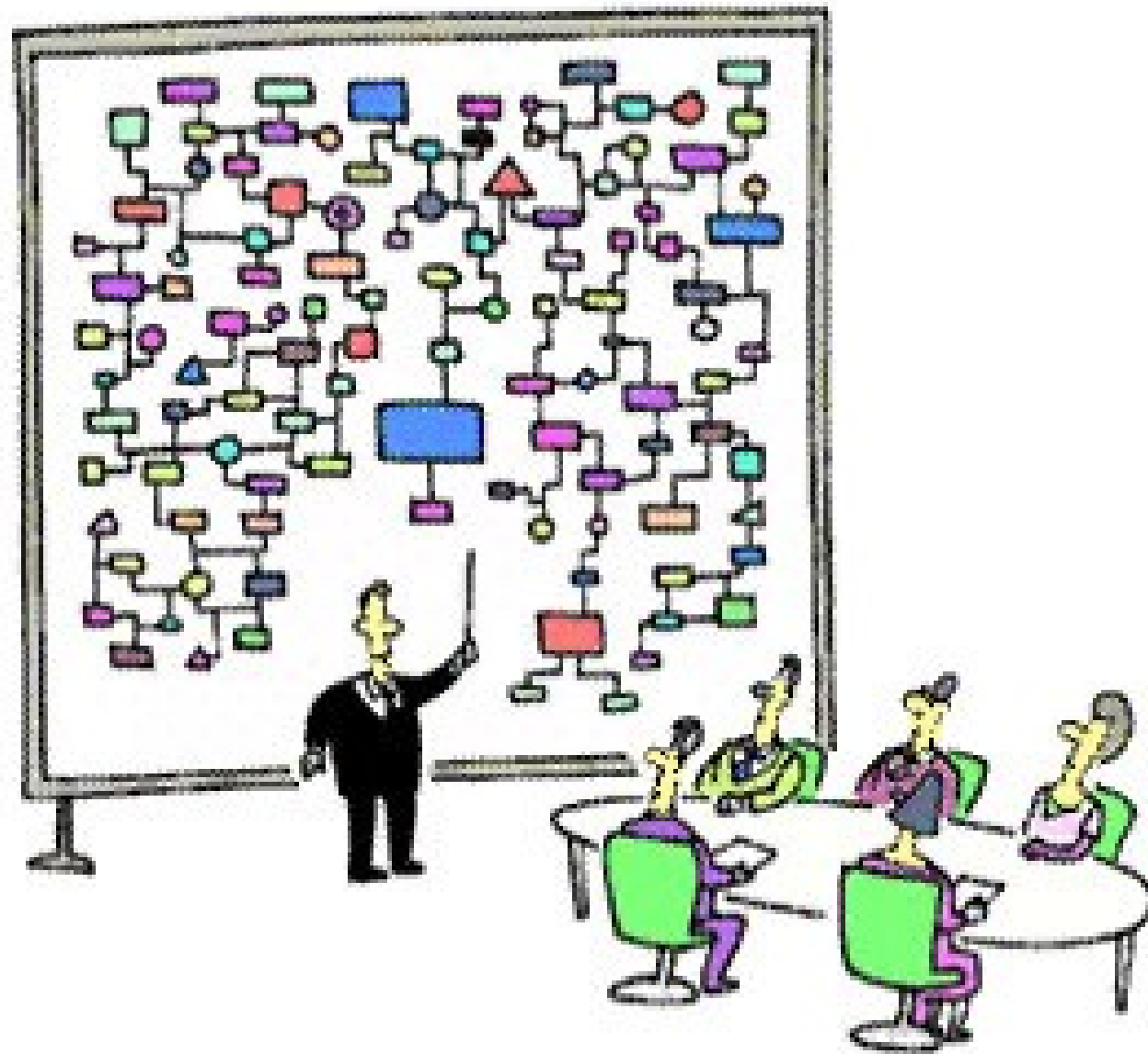
- omics



- OMICS

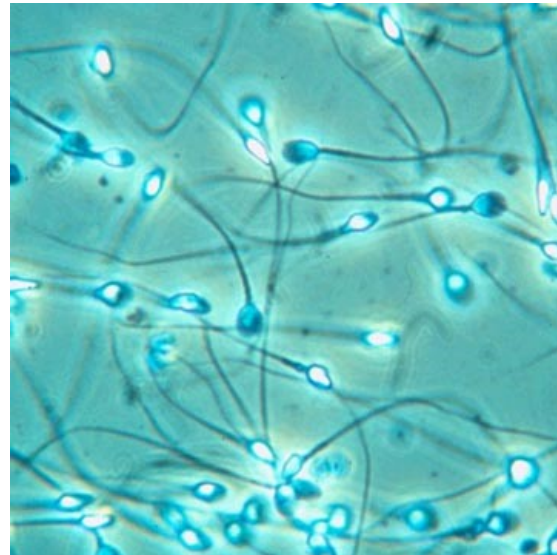
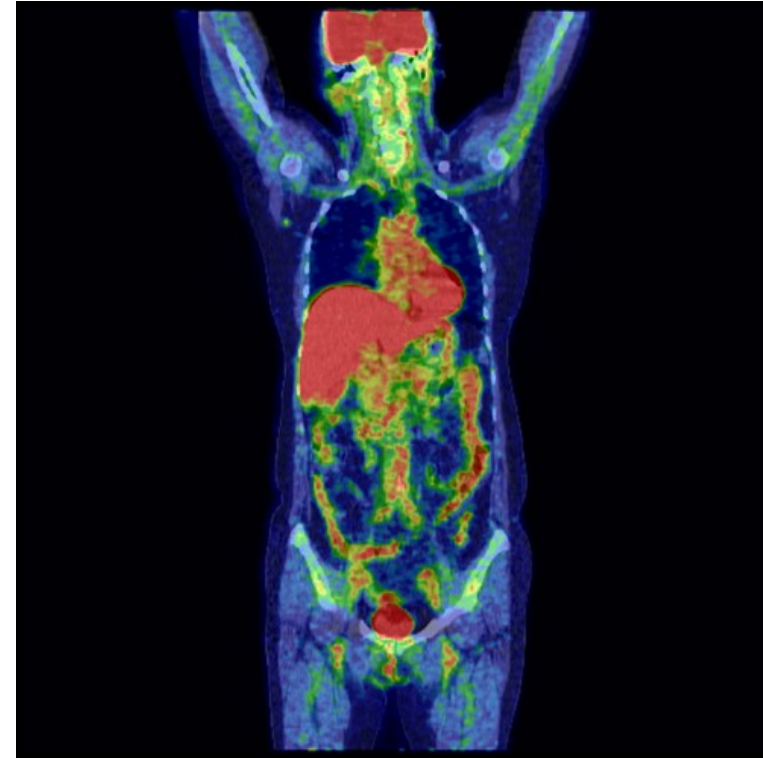
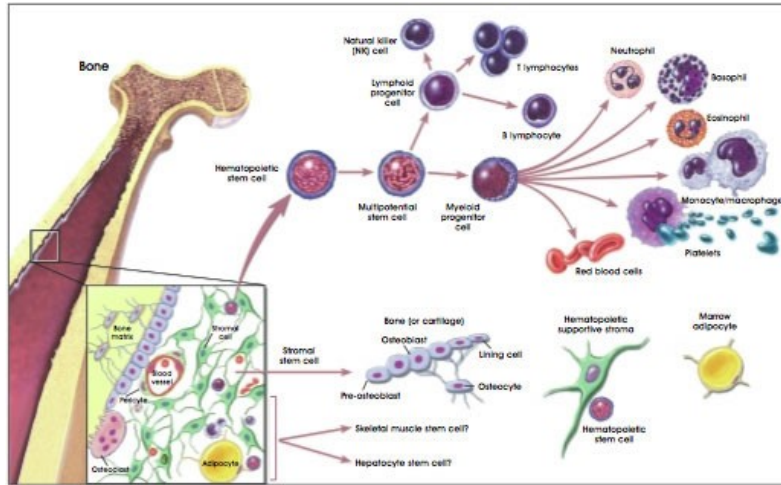






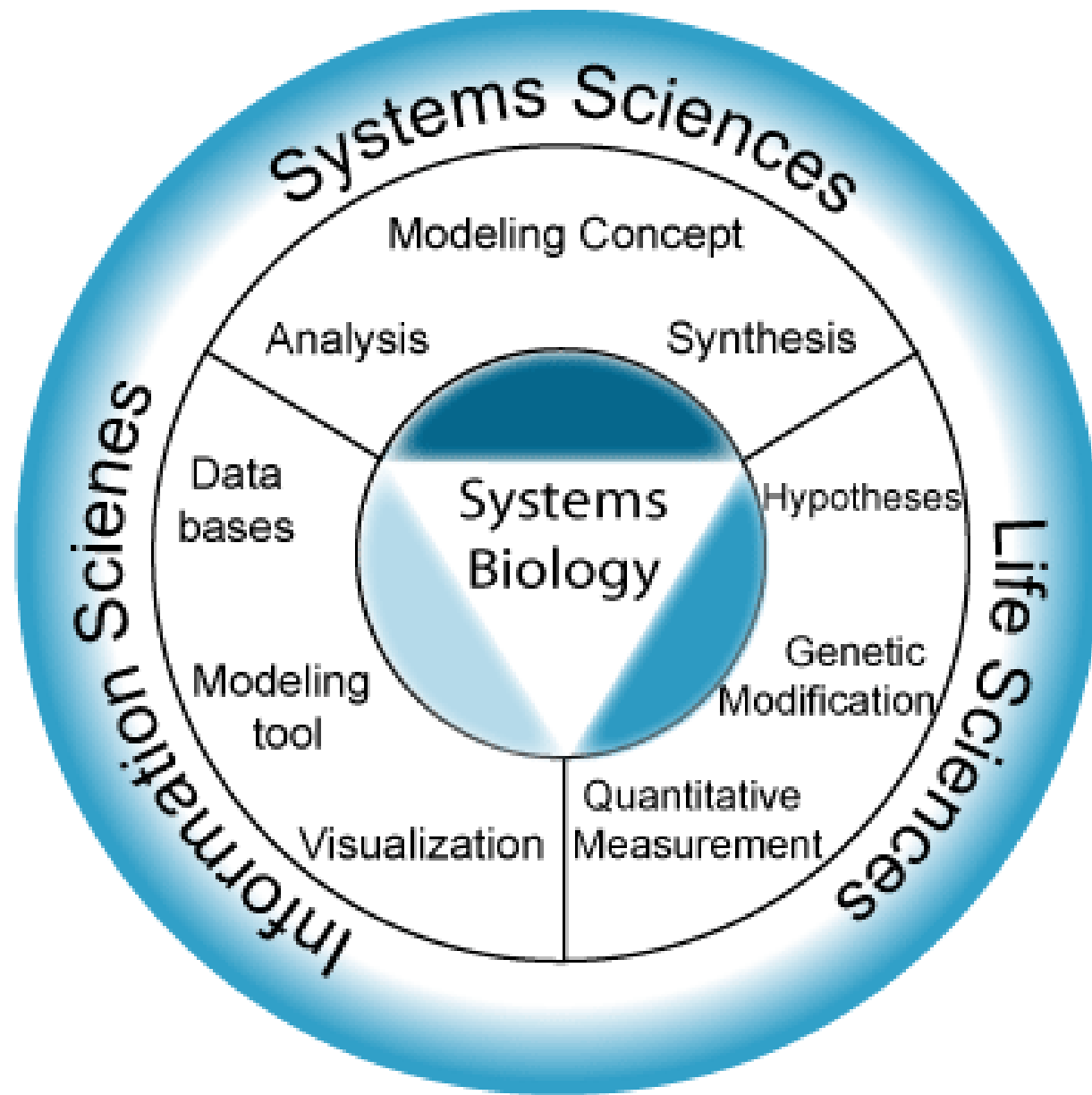
"And that's why we need a computer."

Computational models in biology and medicine



Systems biology

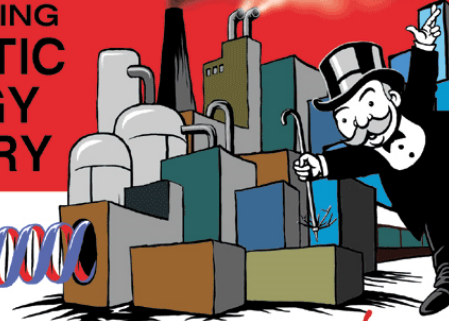




SYNDUSTRY

The news of "Synthia," the world's first human-made species, is just the latest from a rapidly growing artificial life industry. Synthetic biology (or "Syn Bio") aims to profit from the design and construction of industrially useful life-forms.

THE EMERGING SYNTHETIC BIOLOGY INDUSTRY



Syn Bio's Big Shots

Global corporations are investing in synthetic biology labs and partnering with start-up companies.

"Over the next 20 years synthetic genomics is going to become the standard for making anything." - Craig Venter

Cargill
Agribusiness giant. Supports synthetic biology R&D.



BP Energy giant. \$500 million partnership on synthetic biology with University of California Berkeley; holds equity stake in Craig Venter's Synthetic Genomics, Inc.



Du Pont
Chemical giant. Developed first commercial syn bio product with Genencor and sugar giant Tate & Lyle - a fibre called Sorona.



Pfizer
Pharma giant. Conducts in-house syn bio research for drug development.



Virgin Group
Includes Virgin Fuels, investor in synthetic biology. Controlled by celebrity billionaire Richard Branson.

Synthetic Startups

A bevy of 'pure play' syn bio companies is attempting to design synthetic microbes for fuel, chemicals and drugs. Many are university spin-offs.



Gevo
(USA) Developing synthetic biofuels with support from Virgin.



Mascoma
(USA) Developing synthetic biofuels.



Synthetic Genomics
(USA) Constructing synthetic life forms for biofuels and carbon sequestration.



LS9
(USA) Developing synthetic biofuels and industrial chemicals.



Amyris Biotech
(USA) Developing cellular factories to produce drugs, fuels and industrial chemicals.

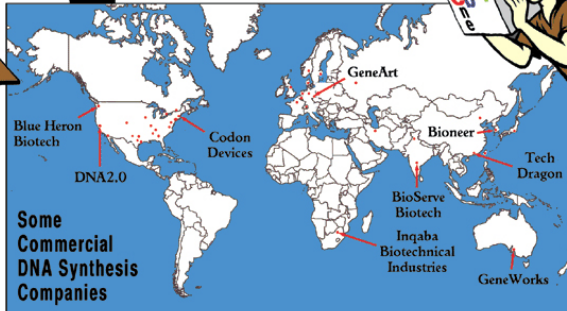


ProtoLife
(Italy) Developing synthetic living systems.

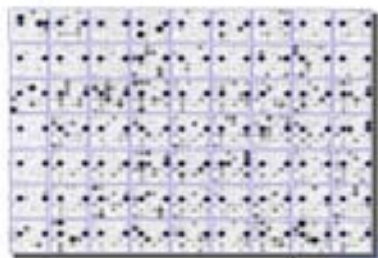
DNA Synthesis Foundries

DNA foundries produce the raw material for creating artificial life: synthetic DNA (sDNA).

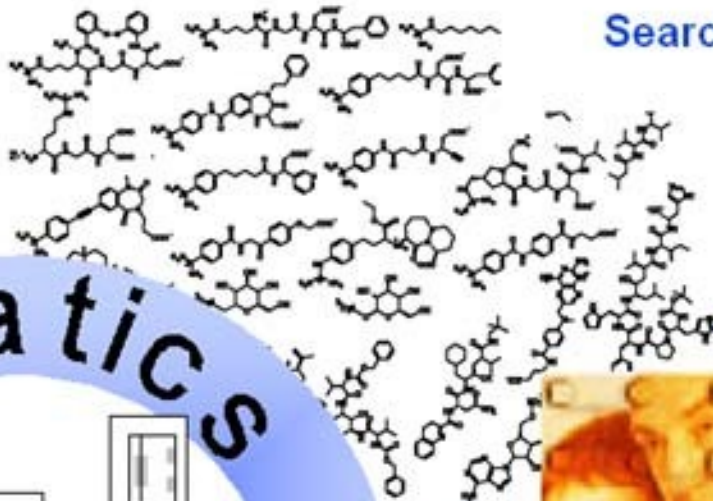
Over 70 DNA foundries worldwide manufacture sDNA for genetic engineers and synthetic biologists. The market for sDNA already exceeds a billion dollars annually. Even long DNA sequences - entire genes, for example - can be ordered over the Internet and delivered within two weeks. The speed of producing accurate DNA sequences is doubling every two years and costs are halving even faster.



DNA chips: comparison of cell states



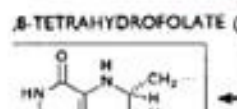
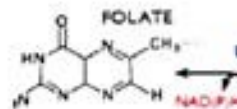
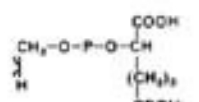
Search for new drugs



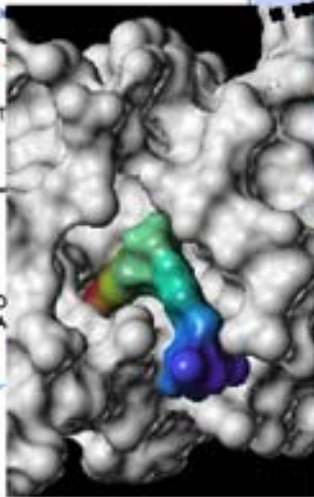
Genetic variations



Biochemical networks



Molecular Interactions



bioinformatics



Data handling, Algorithms
Statistics, Visualisation

Optimizing therapies



Structure prediction

Genomes

```
cactgtggagccaccaccctagggttgcca  
atctactaccaggagcaggga gggcaggag ...
```

Proteins

```
MTNRNFRQINLLDLR VQR RVPVIHQETA  
ECGLACIAMICGHFGKNIDIYLRKFNLS...
```

Sequence analysis

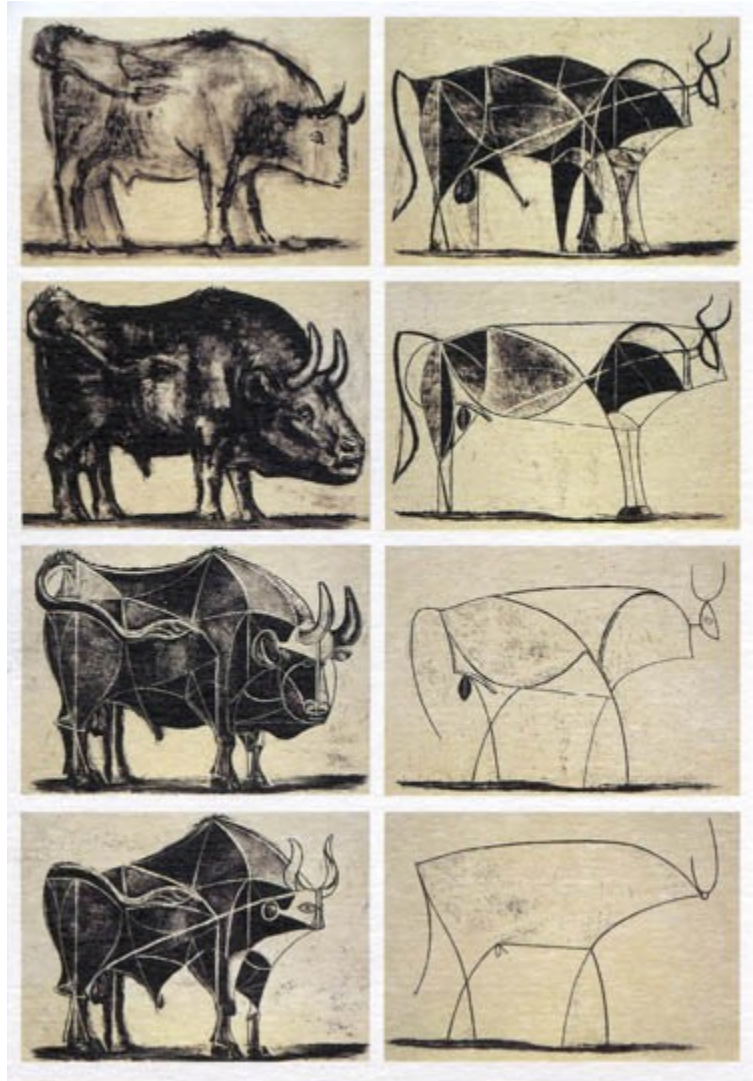
Computational biology



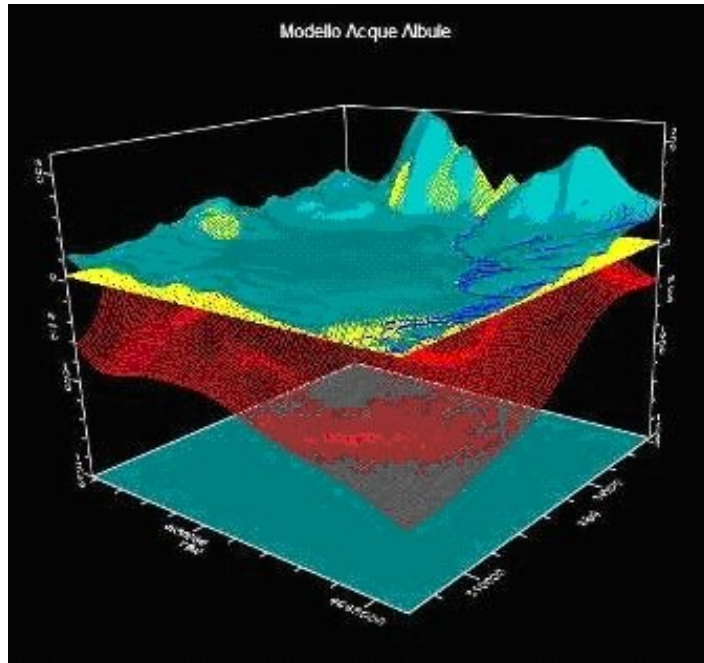
What is a model?



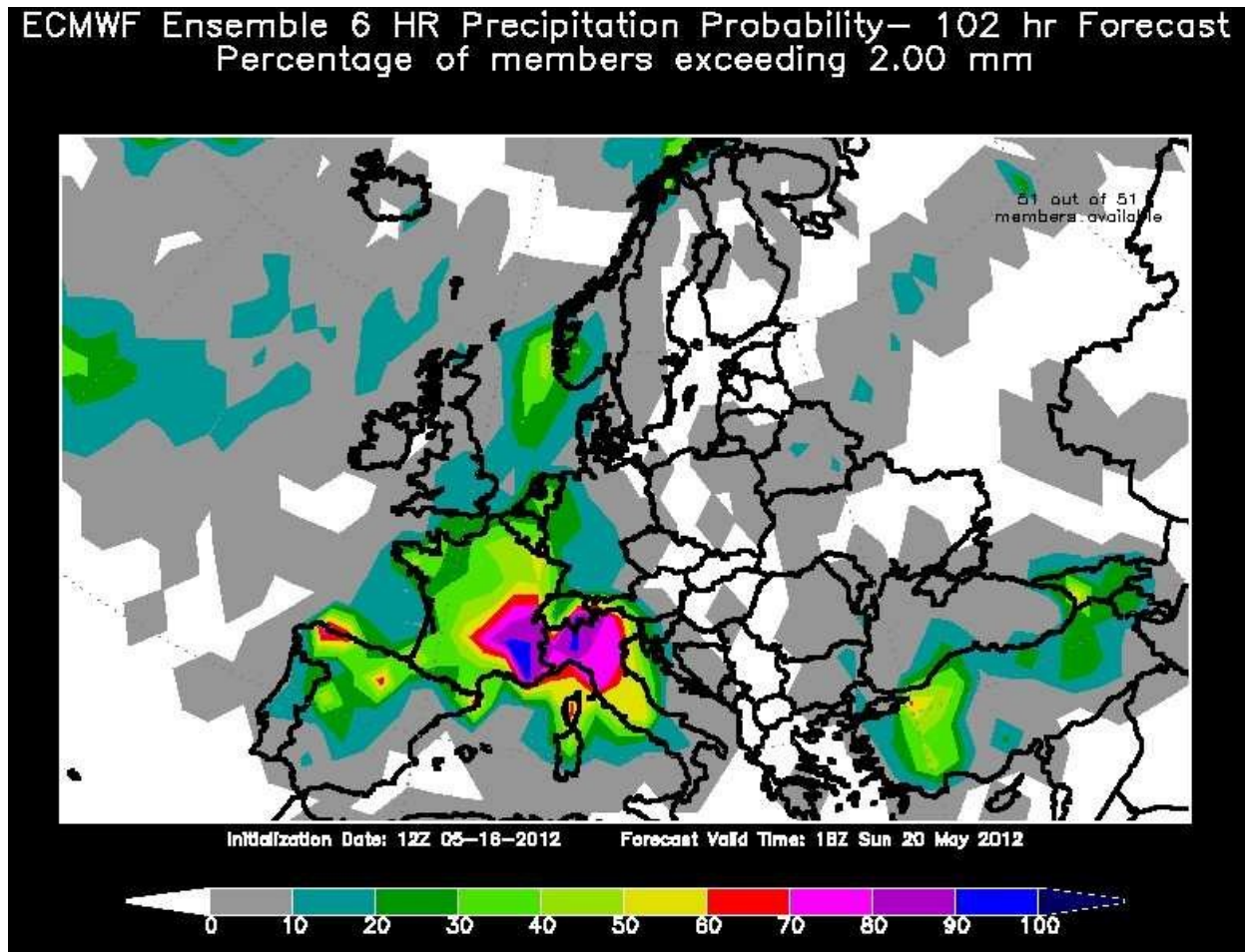
What kind of model we
need?



Numerical models

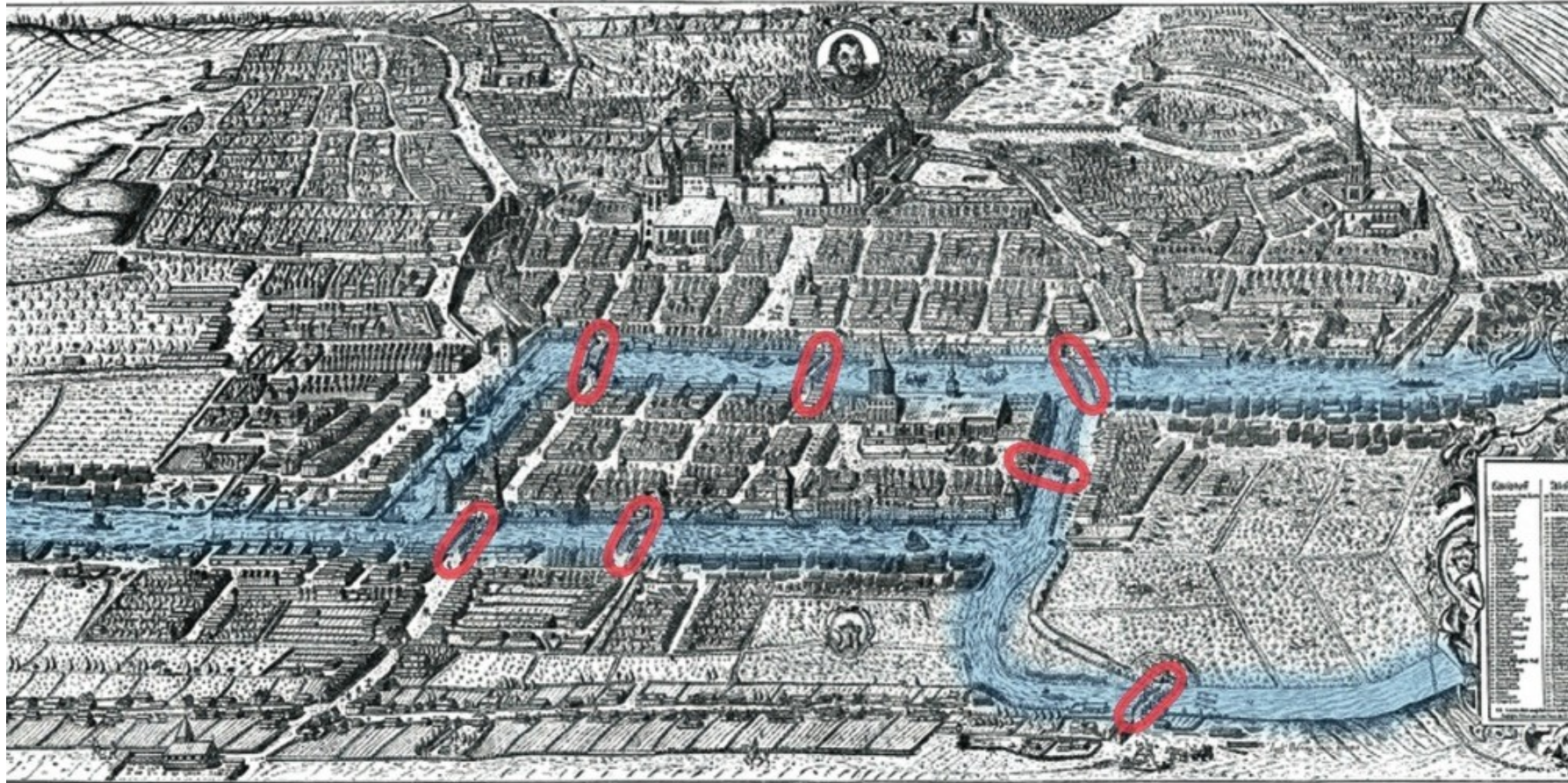


forecasting model



Networks as model

Gedenkblatt zur sechshundert jährigen Jubelfeier der Königlichen Haupt und Residenz-Stadt Königsberg in Preußen.



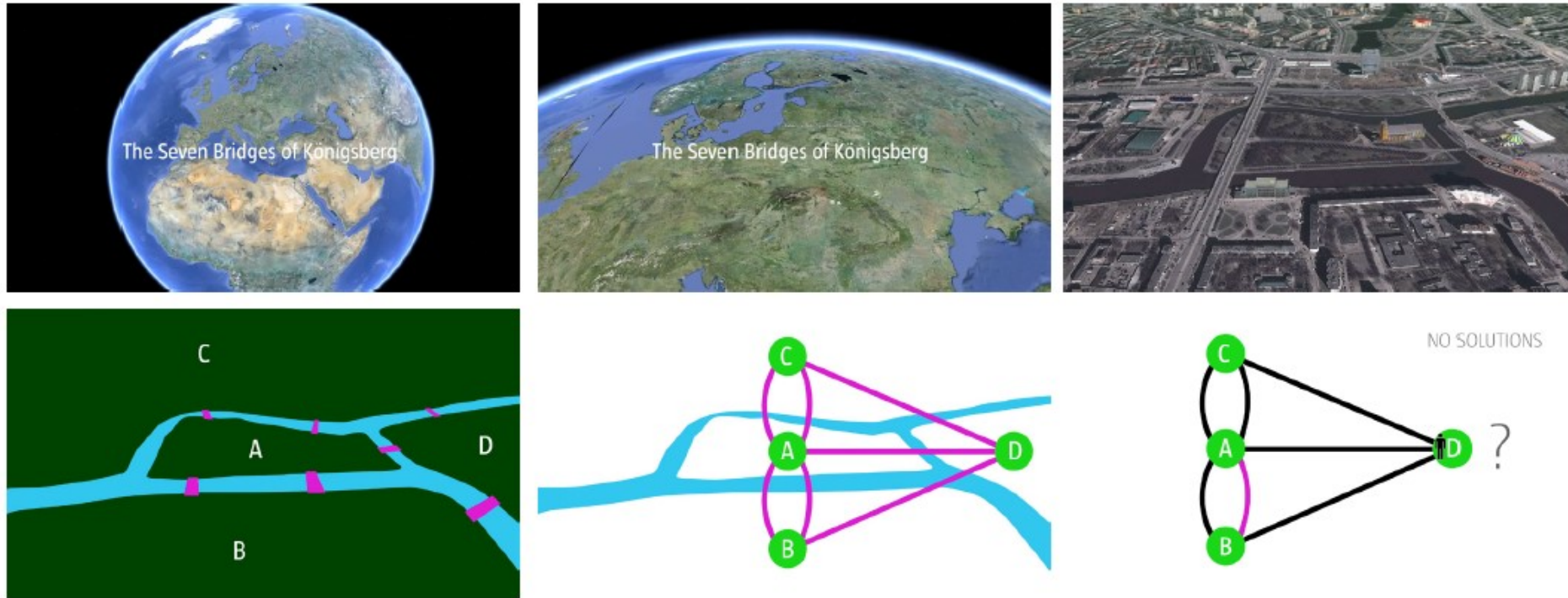


Image 2.1
The bridges of Königsberg.

From the contemporary map of Königsberg (now Kaliningrad, Russia) to Euler's graph. The graph constructed by Euler consists of four nodes (A, B, C, D), each corresponding to a patch of land, and seven links, each corresponding to a bridge. Euler showed in 1736 that there is no continuous path that would cross seven the bridges while never crossing the same bridge twice. The people of Königsberg agreed with him, gave up their fruitless search and in 1875 they built a new bridge between B and C, increasing the number of links of these two nodes to four. Now only one node was left with an odd number of links and it became rather straightforward to find the desired path.

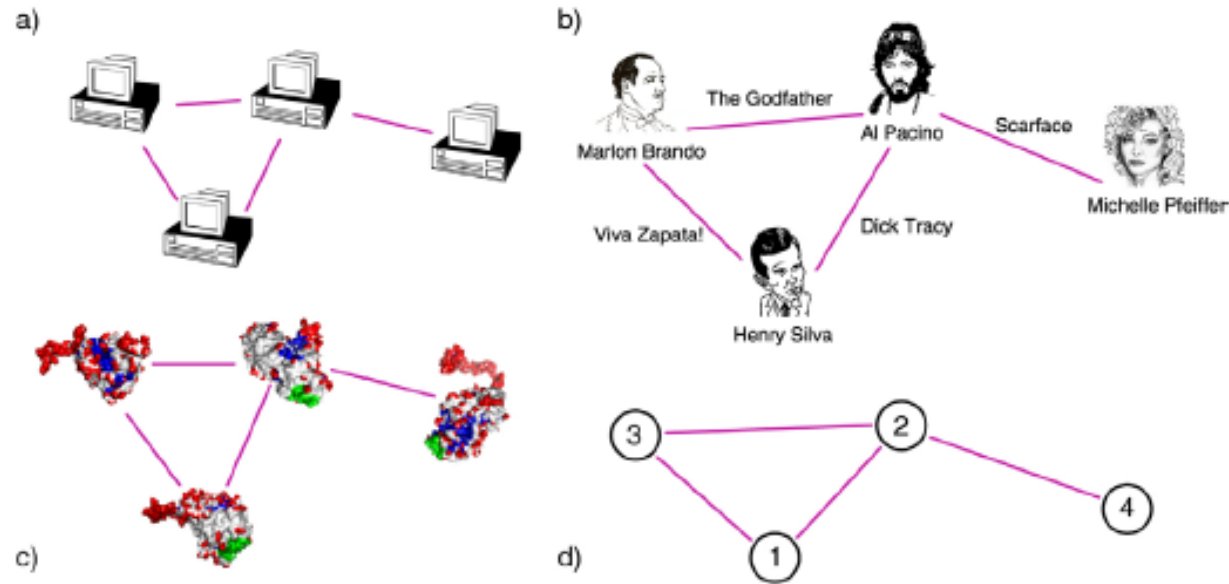
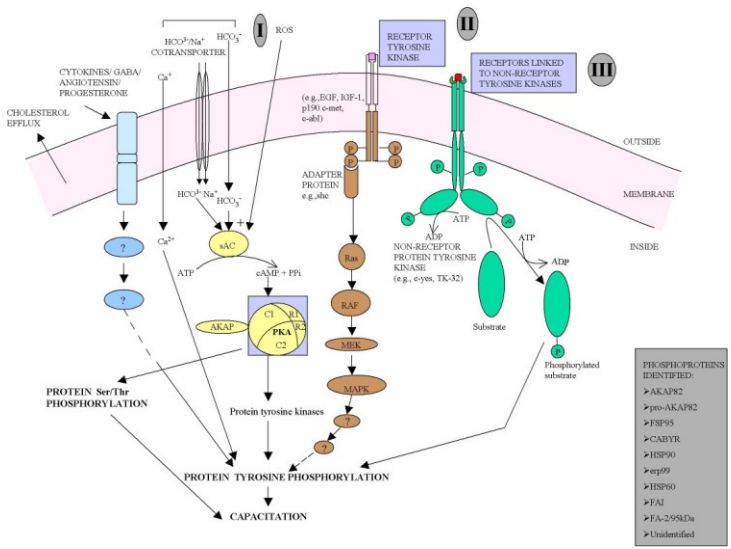


Image 2.3

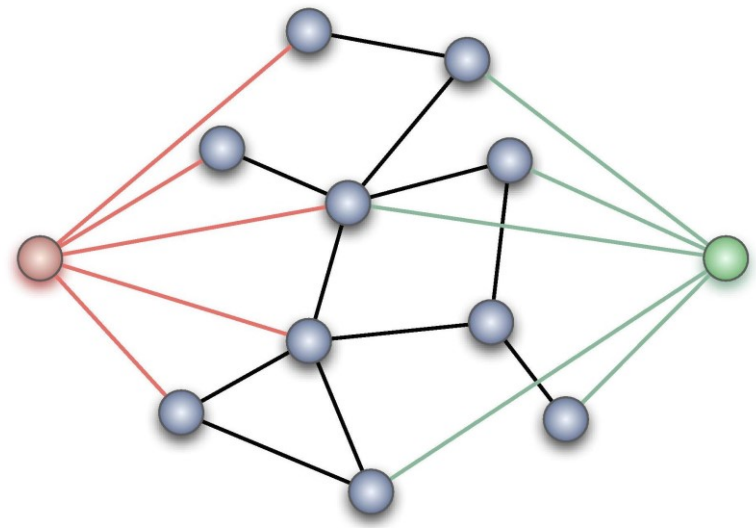
Real systems of quite different nature can have the same network representation.

In the figure we show a small subset of (a) the *Internet*, where routers (specialized computers) are connected to each other; (b) the *Hollywood actor network*, where two actors are connected if they played in the same movie; (c) a *protein-protein interaction network*, where two proteins are connected if there is experimental evidence that they can bind to each other in the cell. While the nature of the nodes and the links differs widely, each network has the same graph representation, consisting of $N = 4$ nodes and $L = 4$ links, shown in (d).

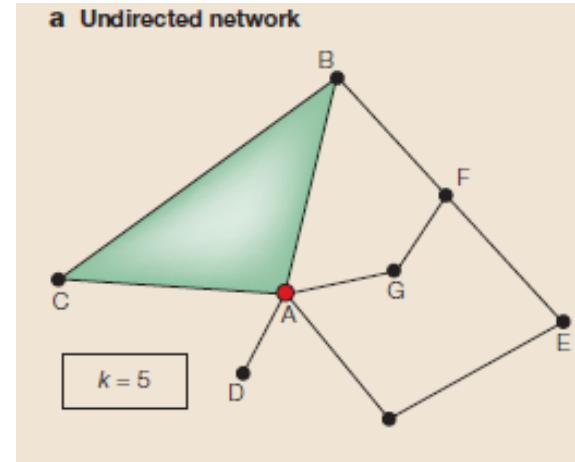
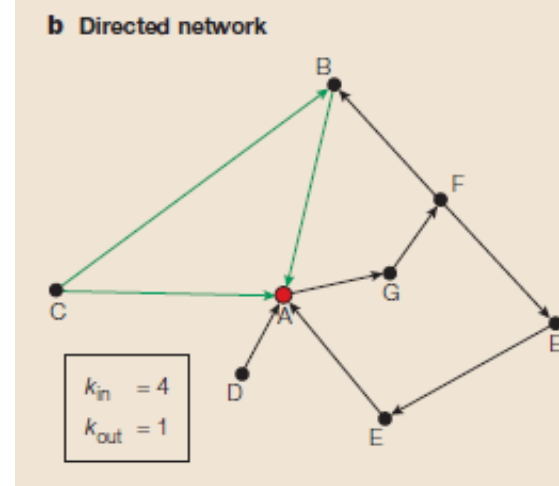
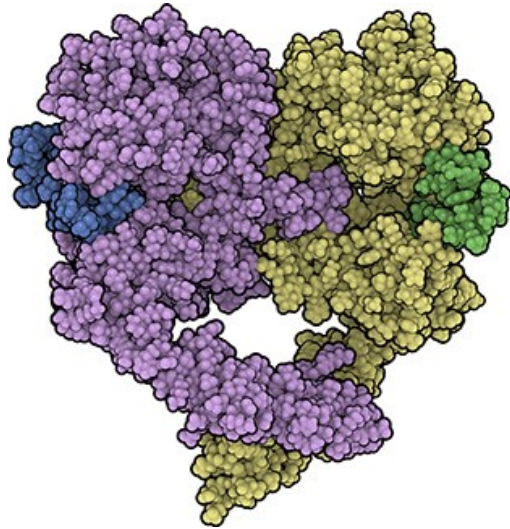
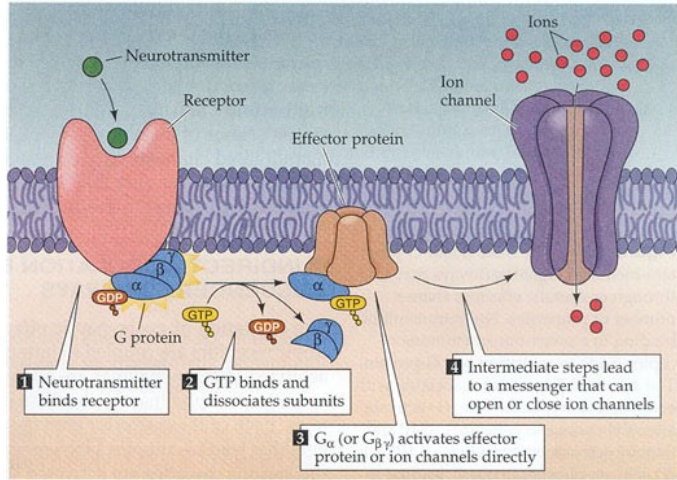
The node



- PHOSPHOPROTEINS IDENTIFIED:
- >AKAP82
 - >pp60-AKAP82
 - >FAP95
 - >CABYR
 - >HSP90
 - >pp59
 - >HSP60
 - >FAI
 - >FA-259kDa
 - >Unidentified



The link

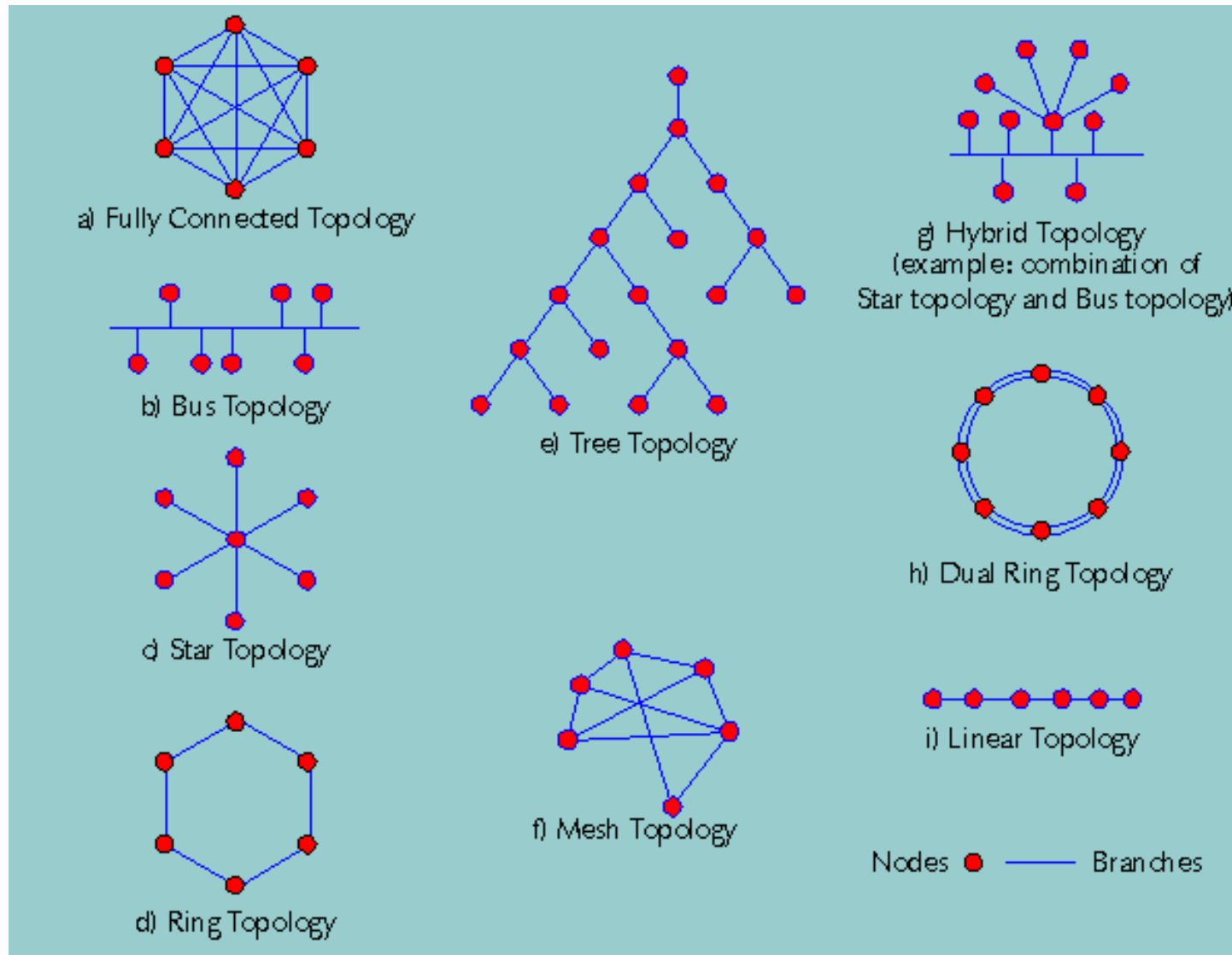


NETWORK NAME	NODES	LINKS	DIRECTED/ UNDIRECTED	N	L	⟨K⟩
Internet	routers	Internet Connections	Undirected	192,244	609,066	2.67
WWW	webpages	links	Directed	325,729	1,497,134	4.60
Power Grid	power plants, transformers	cables	Undirected	4,941	6,594	2.67
Mobile-Phone Calls	subscribers	calls	Directed	36,595	91,826	2.51
Email	email addresses	emails	Directed	57,194	103,731	1.81
Science Collaboration	scientists	co-authorships	Undirected	23,133	186,936	16.16
Actor Network	actors	co-acting	Undirected	212,250	3,054,278	28.78
Citation Network	papers	citations	Directed	449,673	4,707,958	10.47
E. coli Metabolism	metabolites	chemical reactions	Directed	1,039	5,802	5.84
Yeast Protein Interactions	proteins	binding interactions	Undirected	2,018	2,930	2.90

Table 2.1

Network maps and their basic properties.

Network topology



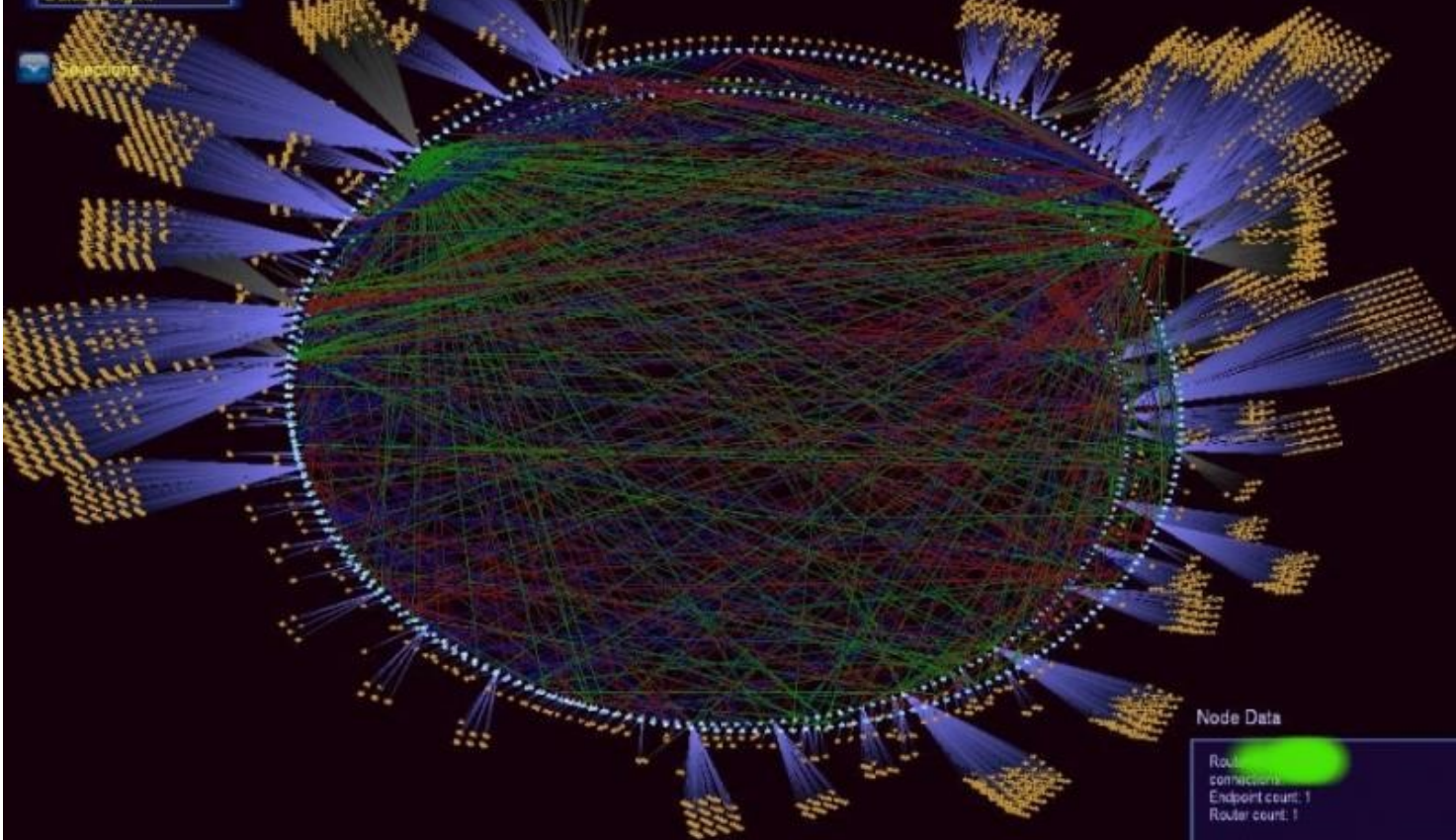
Net Topology Diagram

Main Menu

Vulns/Ps Graph

Ports/Ps Graph

Dataset Mgmt



Node Data

Router
connections
Endpoint count: 1
Router count: 1



Topological parameters

the number of nodes: which represents the total number of molecules involved;

In an undirected network total number of links, L , can be expressed as the sum of the node degrees:

$$L = \frac{1}{2} \sum_{i=1}^N k_i \quad (1)$$

Here the $1/2$ factor corrects for the fact that in the sum (1) each link is counted twice.

the number of edges: which represents the total number of interaction among nodes within the network;

the node degree (or connectivity): which indicates how many links each node has to other nodes;

$$\langle k \rangle \equiv \frac{1}{N} \sum_{i=1}^N k_i = \frac{2L}{N} \quad (7)$$

In directed networks we distinguish between *incoming degree*, k_i^{in} , representing the number of links that point to node i , and *outgoing degree*, k_i^{out} , representing the number of links that point from the node i to other nodes and the *total degree*, k_i , given by

$$k_i = k_i^{in} + k_i^{out} \quad (8)$$

the node degree distribution $P(k)$: which represents the probability that a selected node has exactly k links;

$$p_k = \frac{N_k}{N}$$

the clustering coefficient: it is a measure of how the nodes tend to form clusters: the more the clustering coefficient is higher, the more the presence of clusters will increase;

$$C_I = 2n_I / k(k-1),$$

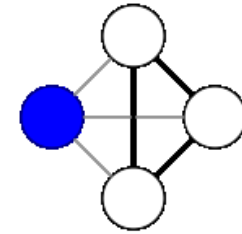
where n_I is the number of links connecting the k_I neighbours of node I to each other

clustering coefficient

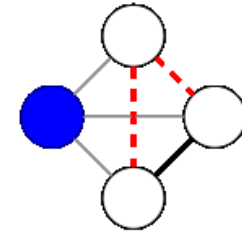
$$c = \frac{3 \times \text{number of triangles}}{\text{number of connected triplets of vertices}}$$

$$= \frac{\text{number of closed triplets}}{\text{number of connected triplets of vertices}}$$

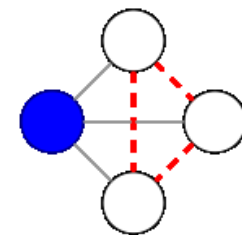
$$\bar{c} = \frac{1}{n} \sum_{i=1}^n C_i.$$



$$c = 1$$



$$c = 1/3$$

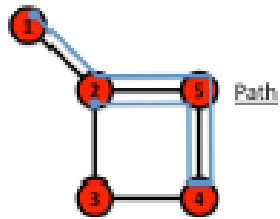


$$c = 0$$

the network diameter: which is the largest distance between two nodes;

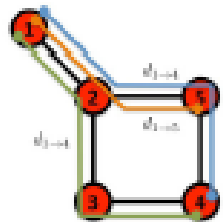
the averaged number of neighbours: which is the mean number of connection of nodes;

the characteristic path length: which is the expected distance between two random individuated connected nodes.



Path

PATH: A sequence of nodes such that each node is connected to the next node along the path by a link. A path always consists of n nodes and $n - 1$ links. The length of a path is defined as the number of its links, counting multiple edges multiple times.



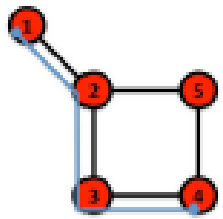
Shortest Path

$$d_{1 \rightarrow 2} = 1$$

$$d_{1 \rightarrow 5} = 2$$

$$d_{2 \rightarrow 5} = 1$$

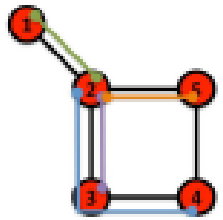
SHORTEST PATH (geodesic path, d): the path with the shortest distance d between two nodes. We will call it the distance between two nodes.



Diameter

$$d_{1 \rightarrow 5} = 3$$

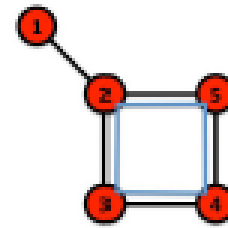
DIAMETER (d_{max}): the longest shortest path in a graph, or the distance between the two furthest away nodes.



Average Path Length

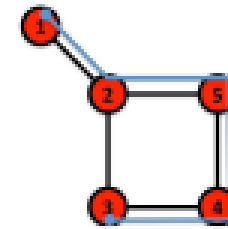
$$(d_{1 \rightarrow 2} + d_{1 \rightarrow 3} + d_{1 \rightarrow 4} + d_{1 \rightarrow 5} + d_{2 \rightarrow 3} + d_{2 \rightarrow 4} + d_{2 \rightarrow 5} + d_{3 \rightarrow 4} + d_{3 \rightarrow 5} + d_{4 \rightarrow 5}) / 10 = 1.6$$

AVERAGE PATH LENGTH ($\langle d \rangle$): the average of the shortest paths between all pairs of nodes.



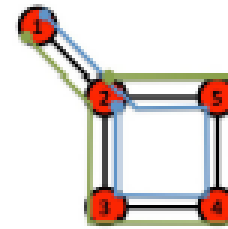
Cycle

CYCLE: a path with the same start and end node.



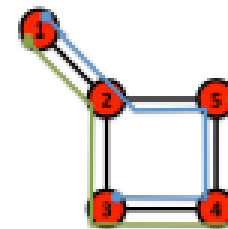
Self-avoiding Path

SELF-AVOIDING PATH: a path that does not intersect itself, i.e. the same node or link does not occur twice along the path.



Eulerian Path

EULERIAN PATH: a path that traverses each link exactly once.



Hamiltonian Path

HAMILTONIAN PATH: a path that visits each node exactly once.

multiple networks

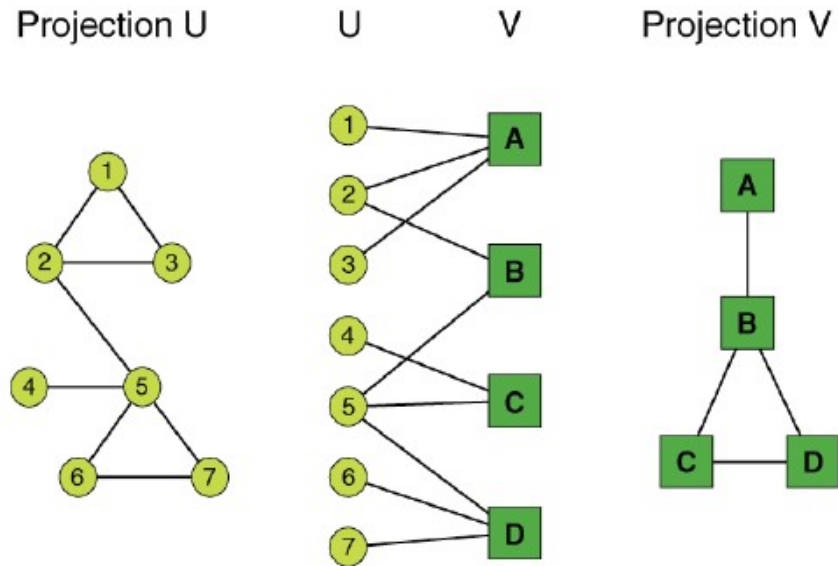


Image 2.9a
Bipartite network.

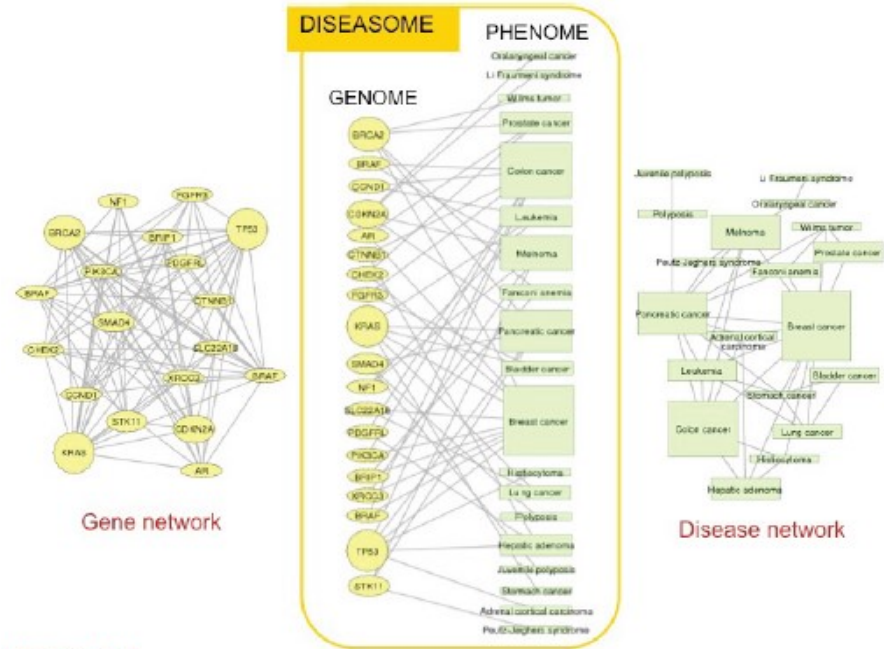


Image 2.9b
Bipartite network.

The tripartite recipe-ingredient network, in which one set of nodes are recipes, like Chicken Marsala, the second set corresponds to the ingredients each recipe has (like flour, sage, chicken, wine, and butter for Chicken Marsala), and the third set captures the flavor compounds, or chemicals that contribute to the taste of a particular ingredient.

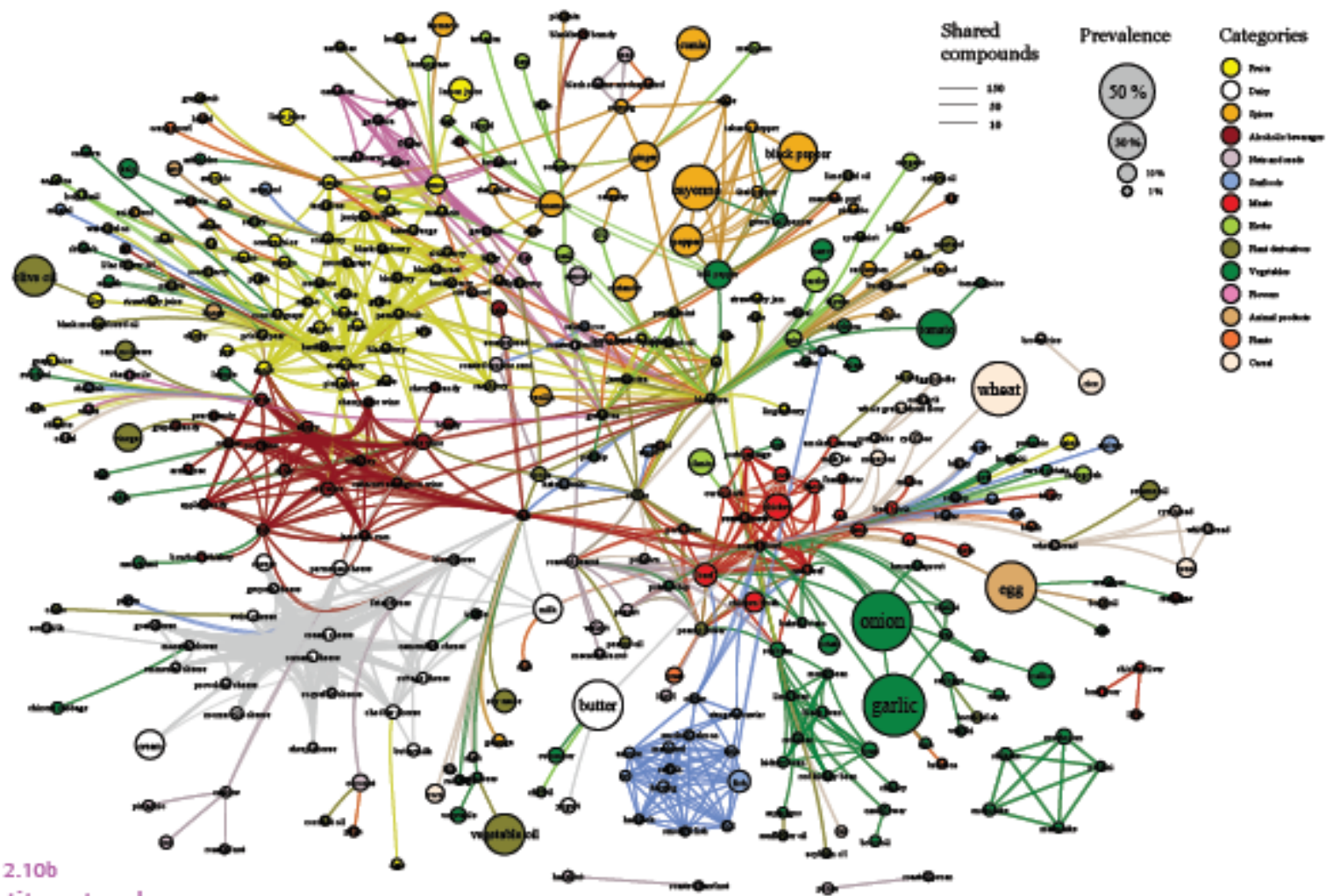


Image 2.10b
Tripartite network.

A projection of the tripartite network, resulting in the ingredient network, often called the flavor network. Each node denotes an ingredient; the node color indicating the food category and node size reflects the ingredient prevalence in recipes. Two ingredients are connected if they share a significant number of flavor compounds, link thickness representing the number of shared compounds between the two ingredients (After [12]).