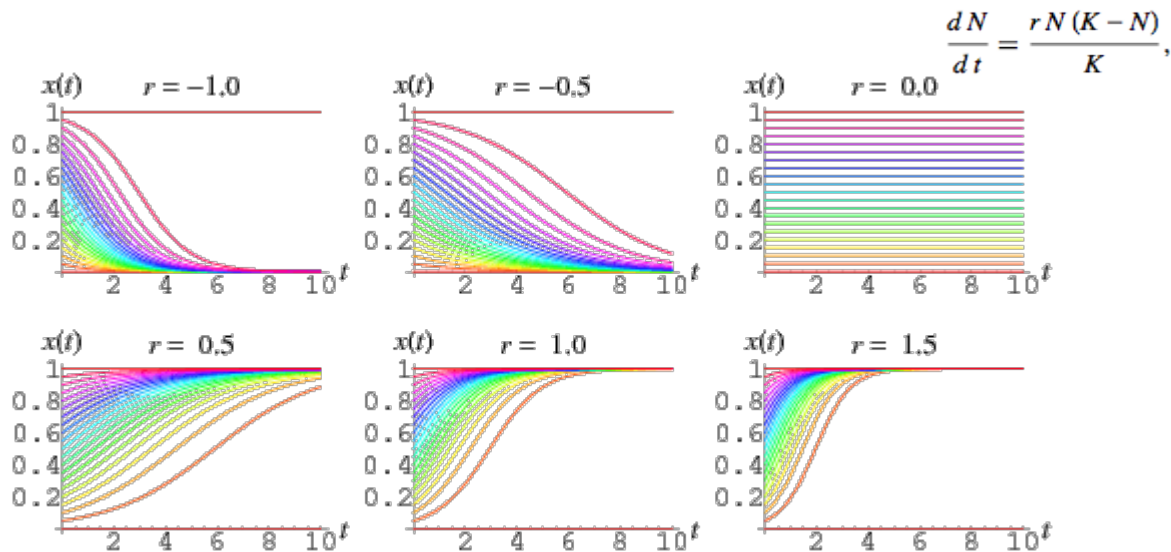# Informatics, algorithmics, and automation: machine learning in 21st century biology
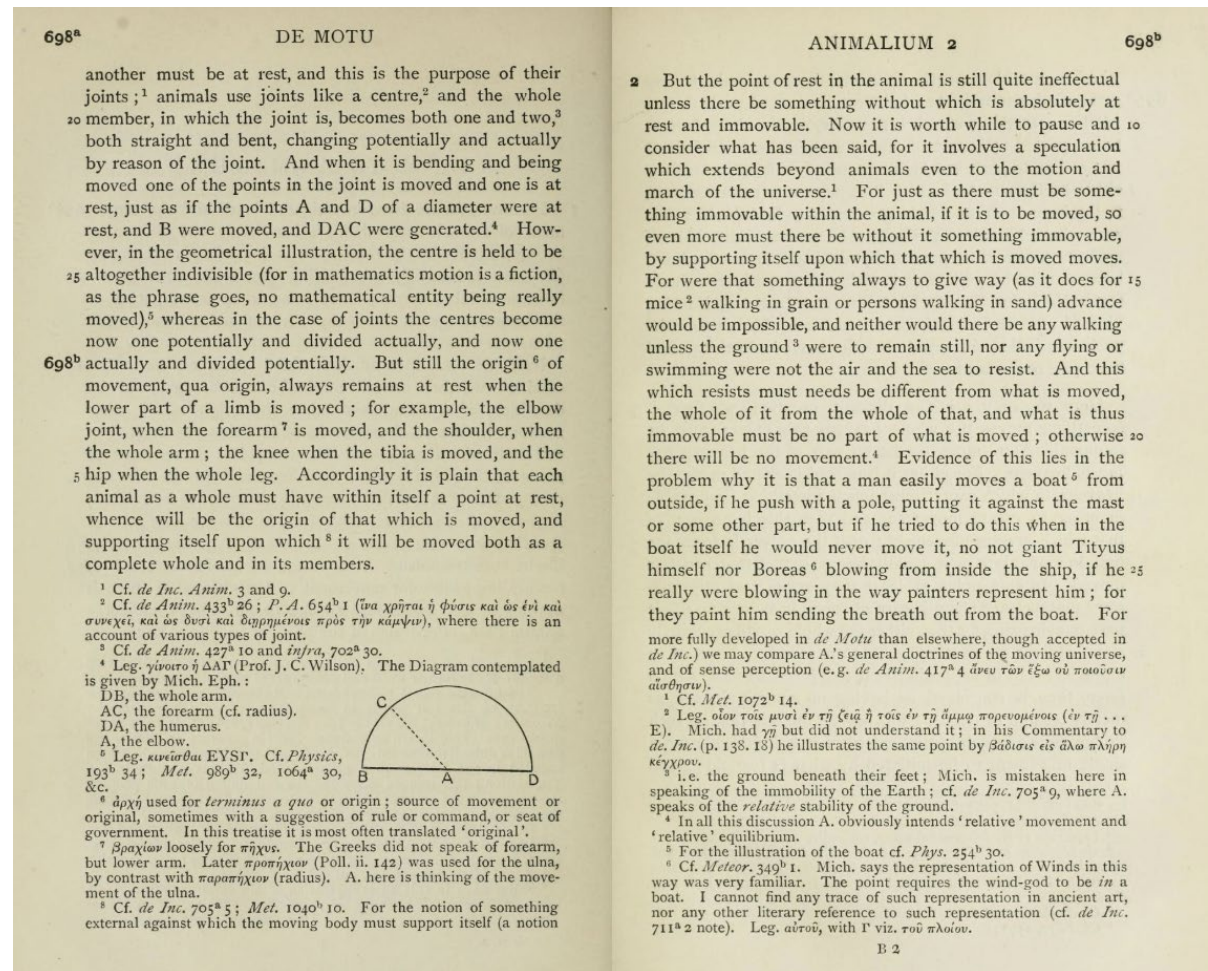
Pradipta Ray

UT Dallas

# Mathematics in biology & medicine

- A very long history – as early as the Hellenic civilization

- All branches of mathematics are harnessed – from classical geometry to calculus

$$\frac{dN}{dt} = \frac{rN(K-N)}{K},$$



Plots of the logistic equation modelling population growth



On the motion of animals, Aristotle, 4th century BC (translated 1912)
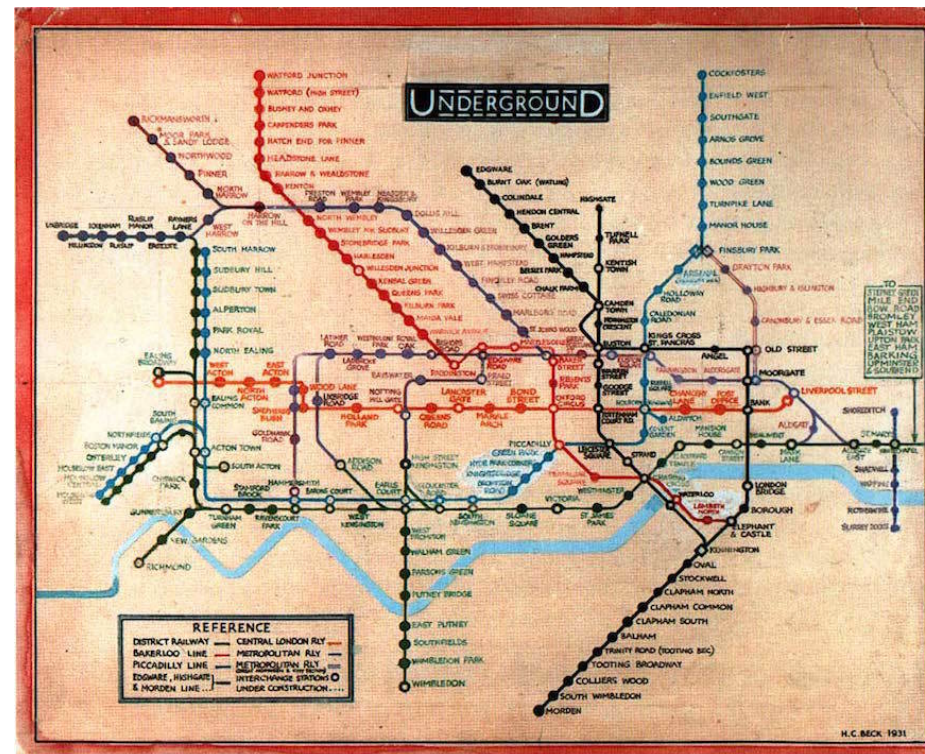
Wolfram

# Informatics : what aspects of mathematics does it deal with ?

- **Mathematics pertaining to data**
  - Automated description, modelling, visualization, prediction (and imputation), etc
- Large overlaps with Computer Science (machine learning) and Statistics community



Term first coined in 1957 though discipline dates back to ancient India / China / Greece

Harry Beck's 1931 visualization of London's tube network

Beck was a draftsman for the underground railway before being fired and voluntarily creating map & selling to ex-employer for 10 GBP !
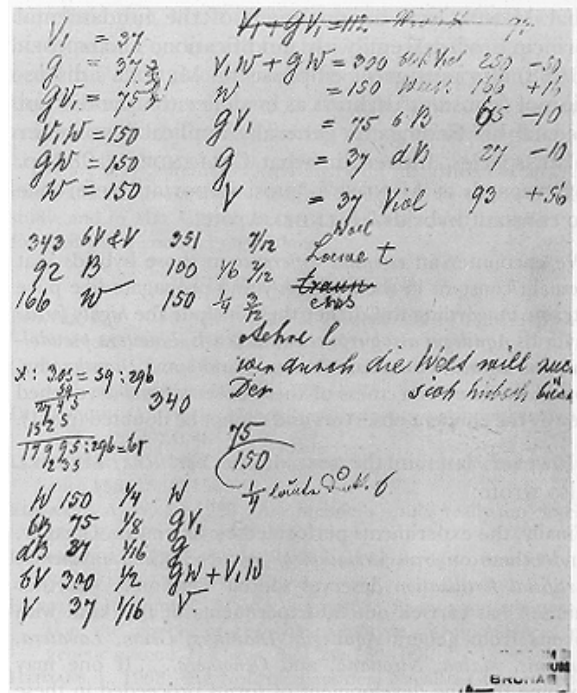
# Algorithms in biology

- Algorithm : a set of instructions for performing a task / computation
  - Set of rules to navigate a car from school to home
  - Computer program for sorting number

**Computational Biology**
Use mathematics & CS to answer questions in biology – which traits of a pea plant are inherited independently of others and at what frequency ?
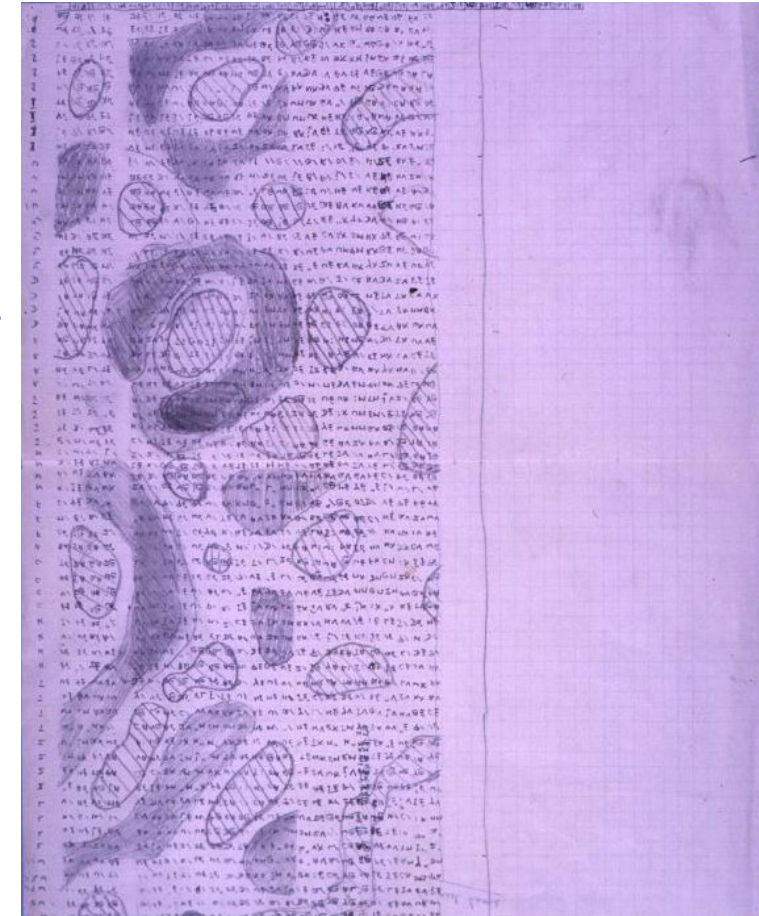


Gregor Mendel's notebook detailing his experiments with breeding peas

**Computation in biology**
What are the underlying instructions driving biological processes – how do molecular gradients give rise to morphological patterns ?



Alan Turing's notebook sketches for modelling how pattern formation occurs in nature

# The dawn of computers

- Algorithm execution and informatics were laborious tasks before computers came about

- The invention of electronic computers and dizzying improvements in speed and memory led to most informatics tasks being performed on computers

Computers + algorithmics + informatics = Machine learning
**Computers + algorithmics + informatics + biology = Bio-informatics**



Wikipedia

Radhanath Sikdar, the human "computer" who first calculated the height of Mt Everest



Computerhistory.org

The Bombe (1942), one of the first large scale electromechanical devices built in Bletchley Park for British WW2 cryptanalysis

# Nucleic acid research : understanding, sequencing, engineering

Timeline | **DNA milestones**
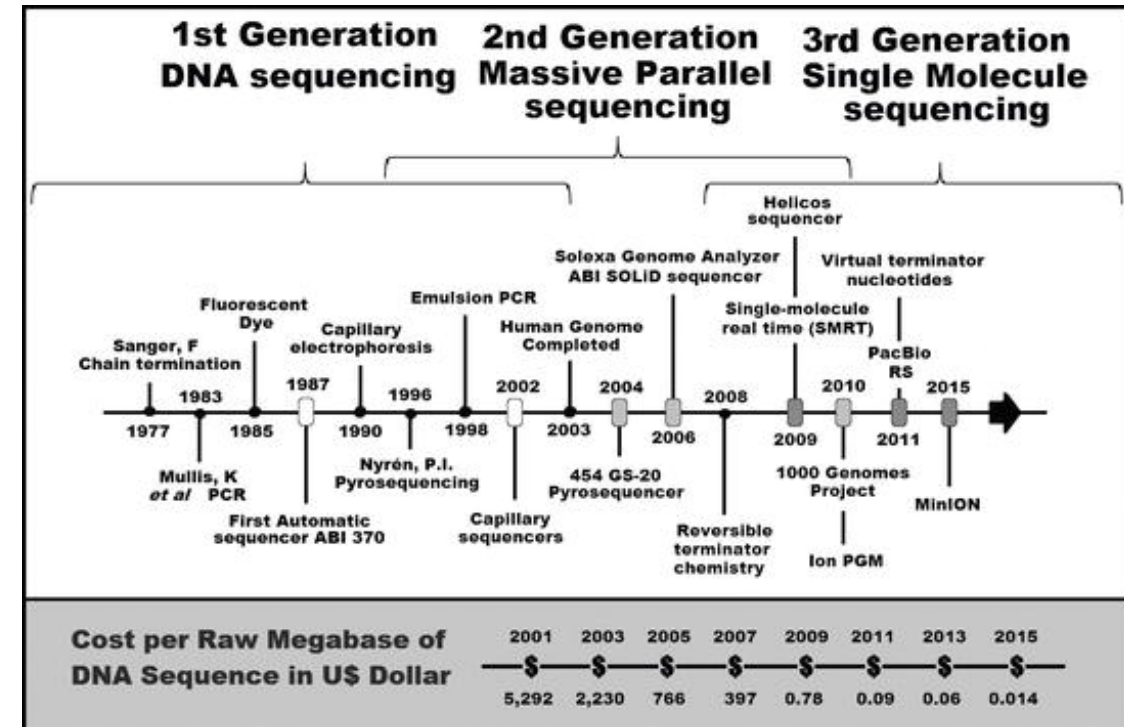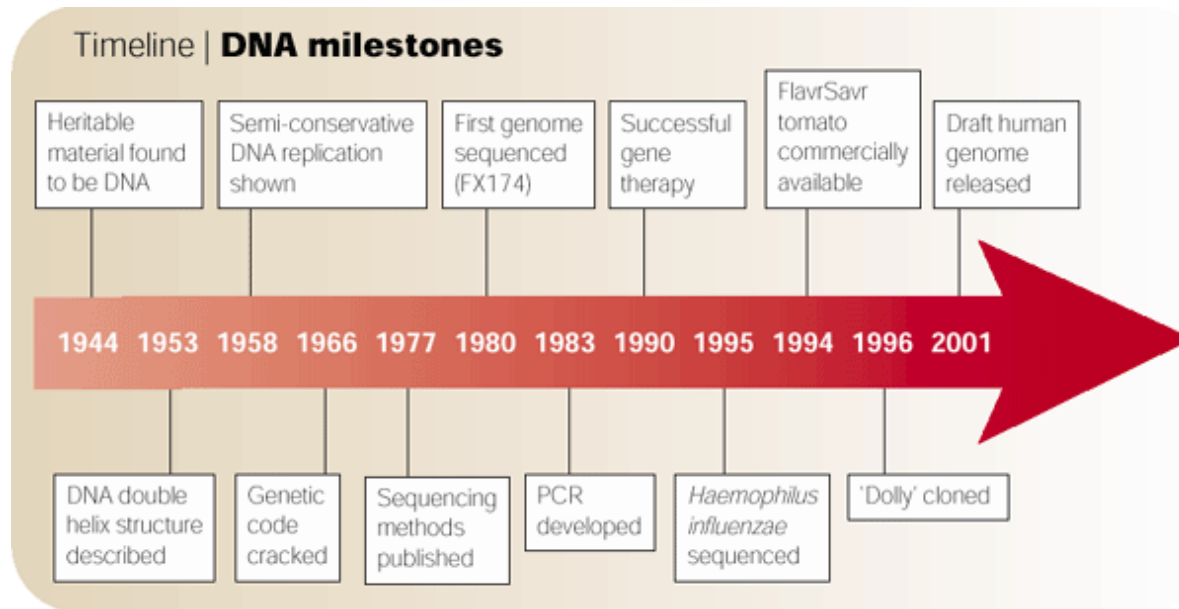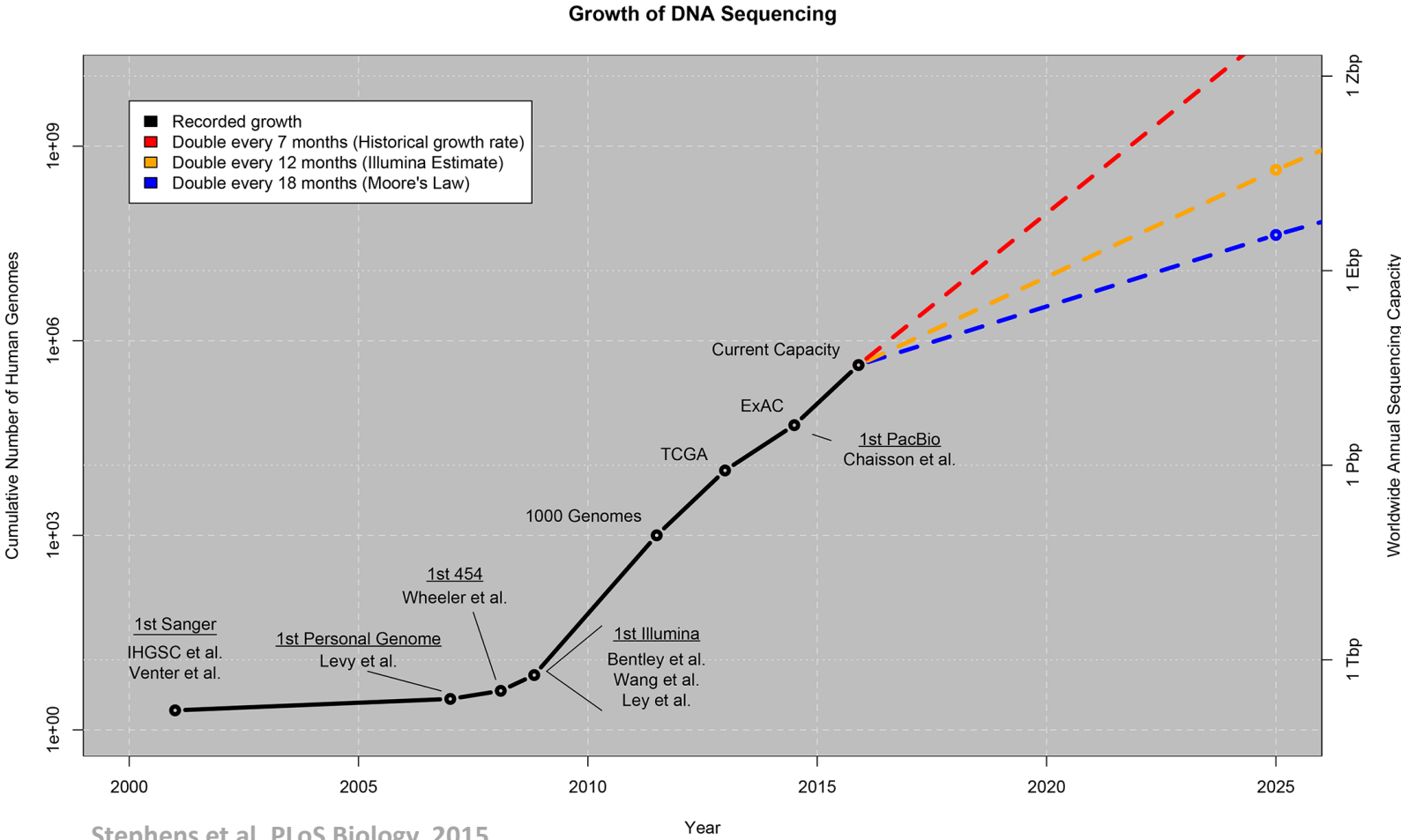
Heritable material found to be DNA · Semi-conservative DNA replication shown · First genome sequenced (FX174) · Successful gene therapy · FlavrSavr tomato commercially available · Draft human genome released

1944 1953 1958 1966 1977 1980 1983 1990 1995 1994 1996 2001

DNA double helix structure described · Genetic code cracked · Sequencing methods published · PCR developed · *Haemophilus influenzae* sequenced · 'Dolly' cloned



**1st Generation DNA sequencing** · **2nd Generation Massive Parallel sequencing** · **3rd Generation Single Molecule sequencing**

Sanger, F Chain termination · Fluorescent Dye · Capillary electrophoresis · Emulsion PCR · Human Genome Completed · Solexa Genome Analyzer ABI SOLiD sequencer · Helicos sequencer · Virtual terminator nucleotides · Single-molecule real time (SMRT) · PacBio RS

1977 · 1983 · 1985 · 1987 · 1990 · 1996 · 1998 · 2002 · 2003 · 2004 · 2006 · 2008 · 2009 · 2010 · 2011 · 2015

Mullis, K et al PCR · First Automatic sequencer ABI 370 · Nyrén, P.I. Pyrosequencing · Capillary sequencers · 454 GS-20 Pyrosequencer · Reversible terminator chemistry · 1000 Genomes Project · Ion PGM · MinION

**Cost per Raw Megabase of DNA Sequence in U$ Dollar**

| | 2001 | 2003 | 2005 | 2007 | 2009 | 2011 | 2013 | 2015 |
|---|---|---|---|---|---|---|---|---|
| $ | 5,292 | 2,230 | 766 | 397 | 0.78 | 0.09 | 0.06 | 0.014 |

# Automation and scale in nucleic acid sequencing



Growth of DNA Sequencing

Stephens et al, PLoS Biology, 2015

# Bioinformatics : 3 decades of explosive growth



Wikipedia

Frederick Sanger : Nobel prize winner (twice!!) for studies on sequence and structure of insulin and for nucleic acid sequencing

- **Rapid improvements in molecular sequencing technologies** for peptides, DNA and RNA produces large amounts of data to be analyzed

- Moore's Law : computers tend to double in computing capacity every year, making **intensive computation feasible, tractable and economical**

- Development of methodology, training a new kind of multidisciplinary scientist, investment in infrastructure – the **Human Genome Project** (1990 – 2003)



A persevering prediction
Number of transistors in CPU*
Log scale

MOORE'S LAW DEFINED

1960  70  80  90  2000  10 14

Source: Intel          *Central processing unit
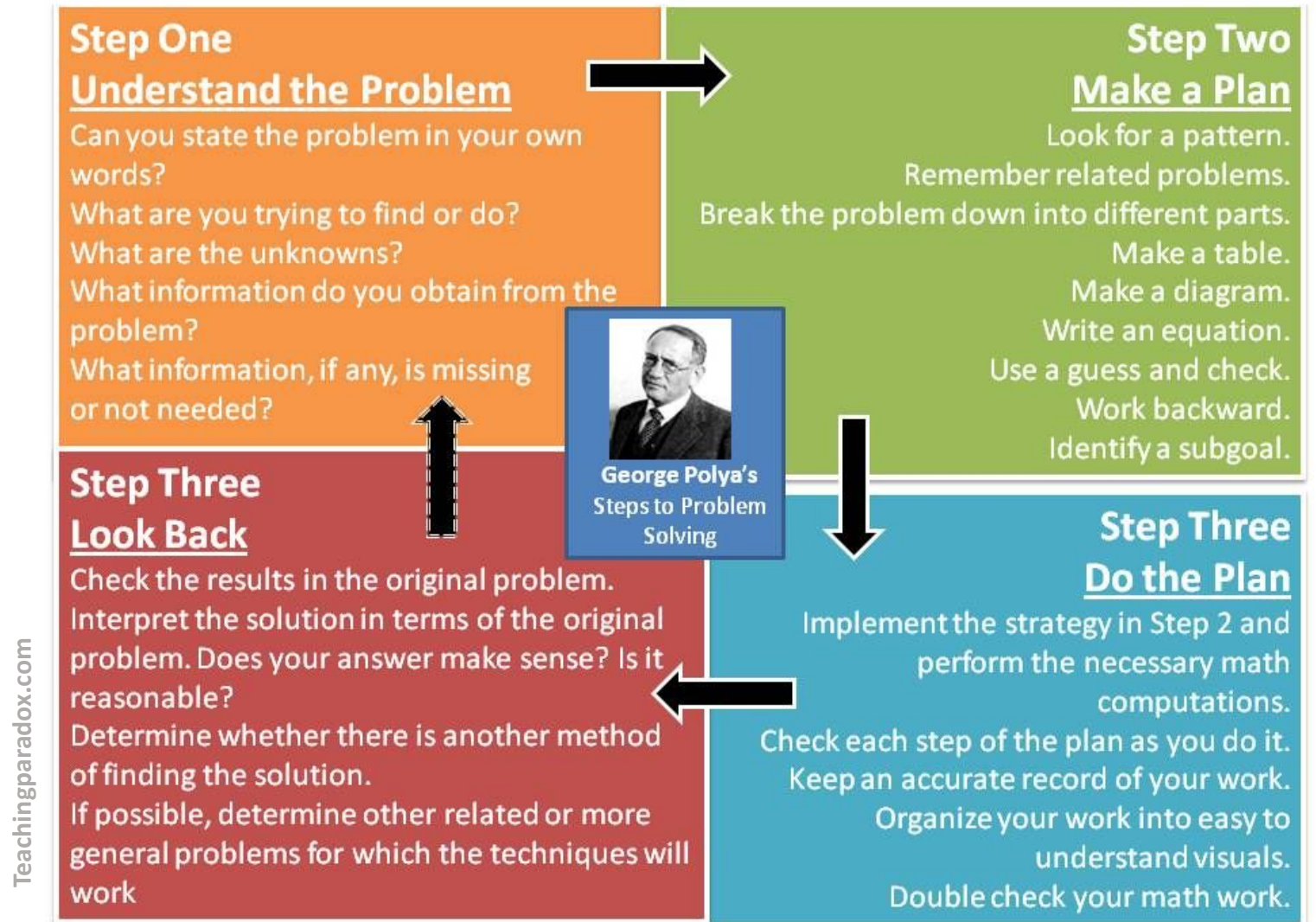
Economist.com

# The mouse & human genome projects

- Sequenced and assembled reference drafts of the human and mouse genome

- Mapped and predicted genes in the reference sequence

- Laid the groundwork for genomics (including computational genomics) as a discipline
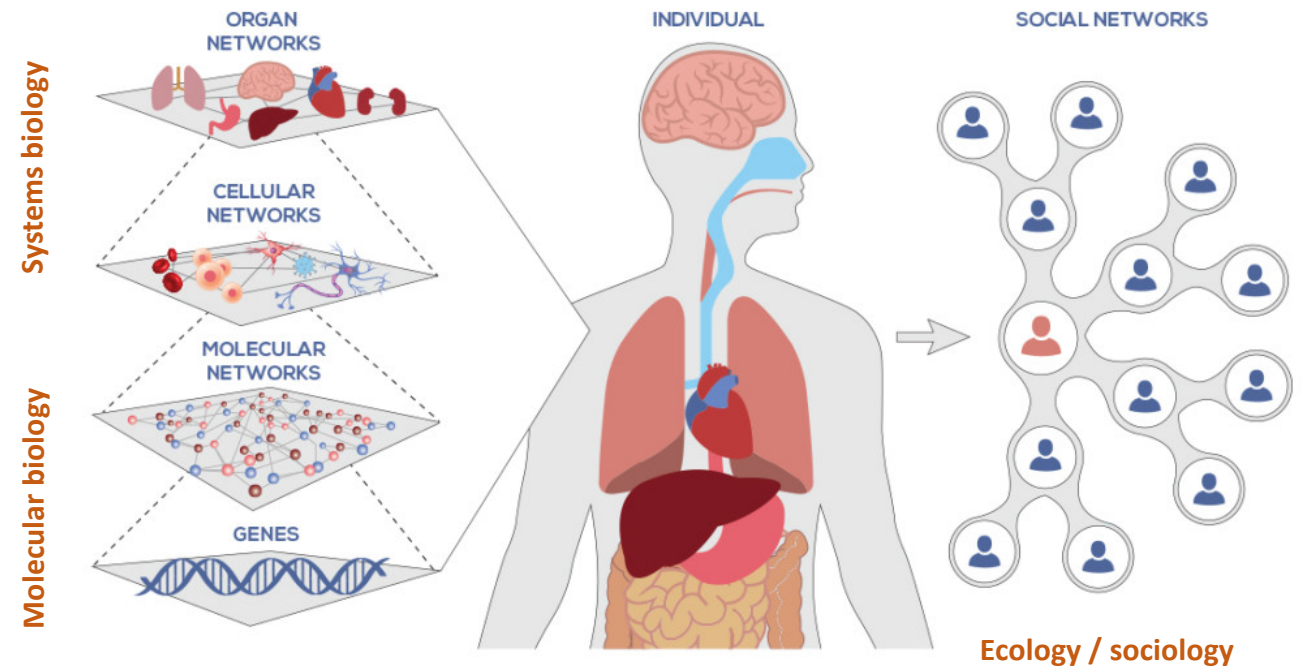
# Asking the right questions

- How to solve it, written by George Polya
  - Strongly encouraged to read this book !

- Problem solving is a cyclic process : true for bioinformatics as well

Teachingparadox.com

**Step One**
**Understand the Problem**
Can you state the problem in your own words?
What are you trying to find or do?
What are the unknowns?
What information do you obtain from the problem?
What information, if any, is missing or not needed?

**Step Two**
**Make a Plan**
Look for a pattern.
Remember related problems.
Break the problem down into different parts.
Make a table.
Make a diagram.
Write an equation.
Use a guess and check.
Work backward.
Identify a subgoal.

**George Polya's Steps to Problem Solving**

**Step Three**
**Look Back**
Check the results in the original problem.
Interpret the solution in terms of the original problem. Does your answer make sense? Is it reasonable?
Determine whether there is another method of finding the solution.
If possible, determine other related or more general problems for which the techniques will work

**Step Three**
**Do the Plan**
Implement the strategy in Step 2 and perform the necessary math computations.
Check each step of the plan as you do it.
Keep an accurate record of your work.
Organize your work into easy to understand visuals.
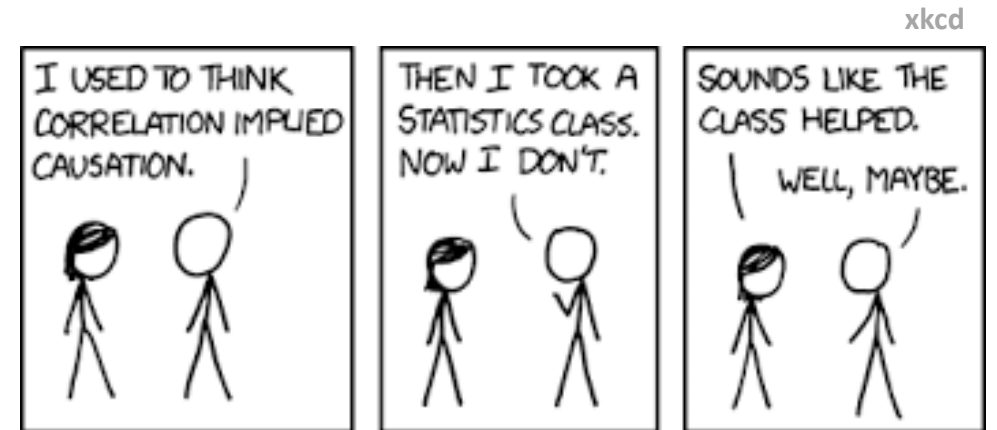Double check your math work.

# Performing analysis at the right granularity

- Given a biological problem, it is important to figure out the right granularity to perform experiments / note observations

- Some problems can be studied at multiple levels : disease - molecular pathology, systems level changes in physiology, transmission of disease
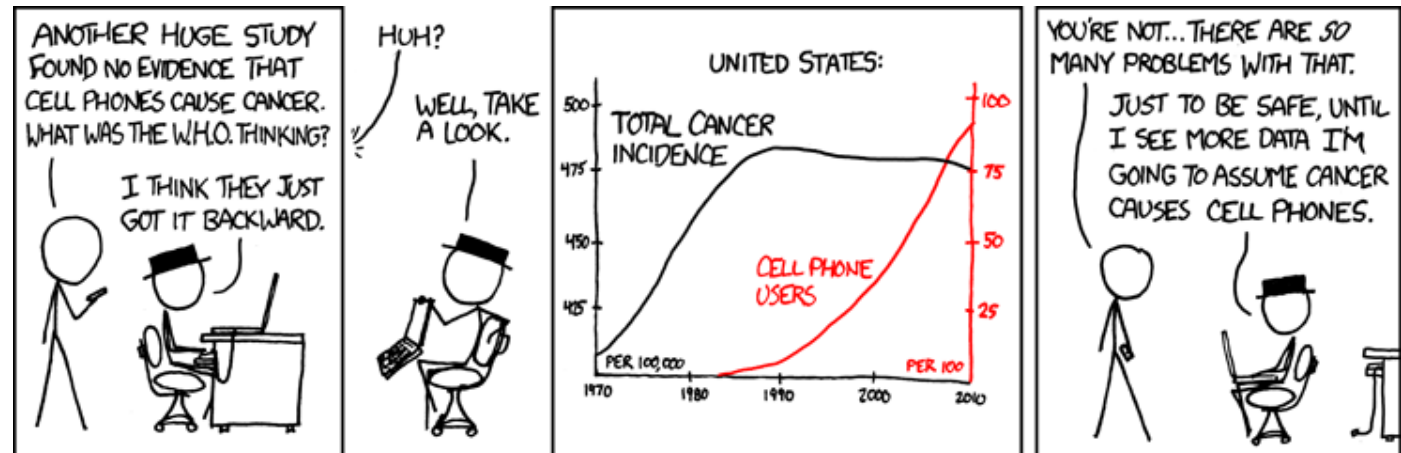


Institute for Systems Biology

# Bioinformatics as an epidemiological science

- A large amount of biology is an empirical science
  - Bioinformatics helps generate, screen and test hypotheses in high throughput fashion to help transition it to a theoretical science
- Since many bioinformatics studies are correlative, follow up studies that de-confound correlation and causation are often required
  - Perturbation studies (eg. gene knockout models)

xkcd

# Pitfalls of correlative studies

- Poorly set up hypotheses, bad inferential mechanisms, and low quality data can all contribute to arriving at the wrong conclusion (artifacts)

- Methodological rigor and domain knowledge are key to avoiding artifacts
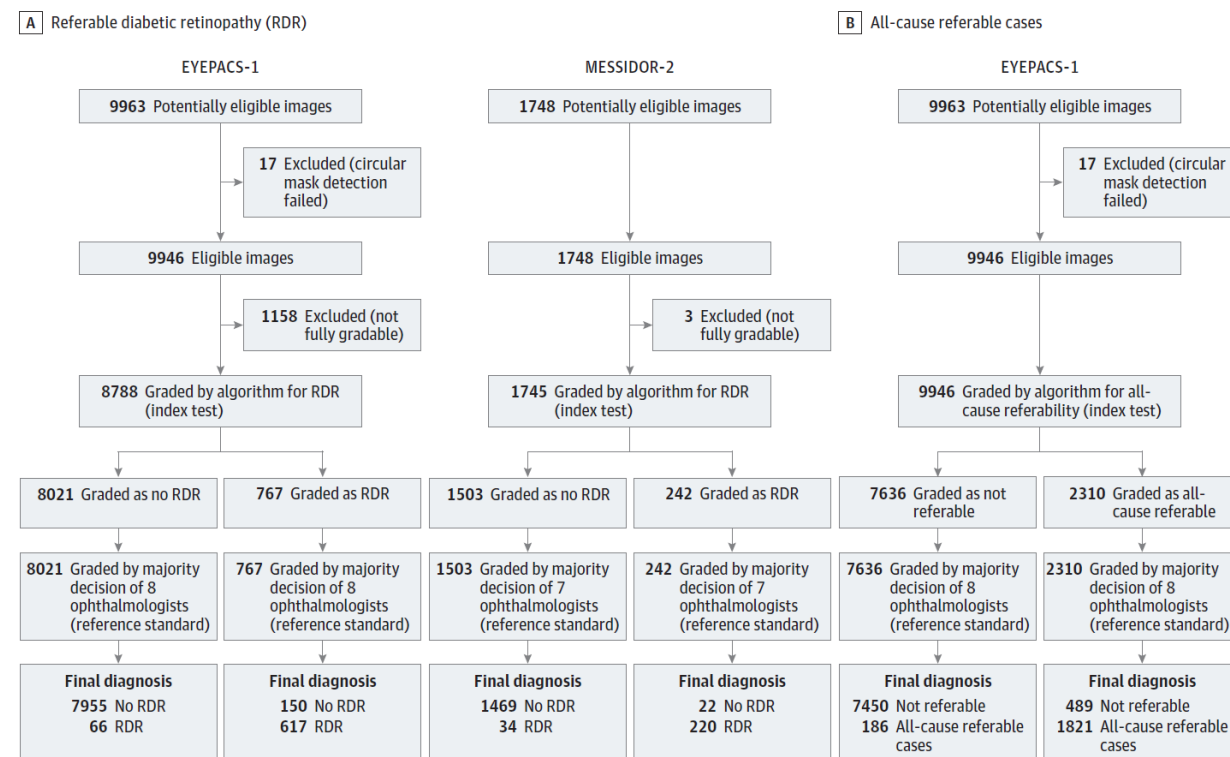


xkcd

# "Artificial intelligence" in biology

The perception of AI has changed over the years :

1970s : expert systems ( eg. early recommender systems in medicine)

1990s : statistical machine learning ( eg. genome wide association studies of diseases )
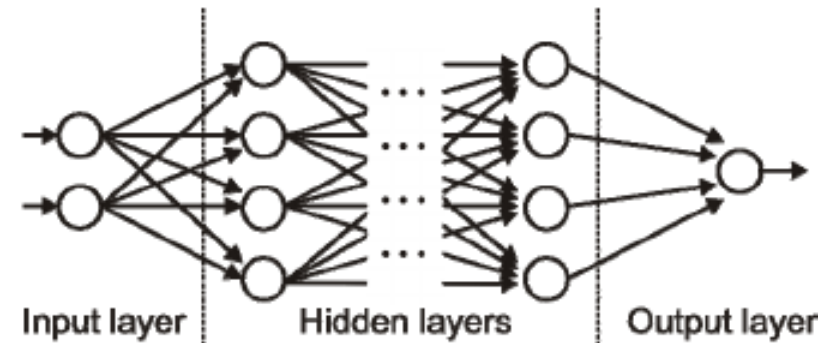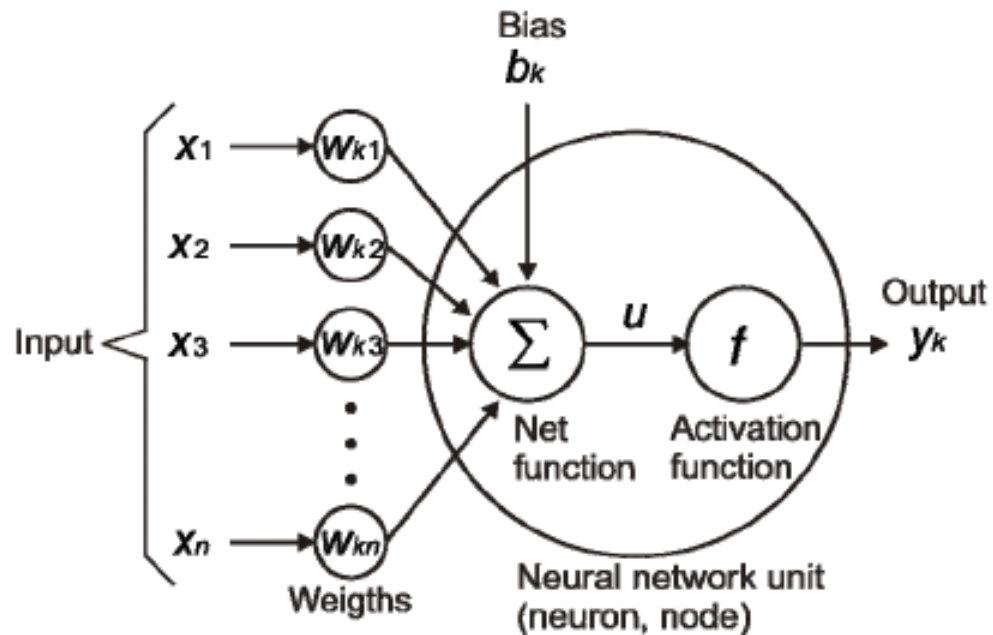
2010s : deep neural networks ( eg. prediction of diabetic retinopathy from imaging )
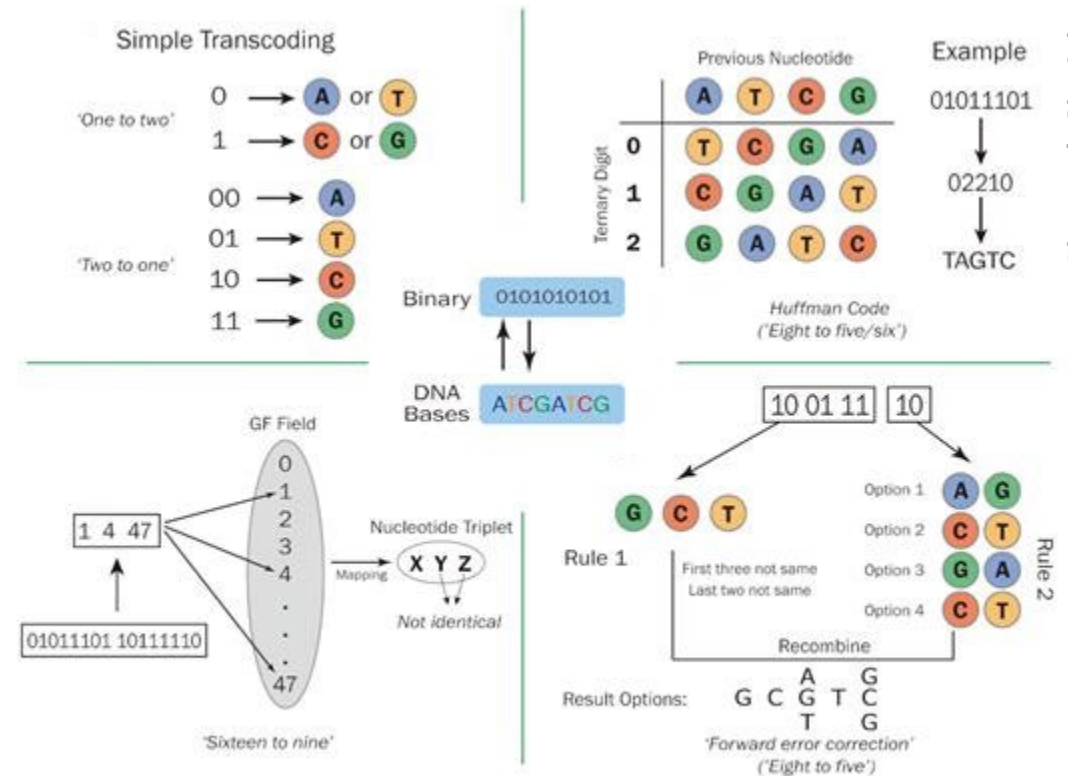


Gulshan et al, JAMA, 2016

# Using neuroscience to model AI

- Neural networks

# Other areas of computer science – biology intersection

- **Biological computation** : how can biological processes be harnessed to improve computing – eg. DNA based data storage

- **Biomedical devices** : electronic devices used for patient monitoring, treatment and therapy

- **Synthetic biology** : creating new biochemical / biological entities



Ping et al, GigaScience

Information encoding strategies for DNA – based storage

# If you'd like to perform research in bioinformatics …

- Knowledge of coding is required

- Send email to [pradiptaray@gmail.com](mailto:pradiptaray@gmail.com)
- Many other wonderful laboratories in UTD working on bioinformatics