

Bphys/Biol E-101 = HST 508 = GEN224

Your grade is based on six problem sets and a course project, with emphasis on collaboration across disciplines.

Open to: upper level undergraduates, and all graduate students. The prerequisites are basic knowledge of molecular biology, statistics, & computing.

Please hand in your questionnaire after this class.
First problem set is due before Lecture 3 starts
via email or paper depending on your section TF.

Bio 101: Genomics & Computational Biology

Week#1 *Intro 1:* **Computing, Statistics, Perl, Mathematica**

Week#2 *Intro 2:* Biology, comparative genomics, models & evidence, applications

Week#3 *DNA 1:* Polymorphisms, populations, statistics, pharmacogenomics, databases

Week#4 *DNA 2:* Dynamic programming, Blast, multi-alignment, **HiddenMarkovModels**

Week#5 *RNA 1:* 3D-structure, microarrays, library sequencing & quantitation concepts

Week#6 *RNA 2:* Clustering by gene or condition, DNA/RNA motifs.

Week#7 *Protein 1:* 3D structural genomics, homology, dynamics, function & drug design

Week#8 *Protein 2:* Mass spectrometry, modifications, quantitation of interactions

Week#9 *Network 1:* Metabolic kinetic & flux balance optimization methods

Week#10 *Network 2:* Molecular computing, self-assembly, genetic algorithms, neural-nets

Week#11 *Network 3:* Cellular, developmental, social, ecological & commercial models

Week#12 Project presentations

Week#13 Project Presentations

Week#14 Project Presentations

Intro 1: Today's story, logic & goals

Life & computers : **Self-assembly** required

Discrete & continuous models

Minimal life & programs

Catalysis & Replication

Differential equations

Directed graphs & pedigrees

Mutation & the Single Molecules models

Bell curve statistics

Selection & optimality

101

101

acgt

1
0

1

1

0

1

1

0

1

1

0

1

00=a

01=c

10=g

11=t

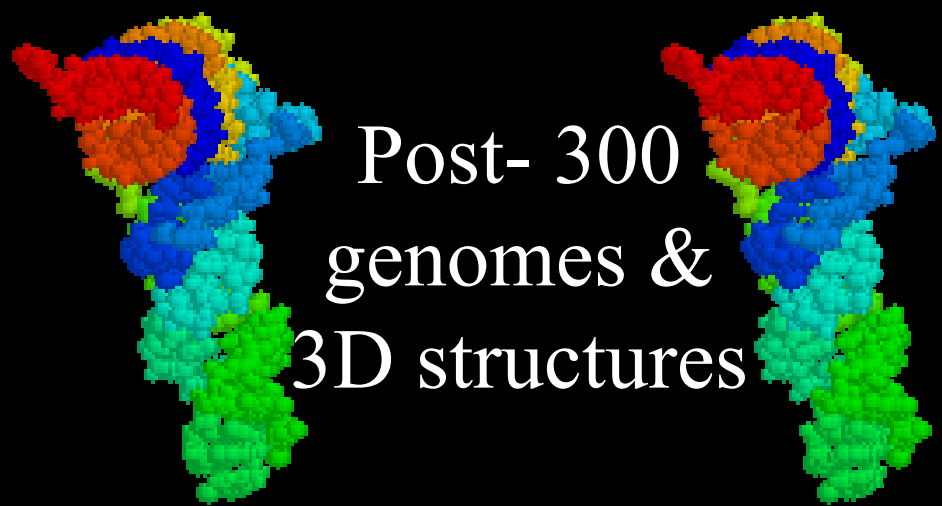
1
0
1

1
0
1

gggatttagctcagtt

gggagagcgcca**gact**

gaa



gat

ttg

gag

gtcctgtgttcgatcc

acagaaattcgacca

Discrete

Continuous

a sequence

lattice

digital

$\Sigma \Delta x$

neural/regulatory on/off

sum of black & white

essential/neutral

alive/not

a weight matrix of sequences

molecular coordinates

analog (16 bit A2D converters)



dx

gradients & graded responses

gray

conditional mutation

probability of replication

Bits (discrete)

bit = binary digit

1 base \geq 2 bits

1 byte = 8 bits

+ Kilo	Mega	Giga	Tera	Peta	Exa	Zetta	Yotta +
3	6	9	12	15	18	21	24
- milli	micro	nano	pico	femto	atto	zepto	yocto -

	Kibi	Mebi	Gibi	Tebi	Pebi	Exbi
1024 = 2^{10}	2^{20}	2^{30}	2^{40}	2^{50}	2^{60}	

<http://physics.nist.gov/cuu/Units/prefixes.html>

Defined quantitative measures

Seven basic (Système International) SI units:
s, m, kg, mol, K, cd, A

(some measures at precision of 14 significant figures)

Quantal: Planck time, length: 10^{-43} seconds, 10^{-35} meters,
mol = 6.0225×10^{23} entities.

Quantitative definition of life?

Historical/Terrestrial Biology vs "General Biology"

Probability of replication ... of complexity from simplicity
(in a specific environment)

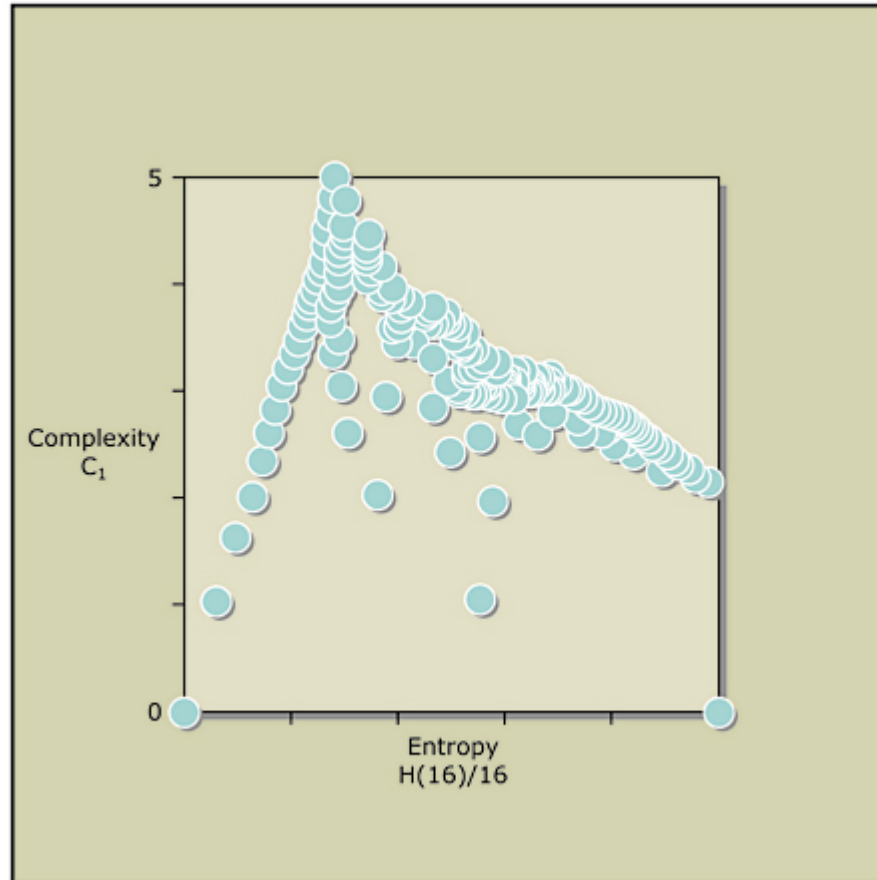
Robustness/Evolvability
(in a variety of environments)

Examples: mules, fires, nucleating crystals,
pollinated flowers, viruses, predators,
molecular ligation, factories, self-assembling machines.

Complexity definitions

1. Computational Complexity = speed/memory scaling P, NP
2. Algorithmic Randomness (Chaitin-Kolmogorov)
3. Entropy/information
4. Physical complexity
(Bernoulli-Turing Machine)

Complexity & Entropy/Information



Why Model?

- To understand biological/chemical data.
(& design useful modifications)
- To share data we need to be able to **search, merge, & check** data via models.
- Integrating diverse data types can reduce random & systematic errors.

Which models will we search, merge & check in this course?

- Sequence: Dynamic programming, assembly, translation & trees.
- 3D structure: motifs, catalysis, complementary surfaces – energy and kinetic optima
- Functional genomics: clustering
- Systems: qualitative & boolean networks
- Systems: differential equations & stochastic
- Network optimization: Linear programming

Intro 1: Today's story, logic & goals

Life & computers : **Self-assembly** required

Discrete & continuous models

Minimal life & programs

Catalysis & Replication

Differential equations

Directed graphs & pedigrees

Mutation & the Single Molecules models

Bell curve statistics

Selection & optimality

Elements of RNA-based life: C,H,N,O,P

Useful for many species:

Na, K, Fe, Cl, Ca, Mg, Mo, Mn, S, Se, Cu, Ni, Co, Si

Group	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
Period																		
1	1 H																	2 He
2	3 Li	4 Be											5 B	6 C	7 N	8 O	9 F	10 Ne
3	11 Na	12 Mg											13 Al	14 Si	15 P	16 S	17 Cl	18 Ar
4	19 K	20 Ca	21 Sc	22 Ti	23 V	24 Cr	25 Mn	26 Fe	27 Co	28 Ni	29 Cu	30 Zn	31 Ga	32 Ge	33 As	34 Se	35 Br	36 Kr
5	37 Rb	38 Sr	39 Y	40 Zr	41 Nb	42 Mo	43 Tc	44 Ru	45 Rh	46 Pd	47 Ag	48 Cd	49 In	50 Sn	51 Sb	52 Te	53 I	54 Xe
6	55 Cs	56 Ba	* 71 Lu	72 Hf	73 Ta	74 W	75 Re	76 Os	77 Ir	78 Pt	79 Au	80 Hg	81 Tl	82 Pb	83 Bi	84 Po	85 At	86 Rn
7	87 Fr	88 Ra	** 103 Lr	104 Rf	105 Db	106 Sg	107 Bh	108 Hs	109 Mt	110 Uun	111 Uuu	112 Uub	113 Uut	114 Uuq	115 Uup	116 Uuh	117 Uus	118 Uuo
*Lanthanoids			* 57 La	58 Ce	59 Pr	60 Nd	61 Pm	62 Sm	63 Eu	64 Gd	65 Tb	66 Dy	67 Ho	68 Er	69 Tm	70 Yb		
**Actinoids			** 89 Ac	90 Th	91 Pa	92 U	93 Np	94 Pu	95 Am	96 Cm	97 Bk	98 Cf	99 Es	100 Fm	101 Md	102 No		

Minimal self-replicating units

Minimal theoretical composition: 5 elements: C,H,N,O,P

Environment = water, NH_4^+ , 4 NTP⁻s, lipids

Johnston et al. [Science 2001 292:1319-1325](#) RNA-catalyzed RNA polymerization: accurate and general RNA-templated primer extension

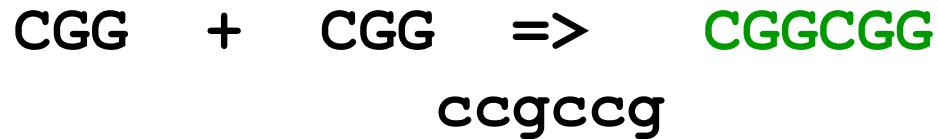
(http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&list_uids=11358999&dopt=Abstract).

Minimal programs

perl -e "print exp(1);"	2.71828182845905
excel : = EXP(1)	2.718281828459050000000000
f77 : print*, exp(1.q0)	2.71828182845904523536028747135266
Mathematica : N[Exp[1],100]	2.71828182845904523536028747135266249775 7247093699959574966967627724076630353547594571382178525166427

- Underlying these are algorithms for arctangent and hardware for RAM and printing.
- Beware of approximations & boundaries.
- Time & memory limitations. E.g. first two above 64 bit floating point:
52 bits for mantissa (= 15 decimal digits), 10 for exponent, 1 for +/- signs. 17

Self-replication of complementary nucleotide-based oligomers



Why Perl & Mathematica?

In the hierarchy of languages, **Perl** is a "high level" language, optimized for easy coding of string searching & string manipulation. It is well suited to web applications and is "open source" (so that it is inexpensive and easily extended). It has a very easy learning curve relative to C/C++ but is similar in a few way to C in syntax.

Mathematica is intrinsically stronger on math (symbolic & numeric) & graphics.

Facts of Life 101

Where do parasites come from?

(computer & biological viral codes)

Over \$12 billion/year
on computer viruses ([ref](http://virus.idg.net/crd_virus_126660.html))
(http://virus.idg.net/crd_virus_126660.html)

20 M dead (worse than black plague
& 1918 Flu)

AIDS - HIV-1 ([download](http://www.ncbi.nlm.nih.gov/htbin-))

([http://www.ncbi.nlm.nih.gov/htbin-
post/Taxonomy/wgetorg?id=11676](http://www.ncbi.nlm.nih.gov/htbin-post/Taxonomy/wgetorg?id=11676))

LoveBug

```
Set dirtemp =3D fso.GetSpecialFolder(2)
Set c =3D fso.GetFile(WScript.ScriptFullName)
c.Copy(dirsystem&"\MSKernel32.vbs")
c.Copy(dirwin&"\Win32DLL.vbs")
c.Copy(dirsystem&"\LOVE-LETTER-FOR-YOU.TXT.vbs")
regruns()
html()
spreadtoemail()
listadriv()
```

Polymerase drug resistance mutations

M41L, D67N, T69D, L210W, T215Y, H208Y

PISPIETVPVKLKPGMDGPK

VKQWPLTEEK

IKALIEICAE **L**EKDGKISKI

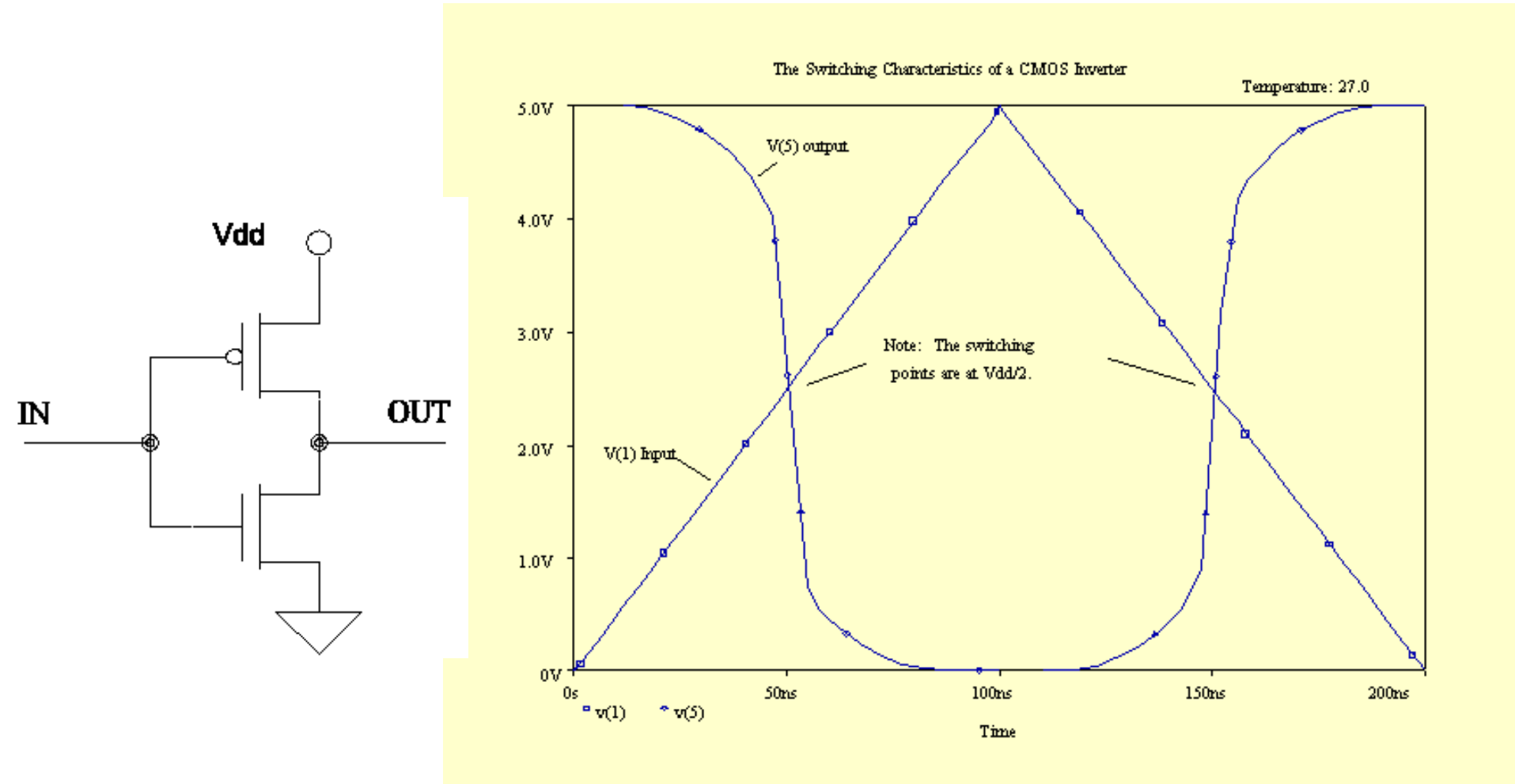
GPVNPYDTPV FAIKKK**NSDK**

WRKLVDFREL NKRTQDFCEV

Conceptual connections

Concept	Computers	Organisms
Instructions	Program	Genome
Bits	0,1	a,c,g,t
Stable memory	Disk,tape	DNA
Active memory	RAM	RNA
Environment	Sockets,people	Water,salts
I/O	AD/DA	proteins
Monomer	Minerals	Nucleotide
Polymer	chip	DNA,RNA,protein
Replication	Factories	1e-15 liter cell sap
Sensor/In	Keys,scanner	Chem/photo receptor
Actuator/Out	Printer,motor	Actomyosin
Communicate	Internet,IR	Pheromones, song

Transistors > inverters > registers > binary
adders > **compilers** > application programs



Self-compiling & self-assembling

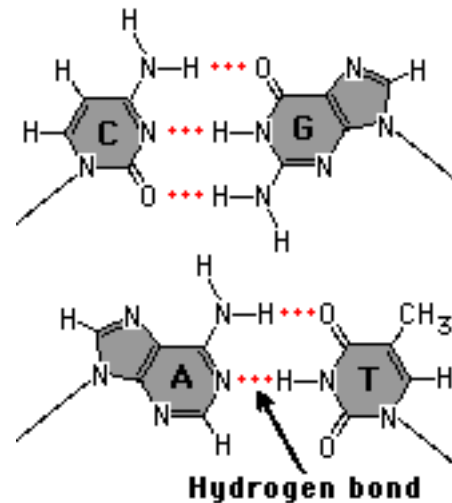


Complementary surfaces

Watson-Crick base pair

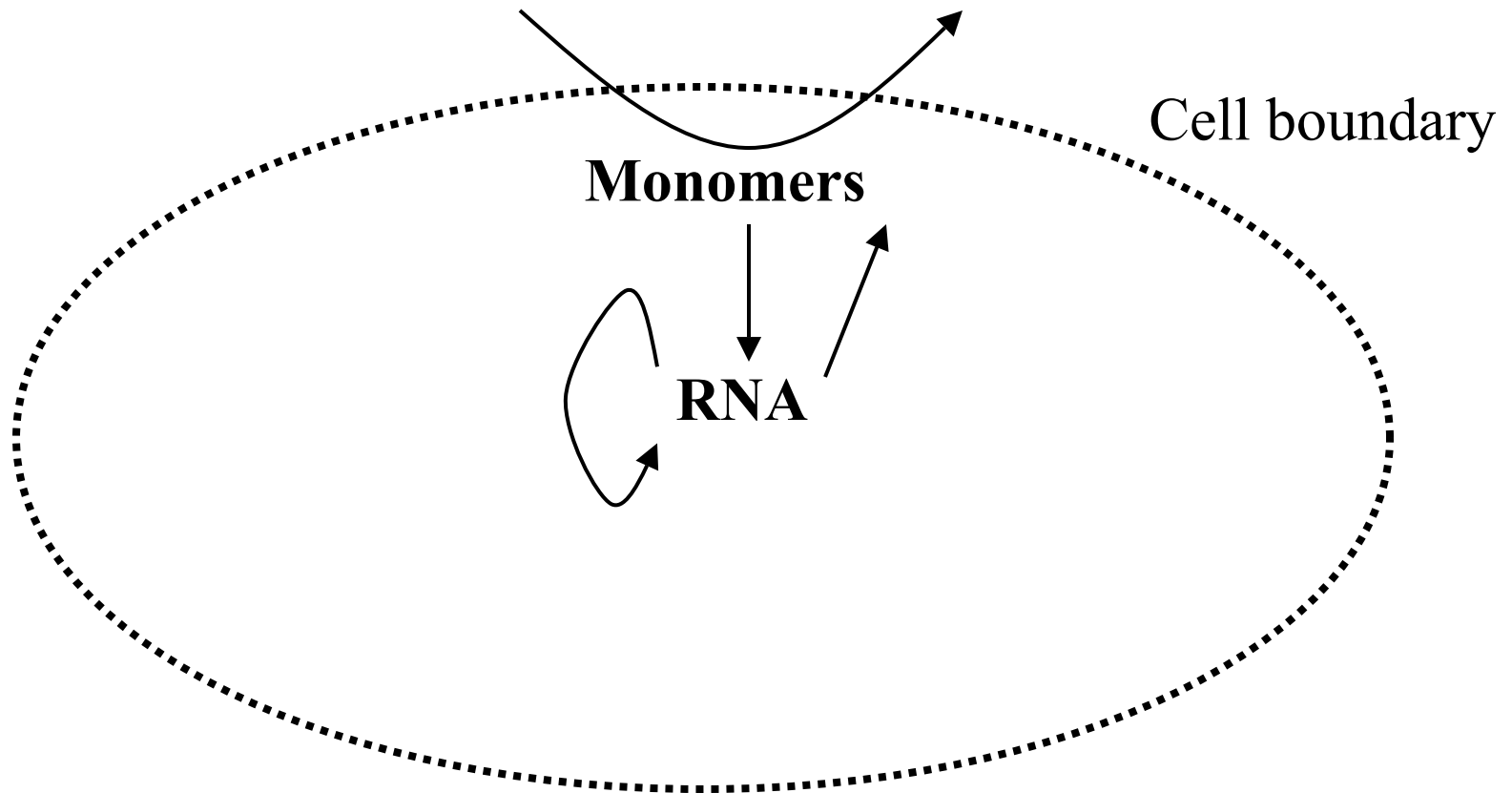
([Nature April 25, 1953](http://www.sil.si.edu/Exhibitions/Science-and-the-Artists-Book/bioc.htm#27))

(<http://www.sil.si.edu/Exhibitions/Science-and-the-Artists-Book/bioc.htm#27>)



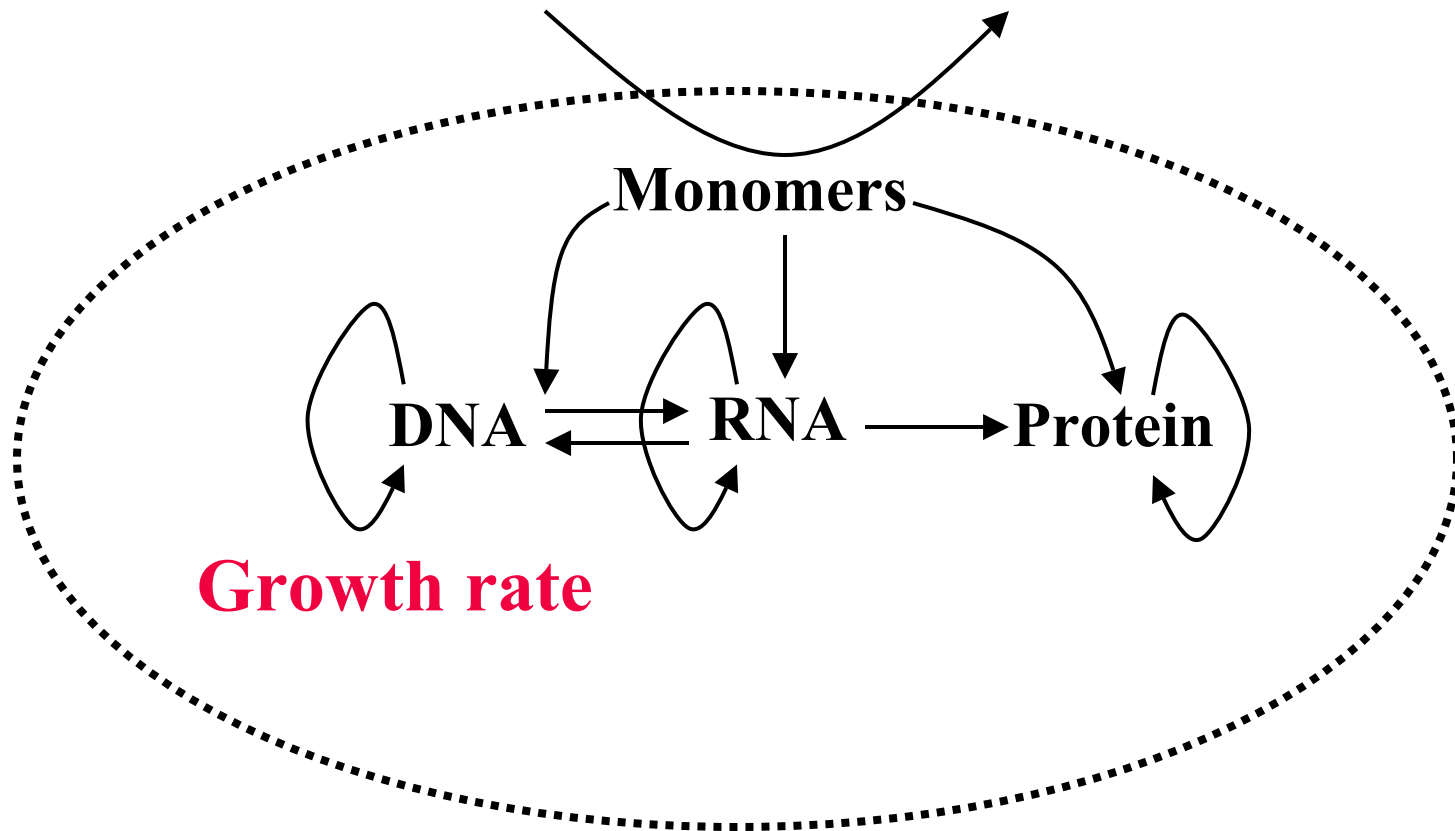
Minimal Life:

Self-assembly, Catalysis, Replication, Mutation, Selection



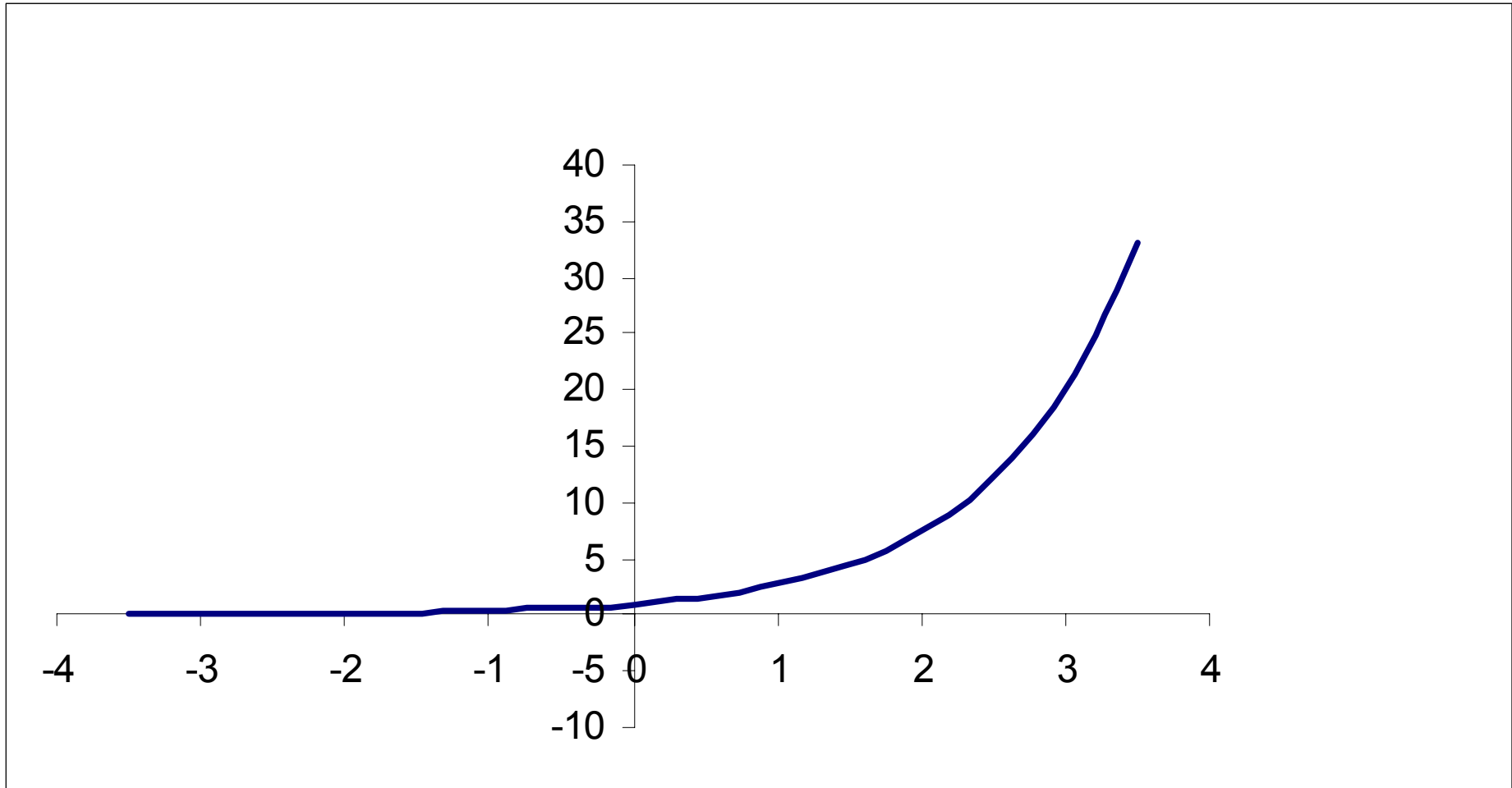
Replicator diversity

Self-assembly, Catalysis, Replication, Mutation, Selection
Polymerization & folding (Revised Central Dogma)



Polymers: Initiate, Elongate, Terminate, Fold, Modify, Localize, Degrade ²⁵

Rorschach Test



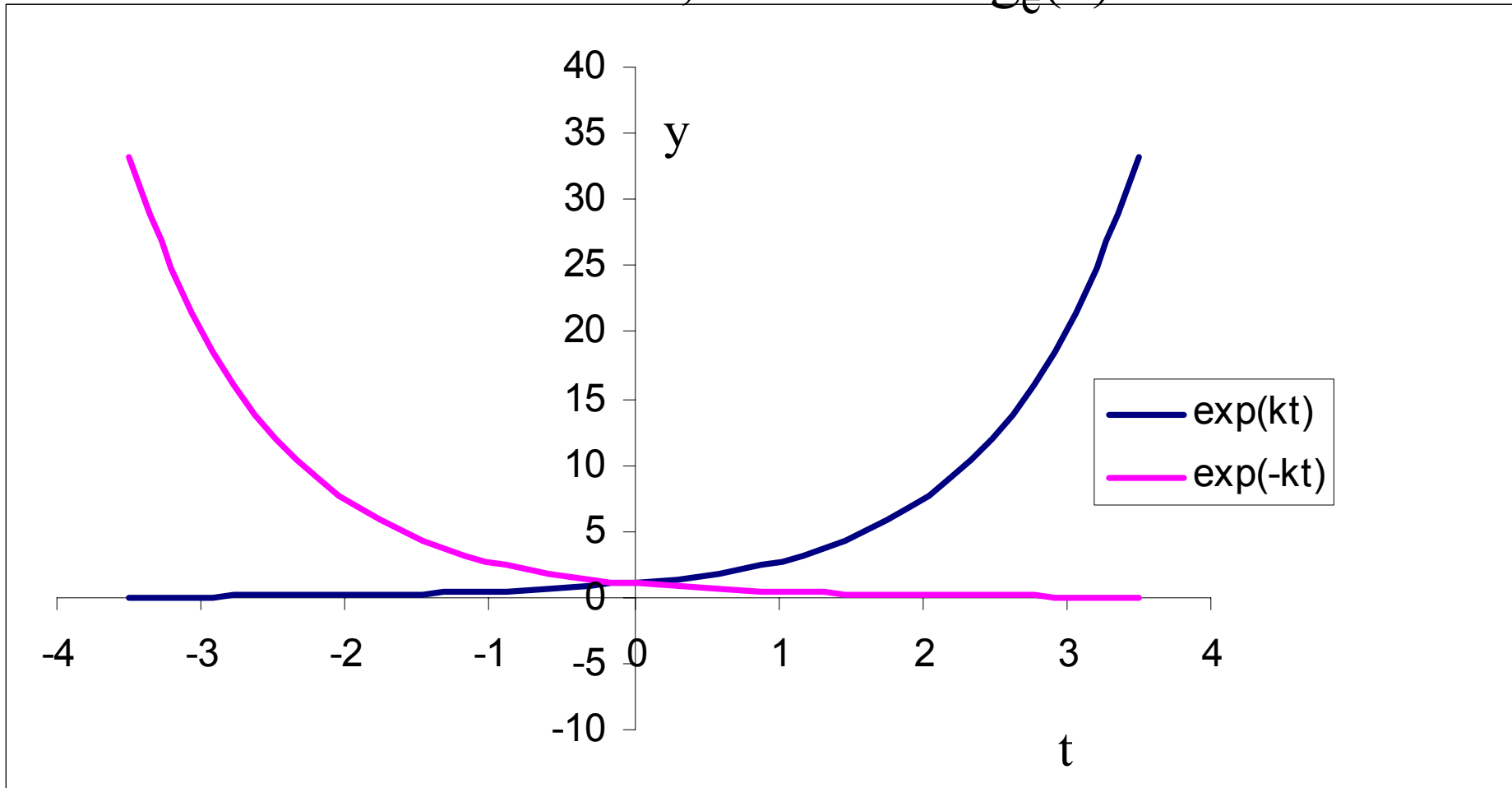
Growth & decay



$$dy/dt = ky$$

$$y = Ae^{kt}; e = 2.71828\dots$$

k =rate constant; half-life= $\log_e(2)/k$



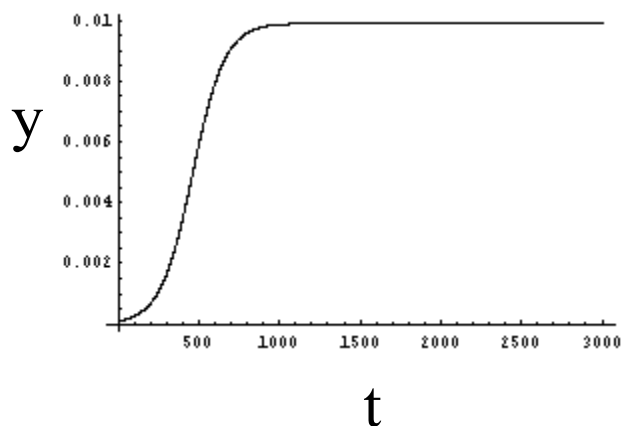
What limits exponential growth?

Exhaustion of resources

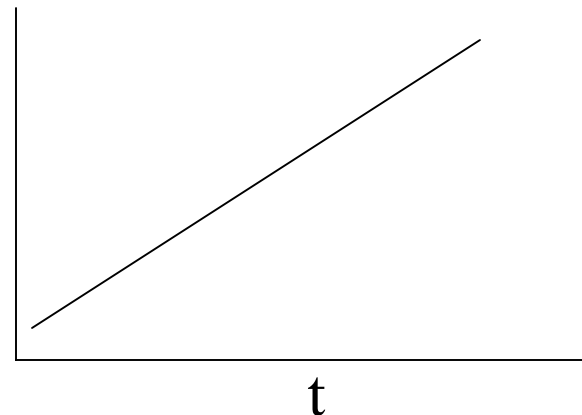
Accumulation of waste products

What limits exponential decay?

Finite particles, stochastic (quantal) limits



Log[y]



Solving differential equations

Mathematica: **Analytical (formal, symbolic)**

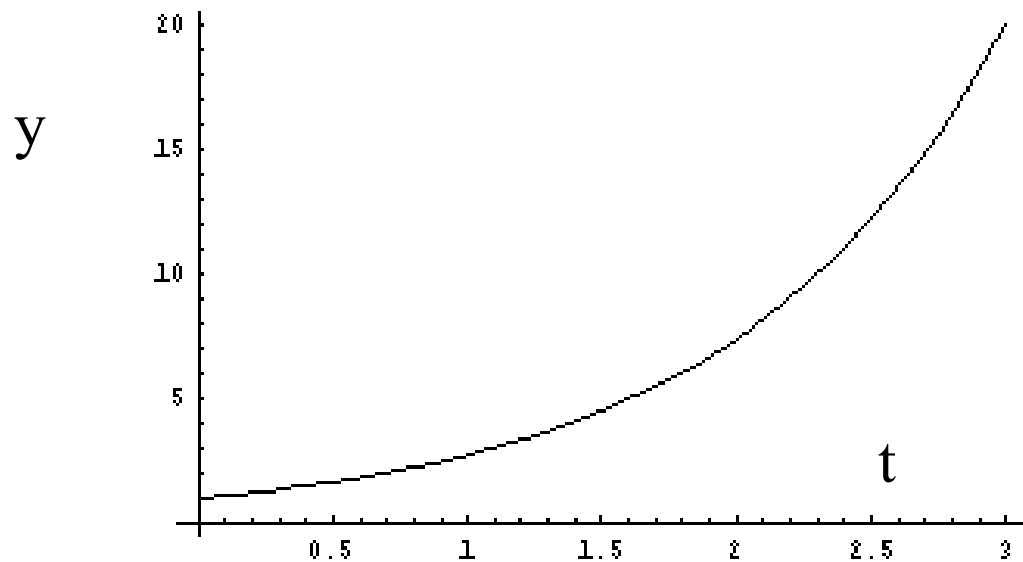
```
In[2]:= DSolve[ {y'[t] == y[t], y[0]==1}, y[t], t ]
```

```
Out[2]= {{y[t]= Et}}
```

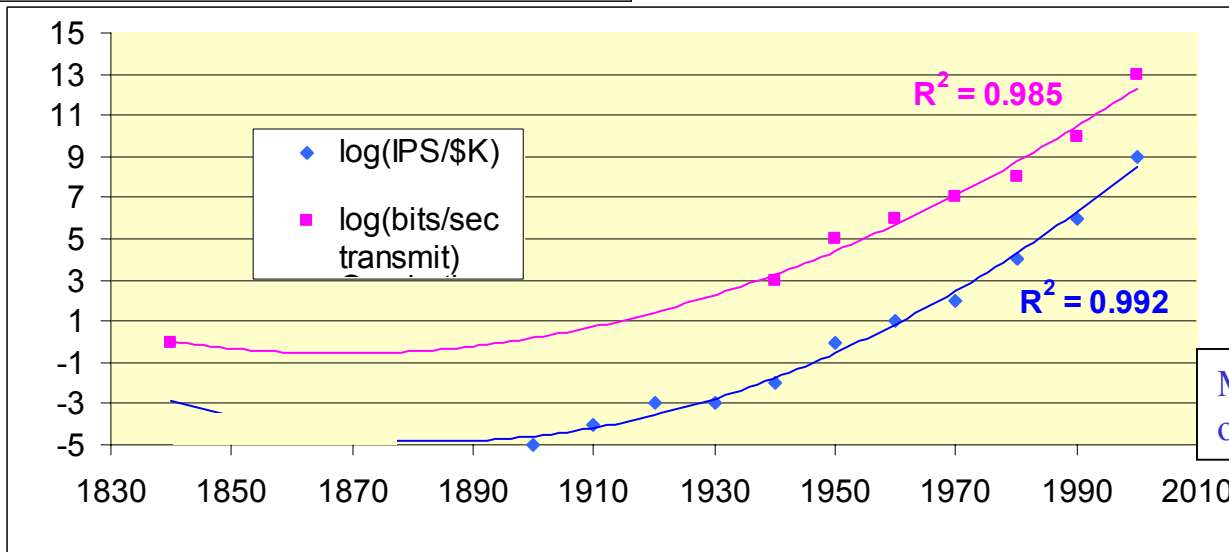
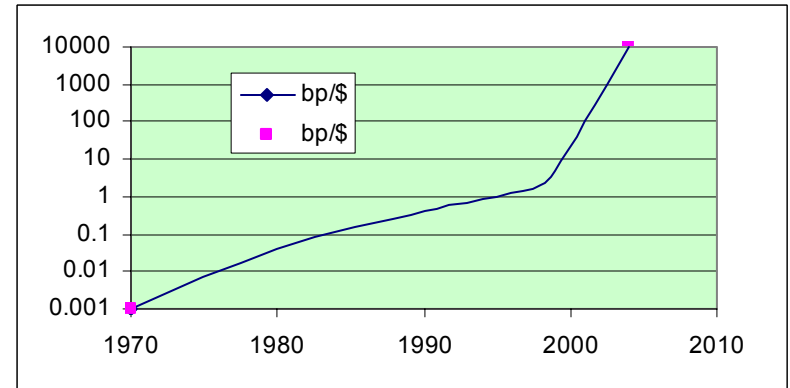
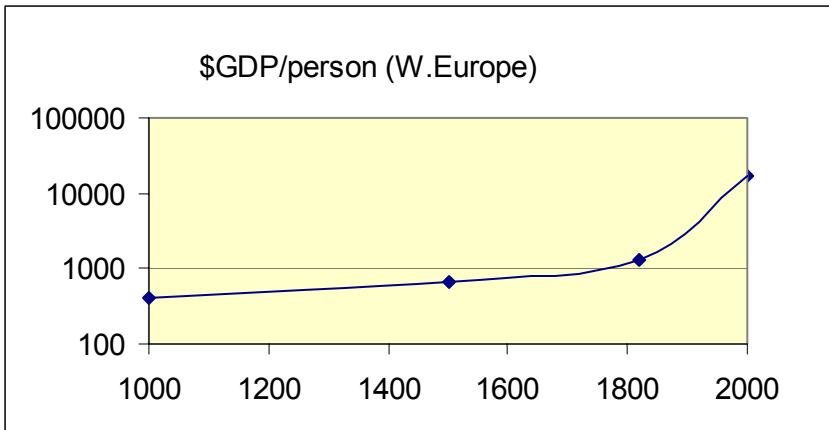
Numerical (&graphical)

```
NDSolve[ {y'[t] == y[t], y[0] == 1}, y, {t, 0, 3}]
```

```
Plot[Evaluate[ y[t] /. % ], {t, 0, 3}]
```



(Hyper)exponential growth



See <http://www.faughnan.com/poverty.html>

See <http://www.kurzweilai.net/meme/frame.html?main=/articles/art0184.html>

Computational power of neural systems

1,000 MIPS (million instructions per second) needed to derive edge or motion detections from video "ten times per second to match the retina ... The 1,500 cubic centimeter human brain is about 100,000 times as large as the retina, suggesting that matching overall human behavior will take about 100 million MIPS of computer power ... The most powerful experimental supercomputers in 1998, costing tens of millions of dollars, can do a few million MIPS."

"The ratio of memory to speed has remained constant during computing history [at Mbyte/MIPS] ... [the human] 100 trillion synapse brain would hold the equivalent 100 million megabytes."

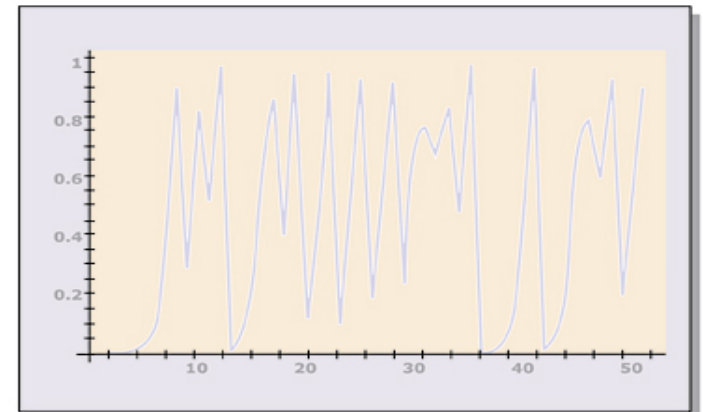
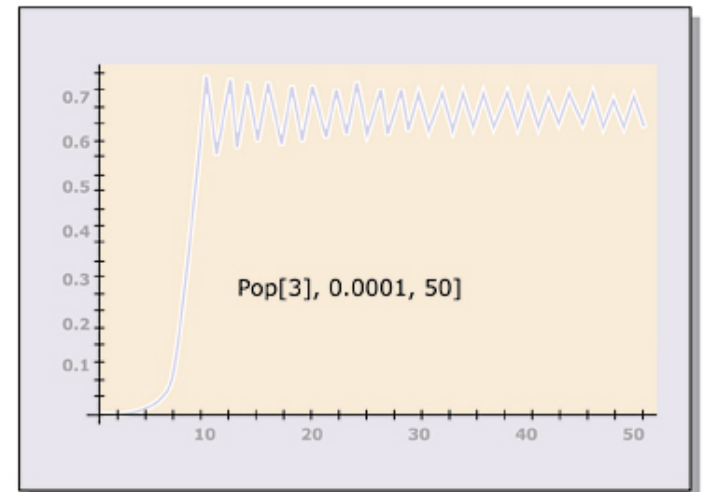
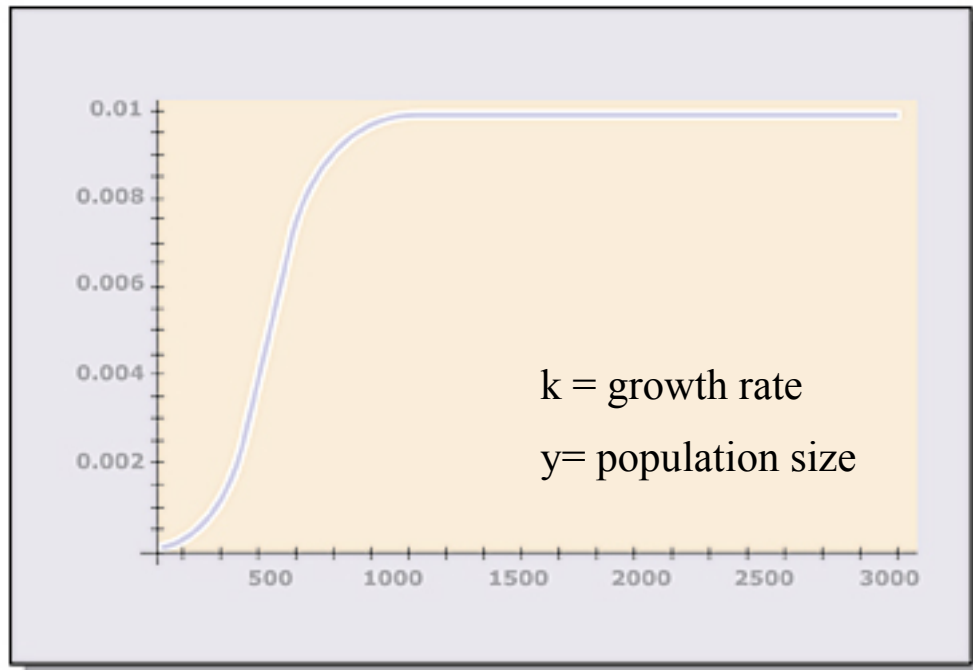
--Hans Moravec <http://www.frc.ri.cmu.edu/~hpm/book97/ch3/retina.comment.html>

2002: the ESC is 35 Tflops & 10Tbytes. <http://www.top500.org/>

Post-exponential growth & chaos

```
Pop[k_][y_] := k y (1 - y);
```

```
ListPlot[NestList[Pop[1.01], 0.0001, 3000], PlotJoined->True];
```



$\text{Pop}[4], 0.0001, 50$

33

Intro 1: Today's story, logic & goals

Life & computers : **Self-assembly** required

Discrete & continuous models

Minimal life & programs

Catalysis & Replication

Differential equations

Directed graphs & pedigrees

Mutation & the Single Molecules models

Bell curve statistics

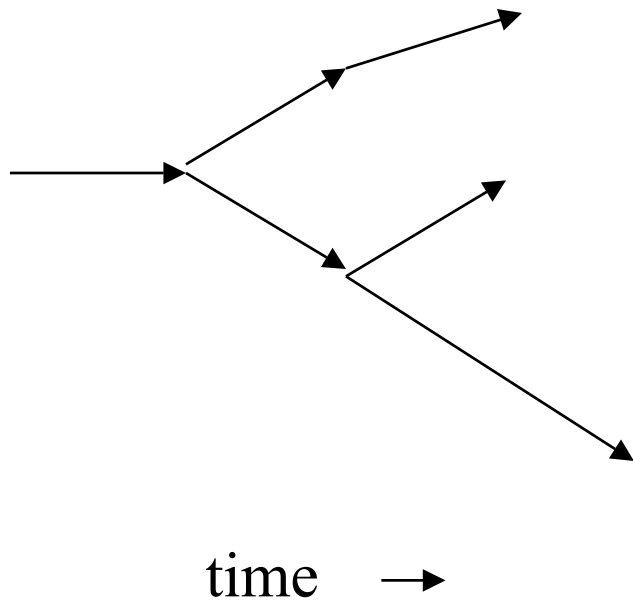
Selection & optimality

Inherited Mutations & Graphs

Directed Acyclic Graph (DAG)

Example: a mutation pedigree

Nodes = an organism, edges = replication with mutation



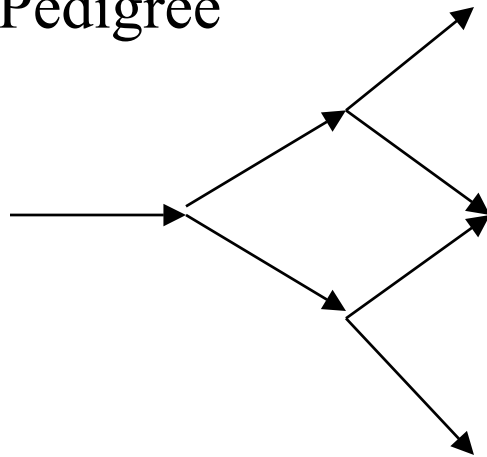
Directed Graphs

Directed Acyclic Graph:

Biopolymer backbone

Phylogeny

Pedigree



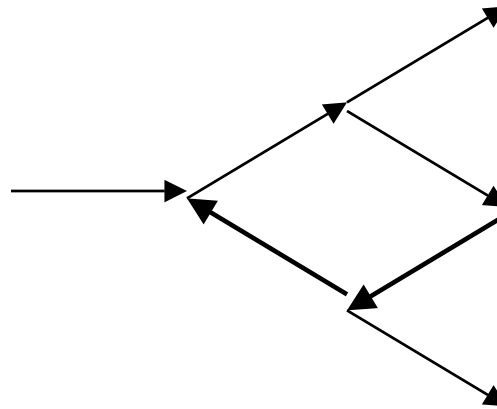
Time →

Cyclic:

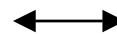
Polymer contact maps

Metabolic &

Regulatory Nets



Time independent or implicit



System models

Feature attractions

E. coli chemotaxis

Red blood cell metabolism

Cell division cycle

Circadian rhythm

Plasmid DNA replication

Phage λ switch

Adaptive, spatial effects

Enzyme kinetics

Checkpoints

Long time delays

Single molecule precision

Stochastic expression

also, all have large genetic & kinetic datasets.

Intro 1: Today's story, logic & goals

Life & computers : **Self-assembly** required

Discrete & continuous models

Minimal life & programs

Catalysis & Replication

Differential equations

Directed graphs & pedigrees

Mutation & the Single Molecules models

Bell curve statistics

Selection & optimality

Bionano-machines

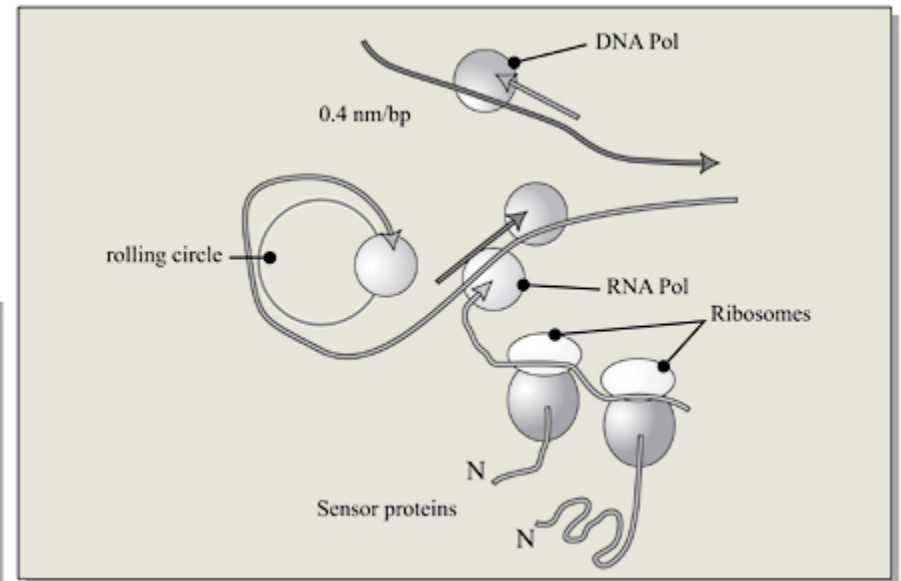
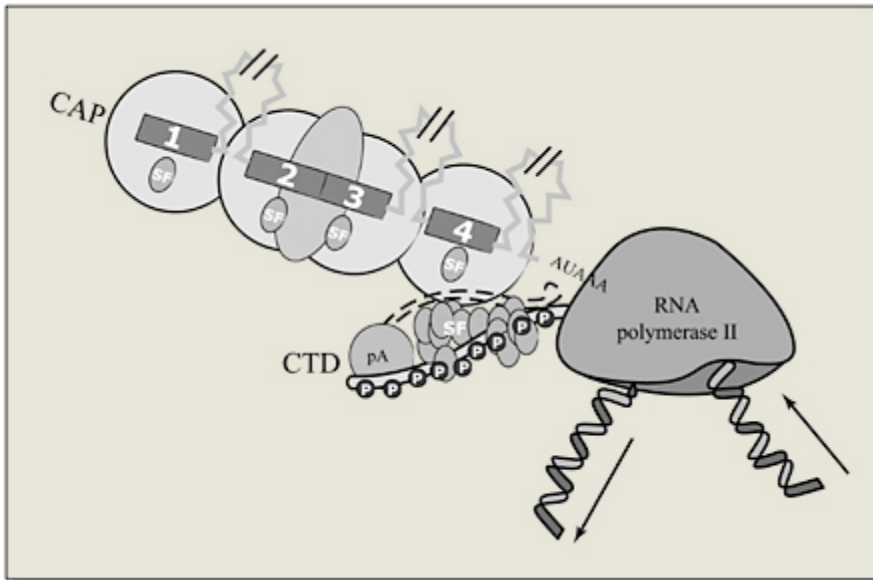
Types of biomodels.

Discrete, e.g. conversion stoichiometry

Rates/probabilities of interactions

Modules vs

“extensively coupled networks”



Maniatis & Reed Nature 416, 499 - 506 (2002)

Types of Systems Interaction Models

Quantum Electrodynamics	subatomic	
Quantum mechanics	electron clouds	
Molecular mechanics	spherical atoms	nm-fs
Master equations	stochastic single molecules	
Fokker-Planck approx.	stochastic	
Macroscopic rates ODE	Concentration & time (C,t)	
Flux Balance Optima	dC_{ik}/dt optimal steady state	
Thermodynamic models	$dC_{ik}/dt = 0$ k reversible reactions	
Steady State	$\sum dC_{ik}/dt = 0$ (sum k reactions)	
Metabolic Control Analysis	$d(dC_{ik}/dt)/dC_j$ (i = chem.species)	
Spatially inhomogenous	dC_i/dx	
Population dynamics	as above	km-yr



Increasing scope, decreasing resolution

Genetic Engineering & Darwinian Selection

Min = 0.1 kg

Teosinte



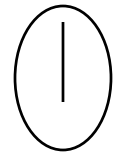
Max = 140 kg

Corn

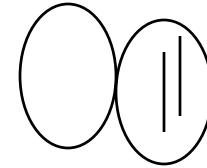
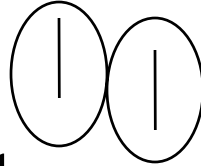


How to do single DNA molecule manipulations?

One DNA molecule per cell



Replicate to two DNAs.



Now segregate to two daughter cells

If totally random, **half** of the cells will have too many or too few.

What about human cells with 46 chromosomes (DNA molecules)?

Dosage & loss of heterozygosity & major sources of mutation in human populations and cancer.

For example, trisomy 21, a 1.5-fold dosage with enormous impact.

Most RNAs < 1 molecule per cell.

See Yeast RNA

25-mer array in

Wodicka, Lockhart, et al. (1997)

Nature Biotech 15:1359-67

(ref)

(http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&list_uids=9415887&dopt=Abstract)

Mean, variance, & linear correlation coefficient

Expectation E (rth moment) of random variables X for any distribution f(X)

First moment= Mean μ ; variance σ^2 and standard deviation σ

$$E(X^r) = \sum X^r f(X) \quad \mu = E(X) \quad \sigma^2 = E[(X-\mu)^2]$$

Pearson correlation coefficient $C = \text{cov}(X, Y) = E[(X-\mu_X)(Y-\mu_Y)] / (\sigma_X \sigma_Y)$

Independent X, Y implies $C = 0$,

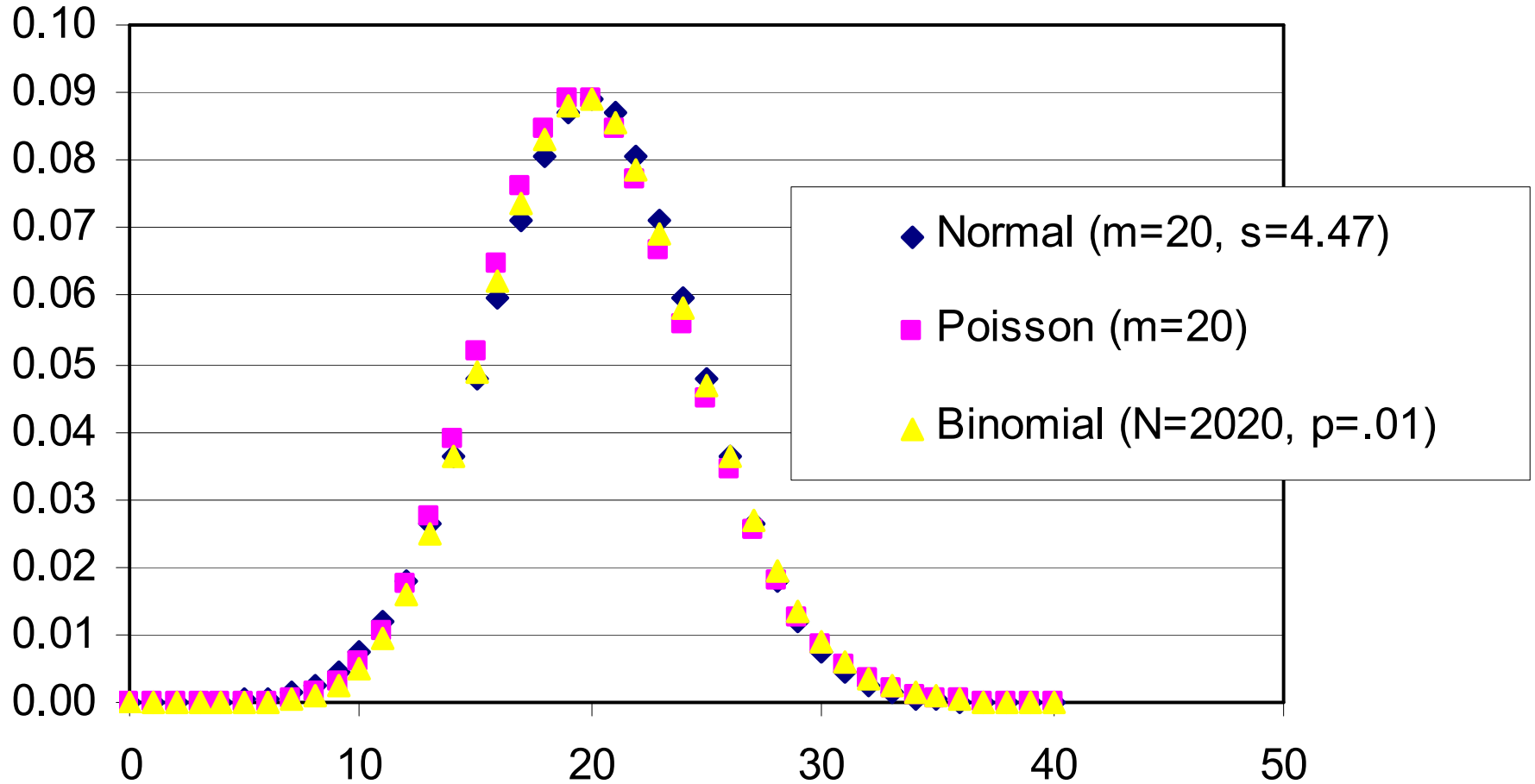
but $C = 0$ does not imply independent X, Y. (e.g. $Y = X^2$)

$P = \text{TDIST}(C * \sqrt{(N-2)/(1-C^2)})$ with dof= N-2 and two tails.

where N is the sample size.



Mutations happen



Binomial frequency distribution as a function of

$$X \in \{\text{int } 0 \dots n\}$$

p and q $0 \leq p \leq q \leq 1$ $q = 1 - p$ two types of object or event.

Factorials $0! = 1$ $n! = n(n-1)!$

Combinatorics ($C = \#$ subsets of size X are possible from a set of total size of n)

$$\frac{n!}{X!(n-X)!} = C(n, X)$$

$$B(X) = C(n, X) p^X q^{n-X} \quad \mu = np \quad \sigma^2 = npq$$

$$(p+q)^n = \sum B(X) = 1$$

$$B(X: 350, n: 700, p: 0.1) = 1.53148 \times 10^{-157}$$

=PDF[BinomialDistribution[700, 0.1], 350] Mathematica

~ = 0.00 =BINOMDIST(350,700,0.1,0) Excel

Poisson frequency distribution as a function of $X \in \{\text{int } 0 \dots \infty\}$

$$P(X) = P(X-1) \mu/X = \mu^x e^{-\mu} / X! \quad \sigma^2 = \mu$$

$$n \text{ large \& } p \text{ small} \rightarrow P(X) \cong B(X) \quad \mu = np$$

For example, estimating the expected number of positives in a given sized library of cDNAs, genomic clones, combinatorial chemistry, etc. $X = \#$ of hits.

$$\text{Zero hit term} = e^{-\mu}$$

Normal frequency distribution as a function of $X \in \{-\infty \dots \infty\}$

$$Z = (X - \mu) / \sigma$$

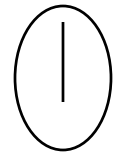
Normalized (standardized) variables

$$N(X) = \exp(-Z^2/2) / (2\pi\sigma)^{1/2}$$

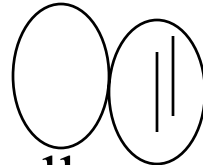
probability density function

$$npq \text{ large} \rightarrow N(X) \cong B(X)$$

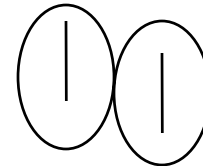
One DNA molecule per cell



Replicate to two DNAs.



Now segregate to two daughter cells



If totally random, half of the cells will have too many or too few.

What about human cells with 46 chromosomes (DNA molecules)?

Exactly 46 chromosomes (but any 46):

$$B(X) = C(n,x) p^x q^{n-x}$$

$$n=46*2; x=46; p=0.5$$

$$B(X) = 0.083$$

$$P(X) = \frac{\mu^x e^{-\mu}}{X!}$$

$$\mu=X=np=46, P(X)=0.058$$

But what about exactly
the correct 46?

$$0.5^{46} = 1.4 \times 10^{-14}$$

Might this select for non random segregation?

What are random numbers good for?

- Simulations.
- Permutation statistics.

Where do random numbers come from?

$$X \in \{0,1\}$$

perl -e "print rand(1);" 0.116790771484375
0.8798828125 0.692291259765625 0.1729736328125

excel: =RAND() 0.4854394999892640 0.6391685278993980
0.1009497853098360

f77: write(*,'(f29.15)') rand(1) 0.513854980468750
0.175720214843750 0.308624267578125

Mathematica: Random[Real, {0,1}] 0.7474293274369694
0.5081794113149011 0.02423389638451016

Where do random numbers come from really?

Monte Carlo.

Uniformly distributed random variates $X_i = \text{remainder}(aX_{i-1} / m)$

For example, $a = 7^5$ $m = 2^{31} - 1$

Given two X_j, X_k such uniform random variates,

Normally distributed random variates can be made

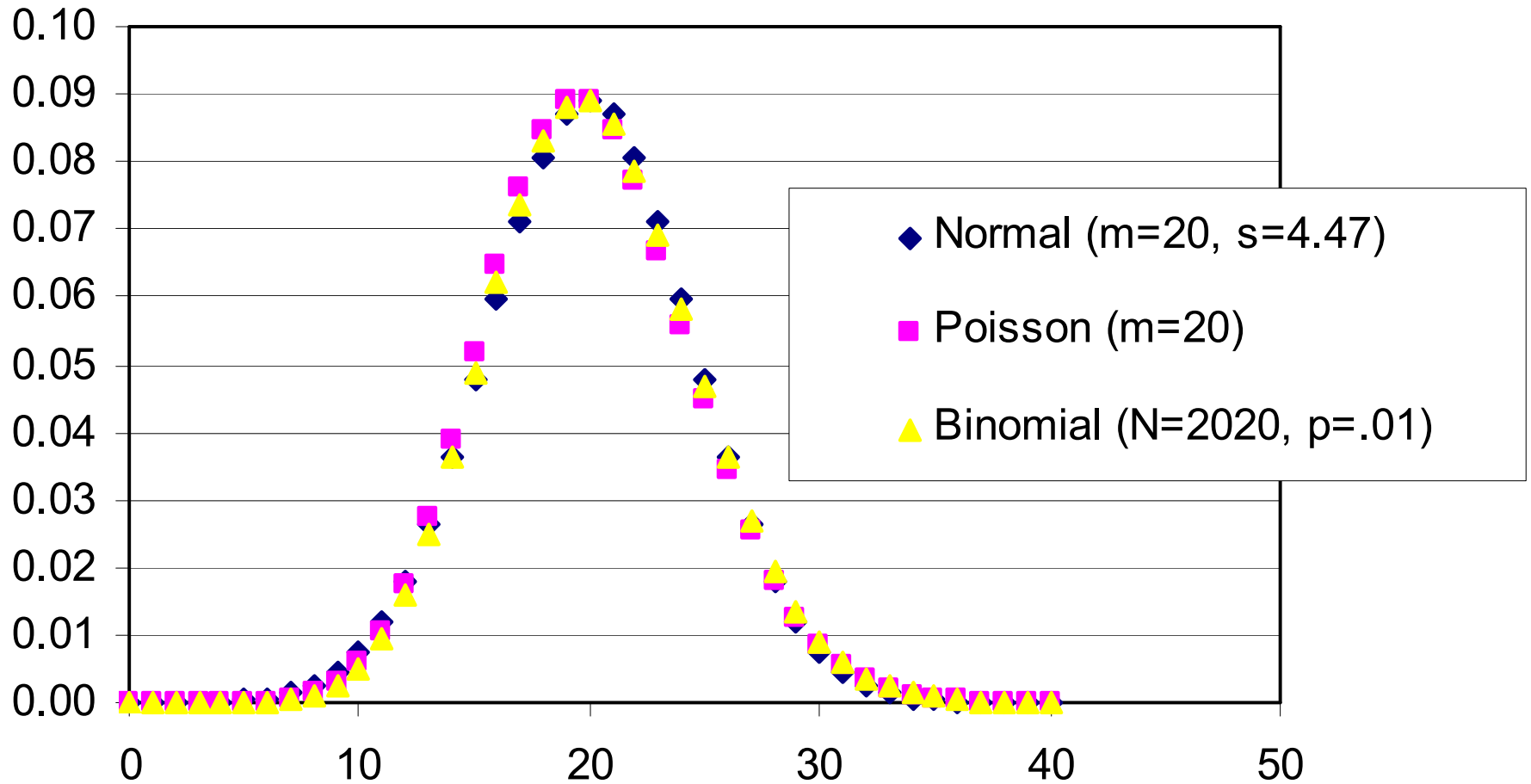
(with $\mu_X = 0$ $\sigma_X = 1$)

$X_i = \text{sqrt}(-2\log(X_j)) \cos(2\pi X_k)$ (NR, Press et al. p. 279-89)

(<http://www.nr.com/>) , (<http://lib-www.lanl.gov/numerical/bookcpdf/c7-1.pdf>).



Mutations happen



Intro 1: Summary

Life & computers : **Self-assembly** required

Discrete & continuous models

Minimal life & programs

Catalysis & Replication

Differential equations

Directed graphs & pedigrees

Mutation & the Single Molecules models

Bell curve statistics

Selection & optimality

Computation and Biology share a common obsession with strings of letters, which are translated into complex 3D and 4D structures. Evolution (biological, technical, and cultural) will probably continue to act via manipulation of symbols (A, C, G, T, 0 & 1, A-Z) plus "selection" at the highest "systems" levels. The power of these systems lies in complexity.

Simple representations of them (fractals, surgery, and drugs) may not be as fruitful as detailed programming of the symbols aided by hierarchical models and highly-parallel testing. Local decisions no longer stay local. Examples are the Internet, computer viruses, genetically modified organisms (GMOs), replicating nanotechnology, bioterrorism, global warming, and biological species transport. Information (& education) is becoming increasingly easy to spread (and hard to control). We are on the verge of begin able to collect data on almost any system at costs of terabytes-per-dollar.

The world is manipulating increasingly complex systems, many at steeper-than-exponential rates. Much of this is happening without much modeling. Some people predict a "singularity" in our lifetime or at least the creation of systems more intelligent (and/or more proliferative) than we are (possibly as little as 100 Teraflops/terabytes). We need to not only teach our students how to cope with this, but start thinking about how to teach these "intelligent" systems as if they were students. As integrated circuits reach their limit soon, the next generation of computers may be based on quantum computing and/or biologically inspired. We need to be able to teach our students about this revolution, and via the Internet teach anyone else listening.