

Exploration Session Week 8: Computational Biology

Melissa Winstanley: mwinst@cs.washington.edu

(based on slides by Martin Tompa, Luca Cardelli)

Exploring DNA Sequences

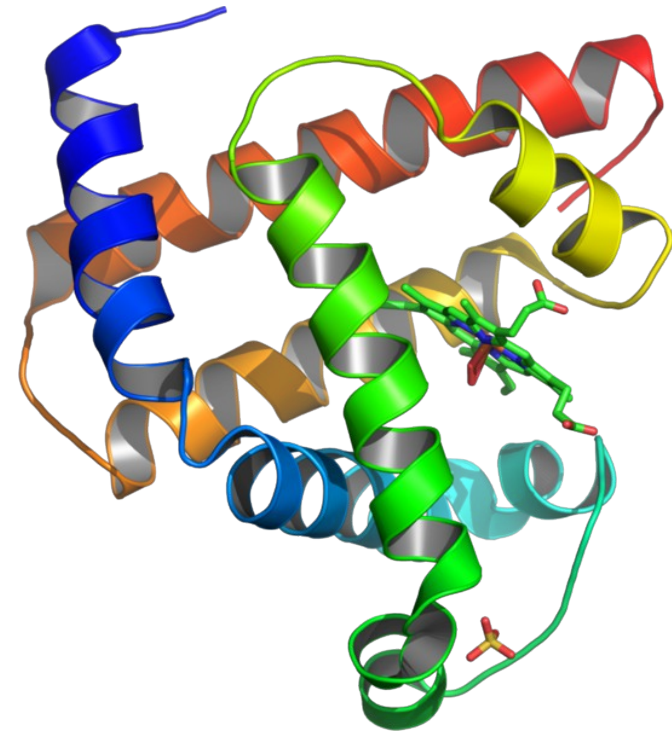
Overview of DNA

- Instructions for cellular function
 - Building proteins
- Composed of *nucleotides*
 - Adenine, thymine, cytosine, guanine
 - A pairs with T, C pairs with G
- Double-stranded: forms a double helix
 - Strands have an *orientation*
 - Pairing of antiparallel strands
- Huge amount of DNA
 - 3 billion base pairs, 2m long in a cell
 - 133 AU long in human
 - 20 million light years long in human population



Overview of Proteins

- Workhorses of cells
- Composed of sequence of *amino acids*
 - 20 to 5000 amino acids in a protein
- 20 possible amino acids
- Proteins fold into complex 3D shapes
 - Fold-It
- Information to make proteins encoded in DNA
 - *Codon*: 3 base pairs
 - Ex. CTA → leucine
 - *Gene*: sequence of DNA for 1 protein



Overall Goals

- Overall
 - Identify key molecules in organisms
 - Identify interactions among molecules
- Computational focus: sequence analysis
 - Identify genes
 - Determine gene function (what protein is produced?)
 - Identify proteins involved in gene expression
 - Identify key functional regions
- Why do we care?
 - Determining function of a new sequence
 - Genetic diseases
 - Evolution

String Alignment

- How to judge how well two strings are aligned?

acbcdb a c - - b c d b
cadbd - c a d b - d -

- Each dash represents an inserted space
- Assign +2 to every exact match, -1 to every mismatch

$$3 * 2 + 5 * (-1) = 1$$

- Higher score indicates a greater match between the strings

BLAST Algorithm

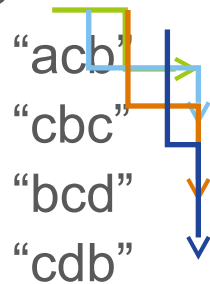
- “Basic Local Alignment Search Tool”
- For comparing biological sequence information
 - Amino acid sequences (proteins) or nucleotide sequences (DNA)
- Inputs
 - A query sequence Q
 - A database D of sequences
- Output
 - Sequences from D that match Q above a certain threshold
- Usefulness
 - Unknown gene in a mouse, so query the human gene database to see if a similar gene exists in humans

BLAST ctd

- Make k-letter subsequences from Q

Ex. $k = 3$:

“acbcd”



- Usually $k = 28$ for DNA, $k = 3$ for proteins

BLAST ctd

- For each subsequence w , find matching subsequences
 - Only consider a matching subsequence if its alignment score is greater than some threshold
 - $\text{Alignment}(\text{seq}) \geq T$

Ex. $T = 2$, $w = \text{"TCG"}$

$\text{seq} = \text{"TCA"} \rightarrow \text{Alignment} = 2 * 2 + 1 * (-1) = 3$
Considered

$\text{seq} = \text{"ACT"} \rightarrow \text{Alignment} = 2 * 1 + 2 * (-1) = 0$
Not considered

BLAST ctd

- Scan the database for exact matches with the high scoring subsequences
- Take each exact match and extend in either direction (no gaps)
 - Until the score decreases below a “dropoff”
 - Forms a “high-scoring segment pair” (HSP)
- Only save match extensions above a certain score threshold S

Query seq:	A	C	T	C	G	G	C
Database:	G	C	T	C	A	G	T
Score	-1	2	2	2	-1	2	-1

Exact match



$$\text{HSP: score} = 2 + 2 + 2 - 1 + 2 = 7$$

BLAST ctd

- For each HSP, do a gapped extension (spaces possible)
- Output each extension that has probability of randomly occurring below a pre-set threshold x

More Complicated Analysis

- Multiple sequence alignment
- Different ways to score subsequences
- Considering context around a sequence
- Predicting 3D structures of proteins

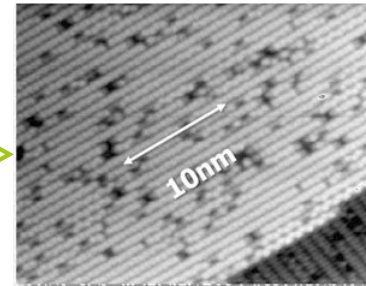
Programming Molecules

Getting Smaller

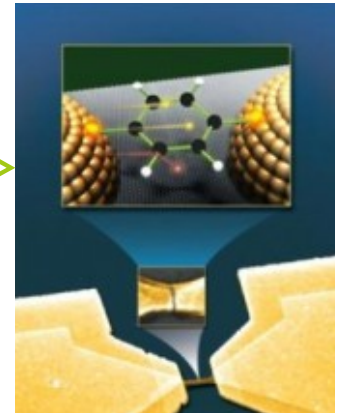
□ First transistor



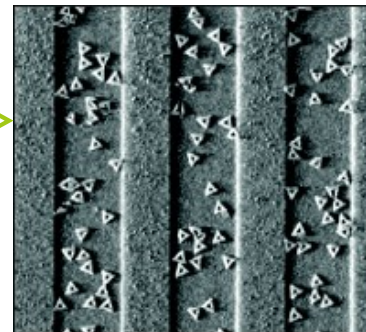
□ 25nm NAND flash



□ Single molecule transistor



□ Molecules on a chip



□ ~10 Moore's Law cycles left

<http://upload.wikimedia.org/wikipedia/commons/thumb/b/bf/Replica-of-first-transistor.jpg/200px-Replica-of-first-transistor.jpg>

<http://www.blogcdn.com/www.engadget.com/media/2010/01/01-30-10intelflash.jpg>

http://www.wired.com/images_blogs/gadgetlab/2009/12/molecular-transistor-264x300.jpg

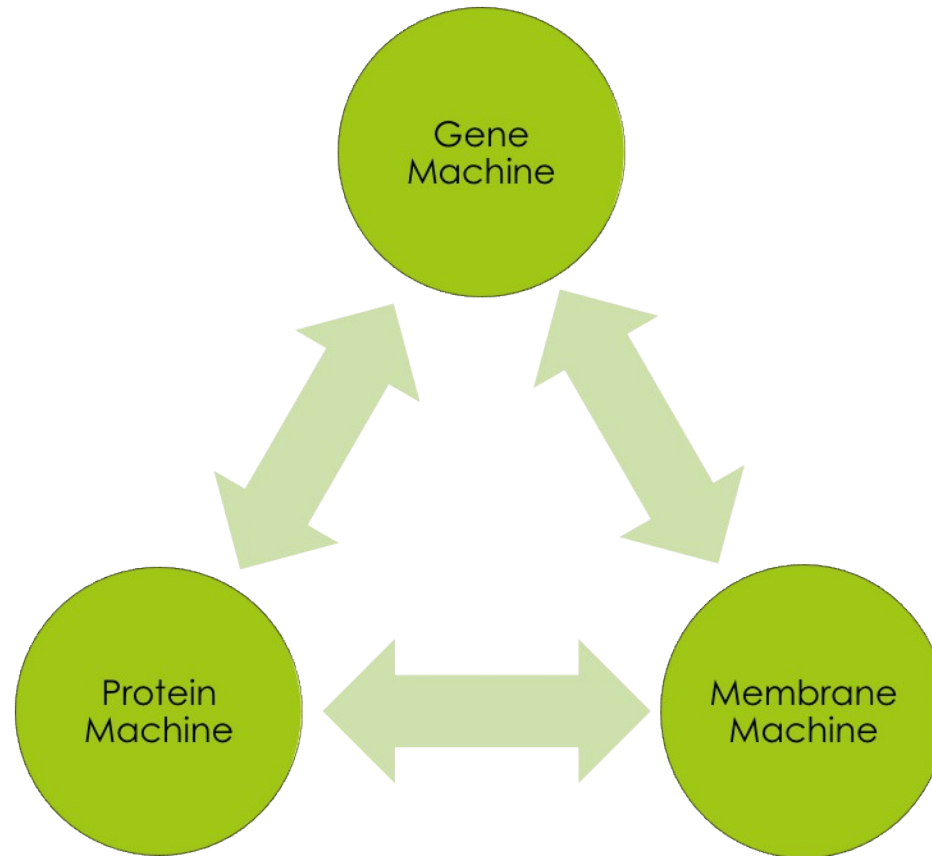
http://www.internetnews.com/img/2009/08/ibm_dna_chips.jpg

Building Smaller

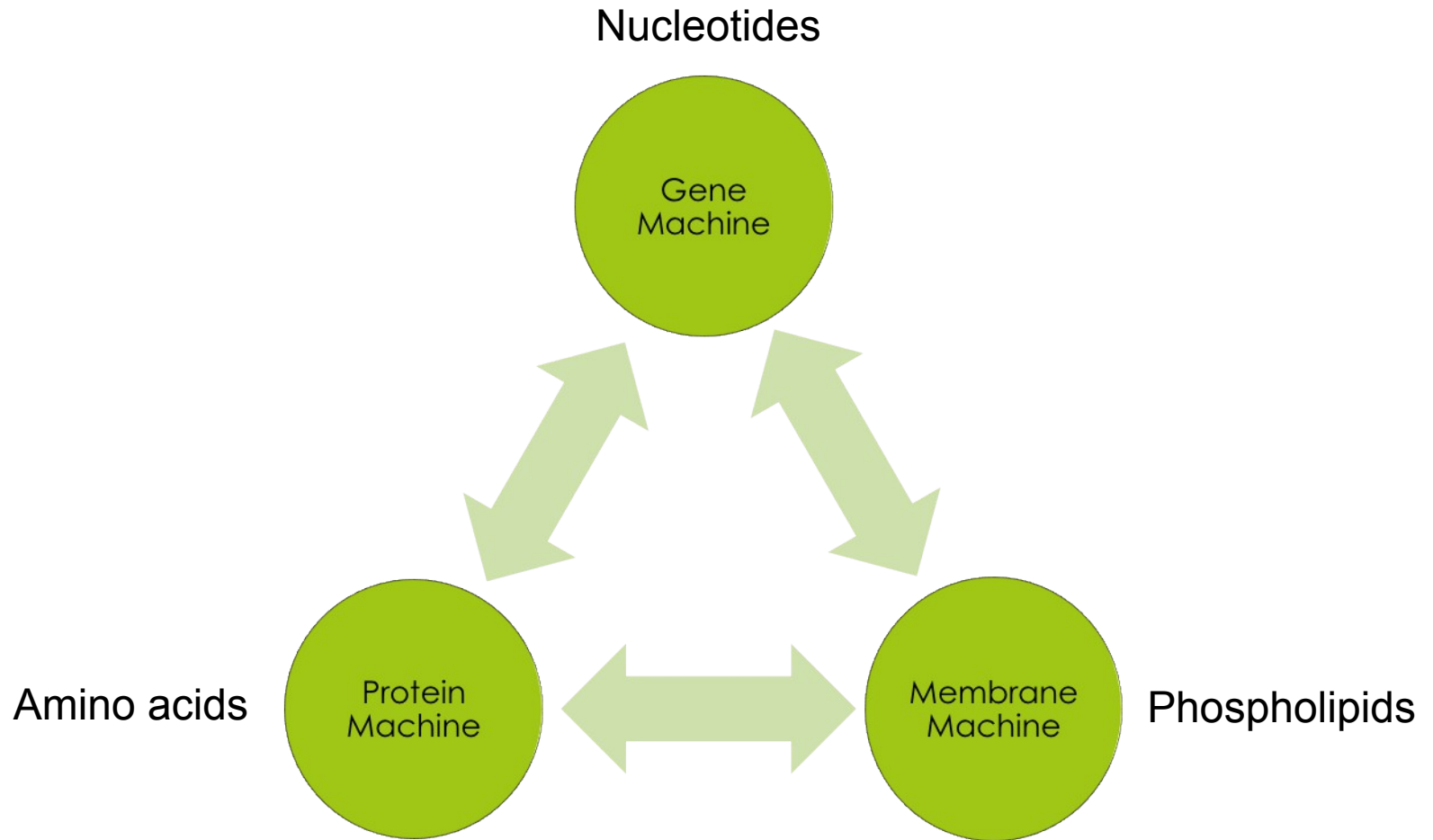
- How to build things smaller than your tools?
 - You can't
 - Solution: self-assembly
 - Molecular IKEA
 - Dear IKEA, please send me a chest of drawers that assembles itself.
 - At a molecular scale, many such materials exist
 - Proteins, DNA/RNA, membranes
 - <http://youtu.be/0N09BIEzDII>



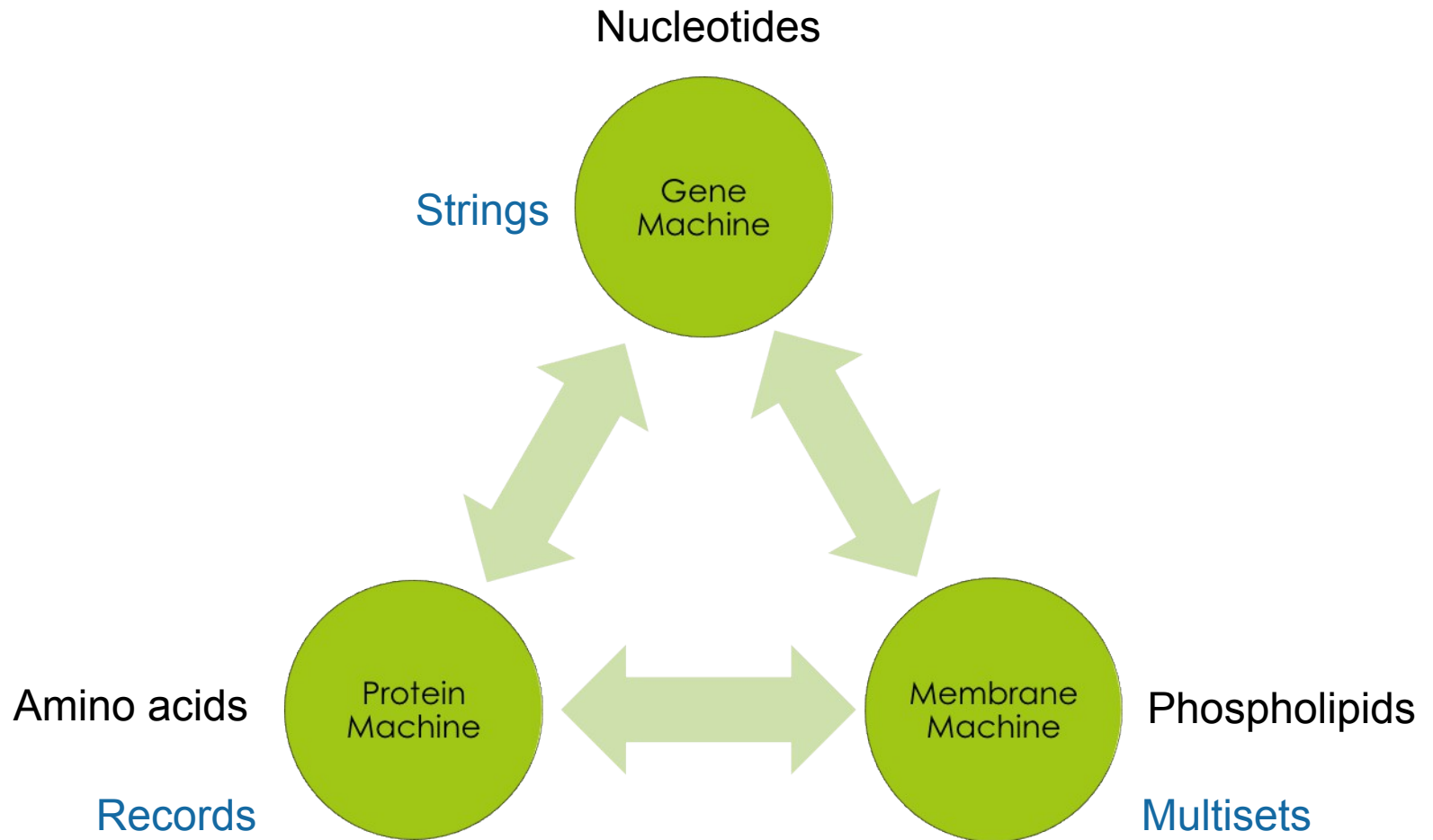
Machines in Biochemistry



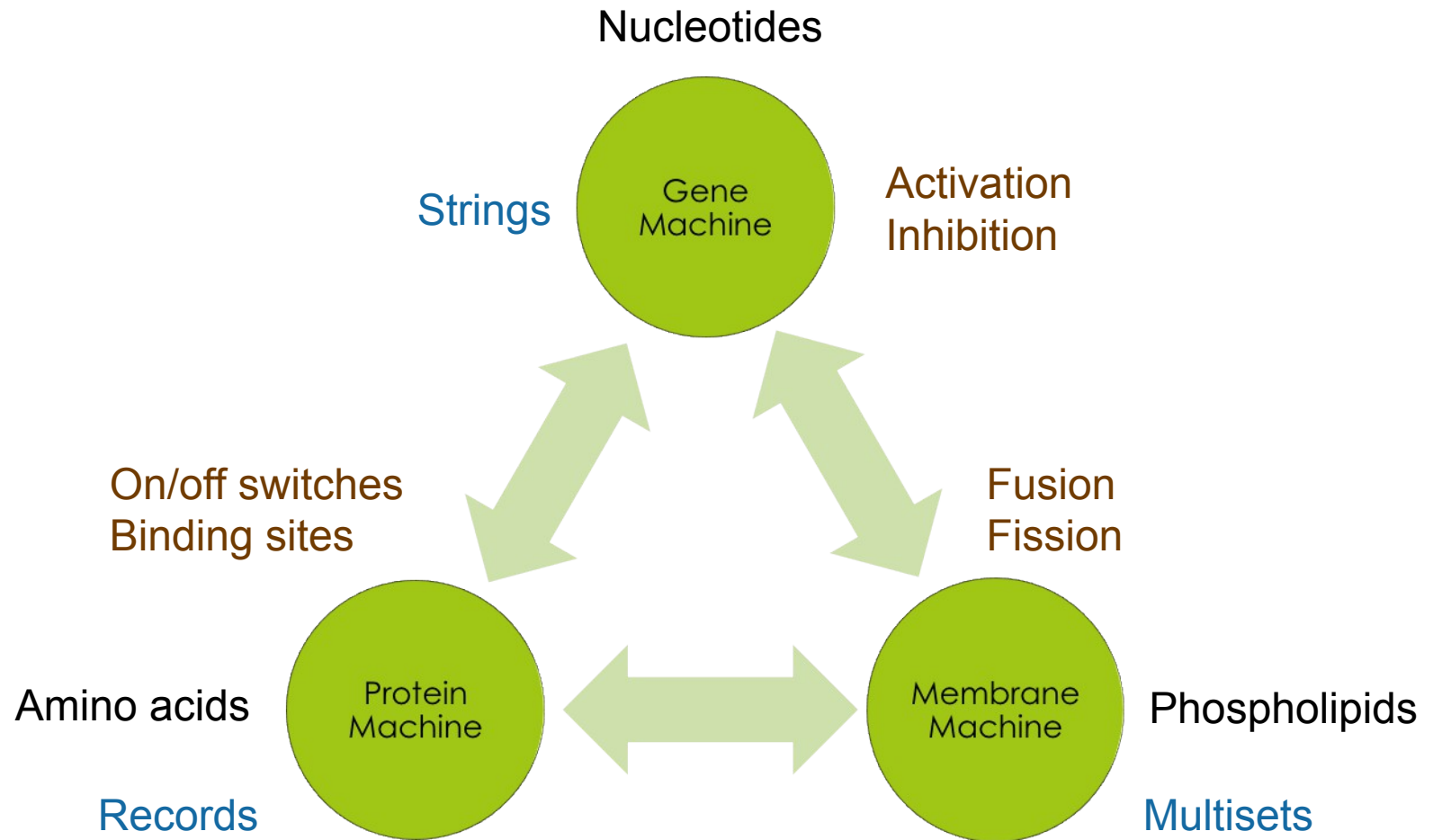
Machines in Biochemistry



Machines in Biochemistry



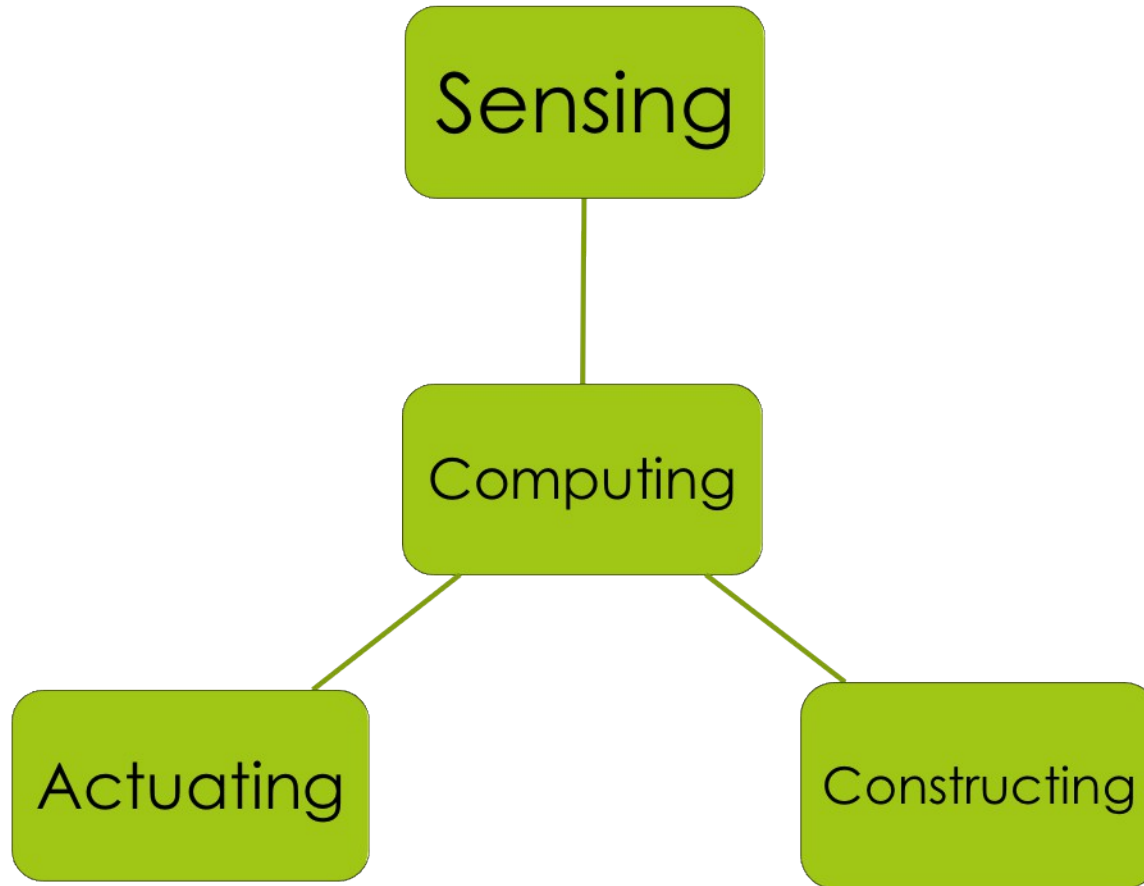
Machines in Biochemistry



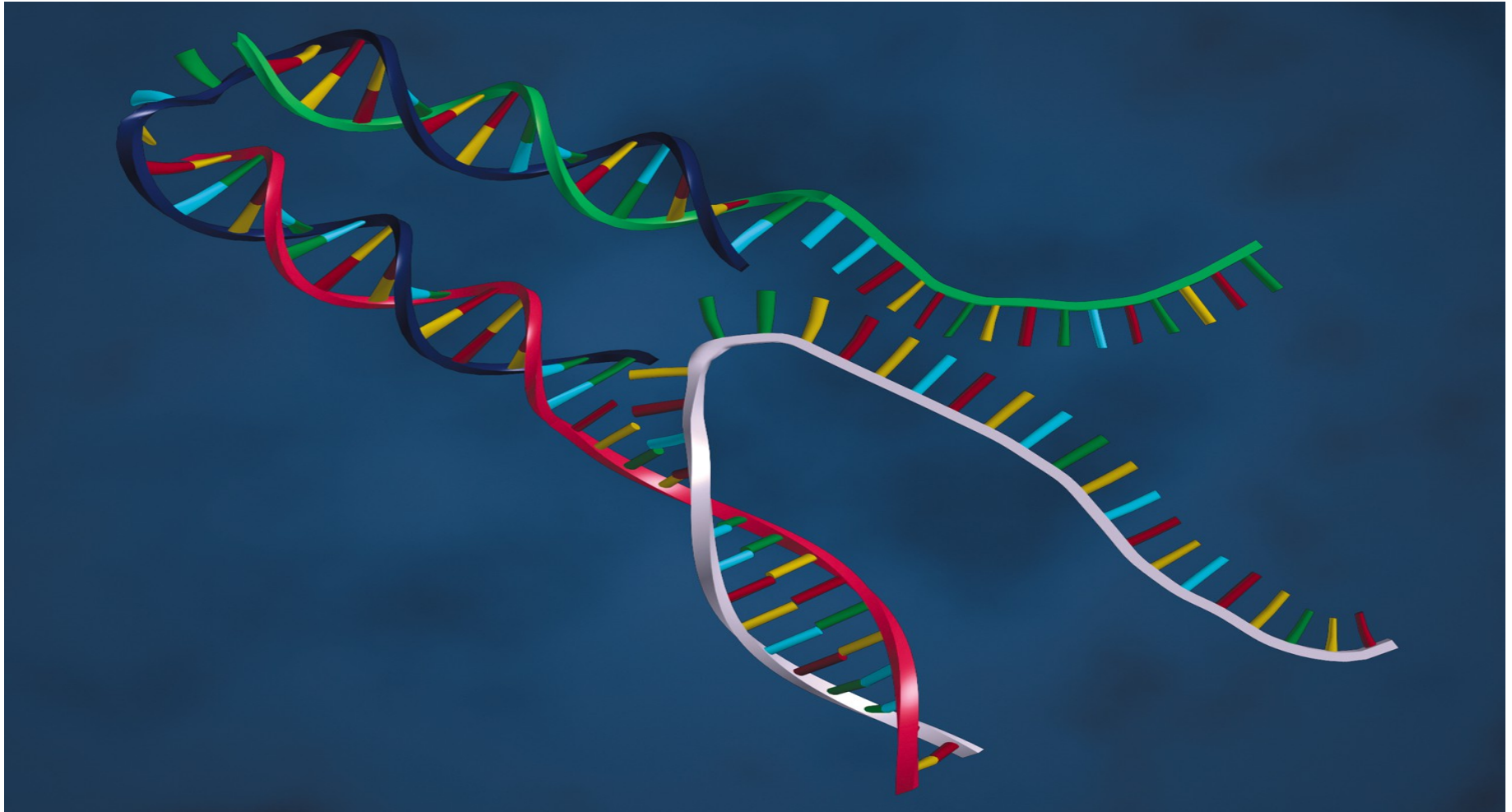
How do we form a “language”?

- Chemical reactions
 - $A + C \rightarrow B + D$
 - Instructions in a “program”
- Problem: combinatorial explosion
 - SO MANY chemical reactions in a cell
- Model reactions as automata – machines that perform a task
- Problem: chemistry is not an executable language
 - Dear Chemist, please execute this arbitrary reaction.

Controlling Systems on a Nanoscale

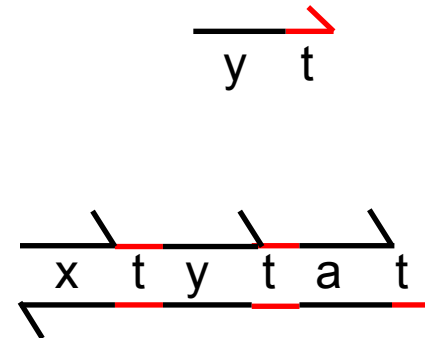
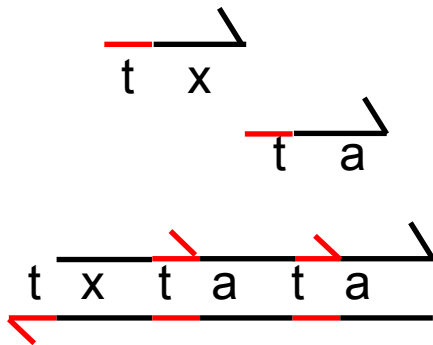


DNA Tweezers



One Approach to Autonomous Computing

- Goal: precisely control organization and dynamics of matter and information at the molecular level
 - Uses DNA, but use is accidental
 - No genes involved



“Gates” and “transducers”

Molecular programming workflow

- First figure out what gates you want to use and signals you want to send
- Signals + gates → structures of DNA
- Structures → sequences of DNA (**NUPACK**)
- Sequences → DNA synthesis (**IDT**)
- DNA synthesis → mail
- Receipt of DNA →_{water} execution
- Florescence is your “print” statement