**Introduction to Bioinformatics**
**Jérôme Waldispühl**
**School of Computer Science, McGill university**
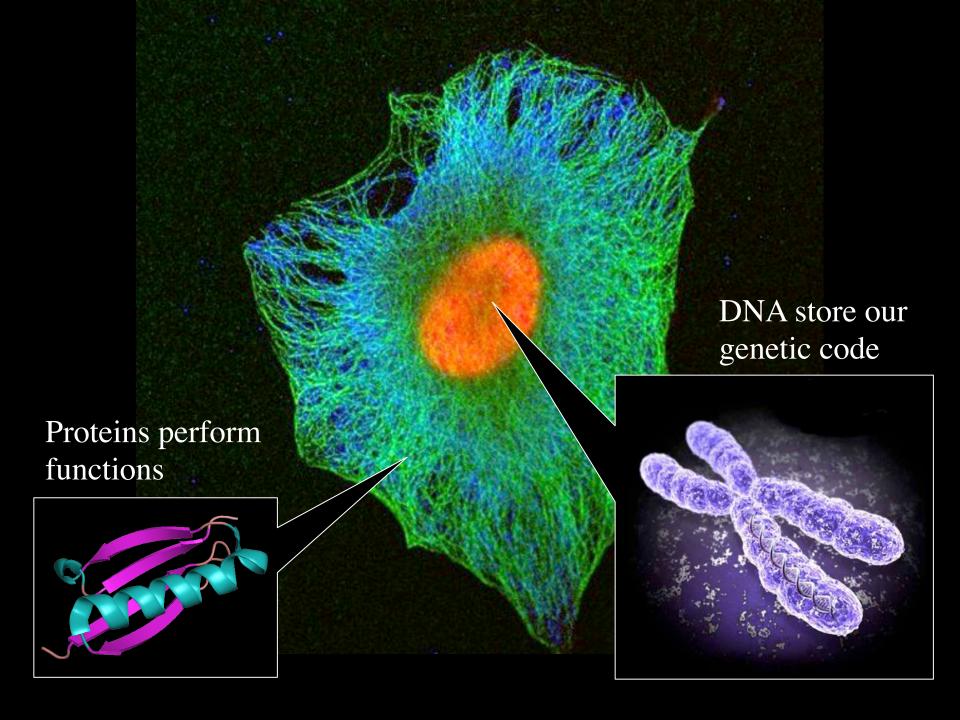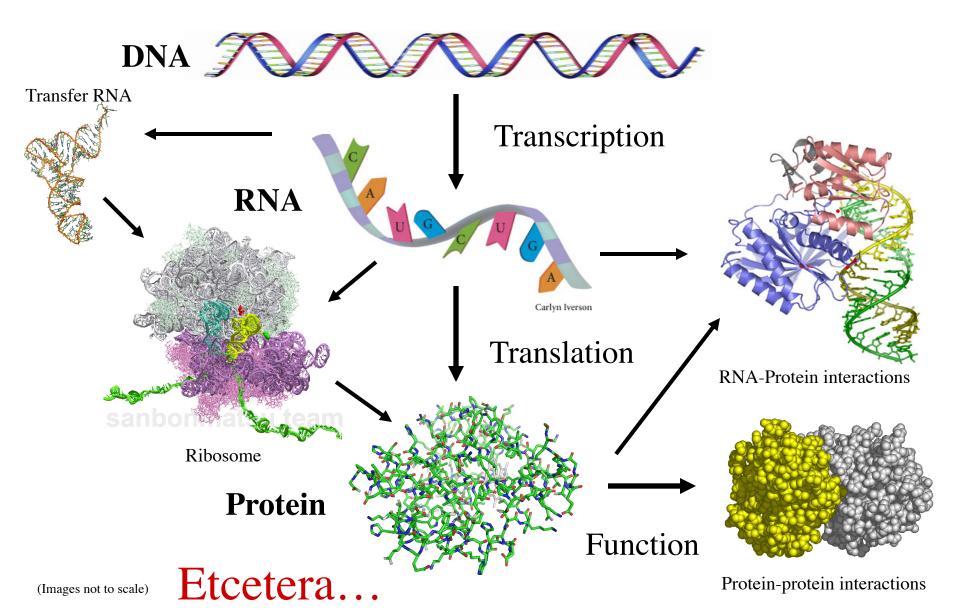(includes slides from Mathieu Blanchette)

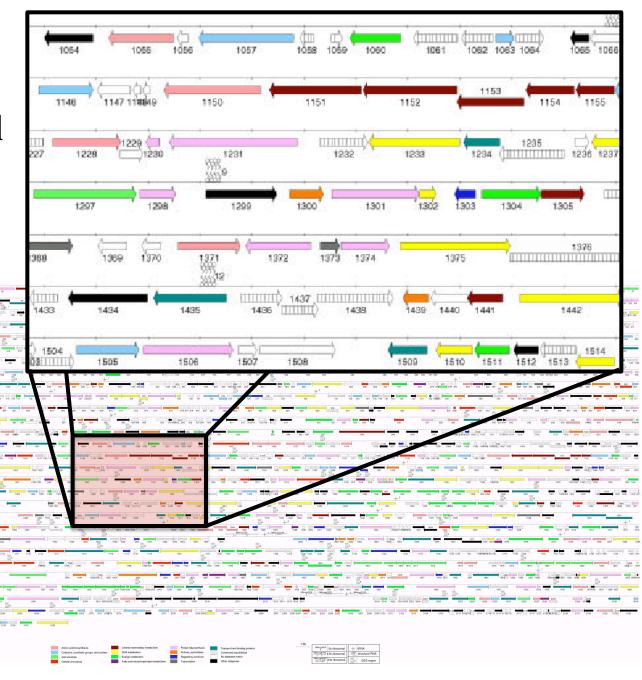How do cells work?

# Definition

**Bioinformatics** is an interdisciplinary field that develops and improves on methods for storing, retrieving, organizing and analyzing biological data. A major activity in bioinformatics is to develop software tools to generate useful biological knowledge.
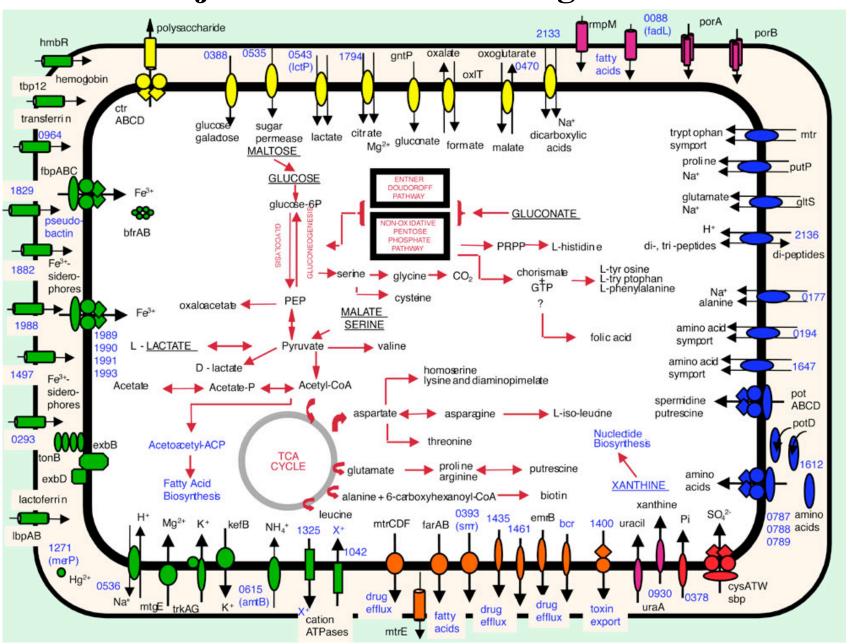
(from Wikipedia)

Proteins perform functions

DNA store our genetic code

# Central dogma of biology

**DNA**

Transfer RNA

Transcription

**RNA**

C

A

U G

C

U

G

A

Carlyn Iverson

Translation

RNA-Protein interactions

Ribosome

sanbonmatsu team

**Protein**

Function

Protein-protein interactions

(Images not to scale)

Etcetera…

**Objective 1: Identify functional regions of the genome**

# Objective 2: Cell Modeling

# The 3 domains of bioinformatics
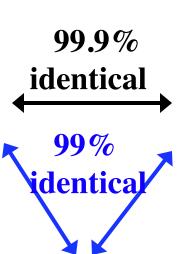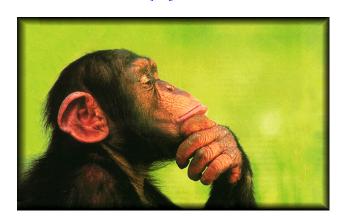
## Structural Bioinformatics



## Genomics

## System Biology

# Genomic

# Genomes

- Human genome: $\{A,C,G,T\}^{3 \times 10^9}$
- Each of your $10^{14}$ cells has two copies

**99.9% identical**

**99% identical**



one helical turn = 3.4 nm

Sugar-phosphate backbone

Base

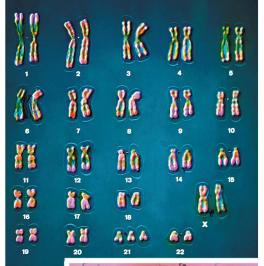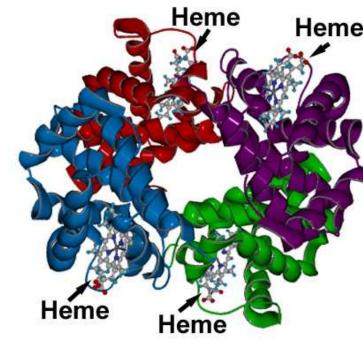Hydrogen bonds

# Roles of the genome



- Genome is a blue print for a cell
- Describes *how* to build proteins
  - 25,000 genes --> 25,000 proteins (+variations)
  - Each protein has its biochemical function
- Describes *when* to build each protein
  - Under which situations should a gene be expressed?
- Proteins allow:
  - Cell administration and maintenance
  - Reaction to stimuli
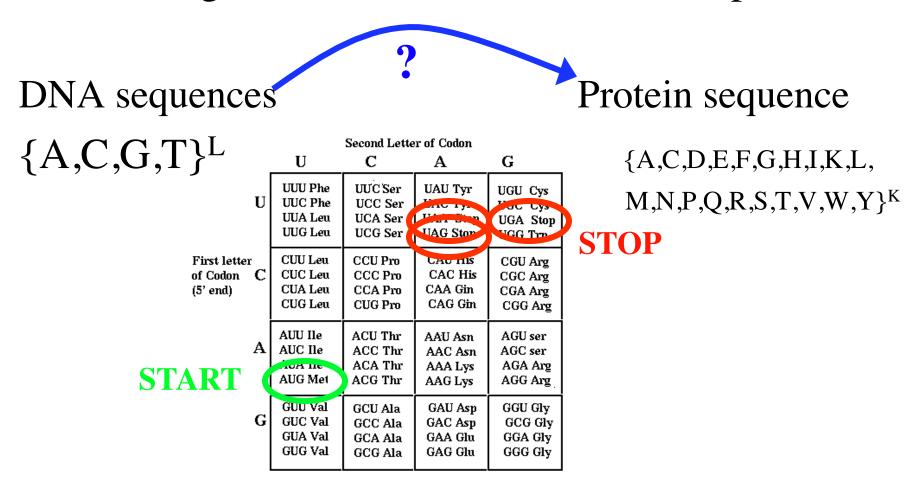  - Protocols for communication between cells

# Roles of the genome: development
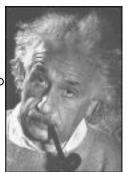
# Content of the genome - Genes

- Gene: region of DNA that encodes one protein

DNA sequences

$\{A,C,G,T\}^L$

Protein sequence

$\{A,C,D,E,F,G,H,I,K,L,$
$M,N,P,Q,R,S,T,V,W,Y\}^K$

**Second Letter of Codon**

|  | U | C | A | G |
|---|---|---|---|---|
| **U** | UUU Phe<br>UUC Phe<br>UUA Leu<br>UUG Leu | UUC Ser<br>UCC Ser<br>UCA Ser<br>UCG Ser | UAU Tyr<br>UAC Tyr<br>UAA Stop<br>UAG Stop | UGU Cys<br>UGC Cys<br>UGA Stop<br>UGG Trp |
| **C** | CUU Leu<br>CUC Leu<br>CUA Leu<br>CUG Leu | CCU Pro<br>CCC Pro<br>CCA Pro<br>CUG Pro | CAU His<br>CAC His<br>CAA Gin<br>CAG Gin | CGU Arg<br>CGC Arg<br>CGA Arg<br>CGG Arg |
| **A** | AUU Ile<br>AUC Ile<br>AUA Ile<br>AUG Met | ACU Thr<br>ACC Thr<br>ACA Thr<br>ACG Thr | AAU Asn<br>AAC Asn<br>AAA Lys<br>AAG Lys | AGU ser<br>AGC ser<br>AGA Arg<br>AGG Arg |
| **G** | GUU Val<br>GUC Val<br>GUA Val<br>GUG Val | GCU Ala<br>GCC Ala<br>GCA Ala<br>GCG Ala | GAU Asp<br>GAC Asp<br>GAA Glu<br>GAG Glu | GGU Gly<br>GCG Gly<br>GGA Gly<br>GGG Gly |

First letter of Codon (5' end)
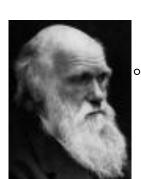
**START**

**STOP**

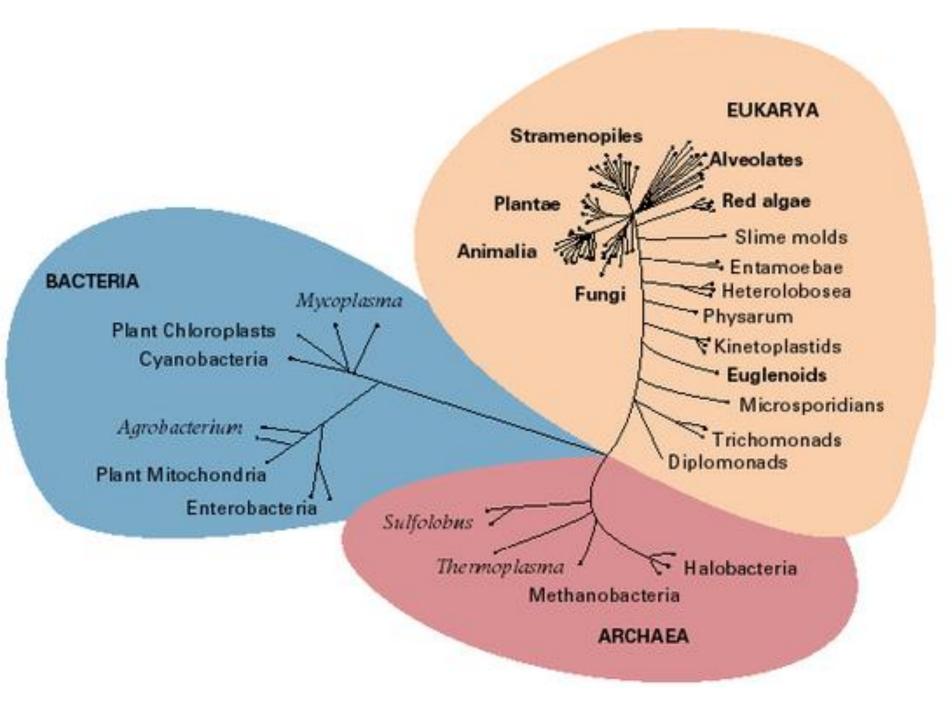# The Programmer - Evolution

- Design principles:
  - Random modifications (variation)
  - Survival of the fittest (natural selection)
- 3 Billion years of evolution
- Today's species are the current solution of the fitness optimization problem
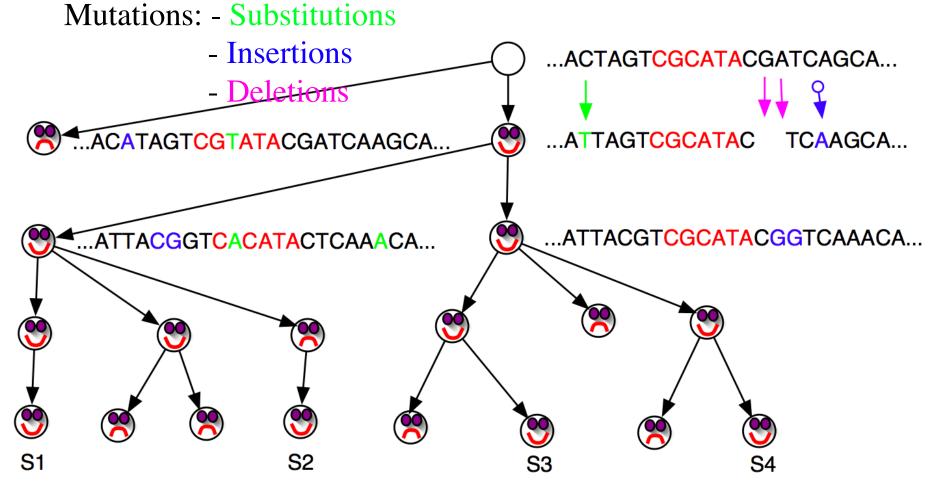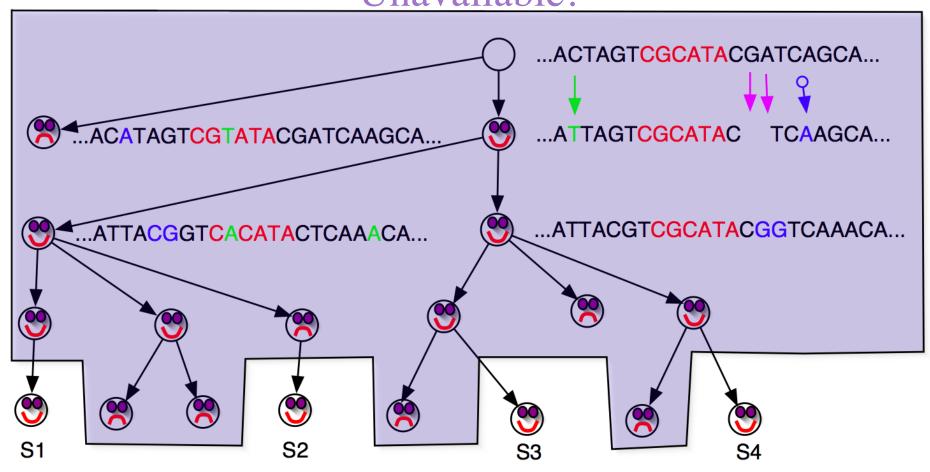
God doesn't play dice

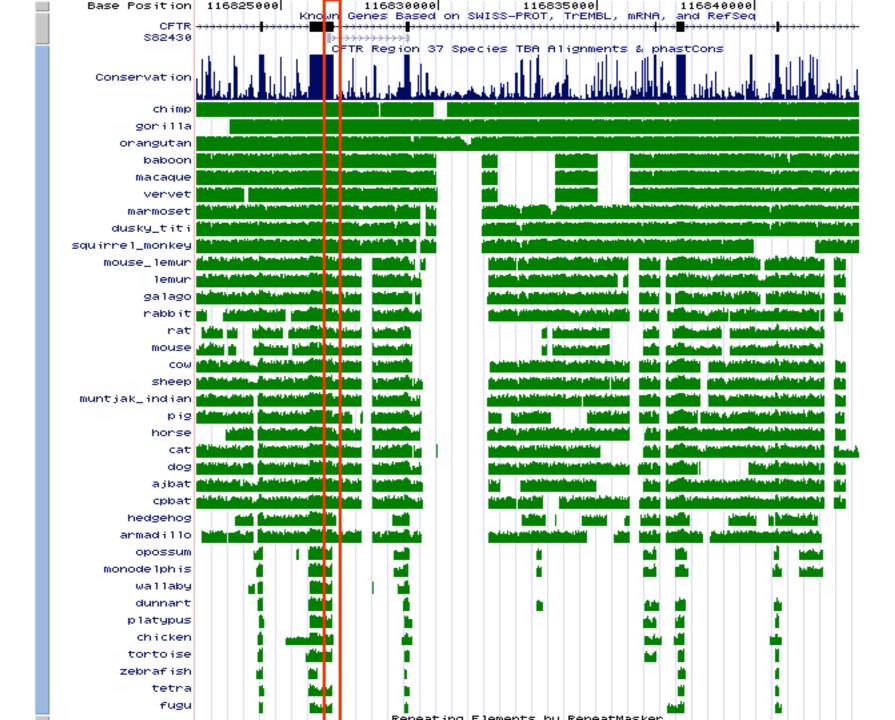Yes he does!

# Central dogma of comparative genomics



Mutations: - Substitutions
- Insertions
- Deletions

...ACTAGTCGCATACGATCAGCA...

...ACATAGTCGTATACGATCAAGCA...

...ATTAGTCGCATAC  TCAAGCA...

...ATTACGGTCACATACTCAAACA...

...ATTACGTCGCATACGGTCAAACA...

S1    ...GTTACGGTCACATACTGAAACA...
S2    ...GTTATGGTCACATACTGAAACTGA...
S3    ...ATTACTCGCATACGGTCTAACA
S4    ...ATTTACTCGCATACGGTCTAGCACT

# Unavailable!



S1:  GTTACGGTCACATACTGAAACA
S2:  GTTATGGTCACATACTGAAACTGA
S3:  ATTACTCGCATACGGTCTAACA
S4:  ATTTACTCGCATACGGTCTAGCAC

Base Position | 116821820 | 116821830 | 116821840 | 116821850 | 116821860 | 116821870

---> GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTAATTATTGGTCTAGCAAGCATTTG

G  Q  R  A  R  I  S  L  A  R  *  I  T  N  Y  W  S  S  K  H  L
V  N  E  Q  E  F  L  *  Q  G  E  *  L  I  I  G  L  A  S  I  C
R  S  T  S  K  N  F  F  S  K  V  N  N  *  L  L  V  *  Q  A  F

Known Genes Based on SWISS-PROT, TrEMBL, mRNA, and RefSeq

CFTR  G  Q  R  A  R  I  S  L  A  R →→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→→

CFTR Region 37 Species TBA Alignments & phastCons

Gaps

human         GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTAATTATTGGTCTAGCAAGCATTTG
chimp         GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTAATTATTGGTCTAGCAAGCATTTG
gorilla       GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTAATTATTGGTCTAGCAAGCATTTG
orangutan     GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTAATTATTGGTCTAGCAAGCATTTG
baboon        GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAA----TAATTATTGGTCTAACAAGCATTTG
macaque       GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAA----TAATTATTGGTCTAACAAGCATTTG
vervet        GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAA----TAATTATTAGTCTAACAAGCATTTA
marmoset      GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAGTAACTGATCATTGGTCTAGCAAGCATTTG
dusky_titi    GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACTGATTATTGGTCTAGCAAGCATTTG
squirrel_monkey GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAGTAACTGATTATTGGTCTAGCAAGCATTTG
mouse_lemur   GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAATGGATTATTGGTCCAGTGAGCATTTG
lemur         GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATAACGGATTATTGGTCCAGTGAGCATTTG
galago        GGTCAGCGAGCAAGAATCTCTTTAGCAAGGTGAATAATGGGGTATGACTCCAGTGGGCGTTTG
rabbit        GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAGTATCTGATTGGTCTACCAAGCATTTG
rat           GGTCAACGTGCAAGAATTTCTTTAGCAAGGTAAACGTTCAACTGTTGGTTTGCTGAGAACTTG
mouse         GGTCAGCGTGCAAGGATTTCTTTAGCAAGGTAAATATTTAACTGTTGGTCTTGTGAGCACTTG
cow           GGTCAGCGAGCGAGAATTTCTTTAGCAAGGTGAATATCTGCCTATGGGTTCAGCAAGCATTTG
sheep         GGTCAACGAGCAAGAATTTCTTTAGCAAGGTAAATATCTGCTTATTGGTCCAGCAAGCATTTG
muntjak_indian GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATATTGGCTTATTGGTCCAGCAAGCATTTG
pig           GGTCAACGAGCAAGAATTTCTTTAGCAAGGTGAATCACTTTTATTGGTCCAGCAAGCGTTTG
horse         GGTCAGCGAGCAAGGATTTCTTTAGCAAGGTGAATAGCTGATTATTGGTCAGTGAGCCTTTG
cat           GGTCAGCGAGCAAGAATTTCCTTAGCAAGGTGAATCTGATTATTGGTCAGTGAGCTTTTG
dog           GGGCAAAGAGCAAGAATTTCCCTAGCAAGGTGAATATCCGACGATTGGTTGGCGAGCATTGG
ajbat         GGGCAGCGAGCAAGAATTTCTTTAGCAAGGTGAATATCTGATTATTGGTCCAGTTAACATTTG
cpbat         GGGCAGCGAGCGAGAATTTGTTTAGCAAGGTGAATATCTGATTATTGGT-CAGTGAGCATTTG
hedgehog      GGTCAACGAGCAAGAATTTCATTAGCAAGGTGAATAT----------------AGCATTTG
armadillo     GGTCAACGAGCAAGAATTTCTTTAGCAAGGT-ATTATCTGACTATTGGACCAGTTGACACTTG
opossum       GGTCAACGAGCAAGAATTTCTTTAGCAAGGTAATATTTTGGTATTTATTCAGTTGAAATTTG
monodelphis   GGTCAACGAGCAAGAATTTCTTTAGCAAGGTAATATTTTGGTATTTGTTCAGTTAAAATTTG
wallaby       GGTCAACGAGCAAGAATTTCTTTAGCAAGGTAATATTTTGGTATTTGTTCAGTTGAAATGTG
dunnart       GGTCAACGAGCAAGAATTTCCTTAGCAAGGTAATATTTTGATATTTGTTCAGCTAAAATTTG
platypus      GGTCAGCGGGCCAGAATTTCATTAGCCAGGTGAGTA-----------TTTCAGGTGGCGTTTG
chicken       GGCCAGCGAGCACGAATCTCACTAGCGAGGTGAGCATTTTGCTATCT----------ATTT
tortoise      GGTCAACGGGCTAGAATCTCACTAGCTAGGTGAATATTTATCATCT----------ATGG
zebrafish     GGTCAGAAGGCACGCGTGGCTCTGGCCAGGT--ATG-TCACACACTT----TTTCACAGCTTC
tetra         GGTCAGAGGGCACGTCTGGGTTTGGCCAGGT--ACT-TCTCTCACAC----CTTCAGCACACC
fugu          GGTCAAAGGGCACGCCTGGGTTTGGCCAGGT--ACT-TCTGCCACAC----TTCAGCGCACC

Conservation
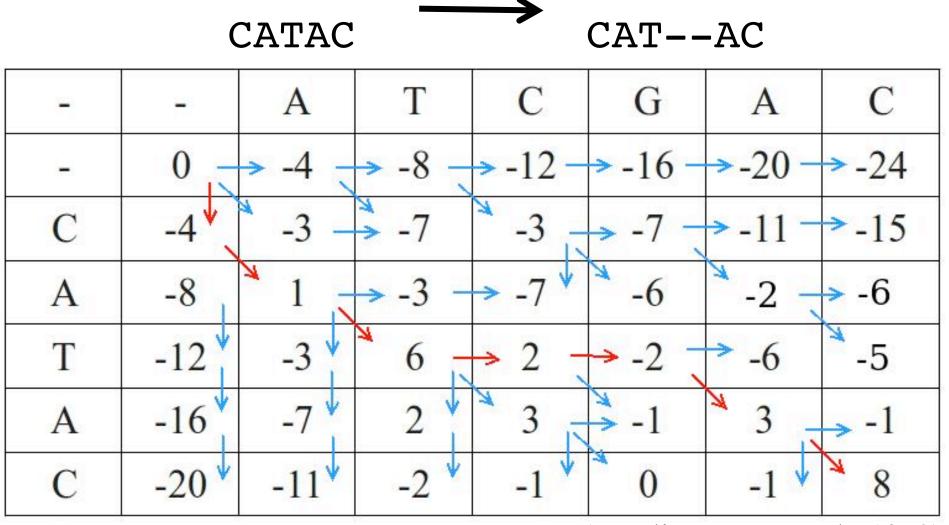
# Algorithms to compare and align genomes

Input

```
ATCGAC
CATAC
```

→

Output

```
-ATCGAC
CAT--AC
```

| - | - | A | T | C | G | A | C |
|---|---|---|---|---|---|---|---|
| - | 0 | -4 | -8 | -12 | -16 | -20 | -24 |
| C | -4 | -3 | -7 | -3 | -7 | -11 | -15 |
| A | -8 | 1 | -3 | -7 | -6 | -2 | -6 |
| T | -12 | -3 | 6 | 2 | -2 | -6 | -5 |
| A | -16 | -7 | 2 | 3 | -1 | 3 | -1 |
| C | -20 | -11 | -2 | -1 | 0 | -1 | 8 |

(Needleman-Wunsch, 1970)

# Mammalian evolution

- Rapid radiation ~75 Myrs ago

  ⬇

- Many nearly independent phyla

- Many "noisy" copies of ancestor

  ⬇

- **Accurate reconstruction of ancestors may be feasible**



Platypus — Monotremata
Opossum — Marsupialia
Tenrec
Elephant Shrew — Afrotheria
Hyrax
Elephant
Armadillo — Xenarthra
Hedgehog
Shrew
Microbat (brown bat)
Megabat (horseshoe bat) — Laurasiatheria
Cow
Cat
Dog
Squirrel
Mouse
Rat
Guinea Pig
Rabbit
Tree Shrew — Euarchontoglires
Lemur
Bushbaby
Macaque
Human
Chimpanzee

Margulies et al., PNAS 2005

# Ancestral mammalian genome reconstruction

Base-by-base reconstruction of complete ancestral genomes
• Including coding, non-coding, repetitive regions

Boreoeutherian ancestor

Expected reconstruction accuracy[*]:

• From ideal choice of extant mammals  99%

• From soon-to-be available genomes:  96%

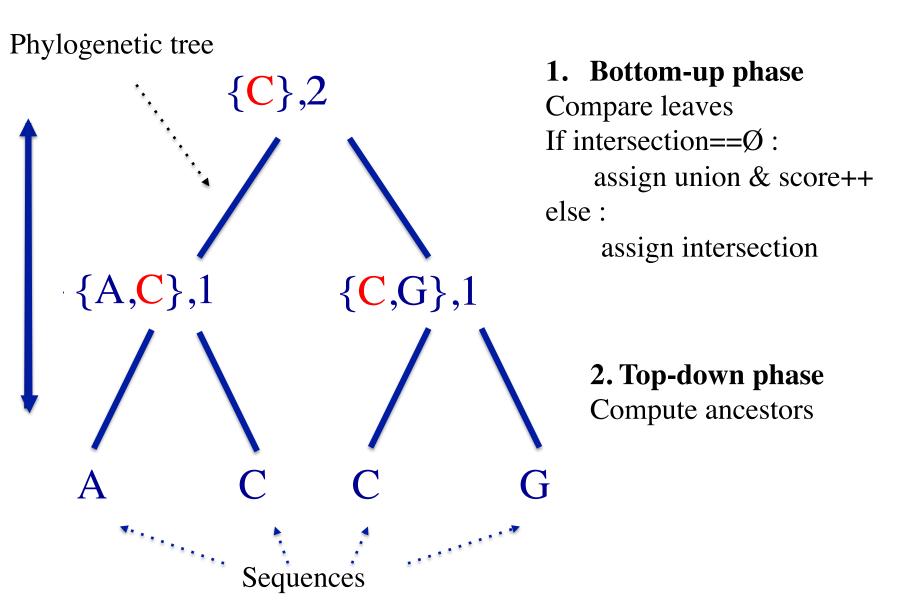• With currently available sequences:  90%
  (full or 2X coverage)

[*] For >90% of euchromatic genome



~70 Myrs

Afrotheria
Xenarthra
Laurasiatheria
Euarchontoglire

Elephant Shrew
Hyrax
Elephant
Armadillo
Hedgehog
Shrew
Microbat (brown bat)
)e bat)
Cat
Dog
Squirrel
Mouse
Rat
Guinea Pig
Rabbit
Tree Shrew
Lemur
Bushbaby
Macaque
Human
Chimpanzee

Tree from Margulies et al., PNAS 200

# The Parsimony Score and Fitch's Algorithm

Phylogenetic tree

$\{C\}, 2$

$\{A,C\}, 1$     $\{C,G\}, 1$

A     C     C     G

Sequences

**1. Bottom-up phase**
Compare leaves
If intersection==∅ :
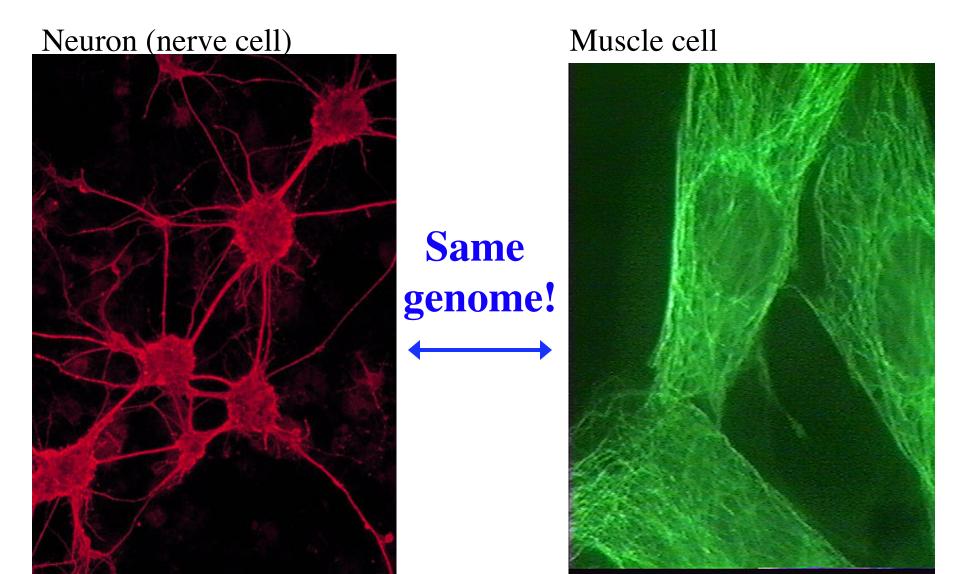    assign union & score++
else :
    assign intersection
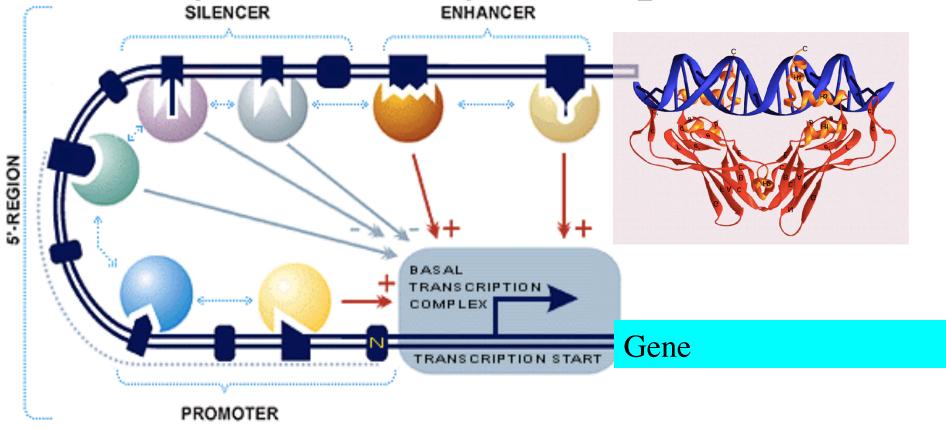
**2. Top-down phase**
Compute ancestors

# Why should we care?

- See genome evolution *happening*, rather than just see its outcome
  - Assign directionality to events (indels, substitutions)
  - Reconstructing the timing of events

- Boreoeutherian ancestor = Ye good olde mammalian stuff
  - BorEut is an archetypical mammalian genome
  - None of the species-specific quirks
  - Human is ~4 times close to BorEut than to mouse
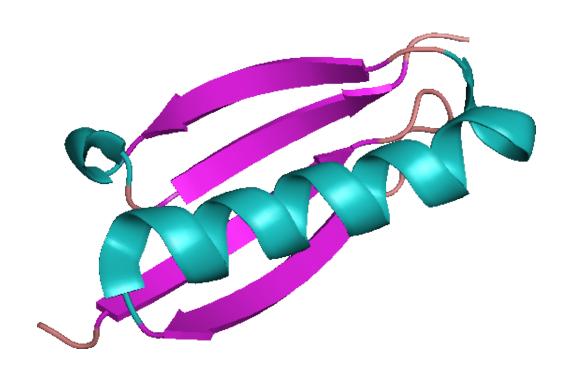
# One program, many functions

Neuron (nerve cell)

Muscle cell



**Same genome!**

# Regulation of gene expression



Transcription factor binding sites:
- ☹ Short: 6 to 20 nucleotides
- ☹ No specific signature; each TF has different binding site
- ☹ Can be up to 1 million nucleotides upstream of gene regulated
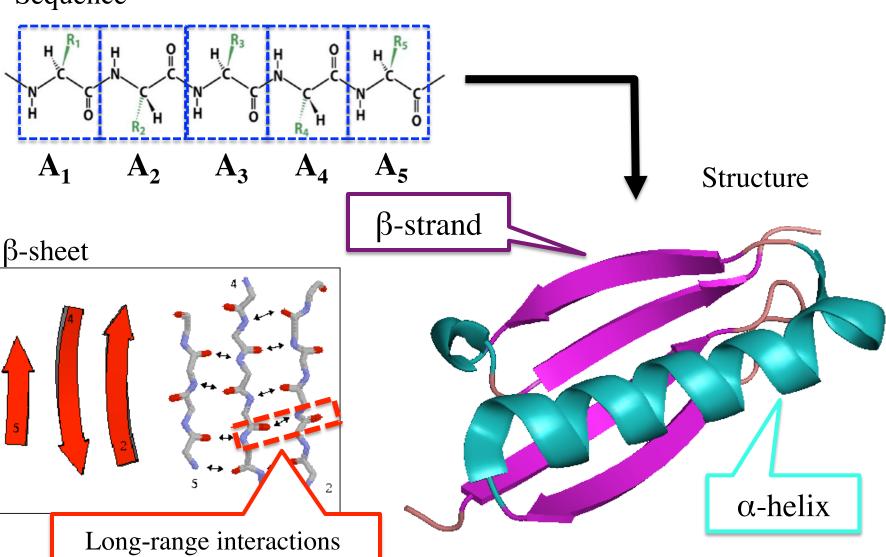- ☺ Often clustered with other binding sites, forming modules

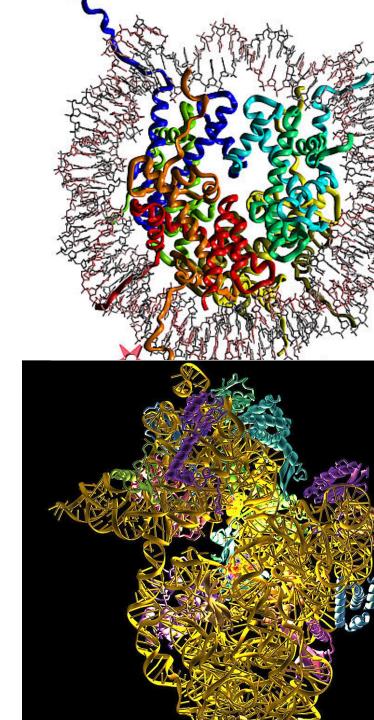# Structural Bioinformatics

# Protein Structure

Sequence



$A_1$  $A_2$  $A_3$  $A_4$  $A_5$

Structure

β-sheet



Long-range interactions
stabilize β–sheet.

β-strand

α-helix

# The hardware - Proteins

- Molecules only obey laws of physics and chemistry.
- Cell organization only relies on interactions between molecules
- Stochastic, dynamic, "chaotic" system
- High error rate in interactions
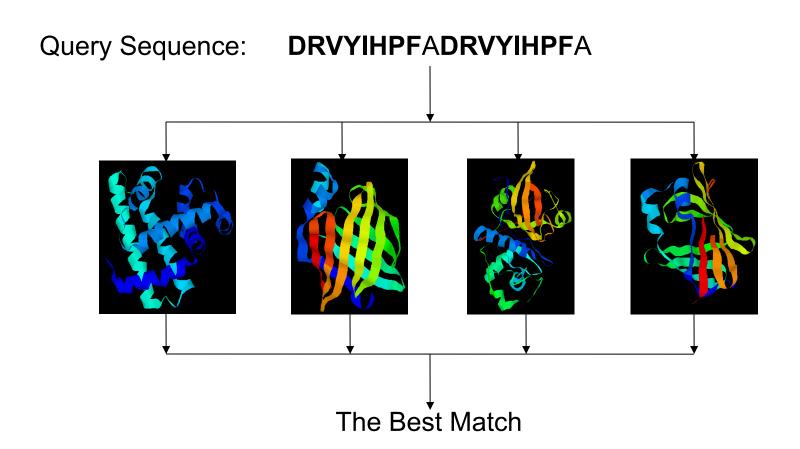- Replicating, self-assembling, self-repairing system
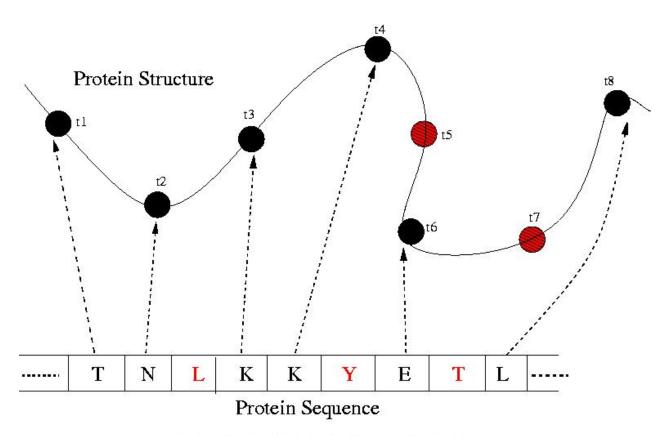
Maybe software engineers have something to learn here…

# Ab-initio vs knowledge-based methods

- Ab inito folding (simulation-based method)
  1998 Duan and Kollman
  36 residues, 1000 ns, 256 processors, 2 months
  Do not find native structure

- Template-based (or knowledge-based) methods
  - Homology modeling: sequence-sequence alignment, works if sequence identity > 25%

  - Protein threading: sequence-structure alignment, can go beyond the 25% limit
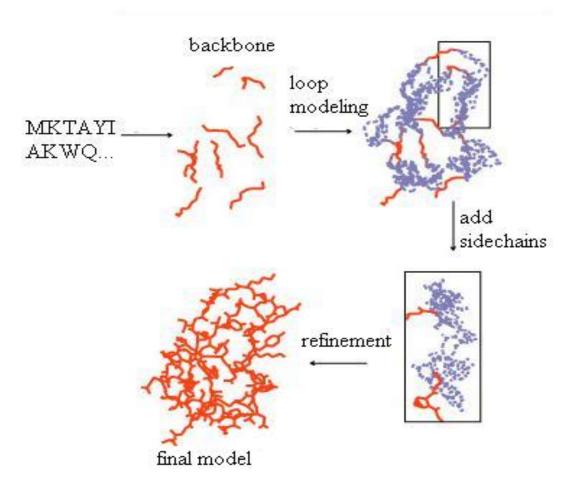
# Protein Threading

Query Sequence:  **DRVYIHPF**A**DRVYIHPF**A



The Best Match

# Threading Example



Protein Structure

Protein Sequence

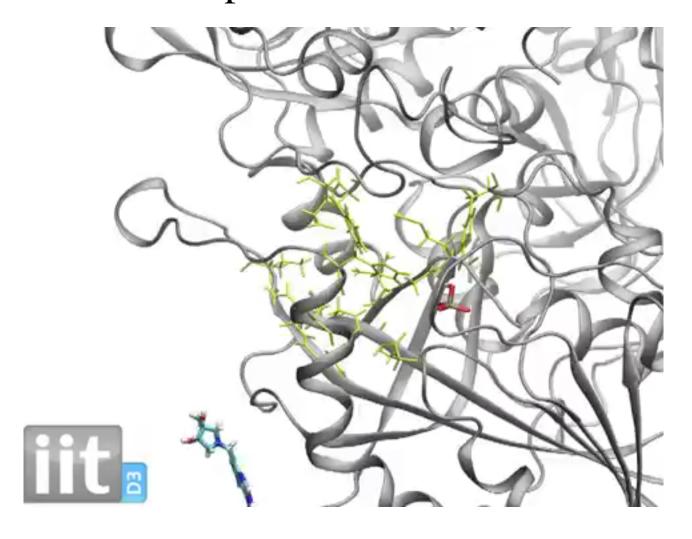Positions or residues in red are gaps

# Protein Structure Prediction



- Stage 1: Backbone Prediction
  - Ab initio folding
  - Homology modeling
  - Protein threading

- Stage 2: Loop Modeling

- Stage 3: Side-Chain Packing
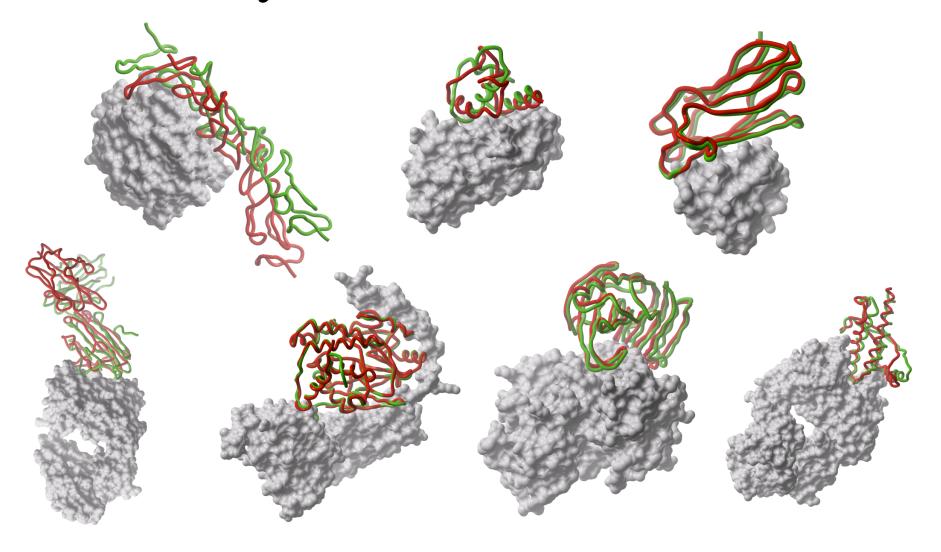
- Stage 4: Structure Refinement

The picture is adapted from http://www.cs.ucdavis.edu/~koehl/ProModel/fillgap.html

# Why are structure useful?
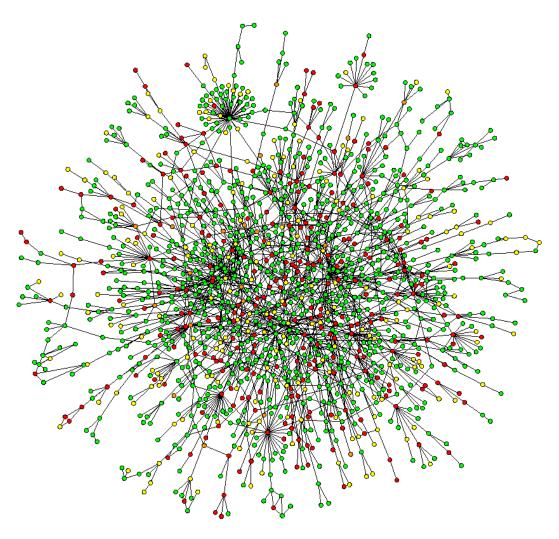## Structure helps to understand function



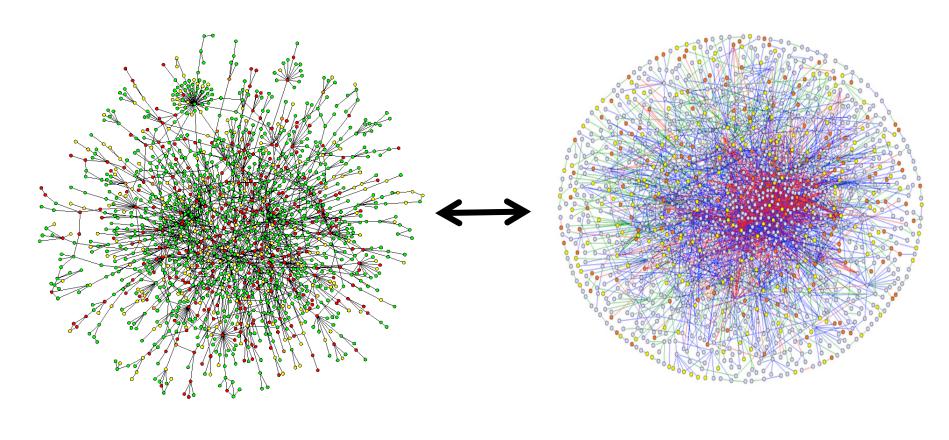Simulation of a drug entering binging site of a protein

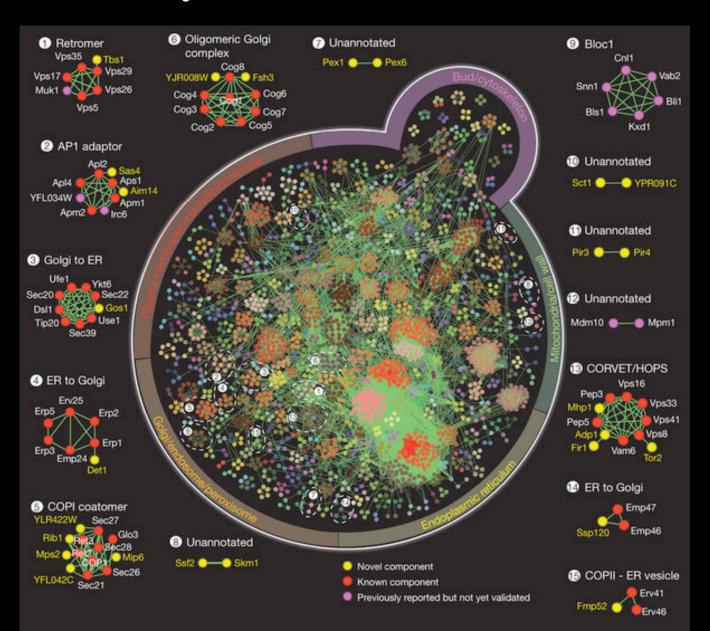# Why are structure useful?

# System Biology
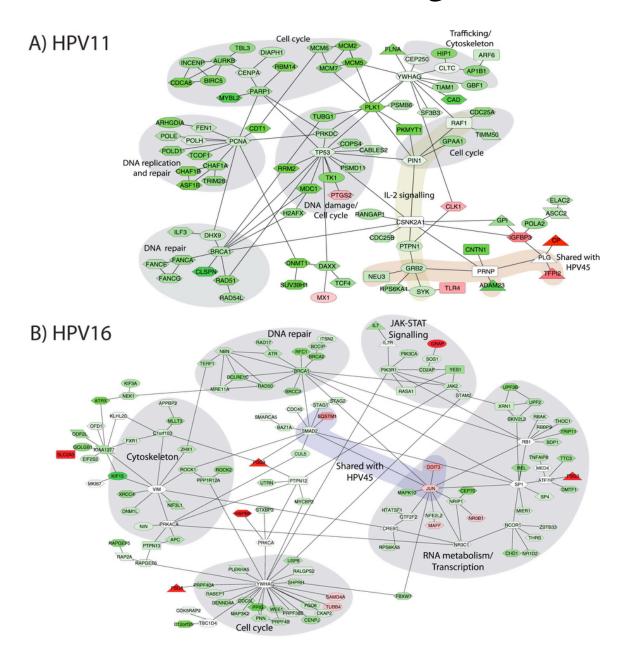
# Interaction network alignment



Yeast         Human

# Discovery of functional motifs in PPI

# How do interaction network change with disease?

# Bioinformatics @ SOCS



McGill Centre for Bioinformatics is located in the Trottier Building

- Mathieu Blanchette
- Michael Hallett (Bellini Bld)
- Derek Ruths
- Jérôme Waldispühl