# CS273A

A Computational Tour of The Human Genome

Lecture 2: Protein Coding Genes

MW  1:30-2:50pm in Clark **S361\*** (behind Peet's)

Profs: Serafim Batzoglou & Gill Bejerano

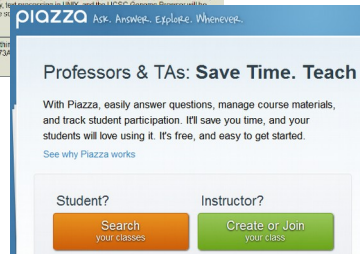CAs: Karthik Jagadeesh & Johannes Birgmeier

\* Handful of lectures/primers elsewhere: track on website/piazza

# Announcements

- http://cs273a.stanford.edu/

  – Course guidelines, office hours, etc.

  – Lecture 1 is posted

  – Problem set 1 rolls out next week

- Course communications via Piazza

  – Auditors please sign up too

- The first tutorial this Friday in Beckman B-302 from 2:00pm-3:30pm. It's the only one some students should consider skipping. While they may be familiar with the first half of the Molecular Biology 101 lecture, we also cover gene regulation and genome rearrangements.

- CAs will be sending out a Doodle poll via Piazza to identify ideal times for office hours. Students can contact them via Piazza for questions.

# Class Goals



- Meet your genome (learn to surf, learn the surf)
- Understand genomic tools (theory, applications)
- DIY (pose questions, write & run tools, understand answers)

# Class Topics

(0) Genome context:
   cells, DNA, central dogma

(1) Genome content / genome function:
   genes, gene regulation, repeats, epigenetics

(2) Genome sequencing:
   technologies, assembly/analysis, technology dependence

(3) Genome evolution:
   evolution = mutation + selection, modes of evolution,
   comparative genomics, ultraconservation, exaptation

(4) Population genomics:
   Tracking human migration patterns via neutral evolution

(5) Genomics of human disease:
   disease susceptibility, cancer genomics, personal genomics

(6) Genome "output" (organism) evolution:
   Evolutionary developmental biology ("evo-devo")

# Organism – Cell - Genome

$10^{13}$ different cells in an adult human. The cell is the basic unit of life.

DNA = linear molecule inside the cell that carries instructions needed throughout the cell's life ~ long string(s) over a small alphabet

Alphabet (nucleotides/bases) {A,C,G,T}

Strings (chromosomes) of length $10^4$-$10^{11}$

Genome:

"instruction"

...ACGTACGACTGACTAGCATCGACTACGACTAGCAC...

# One Cell, One Genome, One Replication

- Every cell holds a copy of all its DNA = its genome.
- The human body is made of ~$10^{13}$ cells.
- All originate from a *single* cell through *repeated* cell divisions.



DNA strings = Chromosomes

cell

genome = all DNA

cell division

egg

egg

chicken

egg

chicken ≈ $10^{13}$ copies (DNA) of egg (DNA)

# What will we study?

The most amazing "Turing tape" in existence, your genome.

# How to Read The Genome

- Genome = DNA.
- Genome is broken up into several strings = chromosomes.
- Humans: Females= (2*chr.1-22)+XX Males= (2*chr.1-22)+XY

cell

DNA strings =
Chromosomes

genome =
all DNA

cell
division

- DNA is double stranded.
- Complementation is rigid.
- Information can be read off of either strand.

5' End        3' End

3' End        5' End

- Every cell contains 2 copies of your genome, one from mom, one from dad.

# The Biggest Challenge in Genomics…

… is <u>computational</u>:

How does this                                                    encode *this*

GCTAGATCGCCTGGTA
GCTTTGCGCCGTCAAA
GTCTTGAAGGCTGTGA
TCAAGCTTCTTGCGAT
CCCGTTTGACCGGAGC
CTTGCCAATGAGTTCT
CAGCTGTCTATATGAA
TCACAAAATACGCAAT



Program                                                              Output

This "coding" question has <u>profound</u> implications for our lives

# Class Topics

(0) Genome context:
   cells, DNA, central dogma

**(1) Genome content / genome function:**
   **genes, gene regulation, repeats, epigenetics**

(2) Genome sequencing:
   technologies, assembly/analysis, technology dependence

(3) Genome evolution:
   evolution = mutation + selection, modes of evolution,
   comparative genomics, ultraconservation, exaptation

(4) Population genomics:
   Tracking human migration patterns via neutral evolution

(5) Genomics of human disease:
   disease susceptibility, cancer genomics, personal genomics

(6) Genome "output" (organism) evolution:
   Evolutionary developmental biology ("evo-devo")

# Genome Content

# Genomes, Genes & Proteins

The most visible instructions in our genome are Genes.

Genes explain exactly HOW to synthesize any protein.

Proteins are the work horses of every living cell.

Genome:

gene

...ACGTACGACTGACTAGCATCGACTACGACTAGCAC...

linear
(folded)
molecule

protein

cell

# Central Dogma of Biology



genome        DNA        {A,C,G,T}

replication (DNA -> DNA)
DNA Polymerase

transcription (DNA -> RNA)
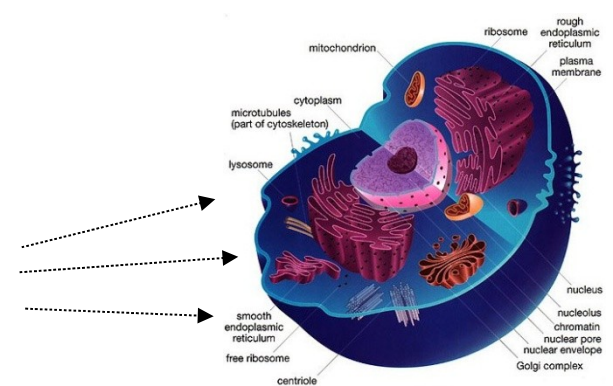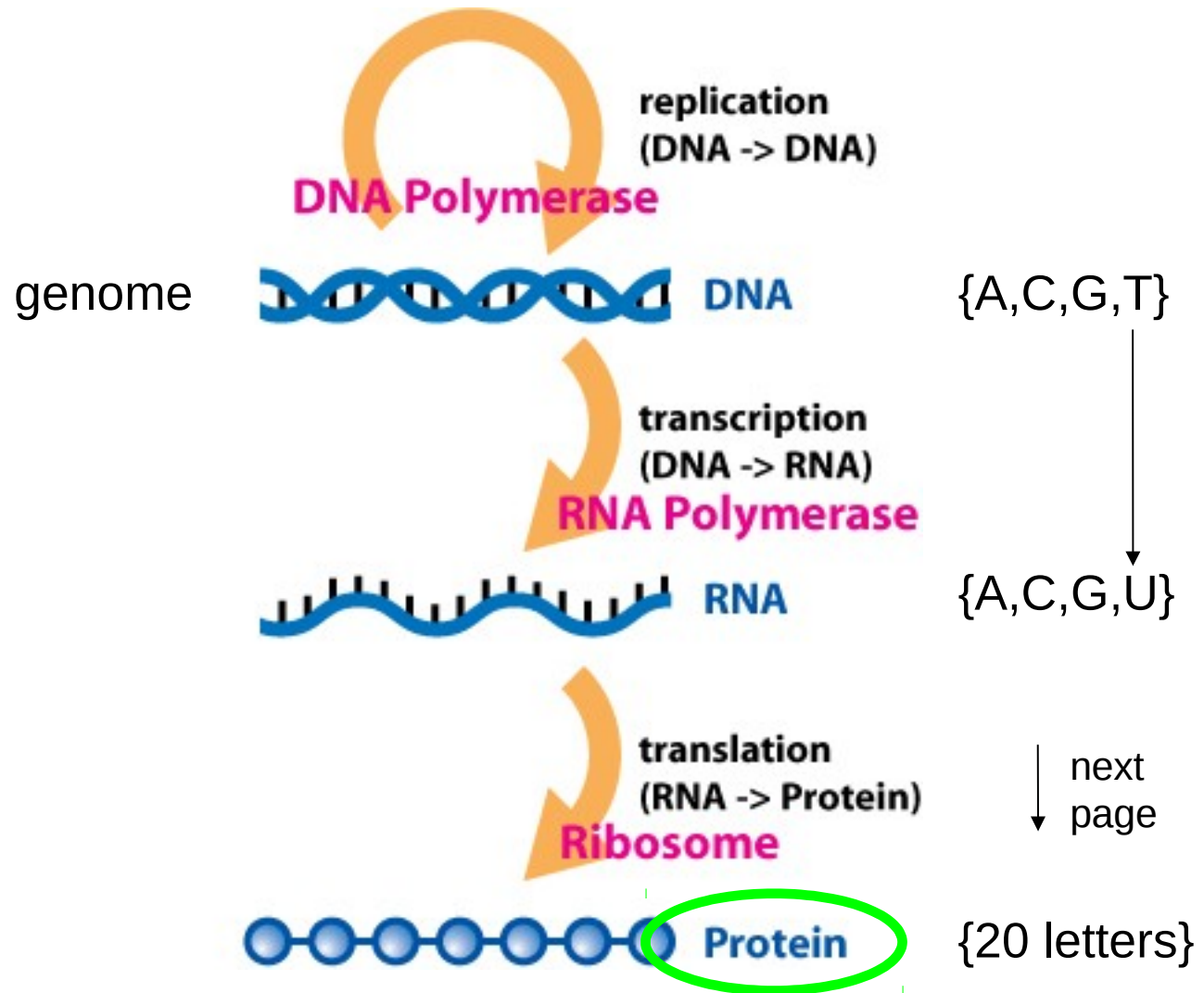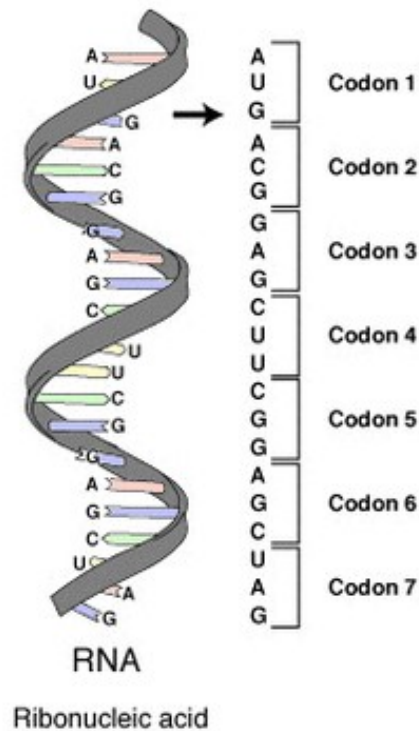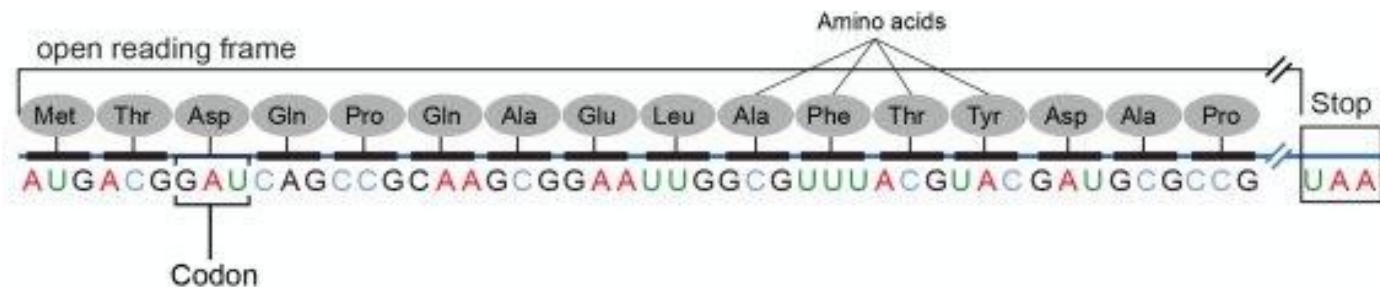RNA Polymerase

RNA        {A,C,G,U}

translation (RNA -> Protein)
Ribosome

next page

Protein        {20 letters}

# Translation: The Genetic Code



RNA
Ribonucleic acid

Codon 1
Codon 2
Codon 3
Codon 4
Codon 5
Codon 6
Codon 7

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

First position (5' end)

Third position (3' end)

Amino acid names:

Ala = alanine
Arg = arginine
Asn = asparagine
Asp = aspartate
Cys = cysteine
Gln = glutamine
Glu = glutamate
Gly = glycine
His = histidine
Ile = Isolevcine
Leu = leucine
Lys = lysine
Met = methionine
Phe = phenylalanine
Pro = proline
Ser = serine
Thr = threonine
Trp = tryptophan
Tyr = Tyrosine
Val = valine



open reading frame

Amino acids

Met Thr Asp Gln Pro Gln Ala Glu Leu Ala Phe Thr Tyr Asp Ala Pro    Stop

AUGACGGAUCAGCCGCAAGCGGAAUUGGCGUUUACGUACGAUGCGCCG    UAA

Codon

# Genes Can Be Encoded on Either Strand

# Gene Structure

# Gene Splicing

# Visualizing Gene Structure

# Genes in the Human Genome



**Intron,** and **direction** of transcription  <<< or >>>



UCSC primer

There are ~20,000 protein coding genes in the human genome.

(Even half way through sequencing the human genome,
Researchers thought there will be well over 100,000 genes).

# Gene Finding

Computational Challenge:

"Find the genes, the whole genes, and nothing but the genes"



Understand Biology  🡒  Write discovery tools

(Our) answer depends on our understanding, data & tools

# Gene prediction approachs

- **Rule-based programs**
  - Use explicit set of rules to make decisions.
  - Example: GeneFinder
- **Neural Network-based programs**
  - Use data set to build rules.
  - Examples: Grail, GrailEXP
- **Hidden Markov Model-based programs**
  - Use probabilities of states and transitions between these states to predict features.
  - Examples: Genscan, GenomeScan

# GenScan States



- N - intergenic region
- P - promoter
- F - 5' untranslated region
- $E_{sngl}$ – single exon (intronless) (translation start -> stop codon)
- $E_{init}$ – initial exon (translation start -> donor splice site)
- $E_k$ – phase k internal exon (acceptor splice site -> donor splice site)
- $E_{term}$ – terminal exon (acceptor splice site -> stop codon)
- $I_k$ – phase k intron: 0 – between codons; 1 – after the first base of a codon; 2 – after the second base of a codon

# Alternative Splicing

# Genes in the Human Genome

When you only show one transcript per gene locus:



If you ask the GUI to show you all well established gene variants:

# Protein Domains



SKSHSEAGSAFIQTQQLHAAMADTFLEHMCRLDIDSAPITARNT
GIICTIGPASRSVETLKEMIKSGMNVARMNFSHGTHEYHAETIK
NVRTATESFASDPILYRPVAVALDTKG**PEIRTGLIKGSGTAEVE**
**LKKGATLKITLDNAYMAACDENILWLDYKNICKVVEVGSKVYVD**
**DGLISLQVKQKGPDFLVTEVENGGFLGSKKGVNLPGAAVDL**PAV
**SEKDIQDLKFGVDEDVDMVFASFIRKAADVHEVRKILGEKGKNI**
**KIISKIENHEGVRRFDEILEASDGIMVARGDLGIEIPAEKVFLA**
**QKMIIGRCNRAGKPVICATQMLESMIKKPRPTRAEGSDVANAVL**
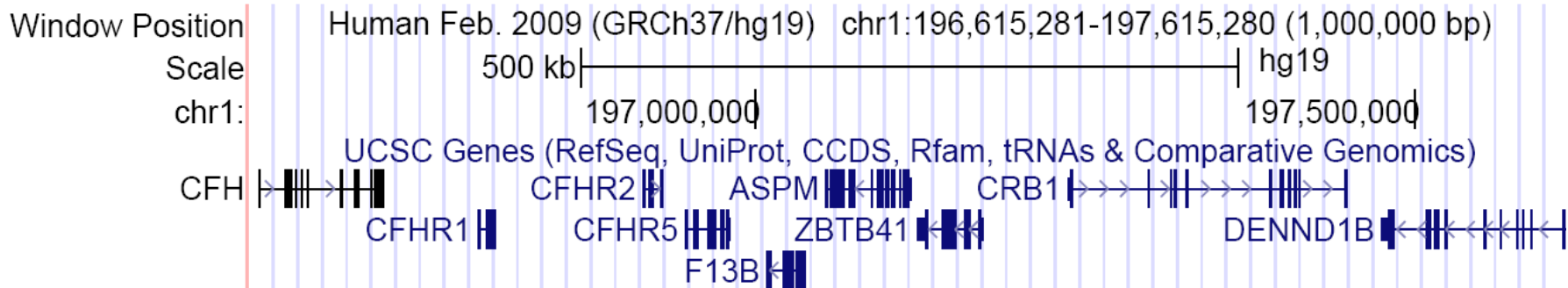**DGADCIMLSGETAKGDYPLEAVRMQHLIAREAEAAMFHRKLFEE**
**L**ARSSSHSTDLMEAMAMGSVEASYKCLAAALIVLTESGRSAHQV
ARYRPRAPIIAVTRNHQTARQAHLYRGIFPVVCKDPVQEAWAED
VDLRVNLAMNVGKAAGFFKKGDVVIVLTGWRPGSGFTNTMRVVP
VP

A protein domain is a subsequence of the protein that folds independently of the other portions of the sequence, and often confers to the protein one or more specific functions.

# Alt. Splicing and Protein Repertoire



Alternative splicing often produces protein variants that have a different domain composition, and thus perform different functions.

What if we want to predict all splice variants that are ever made? Can we even do it from sequence alone?

# Common Problems

- Common problems with gene finders
  - Fusing neighboring genes
  - Spliting a single gene
  - Miss exons or entire genes
  - Overpredict exons or genes

- Other challenges
  - Nested genes
  - Noncanonical splice spites
  - Pseudogenes
  - Different isoforms of same gene

# We can sequence all mRNA of a given cell



(Great, but not all genes/isoforms are expressed in all cells. Some are very exotic).

**AUGGUG** - - - -**GGCCCUUUGGGA** - - - - - **CACUAA**

GTGAGG**ATGGTA**AATA**GGGCAT** - - - **GGA**TTGAG**CACUAA**TAA

# Gene Annotation System



- All Ensembl gene predictions are based on experimental evidence

- Predictions based on manually curated Uniprot/Swissprot/Refseq databases

- UTRs are annotated only if they are supported by EMBL mRNA records

*Val Curwen, et al. The Ensembl Automatic Gene Annotation System Genome Res., (2004)* ***14*** *942 - 950.*

# First full draft of the Human Genome



Human Genome Consortium
(HGC)

Celera

2001

# Everything in Genomics is a *Moving* Target

- The genomes (ie, assemblies)
- Their annotations
- Our understanding of Biology
- The portals

Conclusion: write code that can be run...

and rerun and rerun and rerun and rerun

# Biological Functions of the Human Gene Set



Focus on the X axis:

[HGC, 2001]

# Molecular Functions of the Human Gene Set



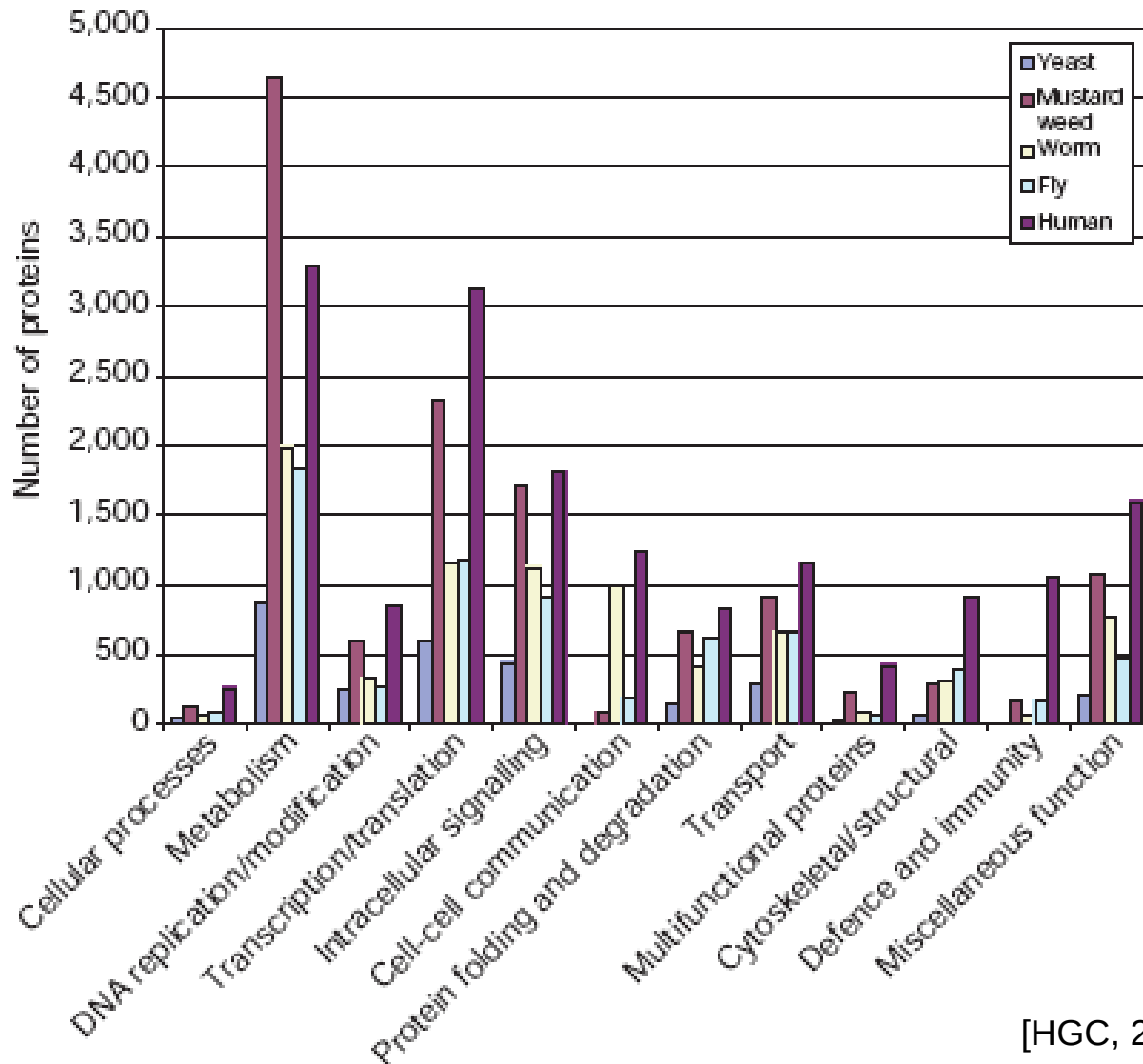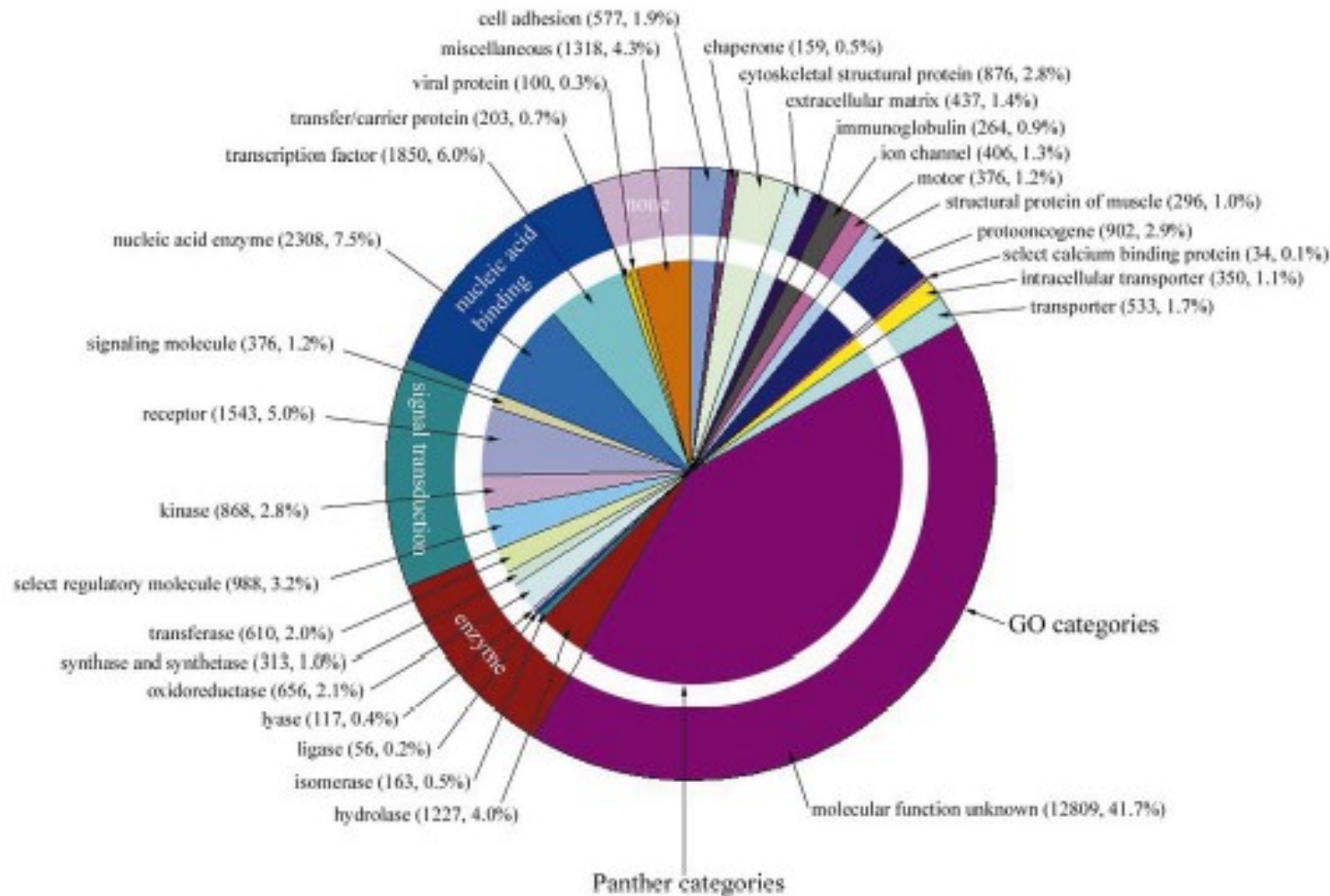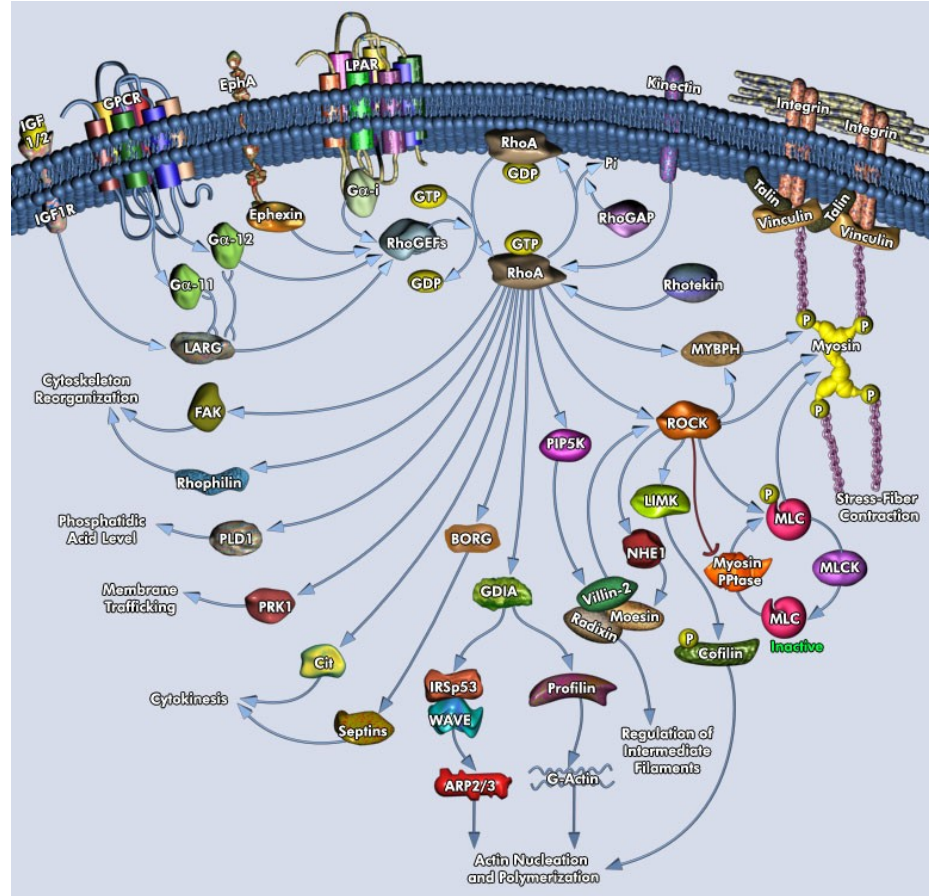Fig. 15. Distribution of the molecular functions of 26,383 human genes. Each slice lists the numbers and percentages (in parentheses) of human gene functions assigned to a given category of molecular function. The outer circle shows the assignment to molecular function categories in the Gene Ontology (GO) (179), and the inner circle shows the assignment to Celera's Panther molecular function categories (116).

cell adhesion (577, 1.9%)
miscellaneous (1318, 4.3%)
viral protein (100, 0.3%)
transfer/carrier protein (203, 0.7%)
transcription factor (1850, 6.0%)
nucleic acid enzyme (2308, 7.5%)
signaling molecule (376, 1.2%)
receptor (1543, 5.0%)
kinase (868, 2.8%)
select regulatory molecule (988, 3.2%)
transferase (610, 2.0%)
synthase and synthetase (313, 1.0%)
oxidoreductase (656, 2.1%)
lyase (117, 0.4%)
ligase (56, 0.2%)
isomerase (163, 0.5%)
hydrolase (1227, 4.0%)

chaperone (159, 0.5%)
cytoskeletal structural protein (876, 2.8%)
extracellular matrix (437, 1.4%)
immunoglobulin (264, 0.9%)
ion channel (406, 1.3%)
motor (376, 1.2%)
structural protein of muscle (296, 1.0%)
protooncogene (902, 2.9%)
select calcium binding protein (34, 0.1%)
intracellular transporter (350, 1.1%)
transporter (533, 1.7%)

nucleic acid binding
signal transduction
enzyme
none

GO categories
molecular function unknown (12809, 41.7%)
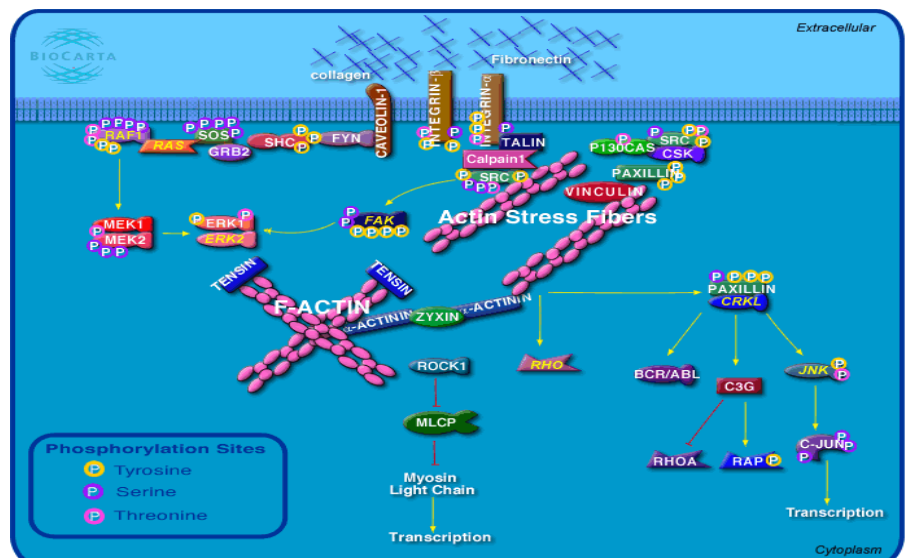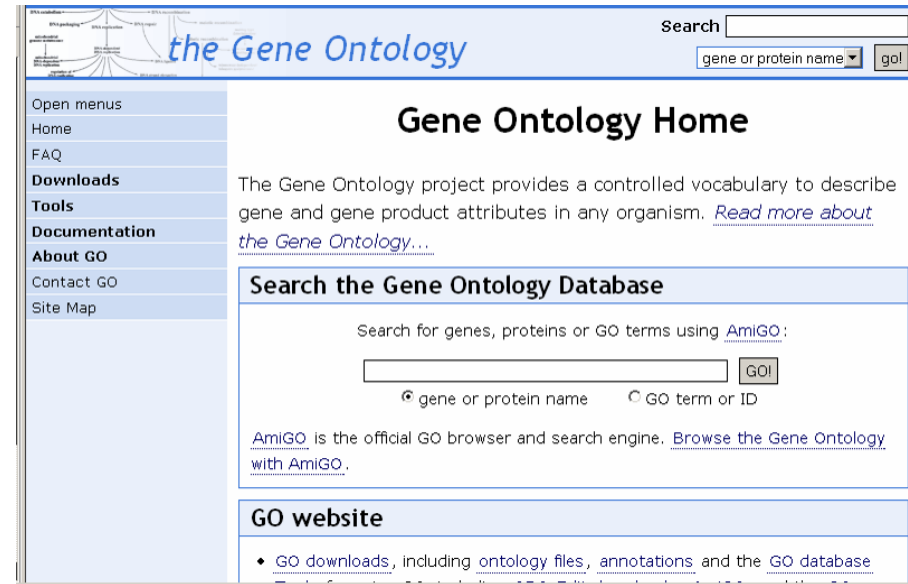Panther categories

[Celera, 2001]

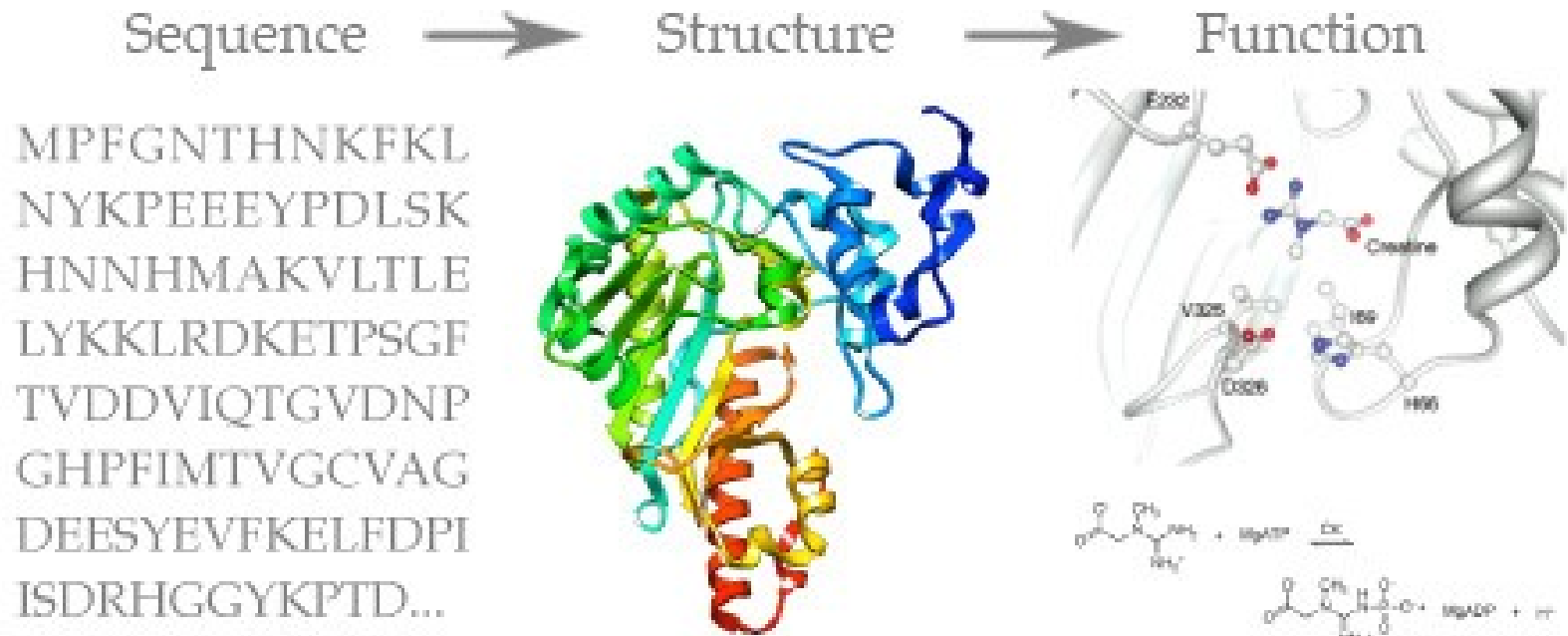# Biological vs. Molecular Function: Pathways



Proteins with very different molecular functions participate to manifest a single biological function, for example: a pathway.

# Gene Sets

- Gene Ontology ("GO")
  - Biological Process
  - Molecular Function
  - Cellular Location
- Pathway Databases
  - KEGG
  - BioCarta
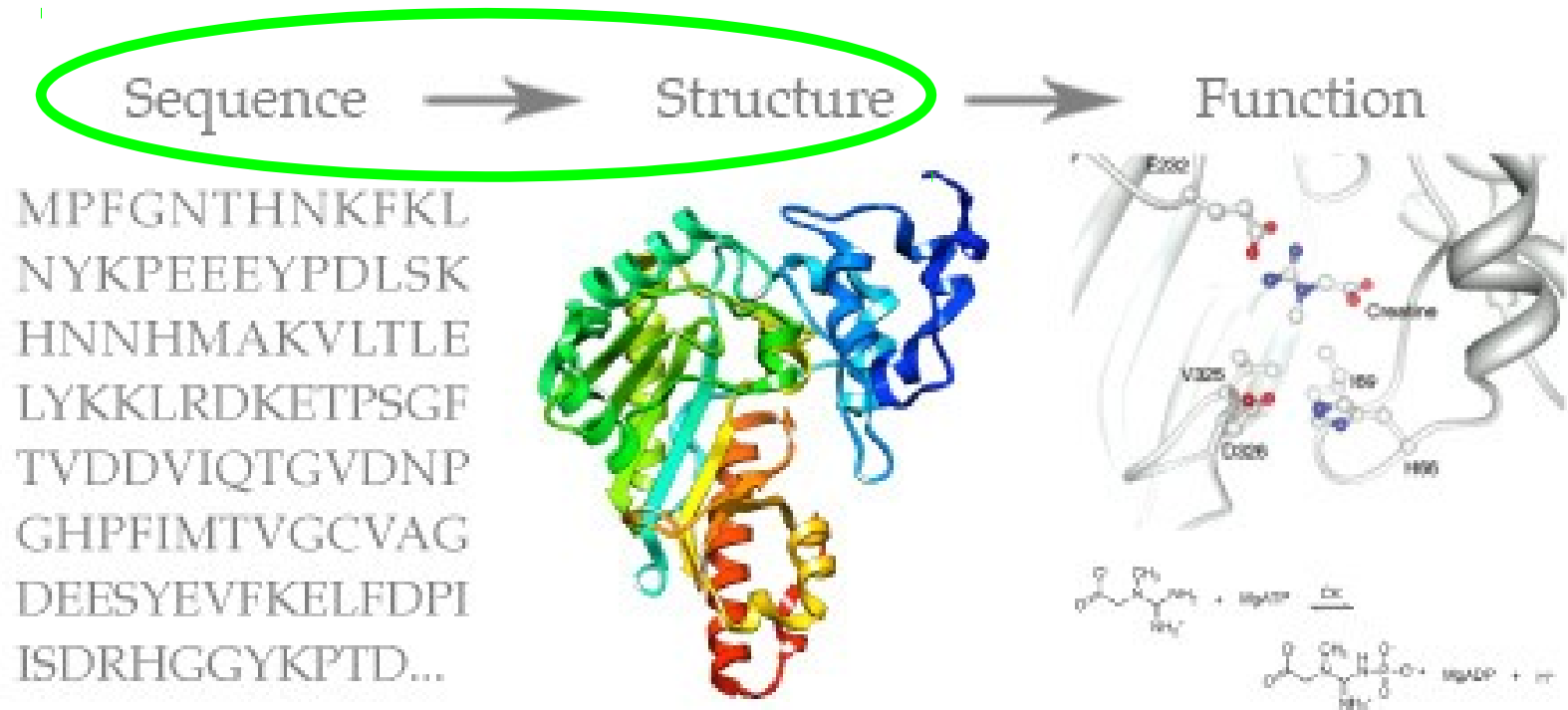  - Broad Institute
- Multiple others

# Genes & Their Functions



Sequence → Structure → Function

```
MPFGNTHNKFKL
NYKPEEEYPDLSK
HNNHMAKVLTLE
LYKKLRDKETPSGF
TVDDVIQTGVDNP
GHPFIMTVGCVAG
DEESYEVFKELFDPI
ISDRHGGYKPTD...
```

Gene (DNA) sequence determines protein (AA) sequence,
  which determines protein (3D) structure,
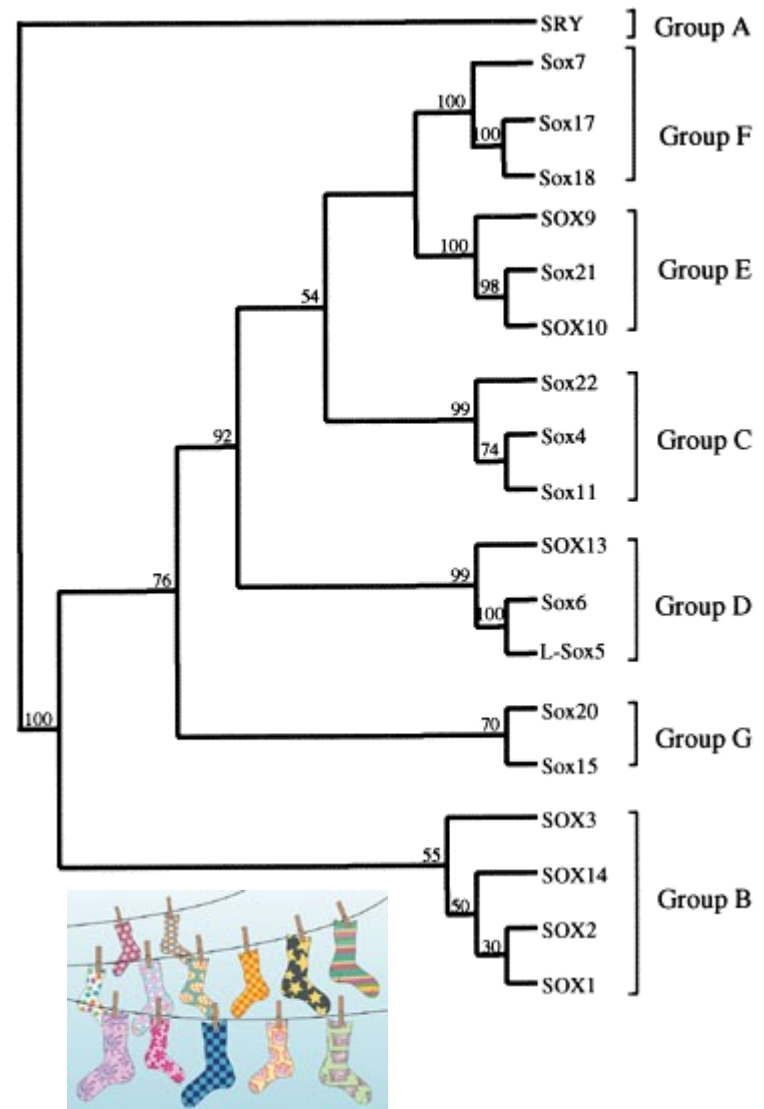  which determines protein's function.

# Protein Folding



Protein folding is the challenge of deducing protein structure from protein sequence.

# Gene Families, Gene Names

Genes (proteins) come in families.

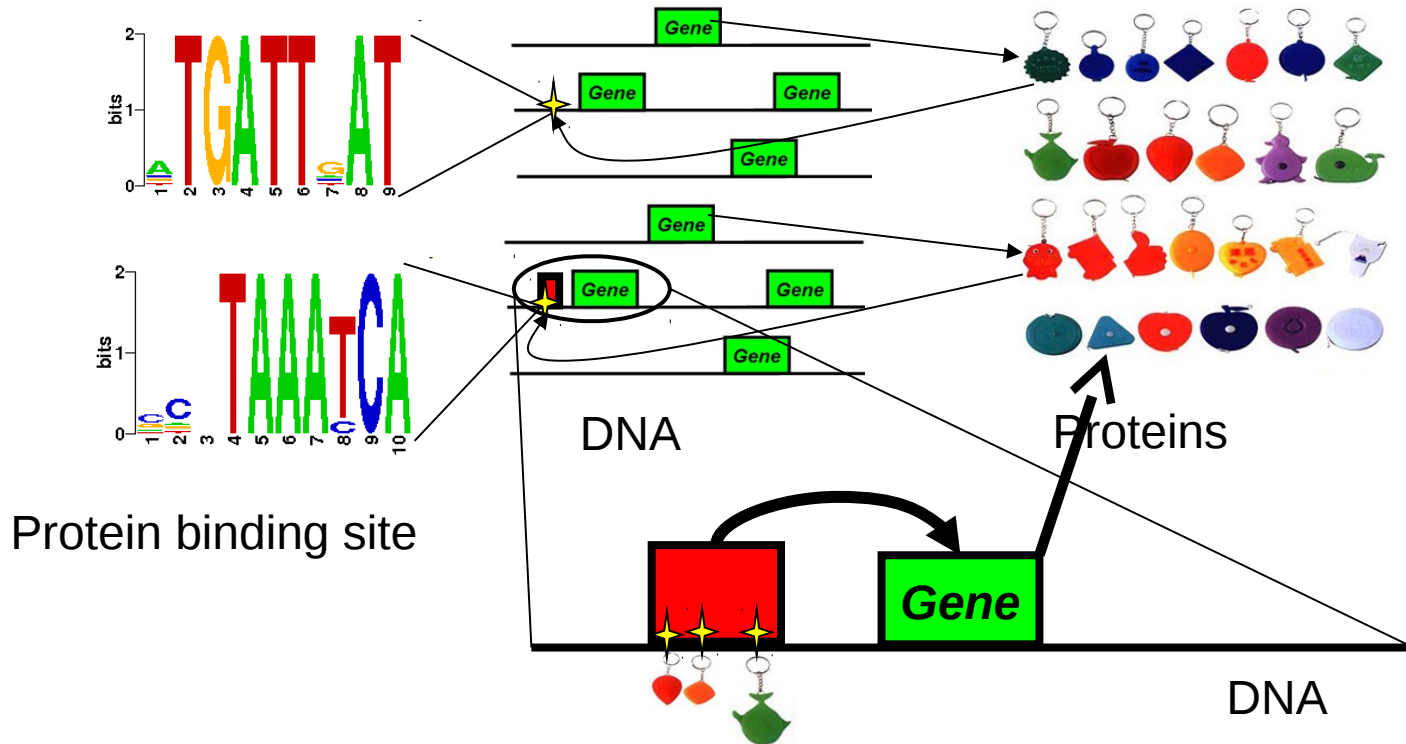Genes of the same family have similar sequences.

Which is why the fold into similar structure and perform similar functions.

Genes of the same family will typically have a "family name" followed by a (sequential) number or "first name".
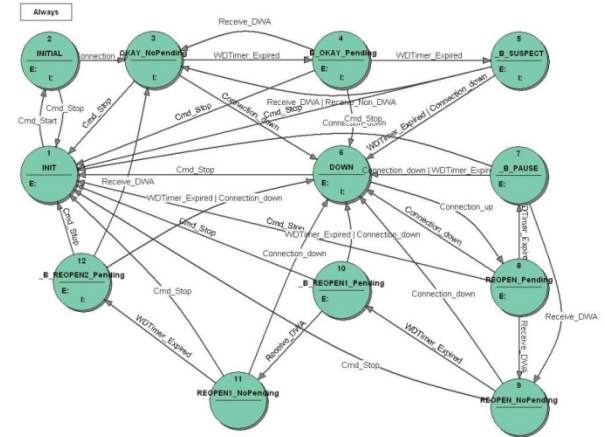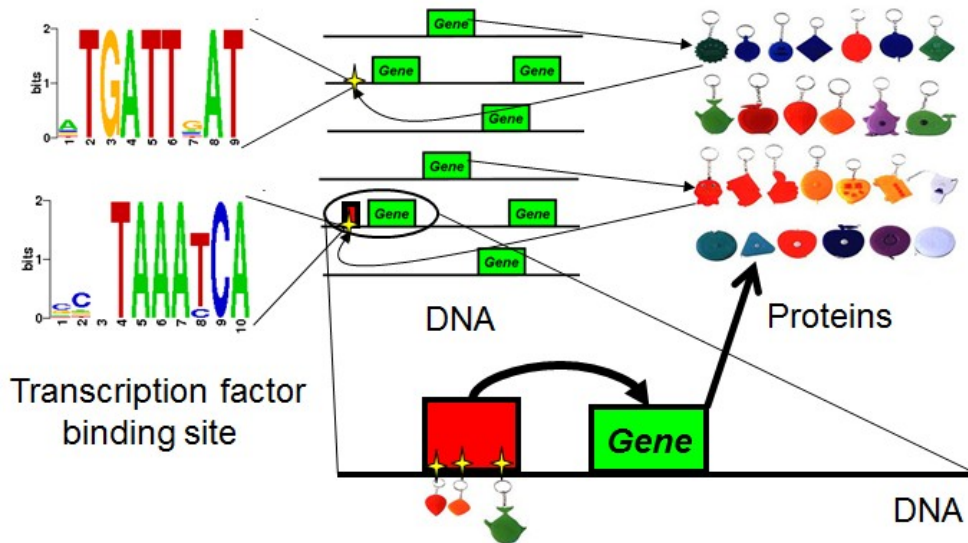
# Some "Special" Functions: Gene Regulation

2,000 different proteins can bind specific DNA sequences.



DNA

Proteins

Protein binding site

*Gene*

DNA

Proteins that regulate the transcription of other proteins
are called <u>transcription factors</u>.

# The Importance of Gene Regulation



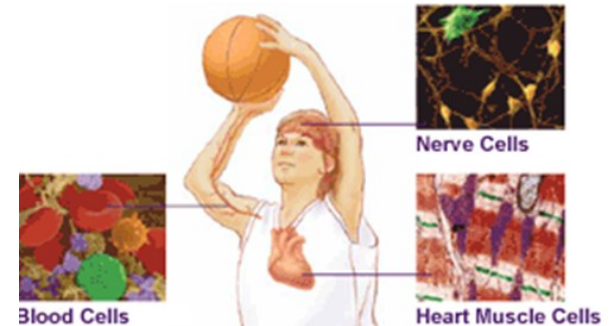Transcription factor binding site

DNA

Proteins

DNA



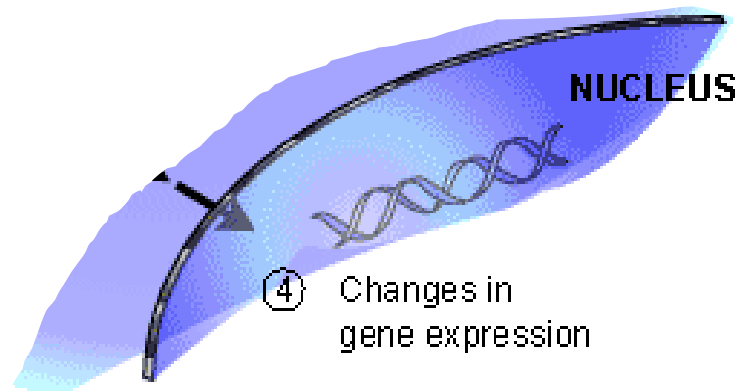The looks & capabilities of different cells are determined by the subset of genes they express.

Different cell types express very different gene repertoires (from the same genome).

To change its behavior a cell can change its transcriptional program.

Think of it as a giant state machine…



Nerve Cells

Blood Cells

Heart Muscle Cells

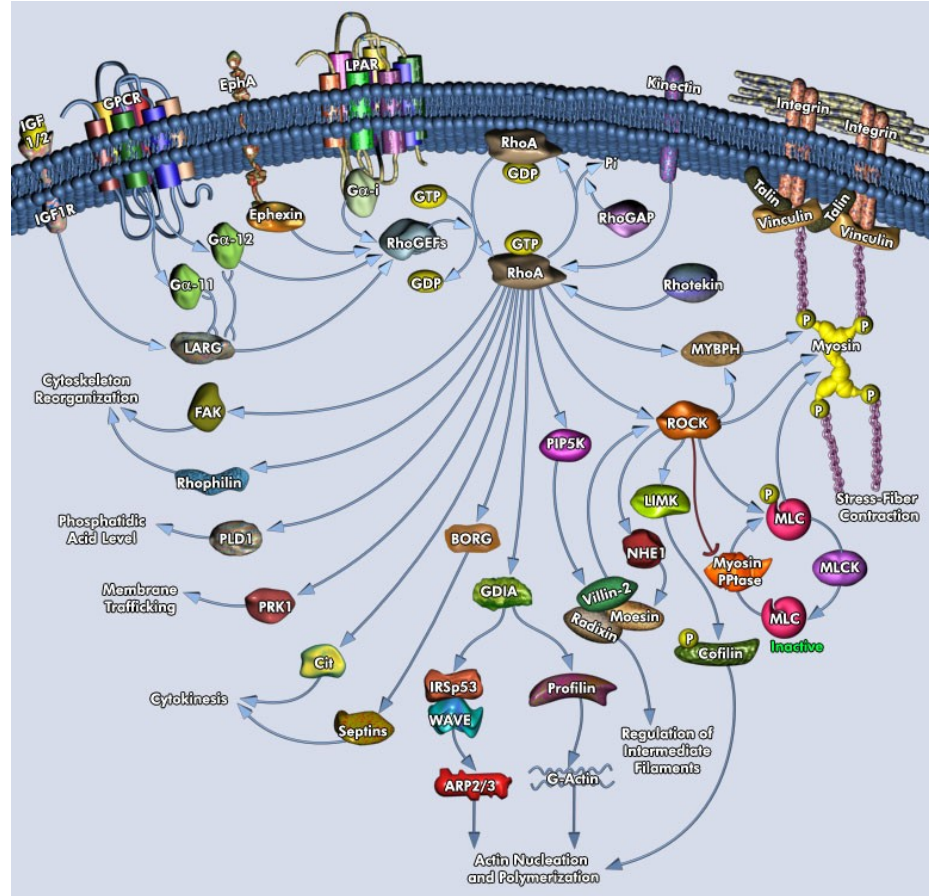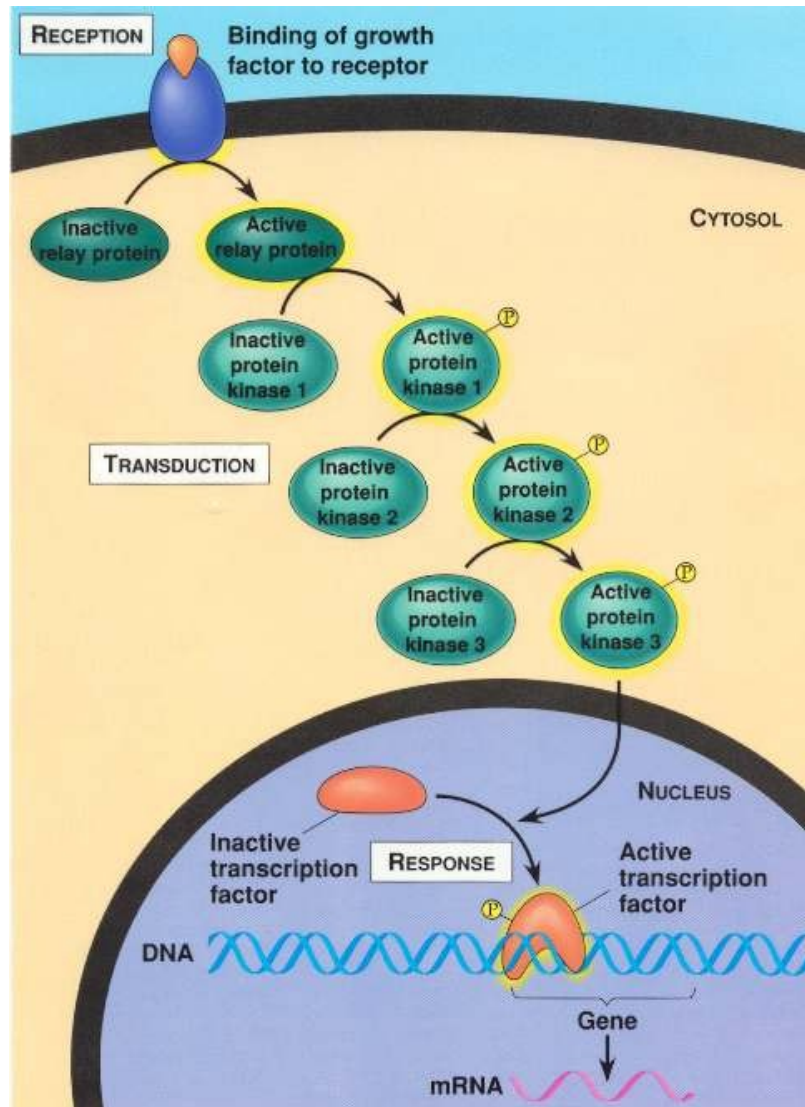# "Special" Function: Cell Signaling



Cells also talk with each other. They send and receive messages, and change their behavior according to messages they receive.

# Biological vs. Molecular Function: Pathways



Proteins with very different molecular functions participate to manifest a single biological function, for example: a pathway.

# Signal Transduction



Now its an even bigger state machine of individual state machines (=cells) talking with each other, orchestrating their individual activities.