

COMPGI10 – Intelligent Systems in Bioinformatics

Lecture 1: Introduction to high-throughput and systems biology

Dr Kevin Bryson

K.Bryson@cs.ucl.ac.uk

This part of the course ...

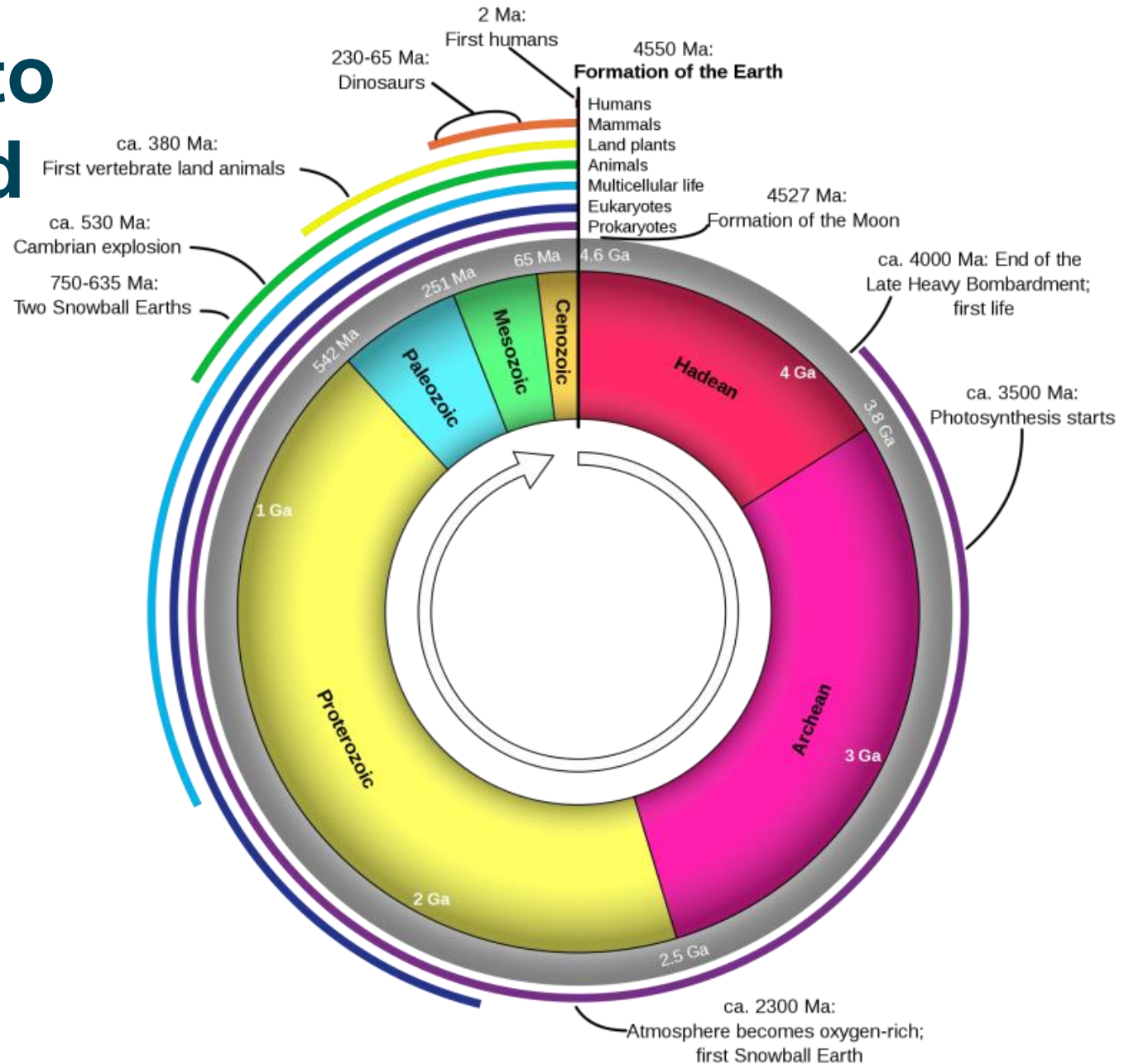
- Up to now you have been looking at ‘individual biological parts’ (DNA, protein, etc.)
- The application of machine learning is valuable in these cases for making predictions about these parts, for instance predicting the secondary structure of a protein from its amino acid sequence.
- But understanding biological systems requires more than just a reductionist understanding of the individual parts that compose them.
- In this part of the course we will look at complete biological systems, techniques that apply to these and how machine learning is useful.

Learning Outcomes

After this lecture, I would hope you:

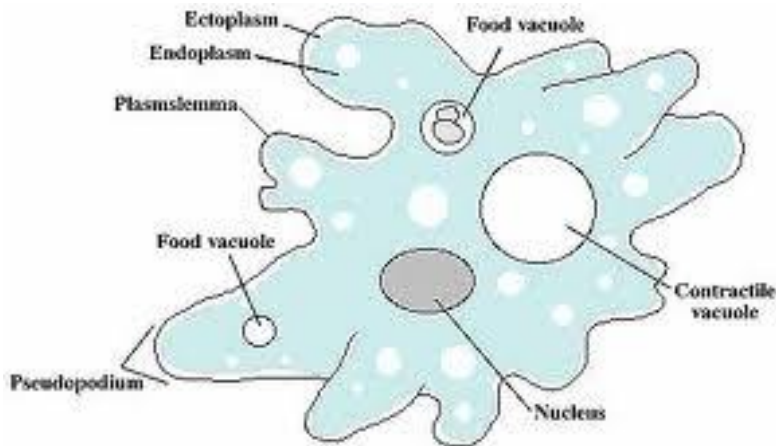
- Understand what we mean by a ‘biological system’ and comprehend the complexity of such a system.
- Appreciate that one way to understand such systems is in terms of their parts and the interactions between these parts.
- Understand different types of high-throughput (HT) experimental techniques to globally characterize the parts.
- Understand different types of biological networks to characterize the interactions between the parts.
- *All this provides a natural problem domain for machine learning.*

Our aim – to understand something that took 4 billion years to evolve!



A simple example from 1.2 billion years ago ... the humble slime mould

Slime moulds are single cell organisms without any nervous system ...



Slime moulds are very simple biological systems - just a bag of chemicals surrounded by a membrane: ... but they can solve mazes!

And they can find the optimal paths between multiple sources of food (like a distribution networks ...)

It's not rocket science ...
but you can see why they
have evolved to have this
relatively complex
behaviour.



Slime moulds are very simple biological systems - just a bag of chemicals surrounded by a membrane:



They have also 'learnt' how to collaborate with each other when required and show 'self-sacrifice' for the common good ...

Aim of Systems Biology and HT technology

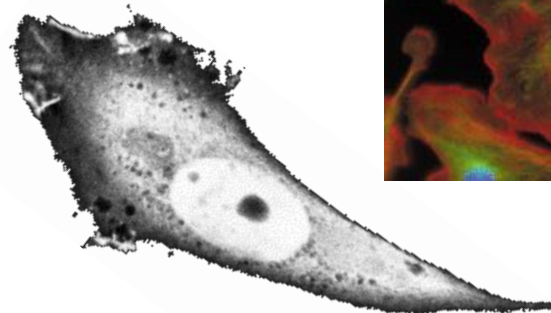
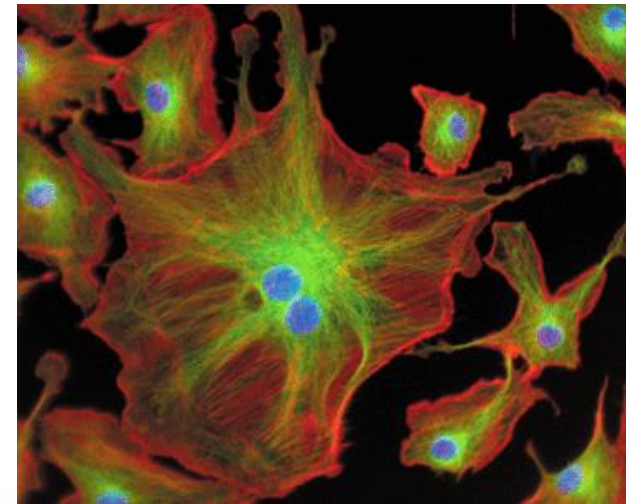
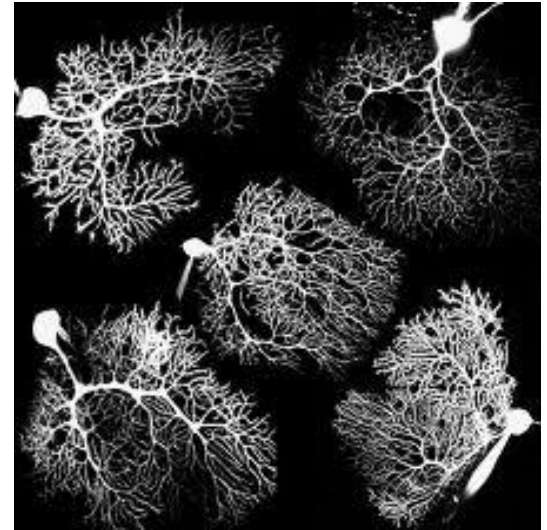
- These relatively complex behaviours are clearly at the ‘systems-level’.
 - They cannot be understood by studying a single gene or single protein structure.
 - They arise from complex interactions between different biological parts (molecules, cells, populations)



Before continuing ...

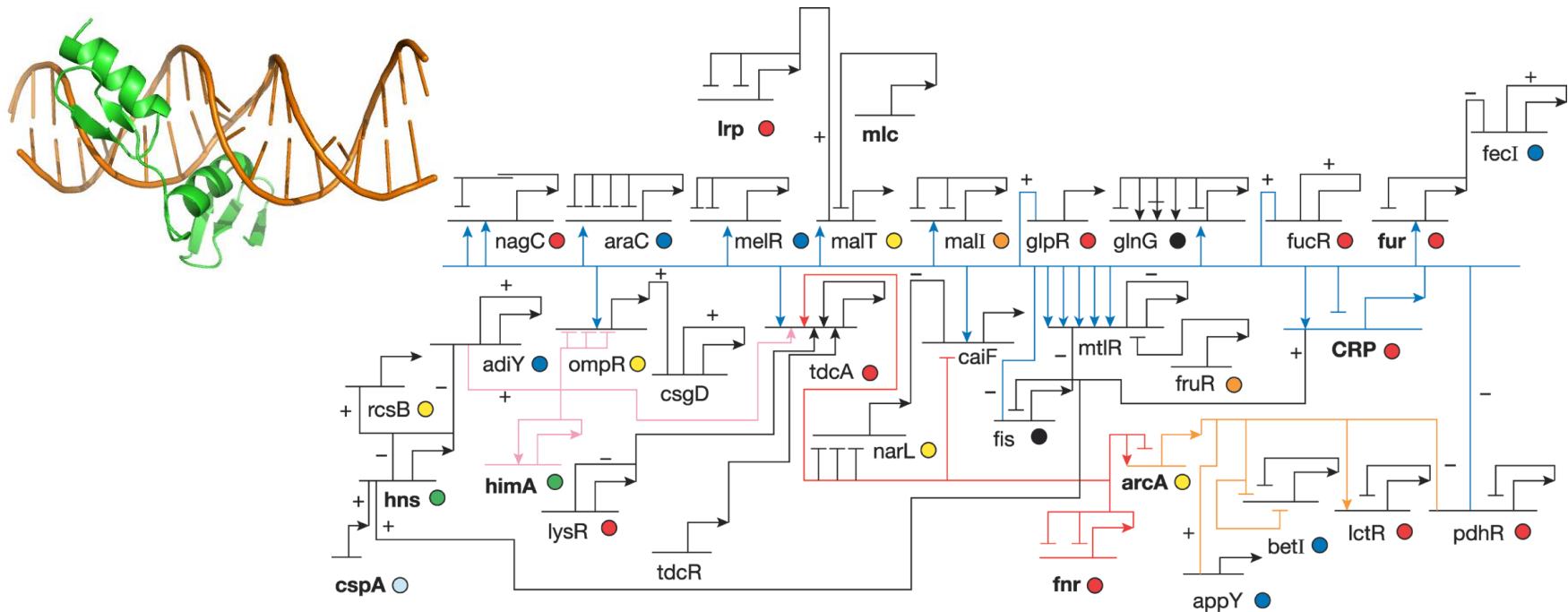
A question ...

- All the cells making up our bodies all share an identical 'blueprint' of instructions (the genome sequence).
- But all our cells look and behave very differently ... how does this happen if they all have the same instructions for proteins?



- Discuss in pairs for a couple of minutes ...

During development each cell receives different ‘signals’ from chemicals released from other cells ... and an intricate ‘genetic program’ where proteins bind to DNA and turn on/off particular genes finally gives all the different cell types.



Biological Network Data –
biological processes
 are accomplished
 by proteins
 working
 together
 in *biological*
networks –
 each protein
 carrying out
 a specific
molecular
function

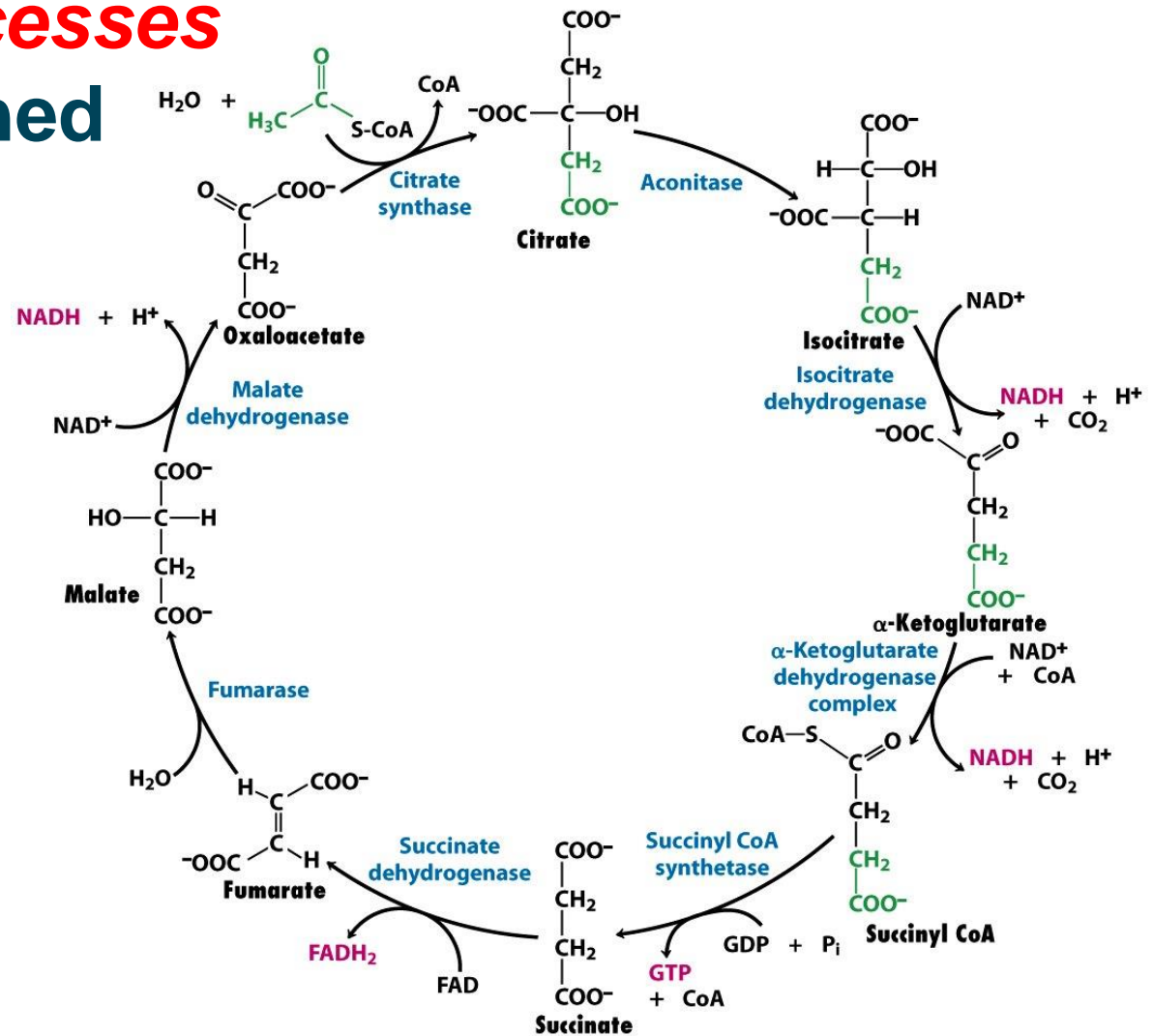
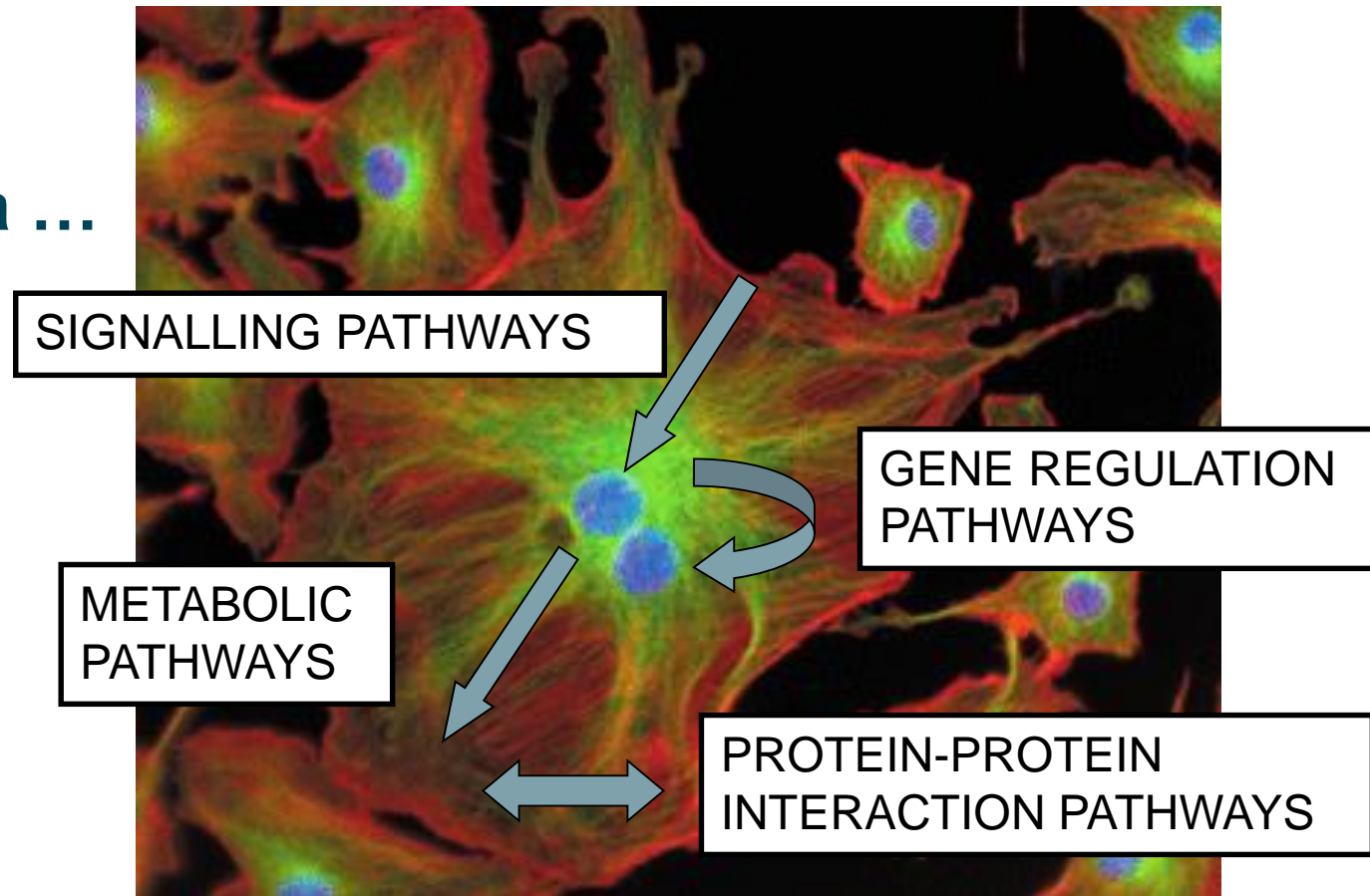


Figure 17-15
 Biochemistry, Sixth Edition
 © 2007 W. H. Freeman and Company

Biological Network Data helps integrate diverse types of data about biological components ...

Key types of biological network data ...

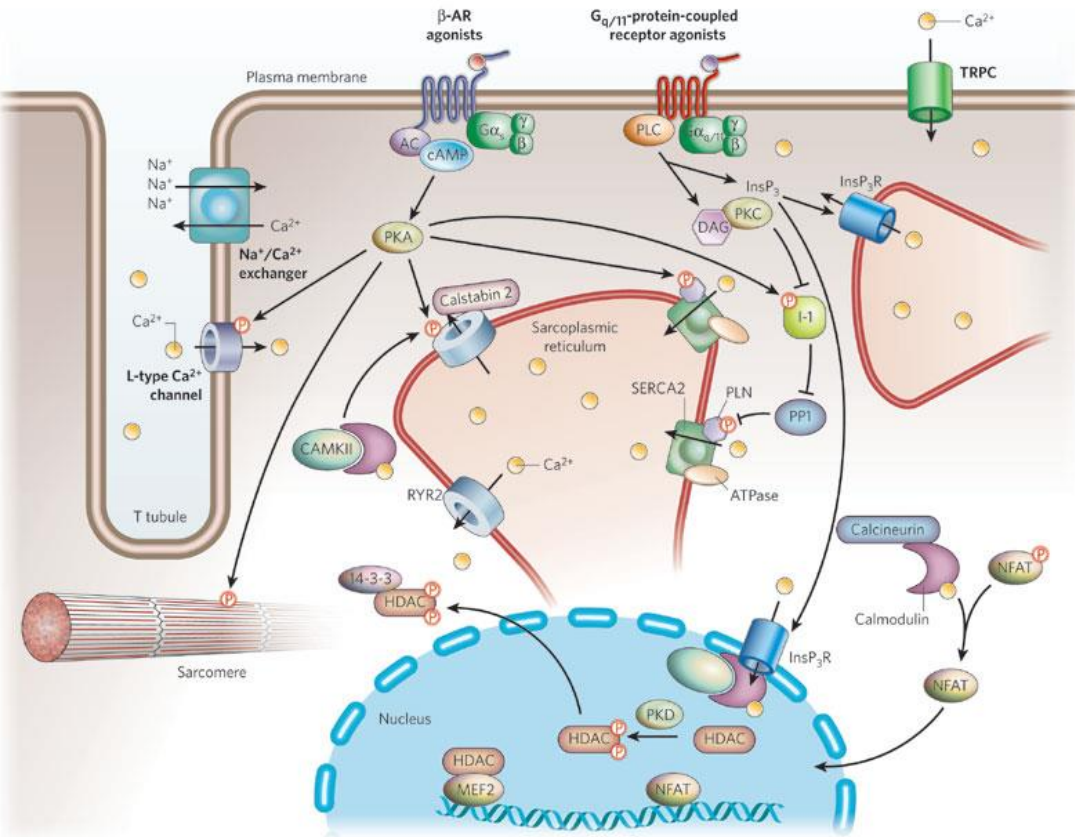


Signalling Pathways – heart disease

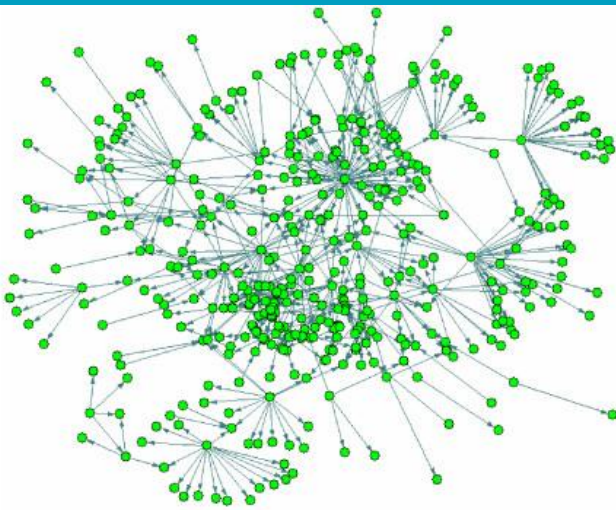
Nodes = Proteins

Edges =
Activation/Inhibition
(e.g., phosphorylation)

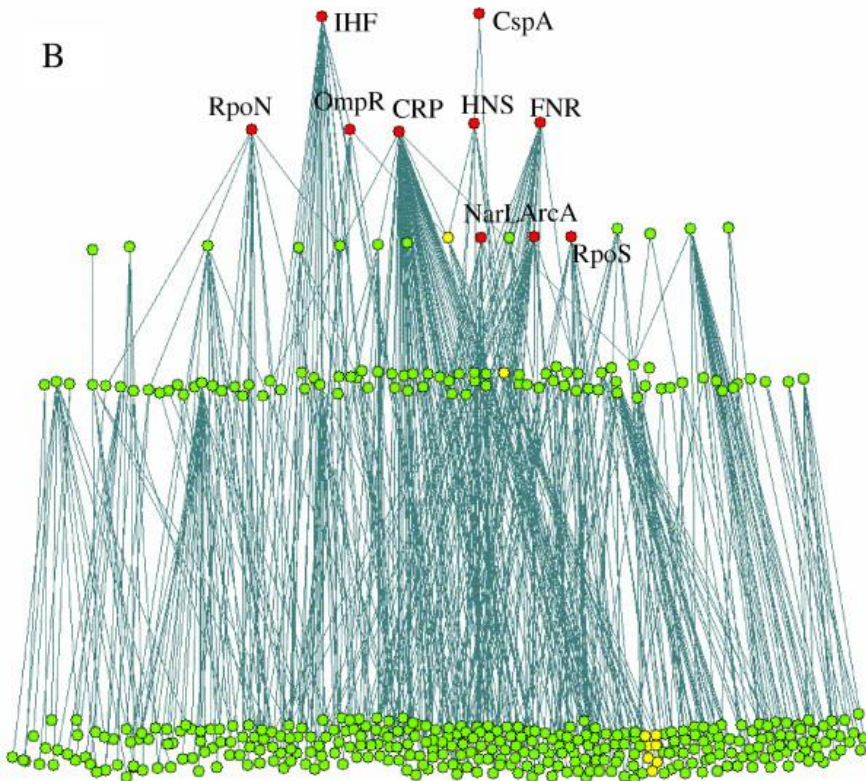
- Signalling pathways complex and 'interwoven'.
- Involve gene regulation and protein/protein interaction, etc.



A



B

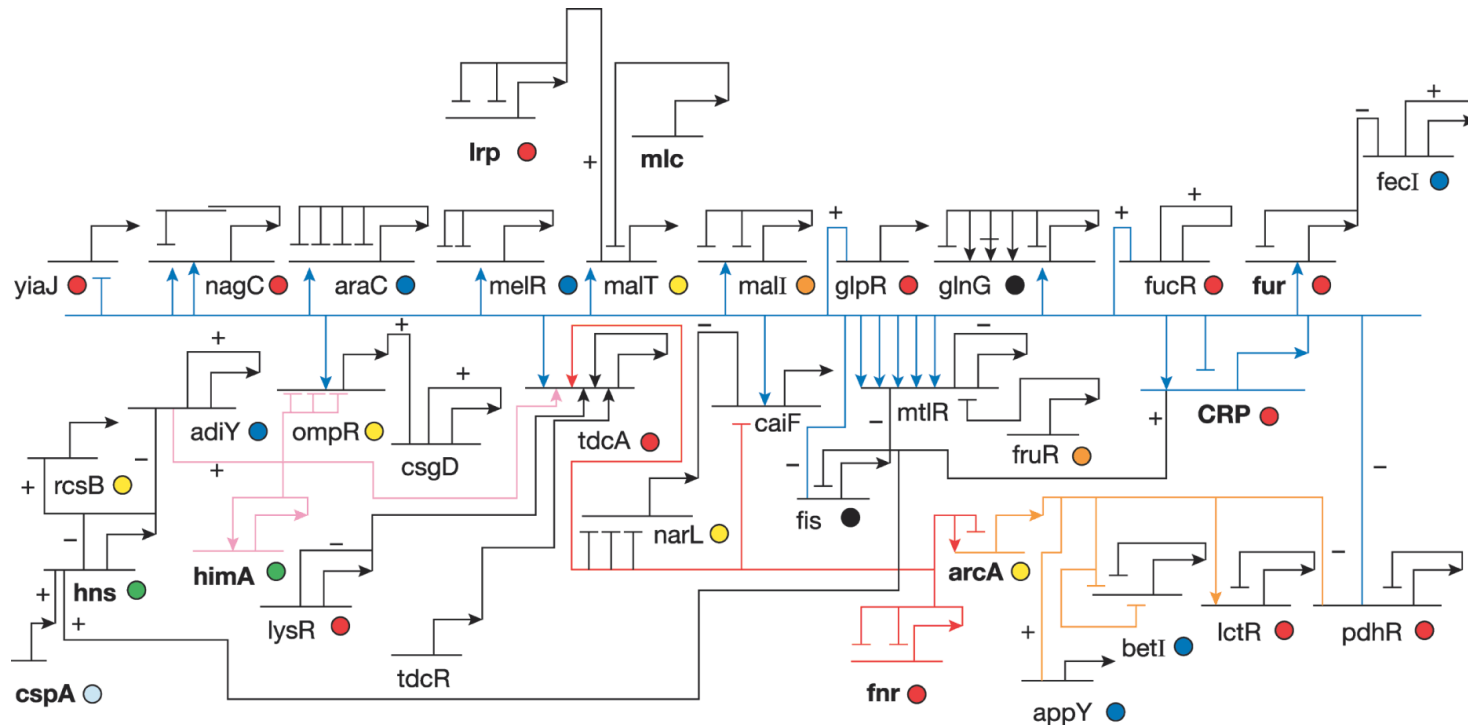


Gene Regulation Networks of *E. coli* ... a very simple organism !

Nodes = Transcription Factors
Edges = Regulatory interaction
Gene regulation is hierarchical.

Hierarchical structure and modules in the Escherichia coli transcriptional regulatory network revealed by a new top-down approach (2004) *BMC Bioinformatics* **5**: 199

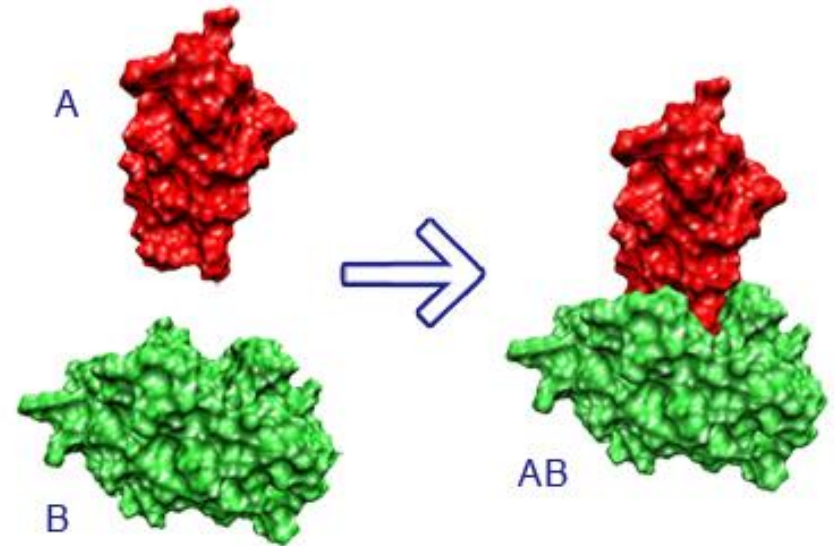
Annotated Gene Regulation Network of *E. coli*



Charting gene regulatory networks: strategies, challenges and perspectives (2004)
Biochem. Journal **381**:1-12.

Protein-Protein (P2P) Interaction Network

- Proteins interact with each other for a number of reasons:
 - Allosteric activation
 - Inhibitory binding (covering the active site)
 - Phosphorylation or modification of residues
 - Physical closeness of enzymes along metabolic pathway.
 - Etc.

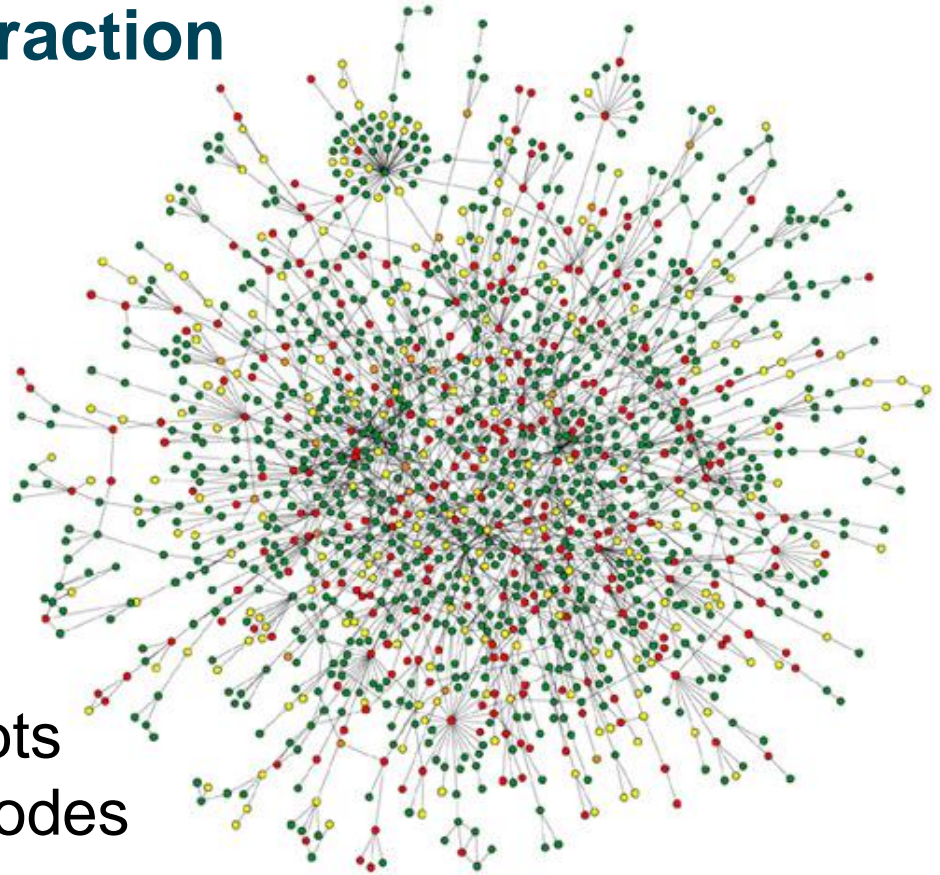


Protein-Protein (P2P) Interaction Networks for Yeast

Nodes = Protein

Edges = Physical interaction between proteins

- Thought to be a 'robust scale-free network' (a few 'fragile' nodes with lots of connections and many nodes with few connections).
- Provides an indication of 'functional modules'



Nature Reviews | Genetics

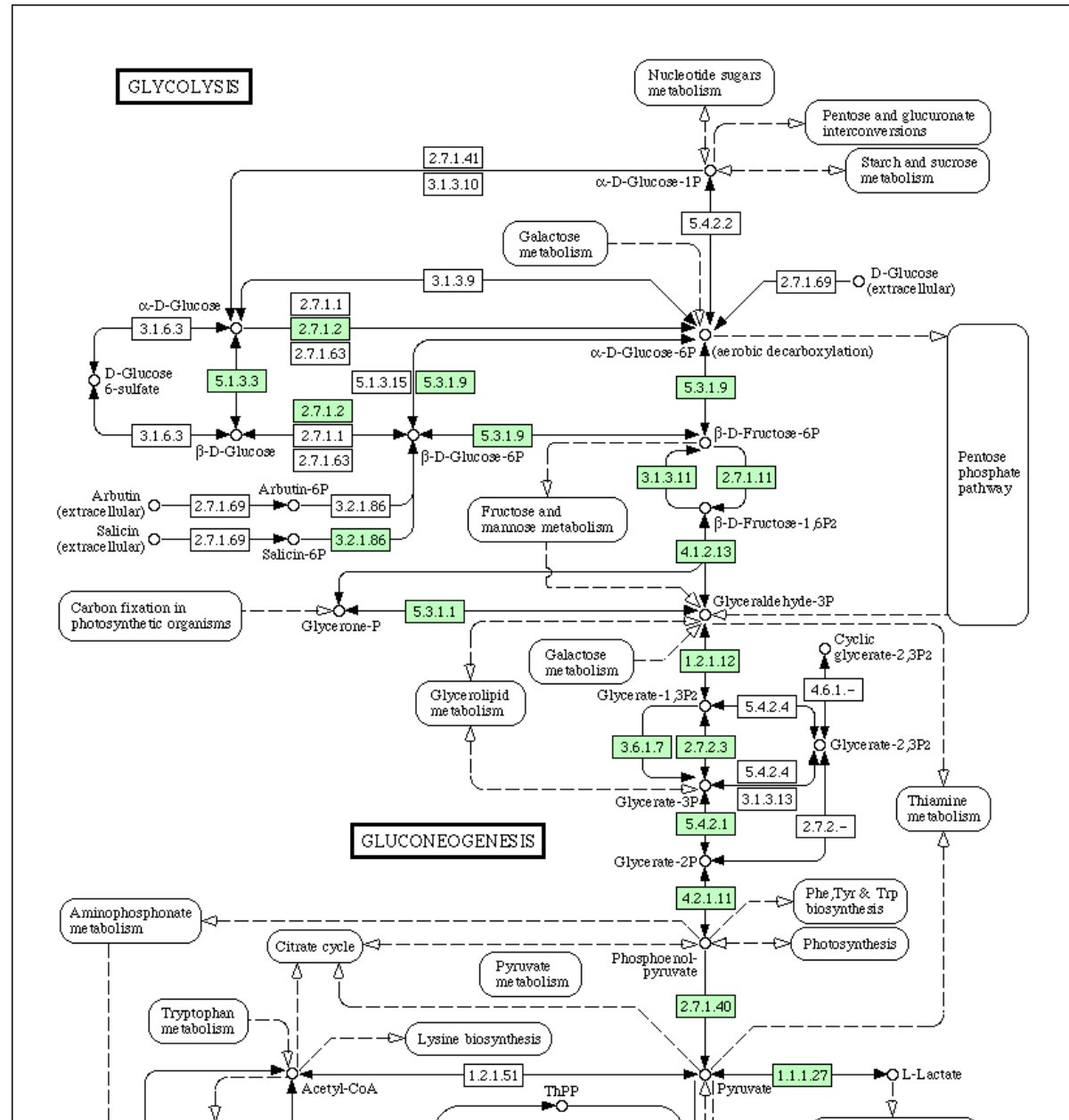
Metabolic Pathways

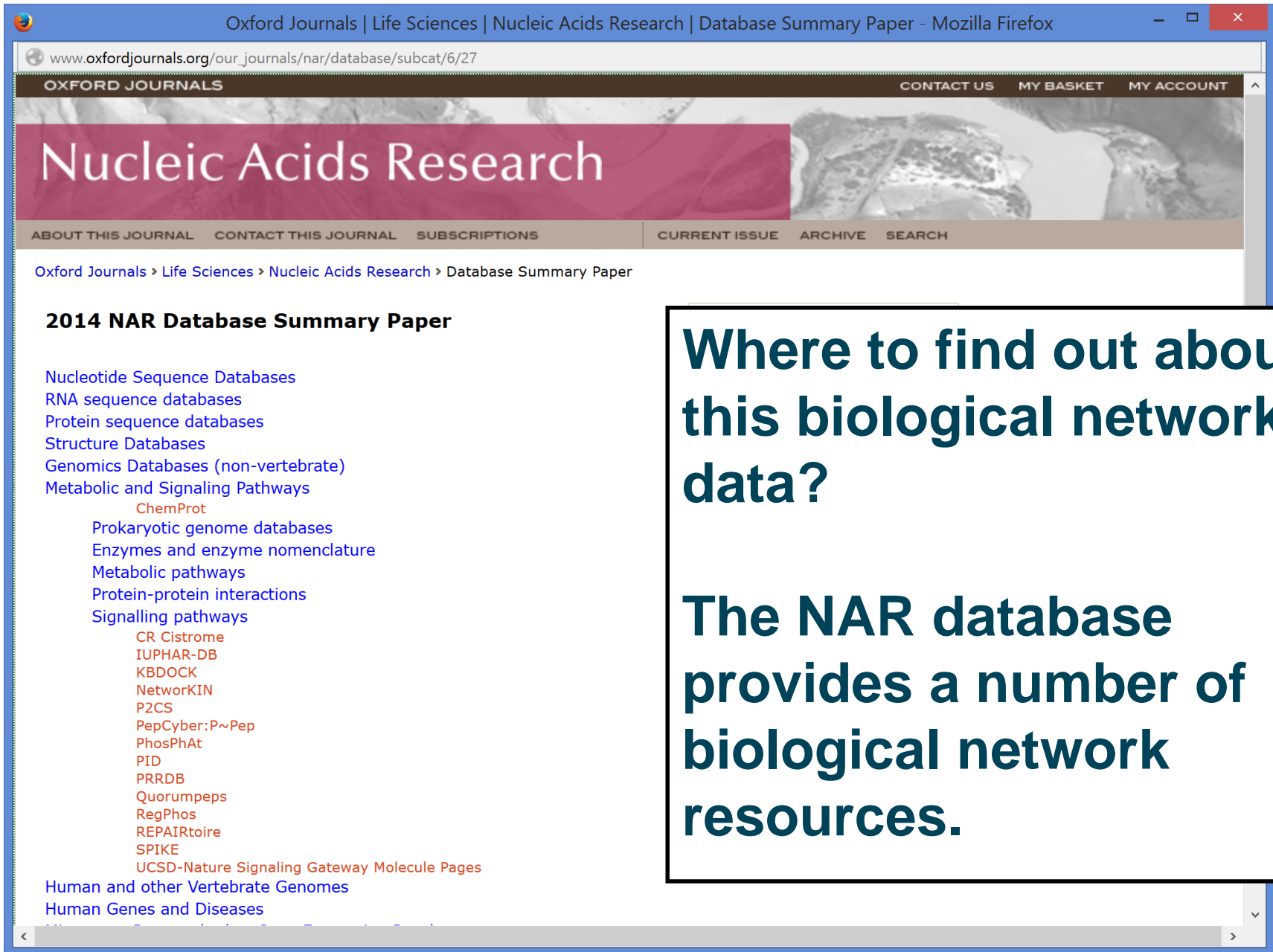
Node = metabolite
(small molecule such as glucose)

Edge = conversion of one metabolite to another via enzyme

This shows a small part of glucose breakdown in a bacteria

The concentration of enzymes are regulated via the gene regulatory network (i.e. all the pathways are linked)





Oxford Journals | Life Sciences | Nucleic Acids Research | Database Summary Paper - Mozilla Firefox

www.oxfordjournals.org/our_journals/nar/database/subcat/6/27

OXFORD JOURNALS CONTACT US MY BASKET MY ACCOUNT

Nucleic Acids Research

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

Oxford Journals > Life Sciences > Nucleic Acids Research > Database Summary Paper

2014 NAR Database Summary Paper

- Nucleotide Sequence Databases
- RNA sequence databases
- Protein sequence databases
- Structure Databases
- Genomics Databases (non-vertebrate)
- Metabolic and Signaling Pathways
 - ChemProt
- Prokaryotic genome databases
- Enzymes and enzyme nomenclature
- Metabolic pathways
- Protein-protein interactions
- Signalling pathways
 - CR Cistrome
 - IUPHAR-DB
 - KBDOCK
 - NetworkIN
 - P2CS
 - PepCyber:P~Pep
 - PhosPhAt
 - PID
 - PRRDB
 - Quorumpeps
 - RegPhos
 - REPAIRtoire
 - SPIKE
 - UCSD-Nature Signaling Gateway Molecule Pages
- Human and other Vertebrate Genomes
- Human Genes and Diseases

Where to find out about this biological network data?

The NAR database provides a number of biological network resources.

Another source of information about biological networks is the PathGuide resource

- Biological pathway resources given at:
<http://www.pathguide.org/>
- Currently 547 resources for:
 - Signalling Pathways
 - Gene Regulation Pathways
 - Protein-Protein Interactions
 - Metabolic Pathways
 - Protein-Compound Interactions

Pathguide: the pathway resource list - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://www.pathguide.org/

Pathguide: the pathway resource list

Home BioPAX cBio MSKCC

Pathguide»the pathway resource list

Navigation

- Protein-Protein Interactions
- Metabolic Pathways
- Signaling Pathways
- Pathway Diagrams
- Transcription Factors / Gene Regulatory Networks
- Protein-Compound Interactions
- Genetic Interaction Networks
- Protein Sequence Focused
- Other

Find: KEGG

Complete Listing of All Pathguide Resources

News

Pathguide: the pathway resource list - Mozilla Firefox

File Edit View History Bookmarks Tools Help

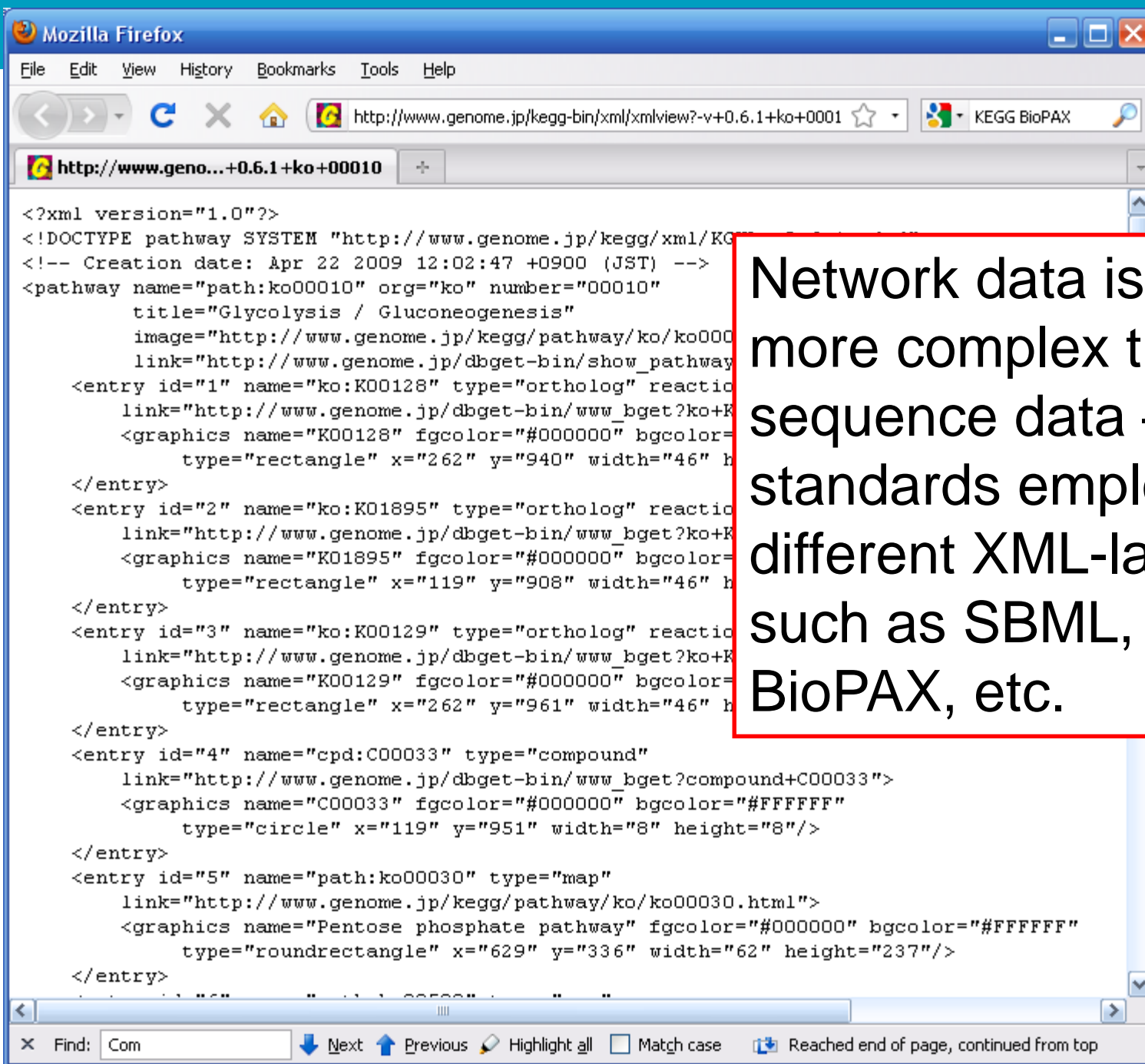
http://www.pathguide.org/

Pathguide: the pathway resource list

GeneNet - Genetic Networks	Details	Free	
GenMAPP - Gene MicroArray Pathway Profiler	Details	Free	
GO - Gene Ontology	Details	Free	
GOLD.db - Genomics of Lipid-associated Disorders	Details	Free	
HMDB - Human Metabolome Database	Details	Free	
IMG - Integrated Microbial Genomes	Details	Free	
Indigo - Gene Neighborhoods and Codon Usage	Details	X	
IntEnz - Integrated relational Enzyme database	Details	Free	
iPath - Invitrogen iPath	Details	\$	
JWS Online - Online Cellular Systems Modelling	Details	Free	SBML
KEGG - Kyoto Encyclopedia of Genes and Genomes	Details	Free	BioPAX
LIGAND - Database of Chemical Compounds and Reactions in Biological Pathways	Details	Free	
Malaria - Malaria Parasite Metabolic Pathways	Details	Free	
MetaCore - MetaCore pathway database	Details	\$	
MetaCyc - Metabolic Pathway Database	Details	Free	BioPAX SBML
MetNetDB - Metabolic Network Exchange	Details	Free	
Millipore Pathways - Your online source for visualizing metabolic and signaling	Details	Free	

Find: KEGG

Next Previous Highlight all Match case



Network data is much more complex than sequence data – many standards employing different XML-languages such as SBML, PSI-MI, BioPAX, etc.

Biological networks/pathways provide ‘global’ views about how the biological parts/components within a cell interact ...

- But what information is available on the biological components themselves?
- Well all the classic data (sequences, structures, function annotation, etc.) is generally available via established databases (UniProt, PDB, etc.) as you would have covered with David Jones.
- In addition, there is **high-throughput ‘-omics’ data** that provides ‘global’ information about particular types of components within a cell (e.g. the concentrations of all mRNAs within a cell in a particular situation).

Introduction to HT –omics data

- *Genomics*
- *Transcriptomics*
- *Proteomics*
- *Metabolomics*

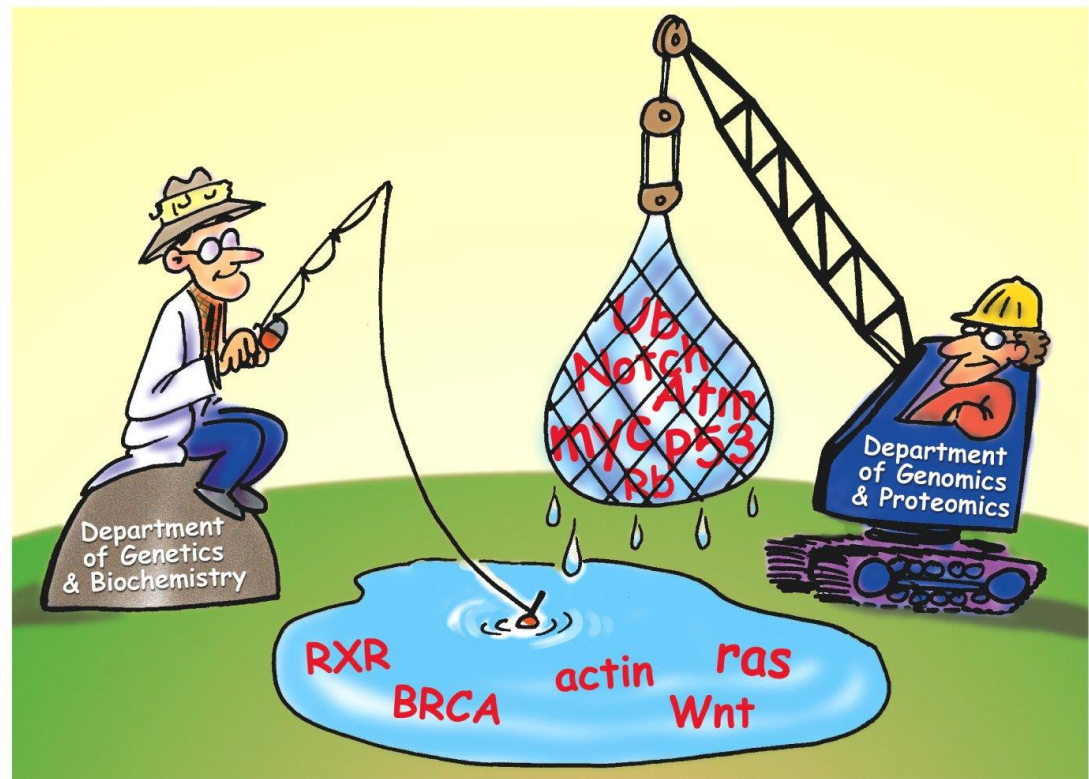


Figure 2-46 Cell and Molecular Biology, 4/e (© 2005 John Wiley & Sons)

Overview ... recall the central dogma of molecular biology ...

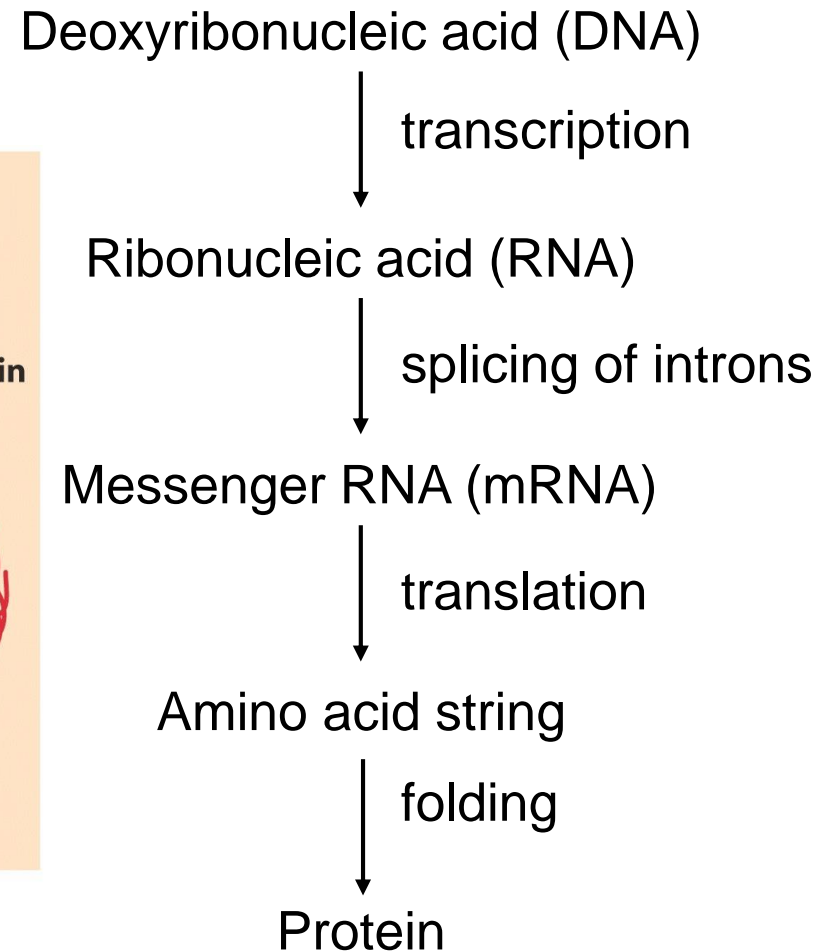
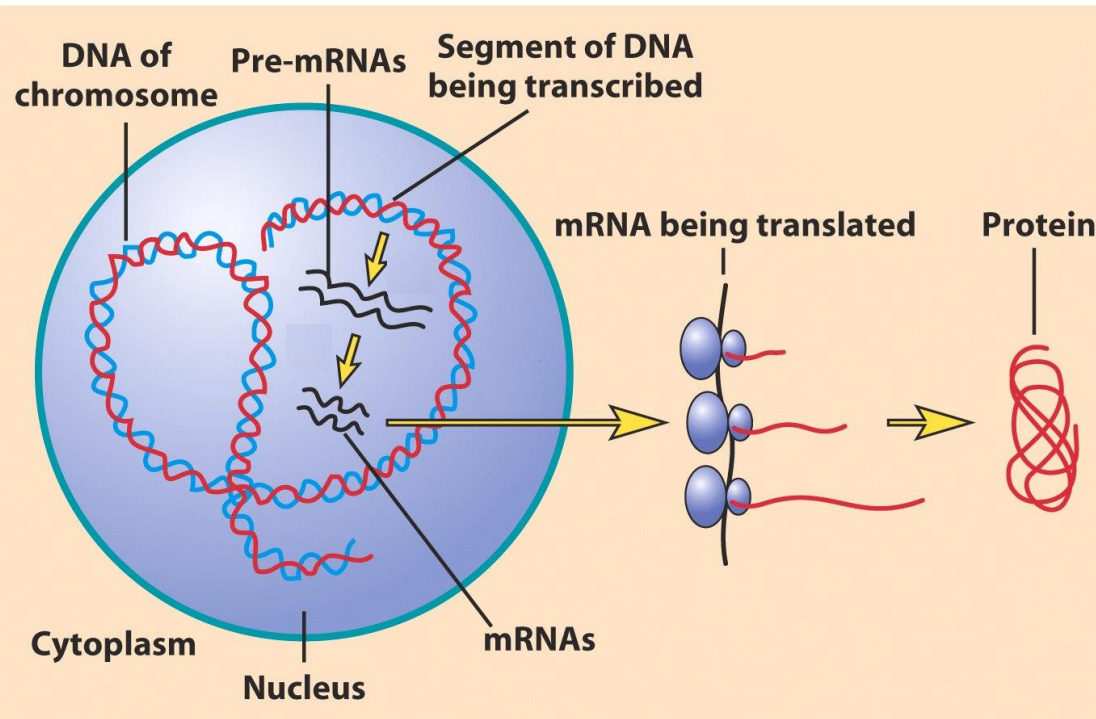


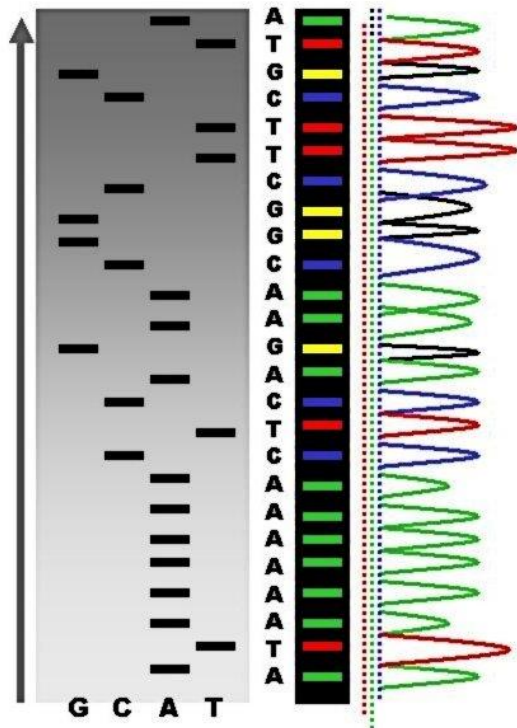
Figure 11-2 Cell and Molecular Biology, 4/e (© 2005 John Wiley & Sons)

Gen-omics

- Genomics is about determining the complete sequence of the DNA genome of an organism.
- The first complete genome ever sequenced was the phage virus Φ -X174 with approx' 5000 bp in 1977 by Fred Sanger.
- The first 'free living' organism was the bacteria *H. Influenza* with approx' 1.8 Mbp in 1995.
- Human genome with approx' 3 Gbp was finally sequenced in April 2003.

The Sanger method has dominated nucleotide sequencing since 1977 ...

Sequence DNA fragments and then assembly these parts into a complete genome.





Firefox

www.oxfordjournals.org/nar/database/subcat/7/28

OUP Oxford Journals | Life Sciences | Nucleic ...

Nucleic Acids Research

ABOUT THIS JOURNAL CONTACT THIS JOURNAL SUBSCRIPTIONS CURRENT ISSUE ARCHIVE SEARCH

Oxford Journals > Life Sciences > Nucleic Acids Research > Database Summary Paper

2012 NAR Database Summary Paper

[Nucleotide Sequence Databases](#)
[RNA sequence databases](#)
[Protein sequence databases](#)
[Structure Databases](#)
[Genomics Databases \(non-vertebrate\)](#)
[Metabolic and Signaling Pathways](#)
[Human and other Vertebrate Genomes](#)
[Model organisms, comparative genomics](#)

[ABA - Ascidian Body Atlas](#)
[ACeDB](#)
[AgBase](#)
[Animal Genome Size Database](#)
[AnimalQTLdb](#)
[ArkDB](#)
[BodyMap](#)
[BodyMap-Xs](#)
[Bovine Genome](#)
[cBARBEL](#)
[CCDS](#)
[ChickVD](#)
[CleanEST](#)
[COG - Eukaryotic Orthologous Groups of proteins](#)
[CORC - A database for Comparative Regulatory Genomics](#)
[Cre Transgenic Database](#)
[DBTGR](#)
[diArk](#)
[Edinburgh Mouse \(EMAP\) Atlas](#)
[EGO - Eukaryotic Gene Orthologs](#)
[EMMA](#)
[Ensembl](#)
[Ensembl Genomes](#)
[Entrez Gene](#)
[euGenes](#)

[Compilation Paper](#)
[Category List](#)
[Alphabetical List](#)
[Category/Paper List](#)
[Search Summary Papers](#)

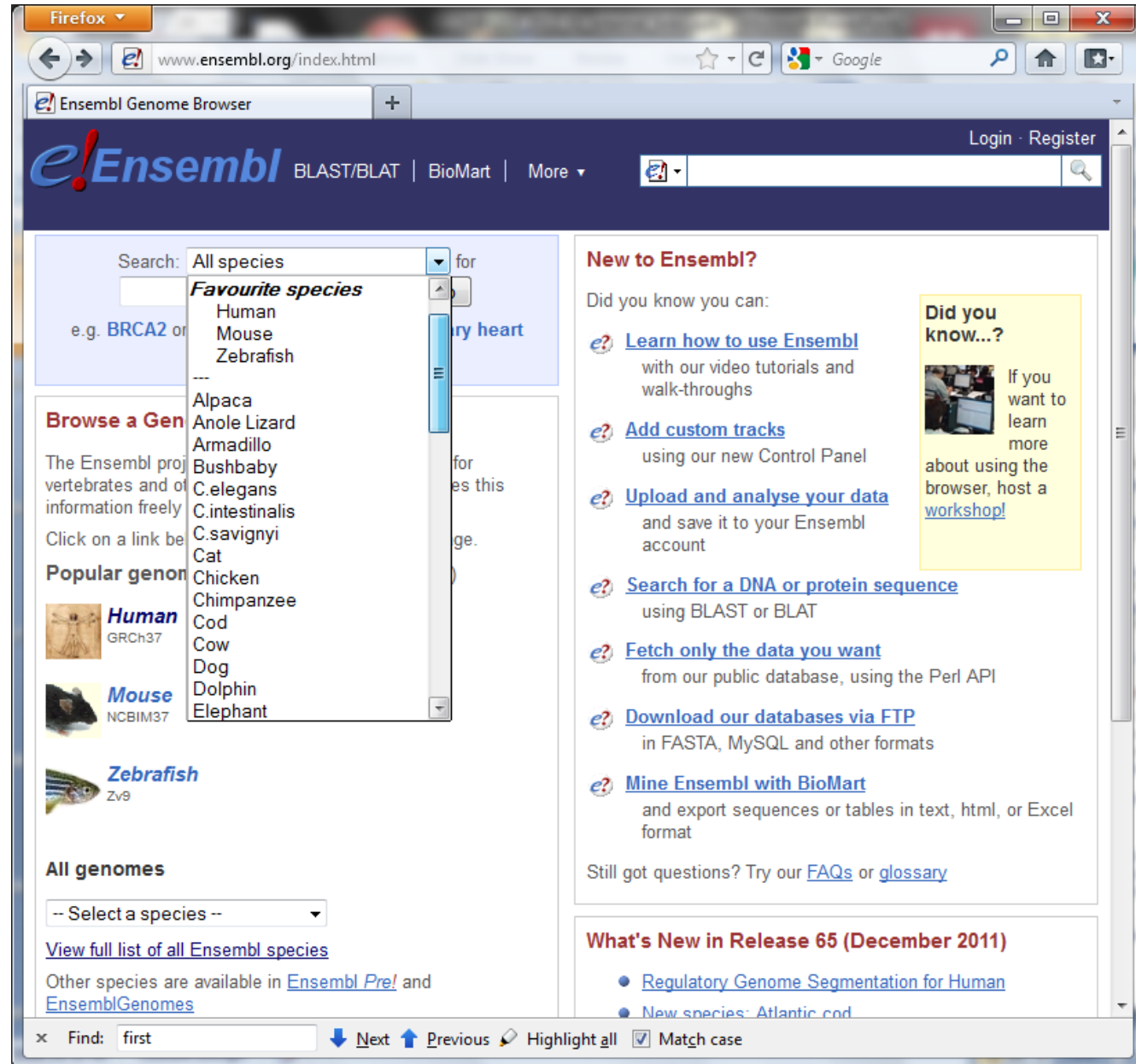
Find: first

Next Previous Highlight all Match case

Looking up genomics data in the NAR list of resources ... many model organism comparative genomics databases including Ensembl

Ensembl genomics resource

- Provides sequence, function and other relevant data organized around the framework of an organism's genome.



The screenshot shows the Ensembl Genome Browser interface. The browser window title is "Firefox" and the address bar shows "www.ensembl.org/index.html". The page header includes the Ensembl logo, navigation links for "BLAST/BLAT", "BioMart", and "More", and a search bar. The main content area is divided into several sections:

- Search:** A dropdown menu for "All species" is open, showing a list of "Favourite species" (Human, Mouse, Zebrafish) and a scrollable list of other species (Alpaca, Anole Lizard, Armadillo, Bushbaby, C.elegans, C.intestinalis, C.savignyi, Cat, Chicken, Chimpanzee, Cod, Cow, Dog, Dolphin, Elephant).
- Browse a Genome:** A section with a description: "The Ensembl project provides vertebrate genome information freely available to all. Click on a link below to browse a genome for this species." It lists "Human" (GRCh37), "Mouse" (NCBIM37), and "Zebrafish" (Zv9) with corresponding icons.
- Popular genomes:** A section with a description: "Click on a link below to browse a genome for this species." It lists "Human" (GRCh37), "Mouse" (NCBIM37), and "Zebrafish" (Zv9) with corresponding icons.
- All genomes:** A section with a dropdown menu for "Select a species" and a link to "View full list of all Ensembl species". It also mentions "Other species are available in Ensembl Pre! and EnsemblGenomes".
- New to Ensembl?:** A section with a list of links: "Learn how to use Ensembl" (with video tutorials and walk-throughs), "Add custom tracks" (using the new Control Panel), "Upload and analyse your data" (and save it to your Ensembl account), "Search for a DNA or protein sequence" (using BLAST or BLAT), "Fetch only the data you want" (from our public database, using the Perl API), "Download our databases via FTP" (in FASTA, MySQL and other formats), and "Mine Ensembl with BioMart" (and export sequences or tables in text, html, or Excel format). It also includes a link to "FAQs or glossary".
- Did you know...?:** A yellow box with a link to "workshop!" and a description: "If you want to learn more about using the browser, host a workshop!".
- What's New in Release 65 (December 2011):** A section with a list of links: "Regulatory Genome Segmentation for Human" and "New species: Atlantic cod".

The footer of the page includes a search bar with "Find: first" and navigation links for "Next", "Previous", "Highlight all", and "Match case".

Firefox

www.ensembl.org/Homo_sapiens/Location/View?g=ENSG00000130203;r=19:45408956-4541261

Ensembl genome browser 65: Homo sap...

e!Ensembl BLAST/BLAT | BioMart | Tools | Downloads | More

Human (GRCh37) Location: 19:45,408,956-45,412,650 Gene: APOE

Login · Register

Location-based displays

- Whole genome
- Chromosome summary
- Region overview
- Region in detail**
- Comparative Genomics
 - Alignments (image) (58)
 - Alignments (text) (58)
 - Multi-species view (54)
 - Synteny (15)
- Genetic Variation
 - Resequencing (20)
 - Linkage Data
- Markers
- Other genome browsers
 - UCSC
 - NCBI
 - Vega

Configure this page

Manage your data

Export data

Bookmark this page

Chromosome 19: 45,408,956-45,412,650

Assembly excepti...

chromosome 19

Assembly excepti...

HSCHR19_1_CT G3

HG730_PATCH HSCHR19_1_CT G3_1

HSCHR19_2_CT G3

HSCHR19_3_CT G3

Export Image

Region in detail [help](#)

Chromosome bands

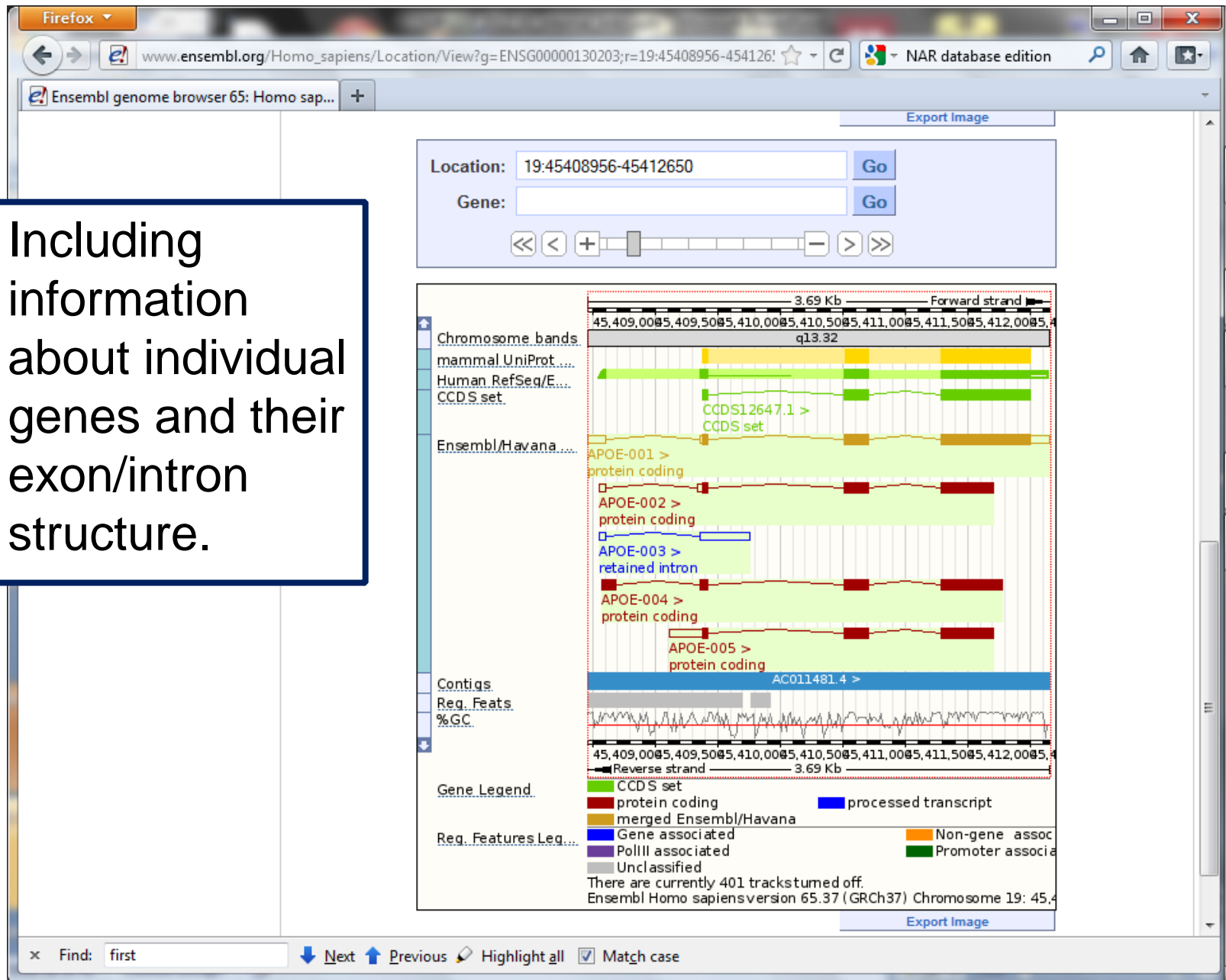
Contigs

Ensembl/Havana...

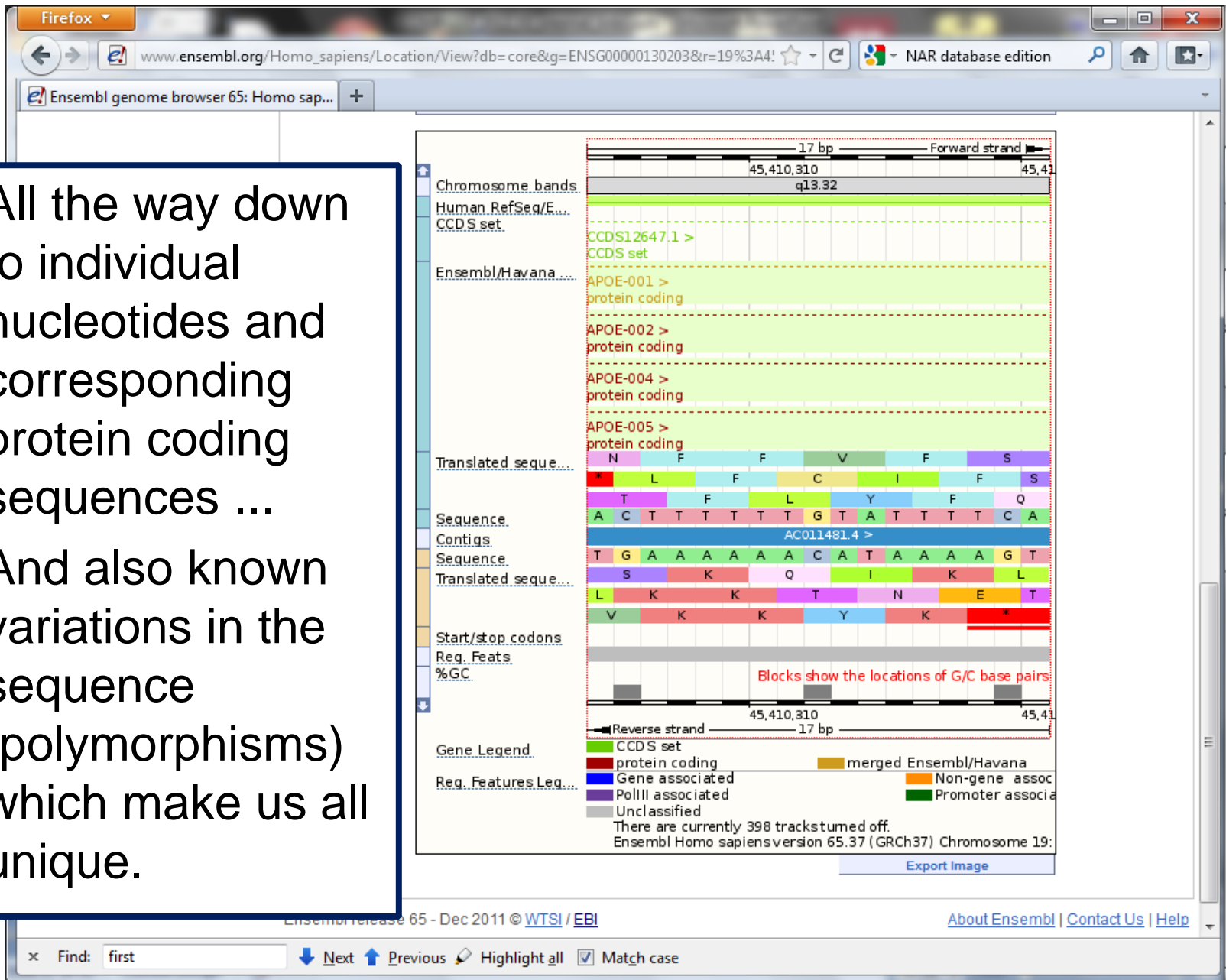
Ensembl/Havana...

Find: first

Next Previous Highlight all Match case



- All the way down to individual nucleotides and corresponding protein coding sequences ...
- And also known variations in the sequence (polymorphisms) which make us all unique.



Genome variation within a species

- Even simple phenotypic traits result from complex interactions between many genes/proteins.
- The dog genome was initiated to try to unravel these relationships between:
genotype \leftrightarrow phenotype



Transcript-omics

- Transcriptomics is about globally measuring the expression level of all the different mRNAs within a tissue.
- Essentially giving you information about the genes being expressed in a tissue.

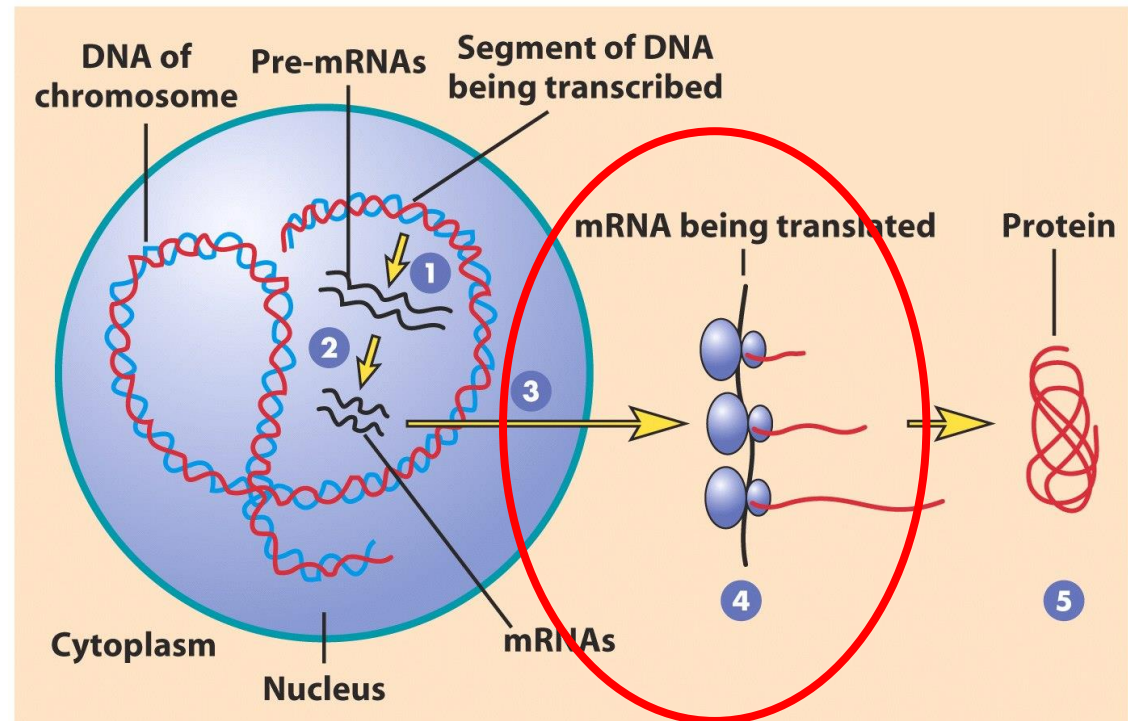
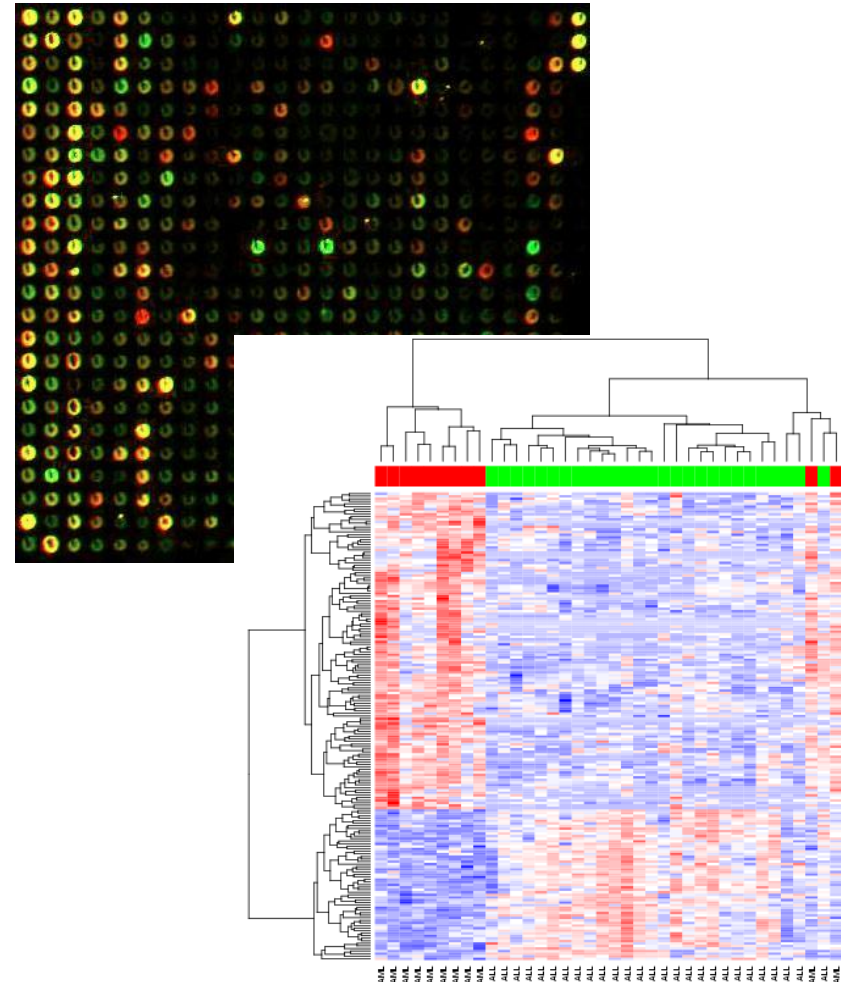


Figure 11-2 Cell and Molecular Biology, 4/e (© 2005 John Wiley & Sons)

Transcriptomics is a very mature field with data standards and analysis techniques

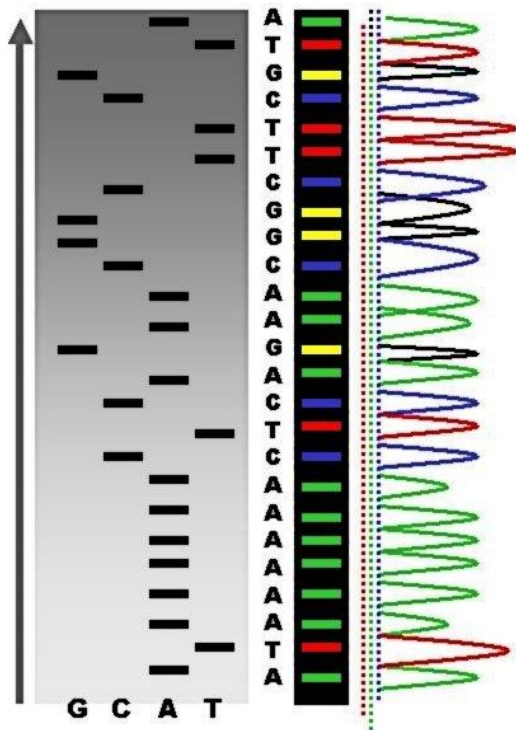
- Provides gene expression levels (mRNA) under different conditions.
- Established databases such as GEO and ArrayExpress.
- Established data formats such as MAGE-ML, etc. <http://www.fged.org/>
- Extensive open-source software such as BioConductor.



Transcriptomics based on *Next Generation Sequencing (NGS)* is currently revolutionizing the field: massively parallel sequencing ... actually counting the mRNA molecules one at a time!



Next Generation Sequencing is like having 160,000,000 Sanger machines all in a row!



Prote-omics

- Proteomics is about globally measuring the concentration/activity of all the different proteins within a tissue.
- Giving you direct info' about what protein activities are present in the tissue.

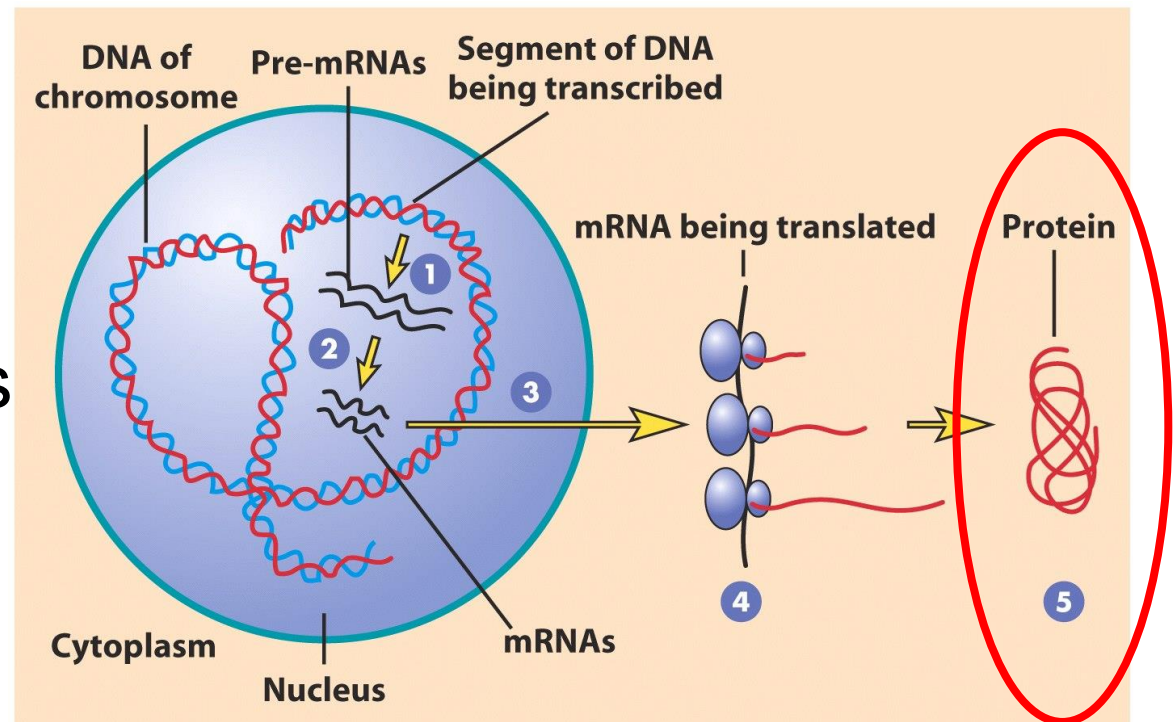
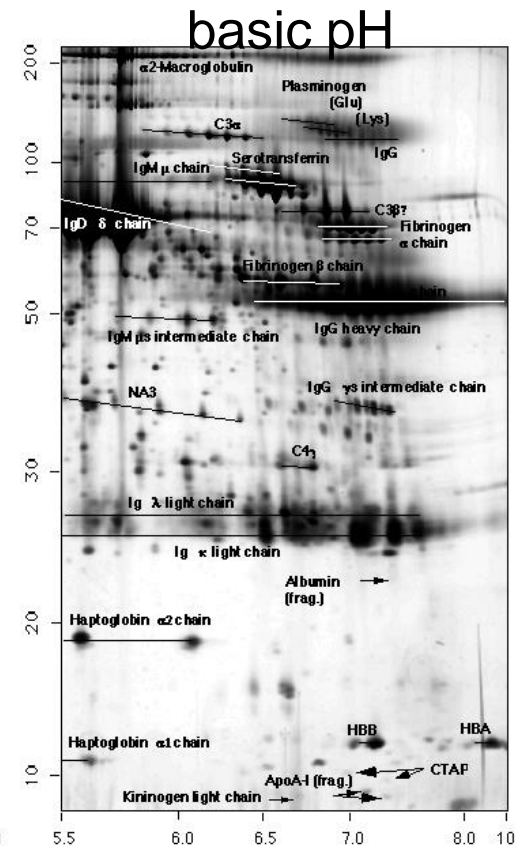
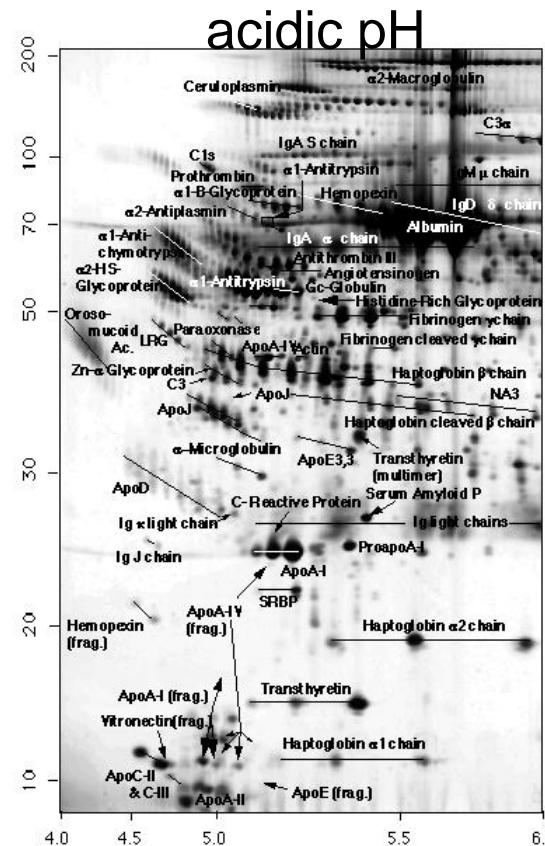


Figure 11-2 Cell and Molecular Biology, 4/e (© 2005 John Wiley & Sons)

Proteomics is starting to mature ...

- Provides information about protein concentration under different conditions.
- Established techniques such as: 2D PAGE, Mass Spec', Protein Chips.
- Proteomics standards being developed:
<http://www.psidev.info/>



Metabol-omics

- Metabolomics is about globally measuring the concentrations of all the metabolites (or small molecules) within a tissue or sample.

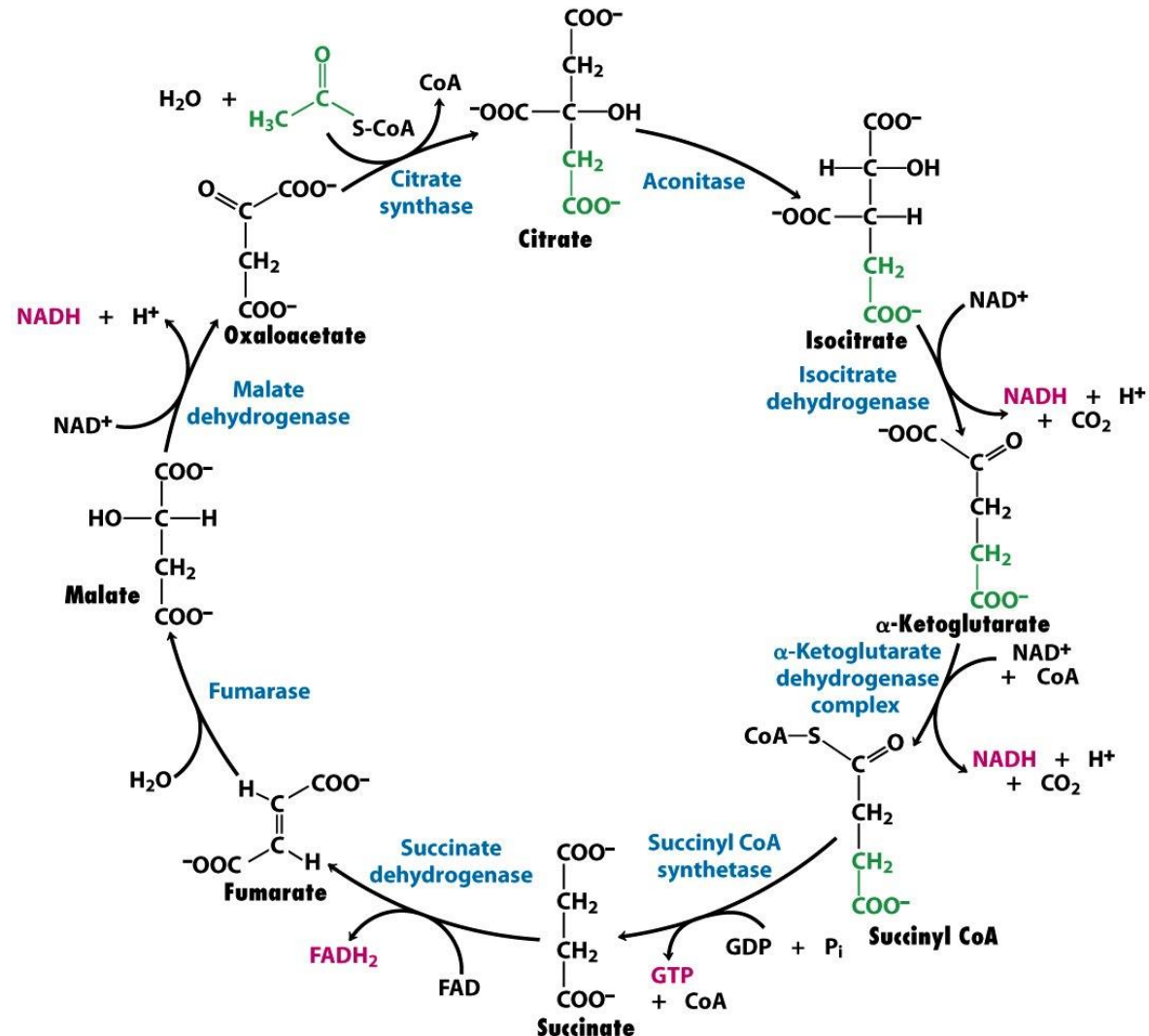
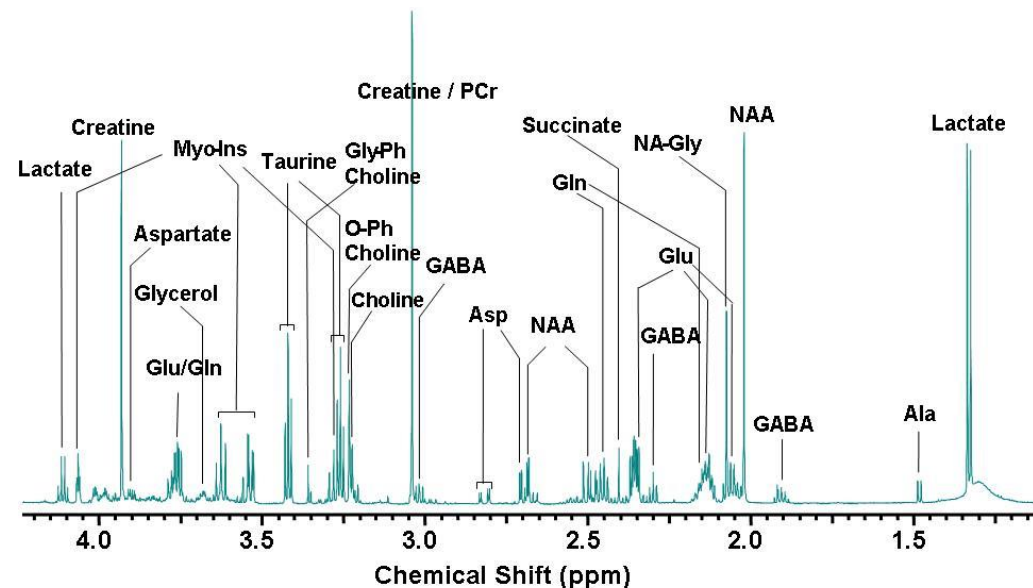


Figure 17-15
Biochemistry, Sixth Edition
 © 2007 W. H. Freeman and Company

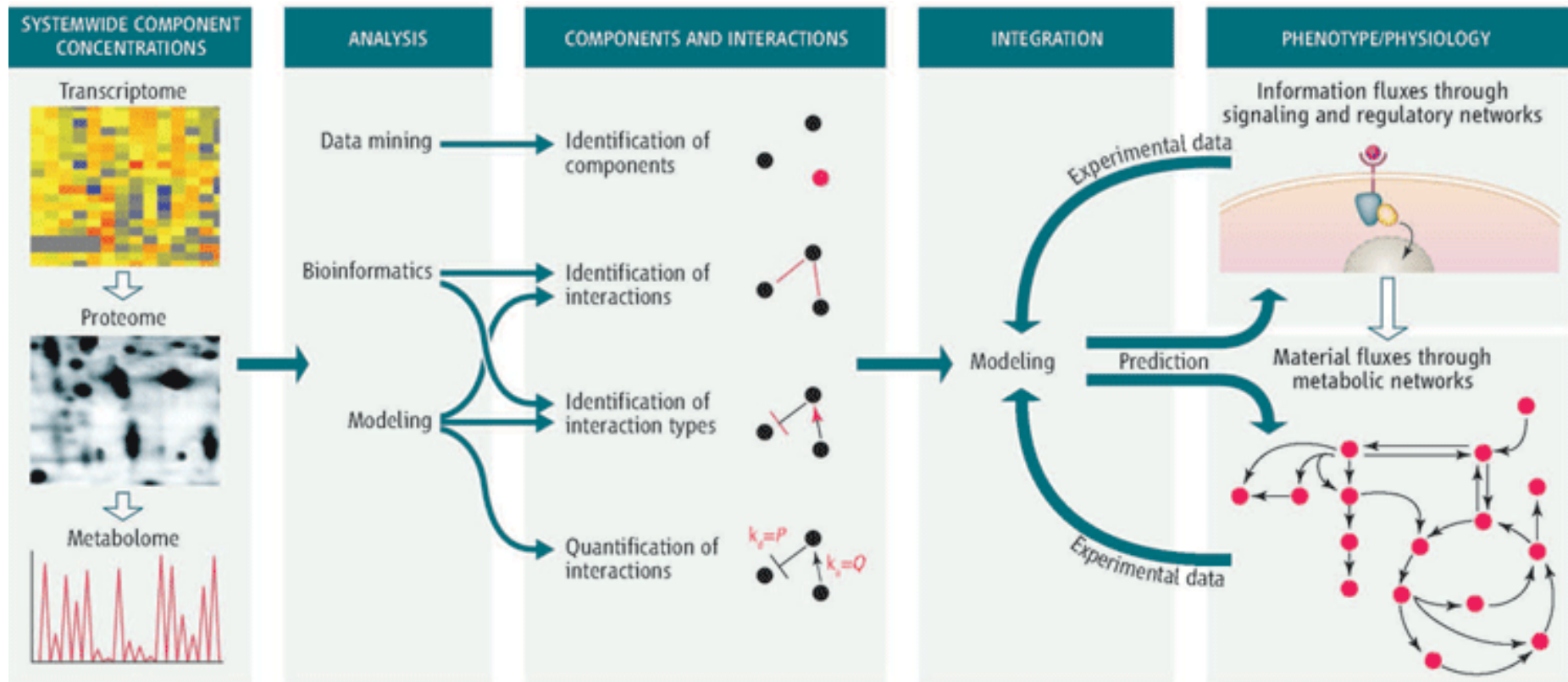
Metabolomics is also maturing rapidly ...

- Provides information about small molecule metabolites under different conditions.
- Key technique is high-resolution NMR.
- Metabolomics data standards currently under development:



<http://msi-workgroups.sourceforge.net/>

Biology is now rich with complex data – machine learning can be used to integrate and unravel complex patterns within this data ...



Getting closer to the whole picture (2007) *Science* **316**(5824), 550 – 551.

Integrating heterogeneous data with Cytoscape

Cytoscape Version 2.6.0 File Edit View Select Layout Plugins Help

Cytoscape Desktop (New Session)

Search: ESP:

Control Panel

Network VizMapper™

Network	Nodes	Edges
tp53 from IntAct	120(0)	261(0)
p53 from NCBI	4615(0)	17256(0)
DNA_Repair_(Re)	290(0)	1766(0)

Import Network From Database

Data Source Pathway Commons Web Service Client About

Search Options

Step 1: Search

TP53 Human Search

Examples: TP53, BRCA1, or SRY.

Step 2: Select

TP53

TP53 regulating kinase

TP53 activated protein 1

TP53 regulating kinase

TP53 regulated inhibitor of ...

Details

TP53

Homo sapiens

Synonyms:

- P53

Step 3: Select Network(s)

Pathway	Data Source
Cell Cycle Ch...	Reactome
G1/S DNA Da...	Reactome
p53-Depend...	Reactome
p53-Depend...	Reactome
Stabilization o...	Reactome
p53-Indepen...	Reactome
G2/M Checkp...	Reactome
G2/M DNA d...	Reactome
DNA Repair	Reactome
Double-Stran...	Reactome
Homologous ...	Reactome
Homologous r...	Reactome

> Double-click pathway to retrieve.

tp53 from IntAct

DNA_Repair_(Reactome)

p53 from NCBI

Results Panel

Node Details

Visual Legend

Filter Edges

Edge Type	Color
COMPONENT_IN_SAME	Yellow
COMPONENT_OF	Orange
CO_CONTROL_DEPENDENT_ANTI	Red
CO_CONTROL_DEPENDENT_SIMILAR	Green
CO_CONTROL_INDEPENDENT_ANTI	Pink

COMPONENT_IN_SAME

Edge is drawn if the first entity is a component of the second entity, which is a complex. This interaction is transient in the sense that A COMPONENT_OF B and B COMPONENT_OF C implies A COMPONENT_OF C. This interaction is undirected.

COMPONENT_OF

Edge is drawn if two entities belong to at least one molecular complex. This does not necessarily mean they interact directly. In a complex with n molecules, this rule will create a clique composed of $n(n-1)/2$ interactions. This interaction is directed.

CO_CONTROL_DEPENDENT_ANTI

Edge is drawn if the first and second entities have control over the same process, their control is dependent, i.e. one of them have effect over control of the other one, and their effect is in different directions (one of them activates, the other inhibits). This interaction is undirected.

CO_CONTROL_DEPENDENT_SIMILAR

Edge is drawn if the first and second entities have control over the same process, their control is dependent, i.e. one of them have effect over control of the other one, and their effect is in the same direction (both activates or both inhibits). This interaction is undirected.

CO_CONTROL_INDEPENDENT_ANTI

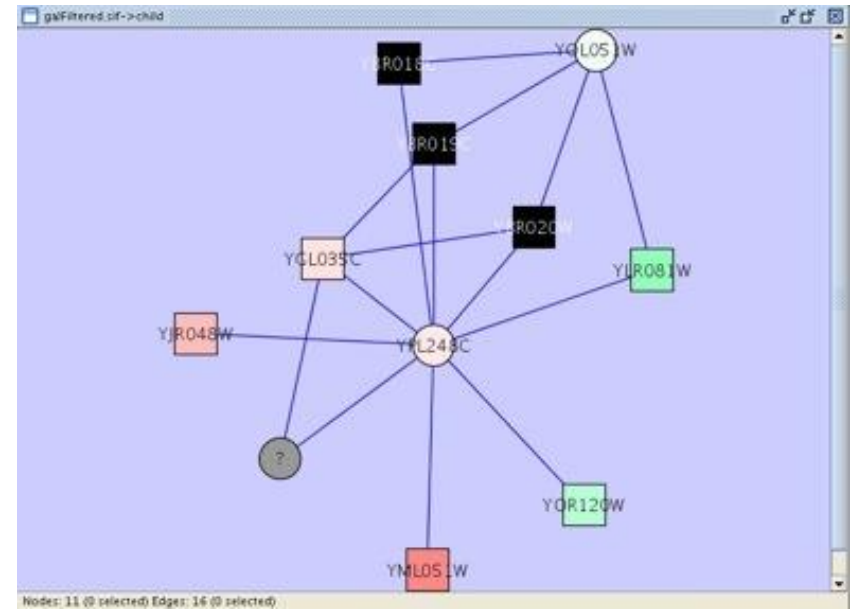
Edge is drawn if the first and second entities have control over the same process, their control is independent, i.e. they act without affecting each other's activity, and their effect

View Details

Middle-click + drag to PAN

Patterns within data (derived by hand by overlaying expression data onto gene/protein networks ...

- YML051W (Gal80) is significantly down-regulated (red square).
 - The black and green squares are significant up-regulation – these seem to be around YPL248C (Gal4).
 - But Gal4 itself is not significantly changed (white circle).
 - However, knowing Gal80 repressed Gal4 ... it can be inferred that Gal80 is down-regulated, thus up-regulating Gal4 which then up-regulates all the genes it activates.
-
- The up-regulation of Gal4 *cannot* be detected via microarrays since most transcription factors have very low expression levels (below noise level).
 - However taking the network information into account does strongly suggests that Gal4 is up-regulated.



Summary

- In this lecture we have introduced systems biology as trying to understand the behaviour of a complete biological system under study (or some aspect of this complete system).
- This requires a good understanding of the components making up a particular biological system (genes, proteins, metabolites, etc.) and how they interact.
- We introduced a number of high-throughput –omics technologies (genomics for getting a ‘parts list’, transcriptomics/proteomics/metabolomics to get a picture of the concentrations of various components).
- We also introduced biological network data (signalling pathways, gene regulation networks, protein-protein interaction networks and metabolite networks).