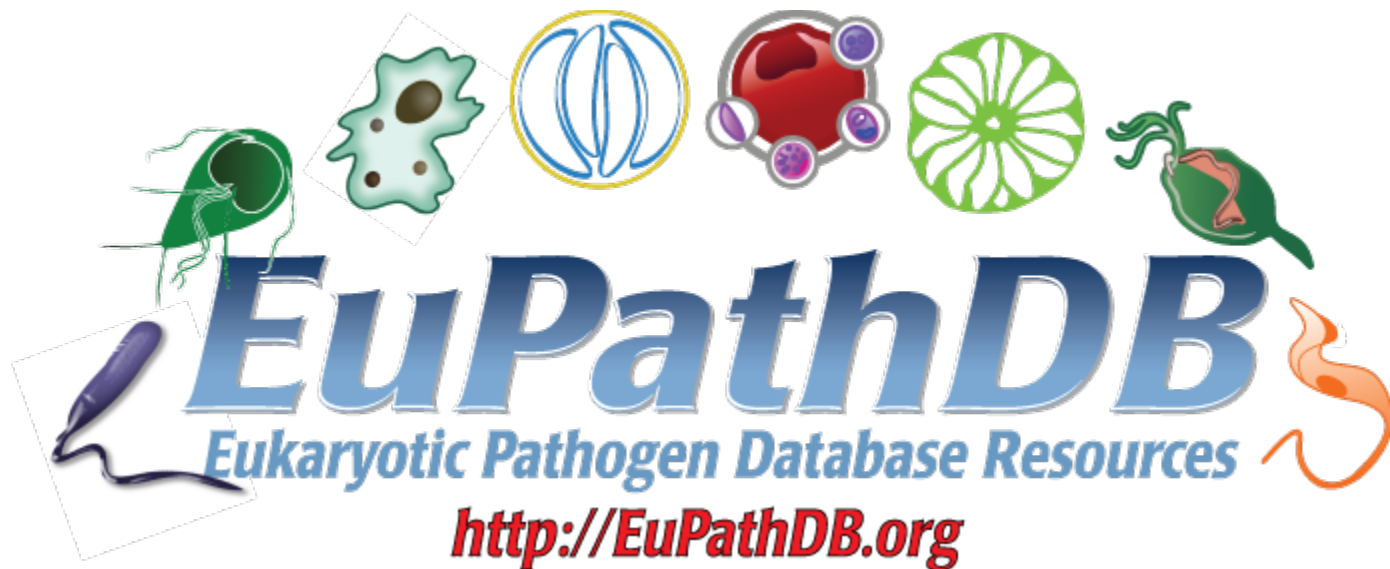


# Welcome!

## EuPathDB Workshop 2013



# Instructors & Friendly Faces

- Betsy Wenthe



- Jessie Kissinger



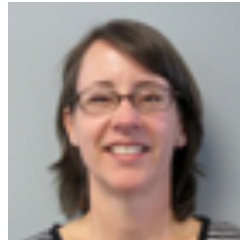
- Brian Brunk



- Omar Harb



- Susanne Warrenfeltz



- Eileen Kraemer



- Cristina Aurrecoechea

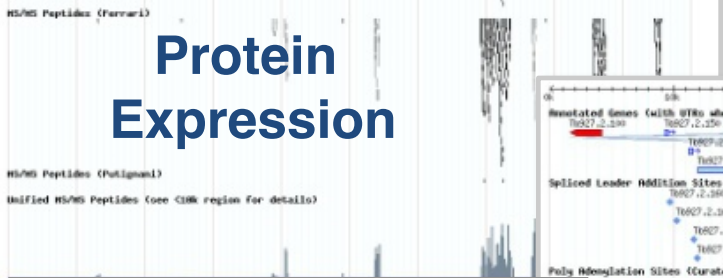


# Crash Course in Omics Terminology, Concepts & Data Types

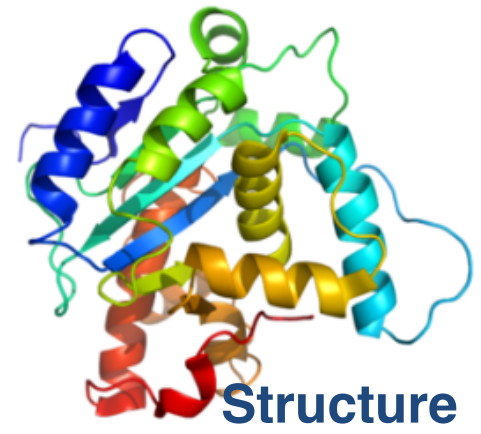
Jessie Kissinger

June 2, 2013

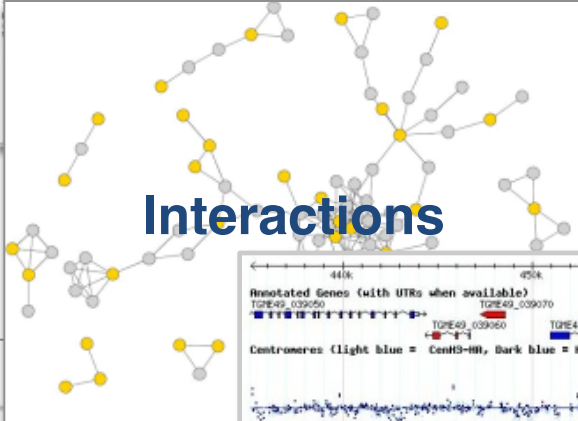
# Protein Expression



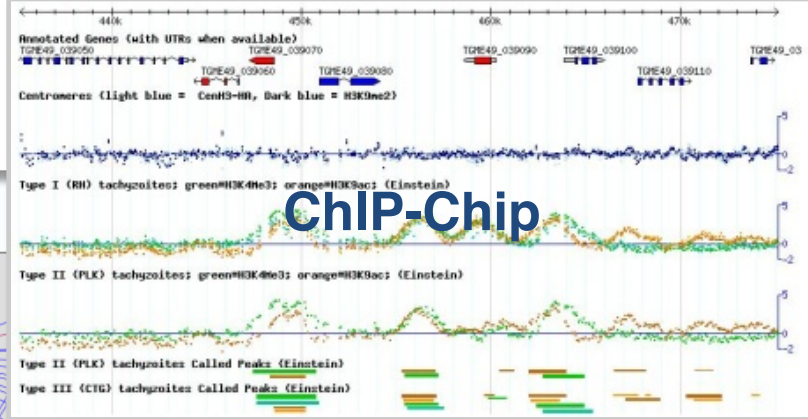
# RNA Sequencing



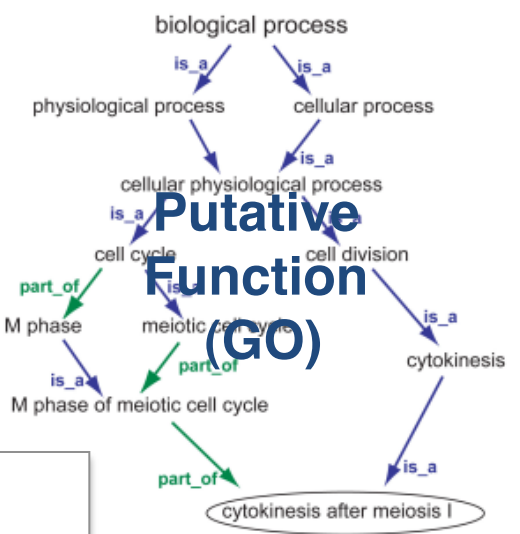
# Interactions



# ChIP-Chip



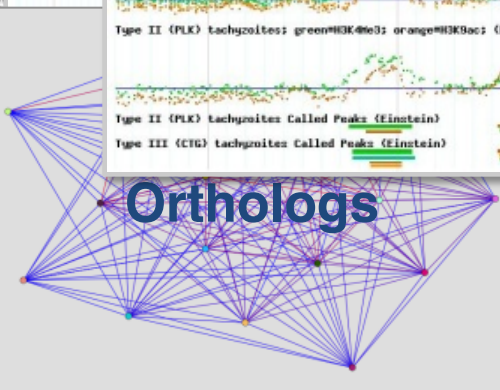
# Ex Prof



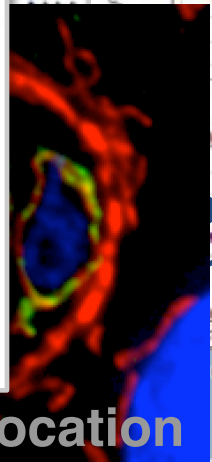
# Putative Function (GO)

# Phy

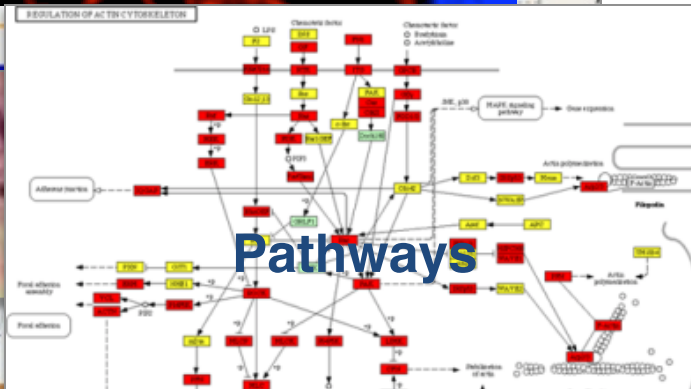
# Orthologs



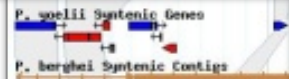
# Subcellular Location



# Pathways



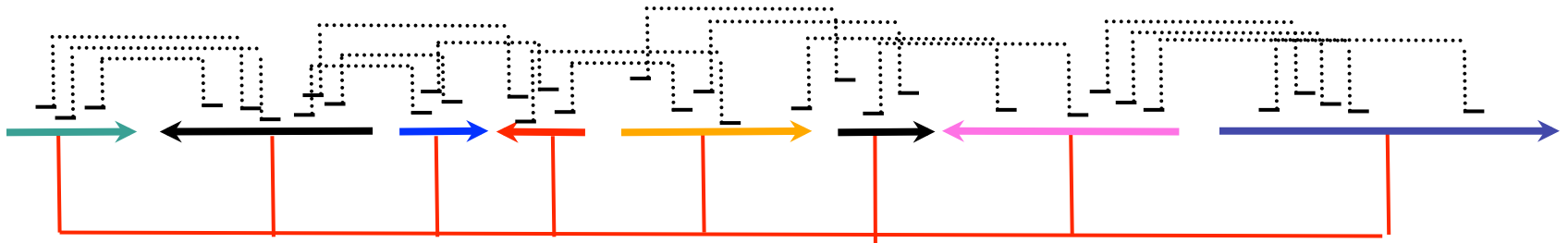
# Isolates



# Genome assembly

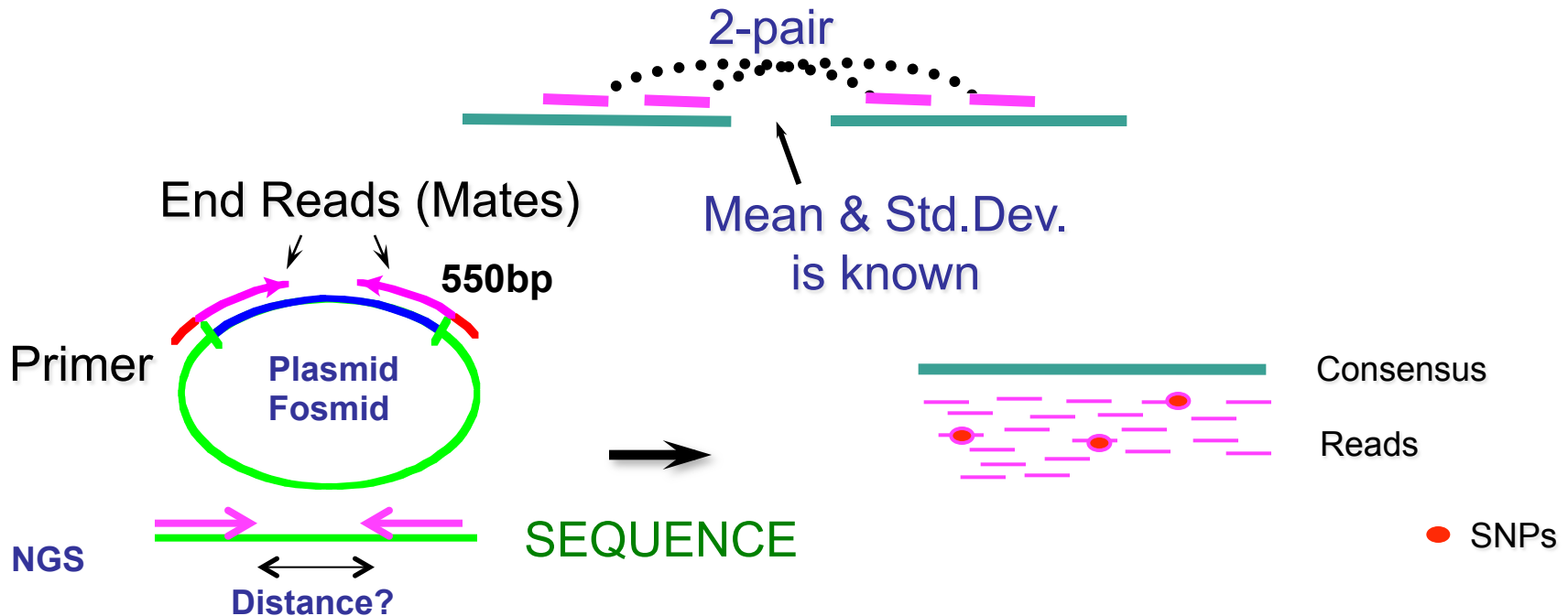
- 5X random genome shotgun
- Library insert size
- Paired-end? (Mated end pairs)
- Contigs
- Scaffolds

# Pairs Give Order & Orientation

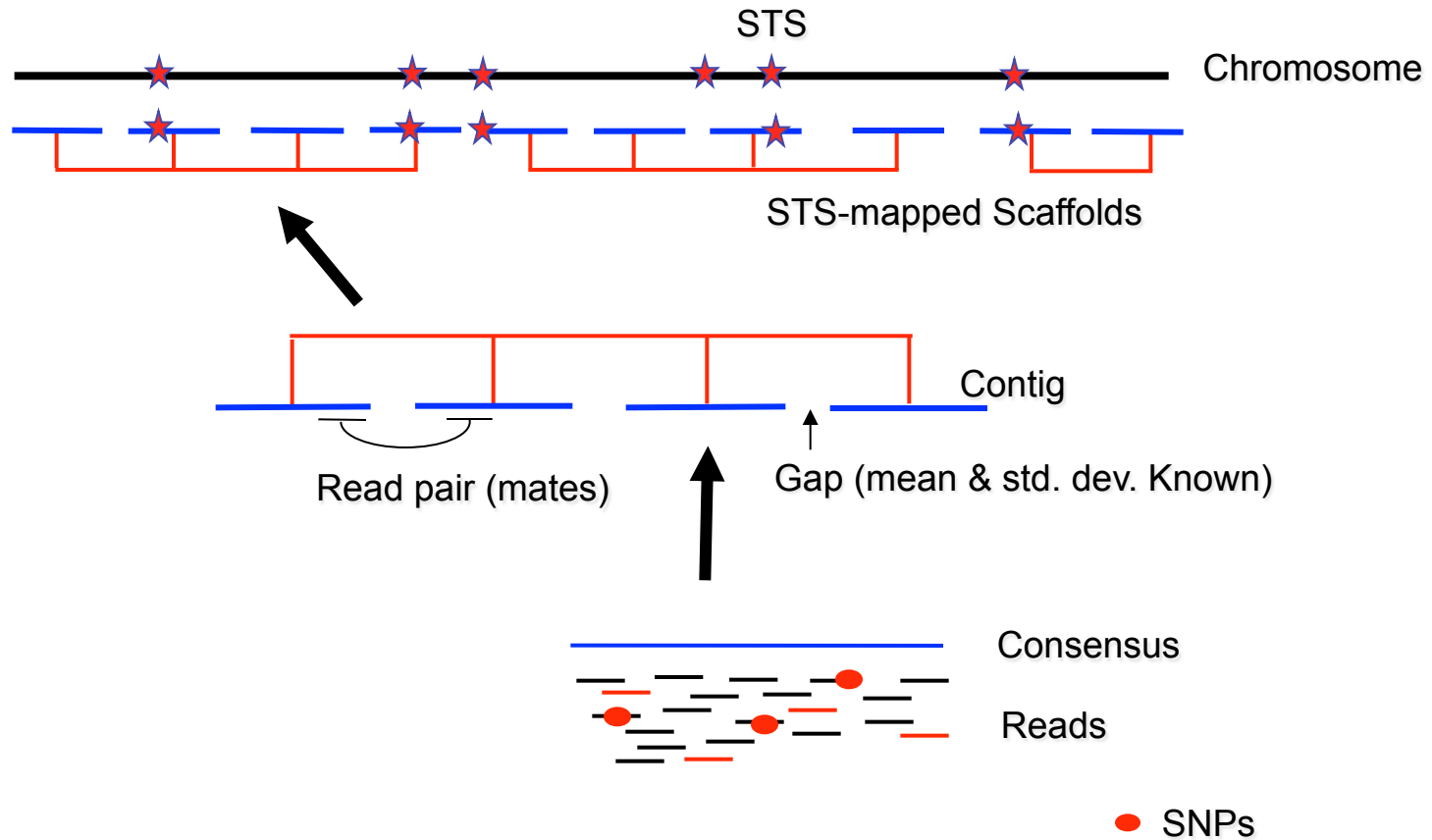


Scaffold

Gaps in scaffolds are traditionally indicated by 100 "N" s



# Anatomy of a WGS Assembly



AAGCTTCGCCAGGCTGTAATCCCGTGAGTCGTCTCAAAAATCATCAAGCAGGTGTCTCAGGGAGACTGCCTGACTGAGTTATGCTAATTCCTTTCTACTTTGGCGTGGTCACGTGTA  
ACCATATCCGAATCATTTCTCTAGCCCTACGAACAGGTAAGAGCGCTAGGGATGTCGGTGGAGTAGTGTGCTTACTCGATAATATTCAGTTGGGACTACCAGCGAGGCGCTCGCTTTGCT  
CACGCAATGCCTGAGACAGTTGCAGAATGAATGGTAACCGACAAACGCGTTTCATATGCGTTTTCAAACCTTAGTAGACGCGTACTGTCTGAAACTGGCGGTACAGGCACCAGATAACGCC  
CTTGGCATCGGCATGTCTCGTACAGAGGTCCGTATGTAGTGCCACGACTTCTAAATCCGGCGACAGGCTGGTCTTTTGTCTTACCACGTATTAGCCCGGTGCGGATTTCTCGGAGCGCAC  
CTGTTCAACACTAGAAAACGGAGTTTCTGATCGAGAAGCCACCACCTTTCCAGAAAGTTGAACGCTAGCATGTCAATTCGATTTTACCCCCCGGTAGTTCTGTGTGTCAATTCGTGTGTC  
GAGACAACTCTGTCCCGCCCCGGTGTGTCCATATGCGTGACTTTCCCGCAATTTTTTTCAGACTTTTCAGGAAAGACAGGCTCCGGAACGATCTCGTCCATGACTGGTAAATCCACGACA  
CCGCAATGGCCCCCAGCACCTCTATCTCTCGTGCCAGGGGACTAACGTTGTATGCGTCTGCGTCTTGTCTTTTTGCATTTCGCTTTCCAAAAAGAGAGCCATCCGTTCCCGCGCACATTC  
AACGCCGAGTCCGGTTTTTTGTCTTTTTGAGTGGTAGGACGCTTTTCATCGCGCAACTACGTGGACATTAAGTTCATTCTTTTTTCGACAGCACGAAACCTTGCAATCAAACCCGC  
CCGCGAAGATCCGATCTTGTCTGCTGTTCGAGTCCCAGTAGCGTCTGTGCGCGCGCTCTGTGTTGGTGGGCGAGCCGCTACACCTGTTATCTGACTGCCGTGCGCAAAATGACGC  
CATTTTTGGGAAAAATCGGGGAACCTCATTCTTTAAAAAGTATGCGGAGGTTTTCTTTTTCTTCTGTTTCGTTTTCTTTTTCTCGGGTTTGATAACCGTGTTCGATGTAAGCACTTTCCGTCTC  
TCCTCCGTGCTTTGTTTCGACATCGAGACCAGGTGTGAGATCCTTCGCTTGTGATCCGGAGACGCGTGTCTCGTAGAACCTTTTCATTTTACCACACGGCAGTGGGAGCACTGCTCTG  
AGTGCAGCAGGGACGGGTGAAGTTTCGCTTTAGTAGTGCCTTTCTGCTCTACGGGGCGTTGTGCTGTCTGGGAAGATGCAGAAACCGGTGTGTCTGGTGTGCGGATGACCCCCAAGAGG  
GGCATCGGCATCAACAACGGCCTCCCGTGGCCCCACTTGACCACAGATTTCAAACACTTTTTCTCGTGTGACAAAAACGACGCCGGAAGAAGCCAGTGCCTGAACGGGTGGCTTCCAGG  
AAATTTGCAAAGACGGGCGACTCTGGACTTCCCTCTCCATCAGTCGGCAAGAGATTC AACGCCGTTGTATGGGACGGAAAACCTGGGAAAGCATGCCTCGAAAGTTTAGACCCCTCGTG  
GACAGATTAACATCGTTCGTTTTCTTCCCTGTGAGCACACAGTAGTGCACACGCTGTTTGGAGCGTCAATCTCAAGAGTGTGGACGCTGTTCCACGTTCTCAAATGTTTTCC  
CAACATCCGTCGTCAGTAGACACACCAACAAAAAGCACACGGCGAATCTGCTCATCGGAGGGAGGAGCCGGGGGCGACACAACATCCTCAACTCTCGAACGAACATATCCGGGGCCGC  
GAAGACGTCAGTCTCTCAAATCCAACCCGGAACGCAAAACATTTCTGCATCAAGTCACGATTGCGCCGGTACCTCCATGTGTAAGCAGTTCCATGAAACCTCCGATATTACACACGACTG  
TGGATATGAATTATATGCAGATGCATATATACTGAGACGCCGATGCAACTATAGTTTCTGCGCCTCCATGGATATTTTCAGACCTTCTCTCACATTTGGTTTTGCCGTACACCTCCGT  
TAGCTTTTTTTTTCTGGCTTTCTTCTTCGTCCTGTATTATCAGCAAAGAAGAAGACATTGCGGGCGAGAAGCCTCAAGCTGAAGGCCAGCAGCGCTCCGAGTCTGTGCTTCACTCCAGC  
AGCTCTCAGCCTTCTGGAGGAAGAGTACAAGGATTTCTGTGACACAGATTTTTTGTGCTGGGTATGTTGTCTTAAACTCCTTGGAACTCCATTTCTGGTCCAGAAACGTAAGTAACTGTATA  
CATGTATATACAGATGTATGGATAATATCTAGAGAAGATACAGGGAAGACTGGCAAGGATGAAAGACATGCAGCTTTAACGAAGCAGAGGGCATTGGCGAGAGGGACCCCGTTATGCT  
GTGTAGTGGCTGTGAATTTTACTTCGCCGTTTACTTGTCTGAGCCTTGTCTTCCACTTGGTACTTGTGTTTCTACCTTCCCCAACGCCCTTCTATTTCCCTTCACTCGCAAGCG  
CGCTCAGTGGGCGCTCACCGAACCCCTTGGTTTTCTGTTTCAGCTGTTGTCTCTTTTTCTCGCTTGTGTTTCTGTTGGCGTCTGTTGCTCGGCTTCTCTCTTTTTCTGTTGGTGCCTCCAG  
ACTATGTGCGCTGTTTCCCCACCCTTCTCGGCTTGTGCTTTTCAGGAGGAGCGGGACTGTACGAGGCAGCGTGTCTCTGGGCGTTGCCTCTCACCTGTACATCACGCGTGTAGCCCGGA  
GTTTTCCGTGCGACGTTTTCTTCCCTGCGTTCCCCGGAGATGACATTTCTTCAAACAATCAACGTCTGCGCAGGCTGCAGCTCTGCGGAGTCTGTGTTGCTTCCCTTTTTGTCCGGAGCT  
CGGAAGAGAGAAGGACAATGAAGCGACGTATCGACCCATCTTCAATTTCCAAGACCTTCTCAGACAACGGGGTACCTACGACTTTGTGGTTCTCGAGAAGAGAAGGAAGACTGACGACGC  
AGCCACTGCGGAACCGGTAAGAGGCAACCGAAGCGCGTAGATAAGAAAAACAACAAAGAGAAGGTGAAACACGAAGAGAAGGGAAAAATGCGGAGAAACCGTGGATTTACAAAGATATCAA  
GAGCAATGCTTTGTGGAGATTTTTTTAATTCAGTAGAGACACCCGCGTGCAGGTTGTGTAGAAATAACTGCGACCCTGGAGACAGAGATCCGCGGATACACCCTTGTGCTTTTTCTC  
TCCTATGTTTATGACGGGTGCTGAACGCTATCGTACTTAATTTGGAGGAGTCTCTCCGAAGCAGCTTTGGCTGGCCATCCGTGTGTTTGCCTTTGTTCCCTGAAAGCCAGAAGGCGCTCC  
ACAGTGGAGCGATATACAGGGACGCTACCGGAGCCCCGTTTTCTGCTTTTGTGACTCTTGCAGAGCAACGCAATGAGCTCCTTGCAGTCCACGAGGGAGACAACCTCCCGTGCACGGGT  
TGCAGGCTCCTTCTCGGCCGACGCAATTGCCCGGTGTTGGCGTGGATGGACGAAGAAGACCGGAAAAACCGGAGCAAAAGGAACGATTCCGGCCGTTCCGATGTTCACTTTAGAG  
GCCATGAAGAATCCAGTACCTTGATCTCATTGCCGACATTATTAACAATGGAAGGACAATGGATGACCGAACGGGTAACGGCGACTGCGAGAAAAGCCACACCGTTTTCTCTGTGAT  
TCTGTCCGCAAGCCCTCTTTTTGCTTCATCCACCTTTGCTATTCTCCGCGCCTTCTTTTTCTGCTCCATGTTCAATTCGTTTCGCTTCTTTCAGTCTTTCCATCTTCCCTGTTACCTCTG  
TCATTCGTTTTCTTGCCTCTATTTAACTGTGTTCTACTCACAGTCTGCATTCGCGGATAGACGAGCTTCCACGCTTTCGCTCTCGACAAGCAACTGTCATTTGTACGCGCTCCCTCCAC  
CGTGAATCGGATTTGCGTTCCGGGTTCTGGGTGAGAAAAGGCTTCCGCGGATTTCTGAATAATACCTTCCGCAATTGTAAGAGGGCAAGGAACAAGAGATATTTCCGGCGCATCT  
TTTTGTGCGGCGCTTTCTCGTGTTCACACCGATGCCCTTCTGTGCATGTTCTGCTTCTCTCTTTTTCCCTGTTTAGGCGTTGGTGTATCTCCAAATTCGGCTGCAC  
TATGCGCTACTCGTGGATCAGGCTTTCCACTTCTCACCAAAAGCGTGTGTTCTGGAAAGGGTAAGGGCGCTTTCAGTGAATGCATATATTTGACTTACAGACTTCTTAACTGTTTGA  
CAACCAACGTACAAATTTGTTTGTCCGTGTGCGTGTTCGACATGTCAAGTATGTGAAGAGTGCCTACTGTAGACTAACGCACGAACAGATTTGTTTATCTGCATGCGCTGTGCACCCGT  
TTCGAGTGTCTGGAGTTTCCGCAACCTTCTTTGAATTTCTGGTTCGTTTTTTATGCGCGCACTGGTTTTGCATGTGGCTGAGAGAGCACAGATCGAAGGTGGGGTGTATGGCGTC  
GCTGCAGAGAACTCCGGCGAAGGCAGAGATAAAGGAGAGTGGAAATCATTGAACAGTGTGGTTCGTTGTTTTCGACAGGTTCCCGAAGAGTTGCTGTGTTTTCATTCCGGCGGACA  
CGAACGCAAACCATCTTCTGAGAAGGGCGTGAAGGCAAGTCTACGTTGTACCTCTTGTCTCTGCCGAAGCTCAGATGTCTCCACGGCGTTGGTTCTTTTTGCTTTTTCGTTGGCA  
TTACCATCGAGTACCACCTCATAGTTGCGTGTGTCTACATGTTTTCTAGAACGTCGGTGTGTTGCTCGTGGCGACCGGCGGAGTGTATGTACCCTGCGCTGTGAGAAGTTGATCCTT



# Six Frame Translation ORF-finding

```

1/1                               31/11                               61/21
M Y A L L I L Y Y I I I R H * S H H A C R G V Y Y I Y
H V R F T D S I L Y Y Y * T L V T S C M * G G L L Y L
A C T L Y * F Y I I L L L D T S H I M H V G G S T I S
GCA TGT ACG CTT TAC TGA TTC TAT ATT ATA TTA TTA TTA GAC ACT AGT CAC ATC ATG CAT GTA GGG GGG TCT ACT ATA TCT
CGT ACA TGC GAA ATG ACT AAG ATA TAA TAT AAT AAT AAT CTG TGA TCA GTG TAG TAC GTA CAT CCC CCC AGA TGA TAT AGA
C T R K V S E I N Y * * * V S T V D H M Y P P R S Y R
M Y A K S I R Y * I I I L C * D C * A H L P T * * I *
H V S * Q N * I I N N N S V L * M M C T P P D V I D I
121/41                             151/51                             181/61
* L E L E R I D L A * L Y N F S D I Y I P A S R G K W
L A R A R T H R L S M T I * F Q R H I Y S R L A G K M
A S S S * N A S T * H D Y I I S A T Y I F P P R G E N
GCT AGC TCG AGC TAG AAC GCA TCG ACT TAG CAT GAC TAT ATA ATT TCA GCG ACA TAT ATA TTC CCG CCT CGC GGG GAA AAT
CGA TCG AGC TCG ATC TTG CGT AGC TGA ATC GTA CTG ATA TAT TAA AGT CGC TGT ATA TAT AAG GGC GGA GCG CCC CTT TTA
S A R A L V C R S L M V I Y N * R C I Y E R R A P F I
* S S S S R M S K A H S Y L K L S M Y I G A E R P F H
L E L * F A D V * C S * I I E A V Y I N G G R P S F P

```

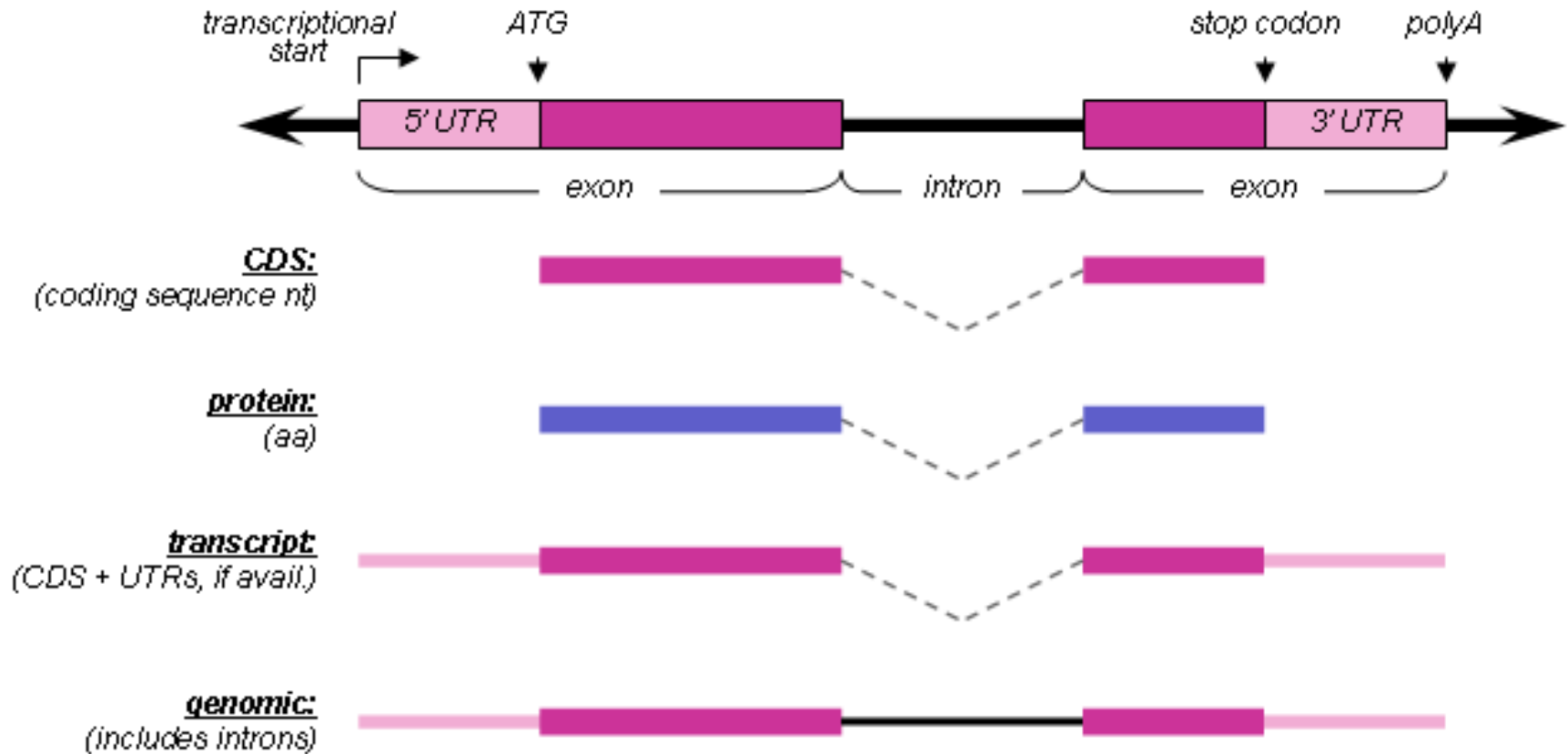
**ORFs ≠ Genes**

ATGCAGAAACCGGTGTGTCTGGTCGTCGCGATGACCCCCAAGAGGGGCATCGGCATCAACAACGGCCTCCCGTGGCCCC  
ACTTGACCACAGATTTCAAACACTTTTCTCGTGTGACAAAACGACGCCCGAAGAAGCCAGTCGCCCTGAACGGGTGGCT  
TCCCAGGAAATTTGCAAAGACGGGCGACTCTGGACTTCCCTCTCCATCAGTCGGCAAGAGATTCAACGCCGTTGTATG  
GGACGGAAAACCTGGGAAAGCATGCCTCGAAAGTTTAGACCCCTCGTGGACAGATTGAACATCGTCGTTTCTCTTCCC  
TCAAAGAAGAAGACATTGCGGCGGAGAAGCCTCAAGCTGAAGGCCAGCAGCGCGTCCGAGTCTGTGCTTCACTCCCAGC  
AGCTCTCAGCCTTCTGGAGGAAGAGTACAAGGATTCTGTGACAGATTTTGTGCGTGGGAGGAGCGGGACTGTACGAG  
GCAGCGCTGTCTCTGGGCGTTGCCTCTCACCTGTACATCACGCGTGTAGCCCGCAGTTTCCGTGCGACGTTTTCTTCC  
CTGCGTTCCCCGAGATGACATCTTTCAAACAATCAACTGCTGCGCAGGCTGCAGCTCCTGCCGAGTCTGTGTTTCGT  
TCCCTTTTGTCCGAGCTCGGAAGAGAGAAGGACAATGAAGCGACGTATCGACCCATCTTCATTTCCAAGACCTTCTCA  
GACAACGGGGTACCCTACGACTTTGTGGTTCTCGAGAAGAGAAGGAAGACTGACGACGACGCCACTGCGGAACCGAGCA  
ACGCAATGAGCTCCTTGACGTCCACGAGGGAGACAACCTCCCGTGACGGGTTGCAGGCTCCTTCTTCGGCCGACGCCAT  
TGCCCCGGTGTGTGCGTGGATGGACGAAGAAGACCGGAAAAACGCGAGCAAAGGAACTGATTCCGGCCGTTCCGCAT  
GTTCACTTTAGAGGCCATGAAGAATCCAGTACCTTGATCTCATTGCCGACATTATTAACAATGGAAGGACAATGGATG  
ACCGAACGG

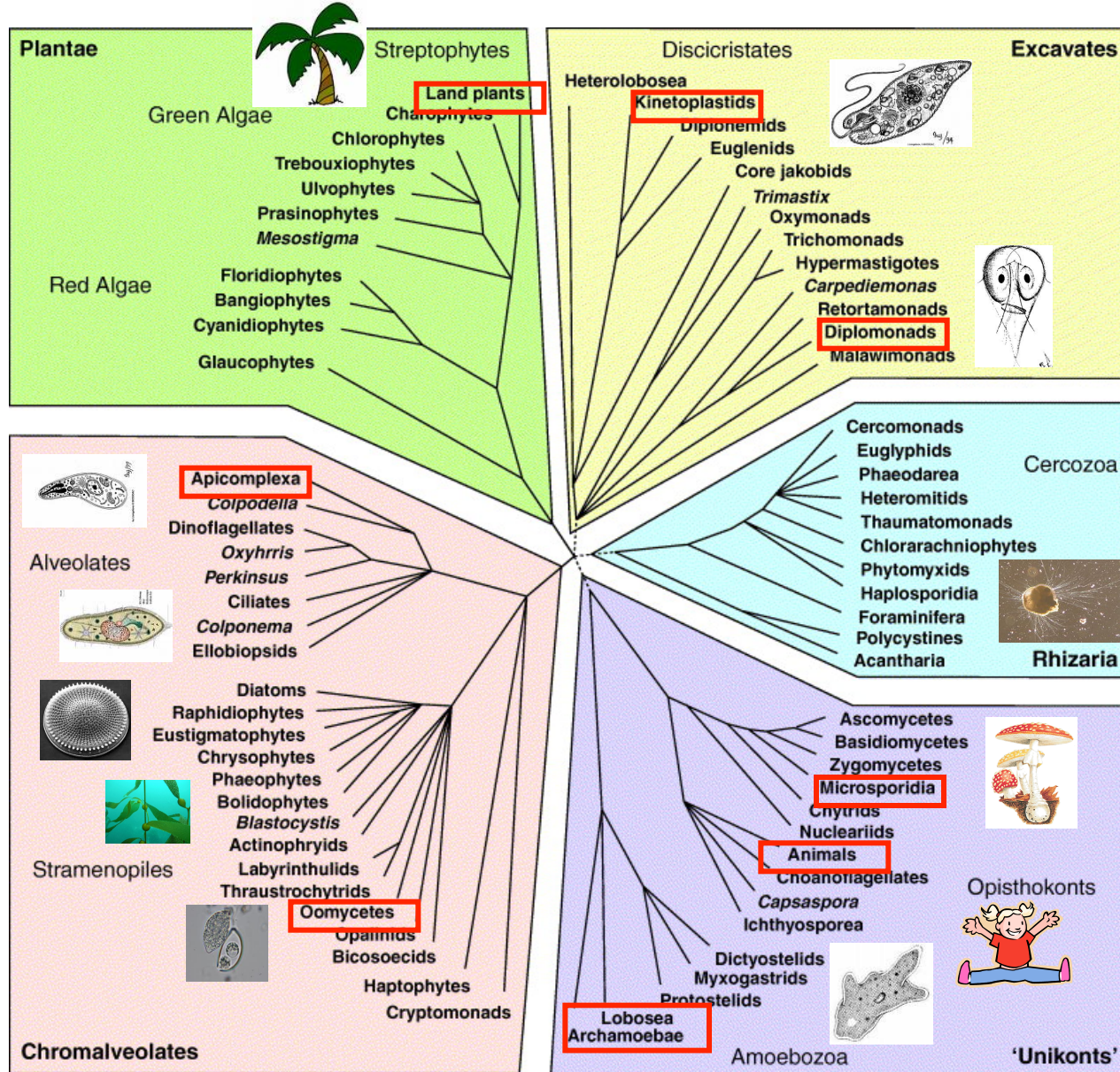
>Translation Frame 1

MQKPVCLVVAMTPKRGIGINNGLPWPHLTTDFKHFSRVTKTTPPEASRLN  
GWLPRKFAKTGDSGLSPSPVGRFNAVVMGRKTWESMPRKFRPLVDRLNI  
VVSSSLKEEDIAAEKPQAEQQQRVRVCASLPAALSLLLEEEYKDSVDQIFV  
VGGAGLYEAALS LGVASHLYITRVAREFFPCDVFFPAFFGDDILSNKSTAA  
QAAAPAESVFPFCPELGREKDNEATYRPIFISKTFSDNGVPYDFVPLEK  
RRKTDDAATAEPSNAMSSLTSTRETTPVHGLQAPSSAAAIAPVLAWMDEE  
DRKKREQKELIRAVPHVHFRGHEEFQYLDLIADIINNGRTMDDRT

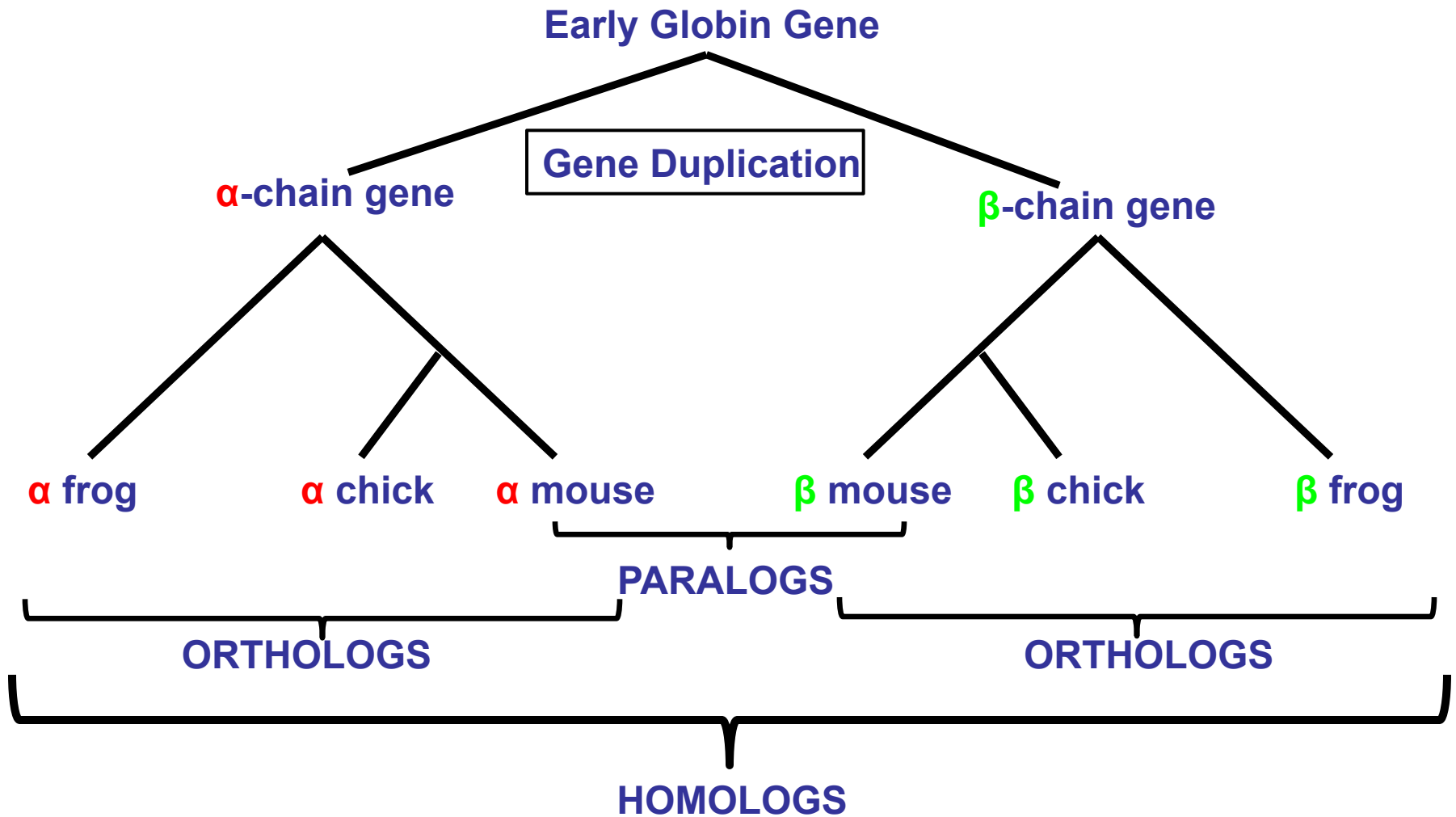
# Terminology



# A Tree of Eukaryotes (Keeling, 2005)



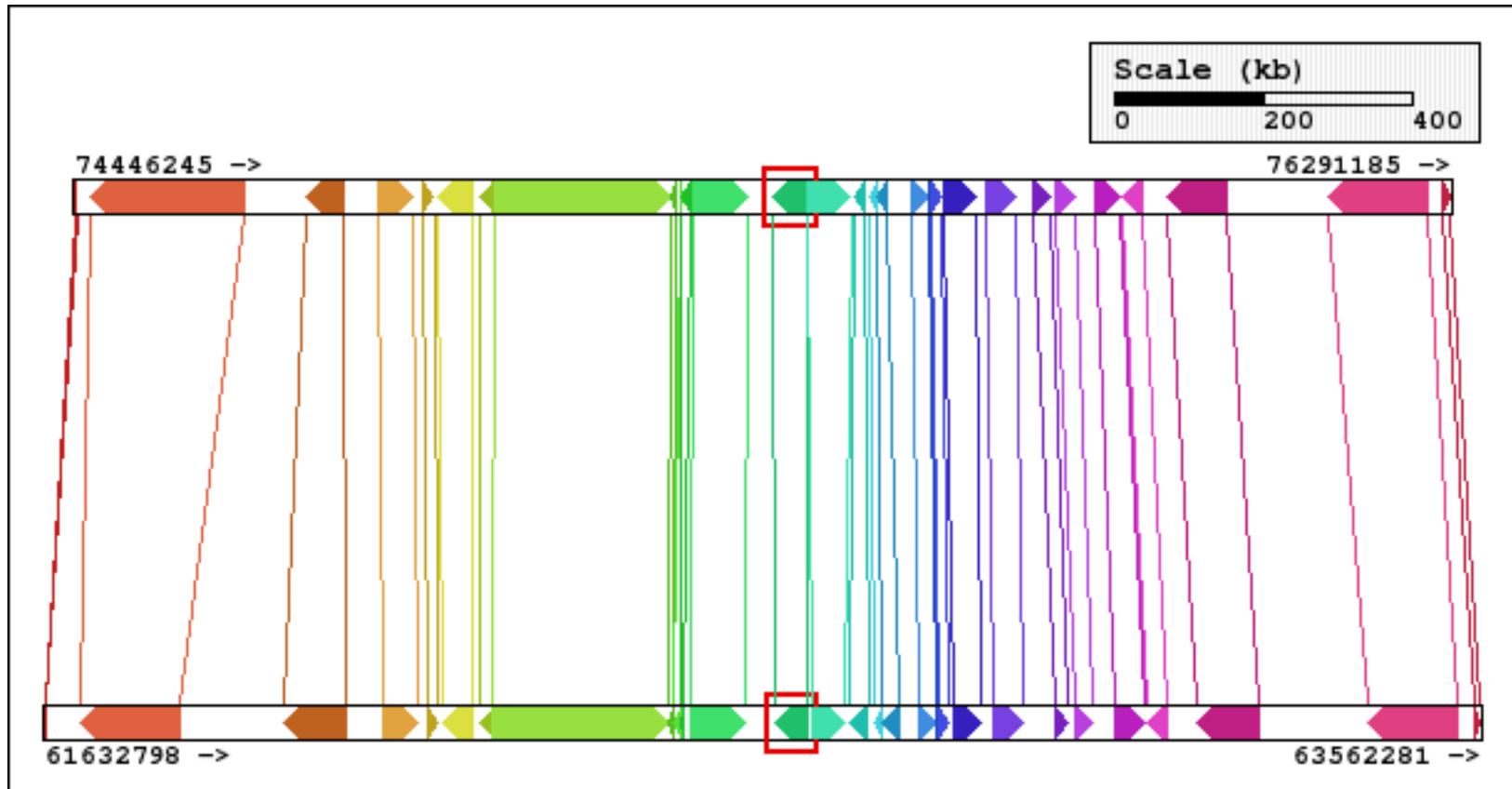
# Homology



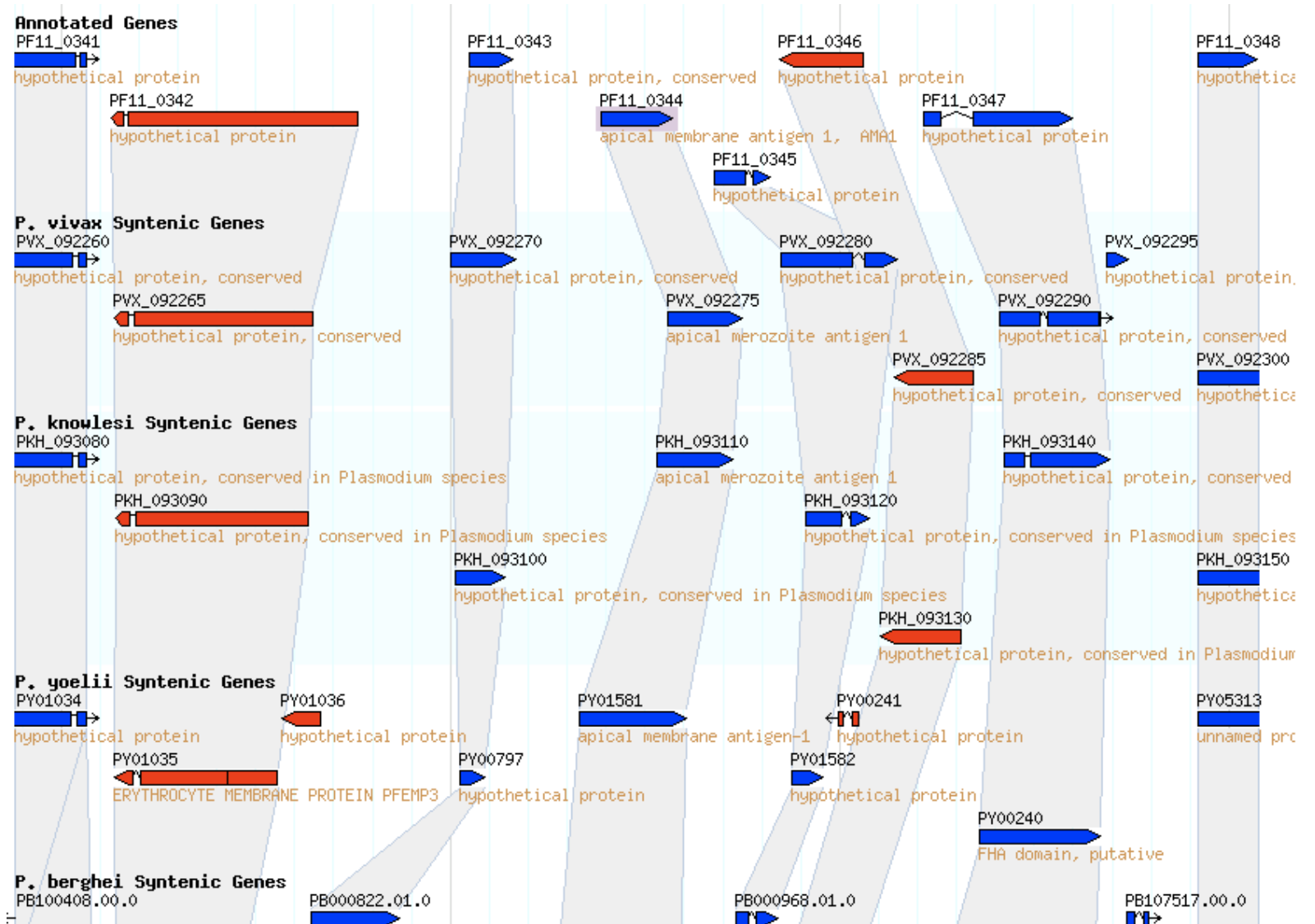
# Evolutionary relationships

- Homology - related by evolutionary descent not equivalent to similarity
- Orthology - same gene in different organisms, e.g. alpha hemoglobin in humans and chimps
- Paralogy - genes within an organism related by gene duplication, e.g. alpha and beta hemoglobin in humans
- Xenology - genes related by gene transfer

Synteny = large regions of chromosomes containing the same genes



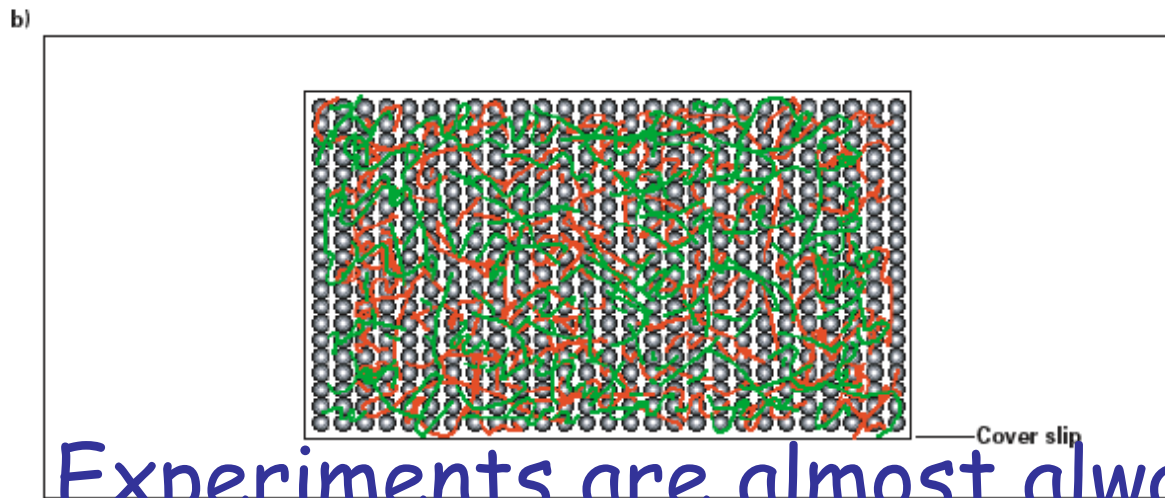
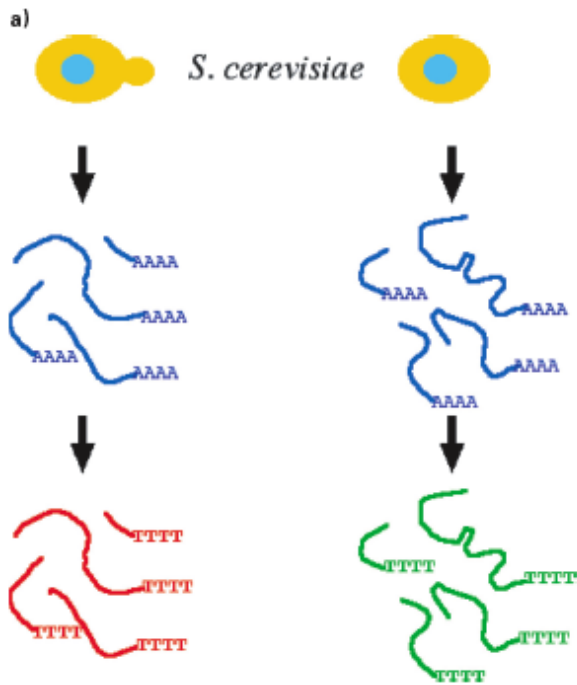
# Synteny among Plasmodia



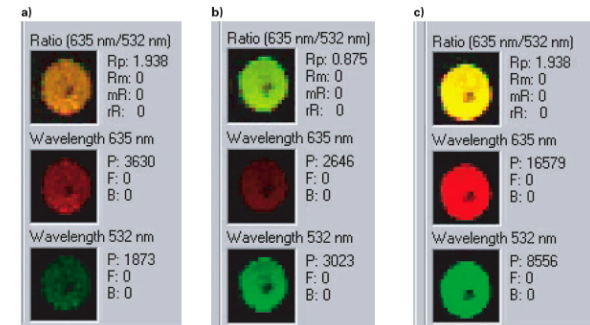


# Expression Profiles

- The pattern of expression of one or more genes over time or a set of experimental conditions, e.g. during development or a drug treatment or in a genetic mutant such as a gene knock-out.
- Always... has a time and location component



The RNA samples from the test and the control are labeled with different colors in a reverse-transcription reaction and then hybridized, together, competitively to a slide or chip containing gene sequences in multiple copies.



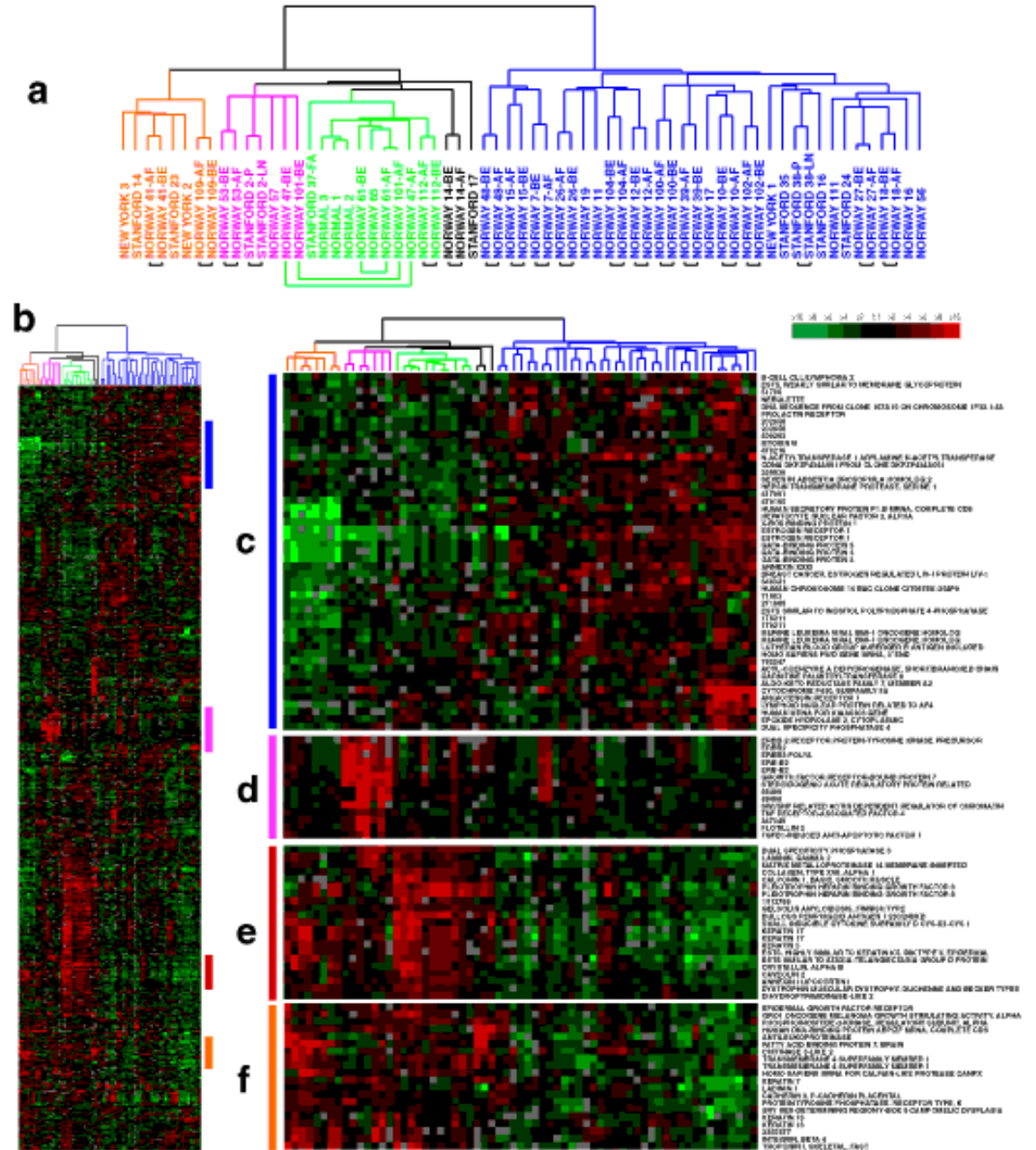
Experiments are almost always  
Competitions between conditions or stages

Ratios of experimental to control expression are often expressed as colors rather than numbers



Figure 2

Clustered  
Microarray  
Data  
Genes with  
Similar  
Expression  
Profiles are  
Grouped  
together



# Other RNA expression

- Expressed Sequence Tags, ESTs
  - Usually represent partial cDNA
  - Often clustered
  - Come from libraries that may, or may not be normalized
  - Often used to identify genes in genomes and locations of introns
- SAGE-tags (Serial Analysis of Gene Expression)
  - Primary purpose is relative levels of gene expression
- RNA-Seq (NGS)
  - Little sequence bias
  - Quantitative
  - Can be strand-specific

# Genes can be located on either DNA strand

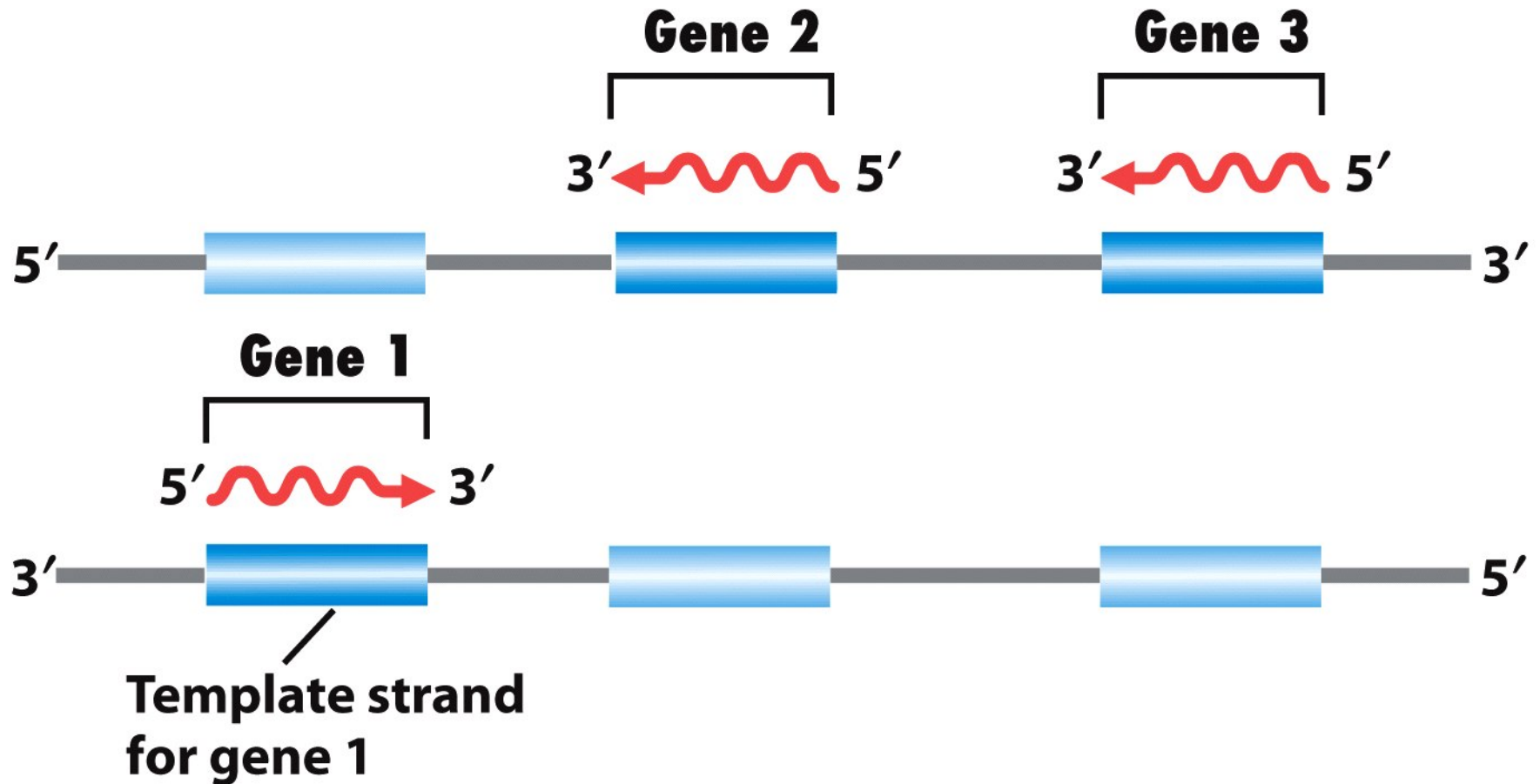


Figure 8-3  
*Introduction to Genetic Analysis, Ninth Edition*  
© 2008 W. H. Freeman and Company

Overview of transcription: Either strand can serve as a template for a gene

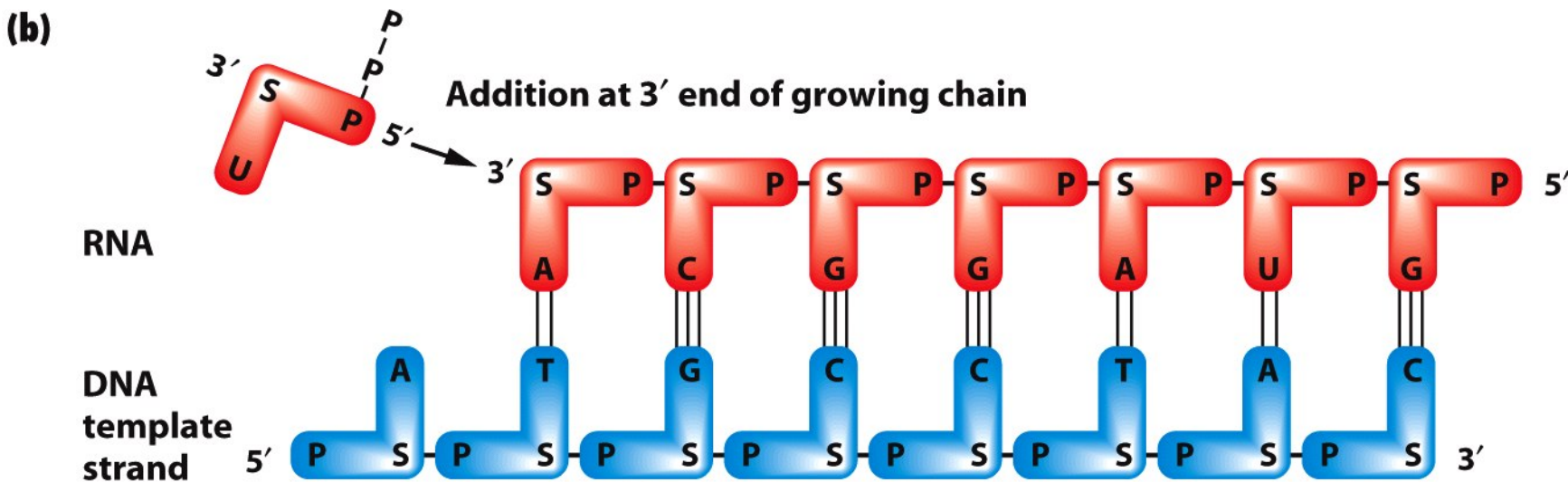
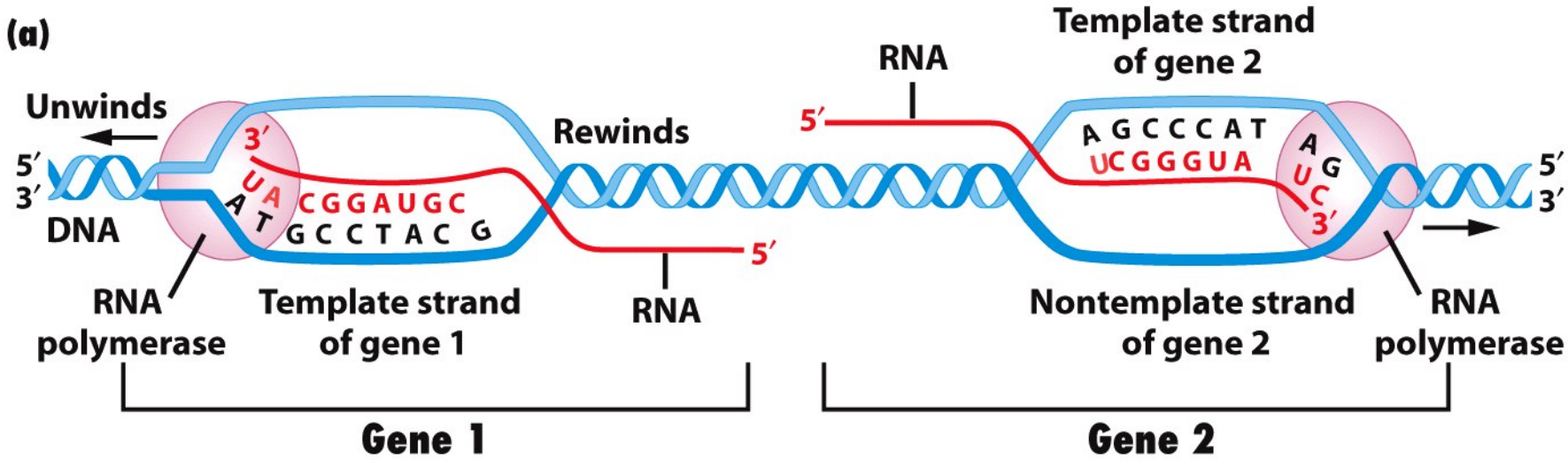


Figure 8-4  
*Introduction to Genetic Analysis, Ninth Edition*  
 © 2008 W. H. Freeman and Company

# Convention

Gene location = non-template strand,  
i.e. same as the mRNA



Figure 8-6  
*Introduction to Genetic Analysis, Ninth Edition*  
© 2008 W. H. Freeman and Company



# Complex patterns of eukaryotic mRNA splicing: What is a Gene?

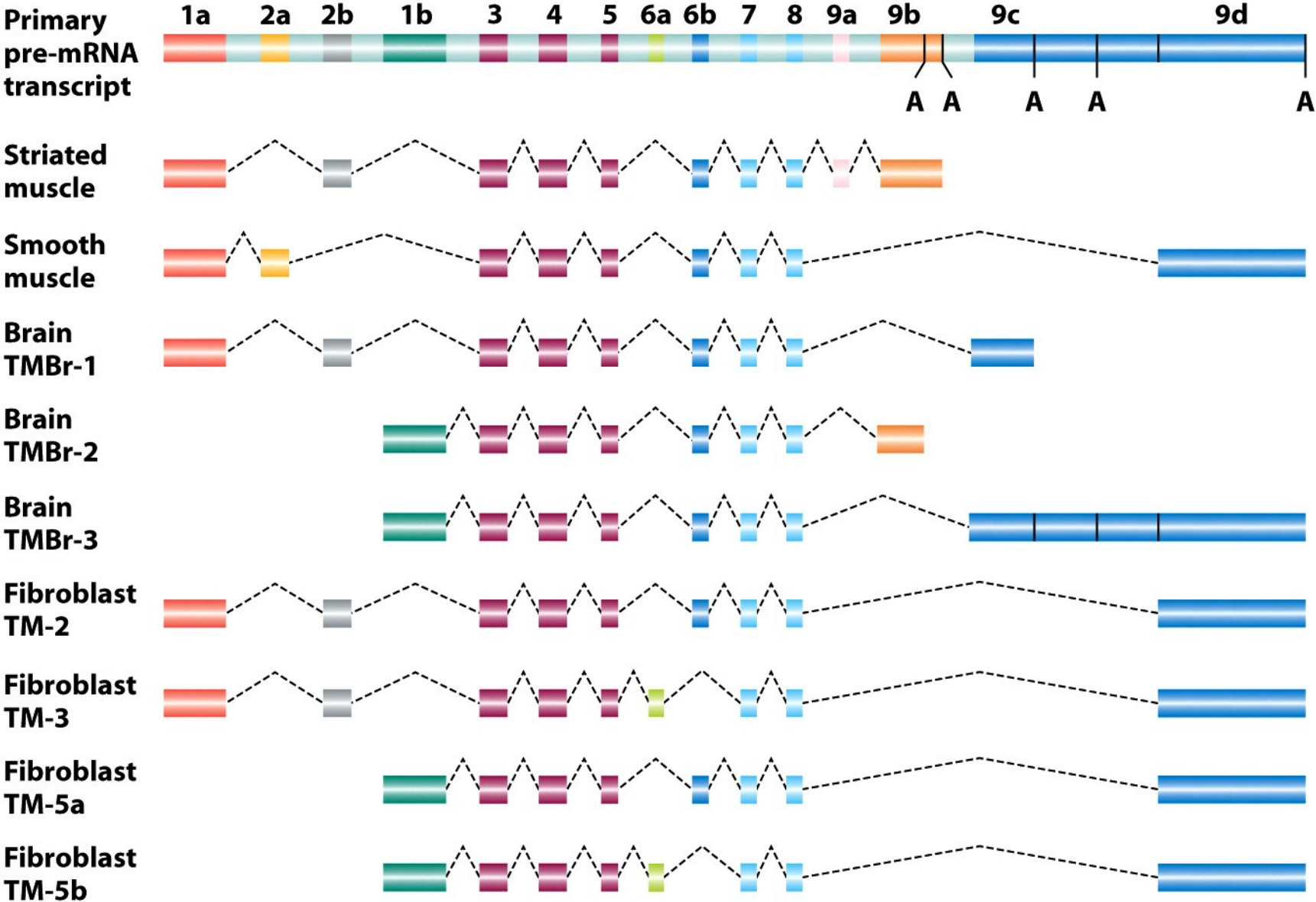


Figure 8-14  
*Introduction to Genetic Analysis, Ninth Edition*  
 © 2008 W. H. Freeman and Company

# Bioinformatics uses algorithms

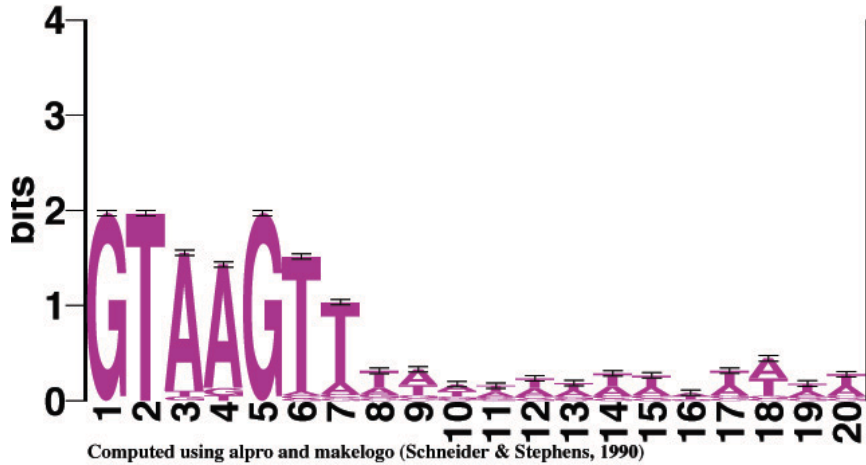
- Algorithms are sets of rules for solving problems or identifying patterns
- Algorithms can be general or case specific and often need to be trained
- Computational analysis, like wet-bench analyses are only as good as the tools, techniques and material allow, and all interpretations come with caveats (like the experimental conditions, often call parameters in bioinformatics).

# How to find an intron

- Usually begins with GT and end with AG
- Must be longer than 19 nucleotides
- Must contain a branchpoint “A”
- Donor GT often followed by a sequence pattern. This pattern is species-specific
- Acceptor AG often preceded by pyrimidine stretch
- Has a mean length of “X” as is observed in this species

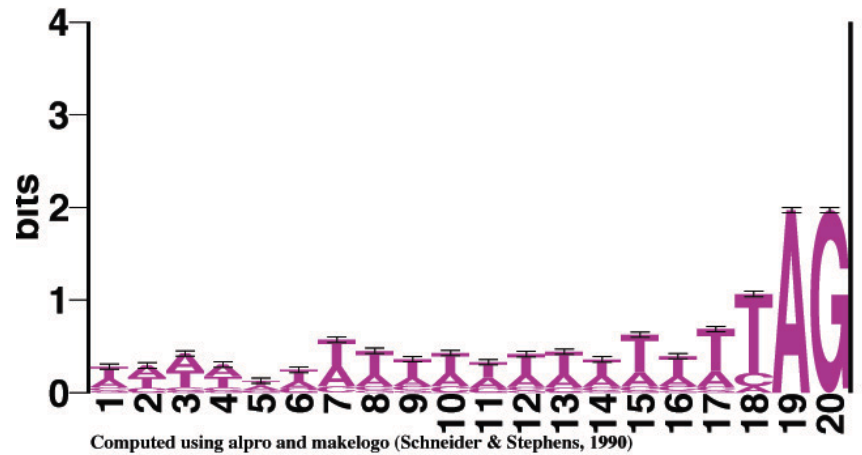
# Donor Site

Generated by <http://www.bio.cam.ac.uk/seqlogo/logo>

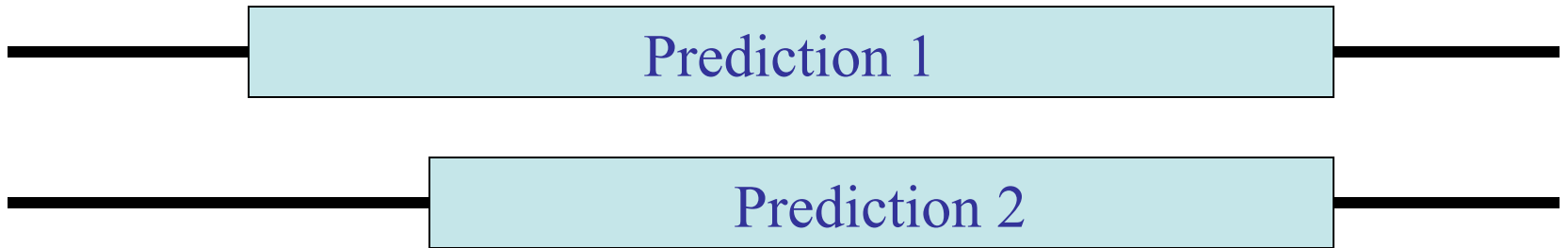


# Acceptor site

Generated by <http://www.bio.cam.ac.uk/seqlogo/logo>



Different prediction methods  
often generate different  
results



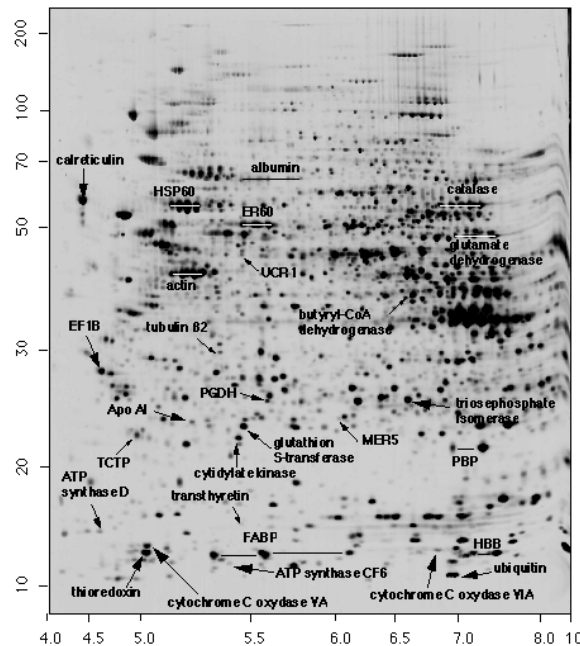
# Protein Expression/Sequence

## Data

- MW-Isoelectric point
- MW
- Sequence/spans

## Technology

- 2D gel electrophoresis
- Mass spectrometry
- Tandem MS (MS-MS, LC MS-MS etc)

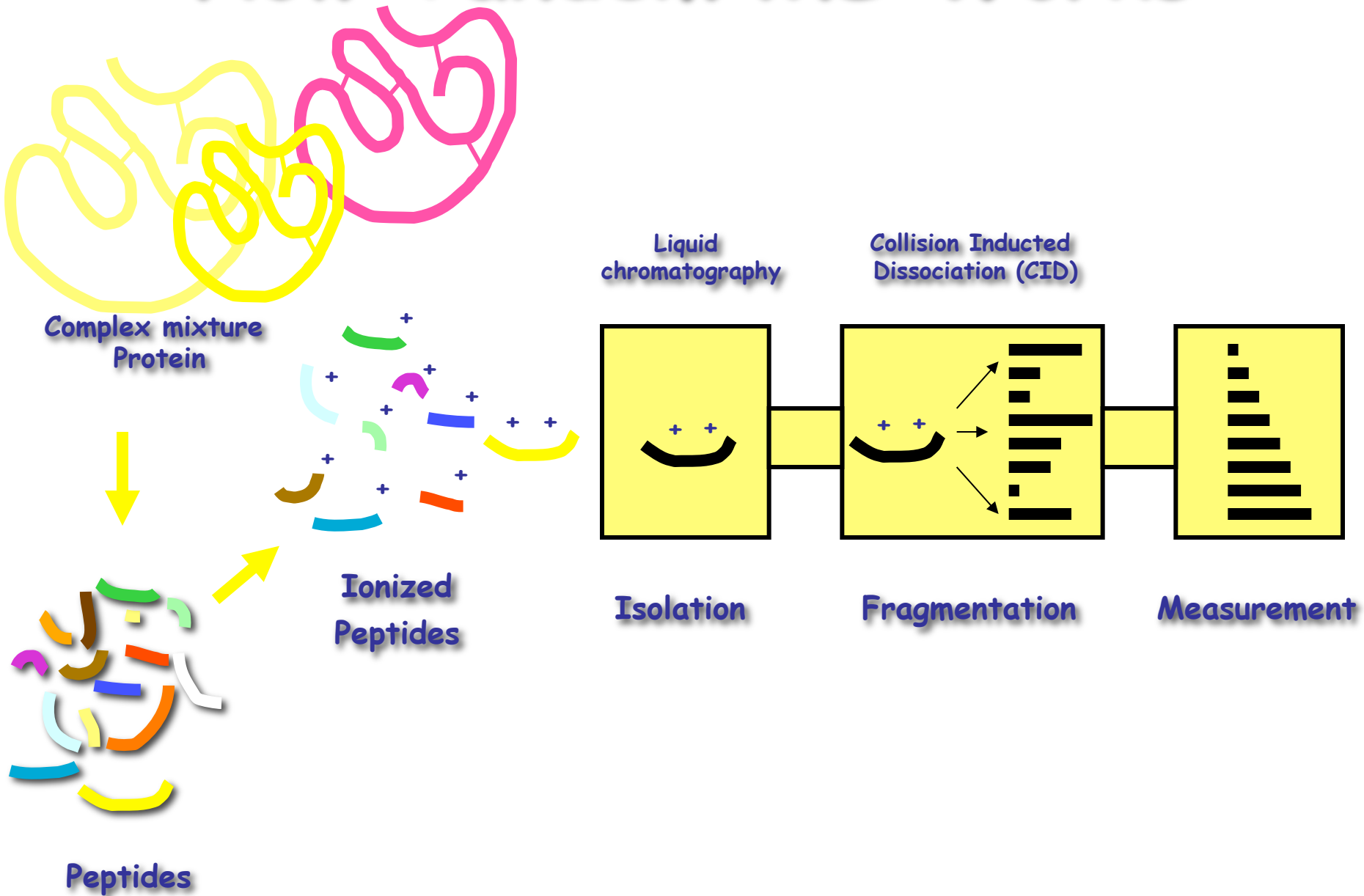


Typical 2 D gel

# High throughput mass spectrometry

- Direct identification of proteins from biological sample.
- Capillary liquid chromatography apparatus (LC) coupled with...
- Electrospray tandem mass spectroscopy (MS/MS)
- “Sequest”, Mascot, or other software links mass spectra with genomic sequence database.

# How Tandem MS Works





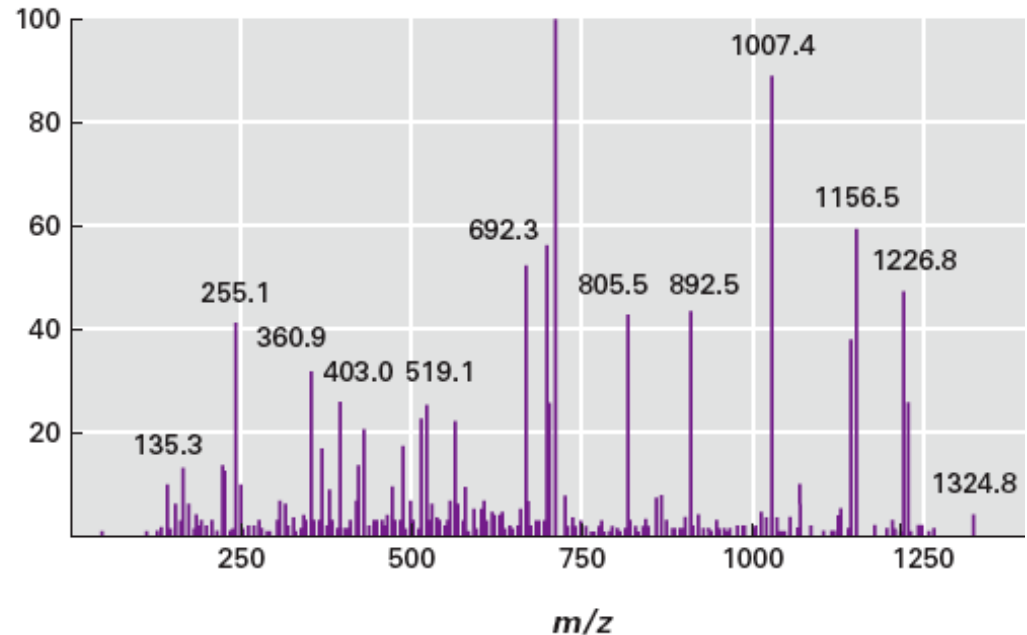
# Tandem MS protein data

a)

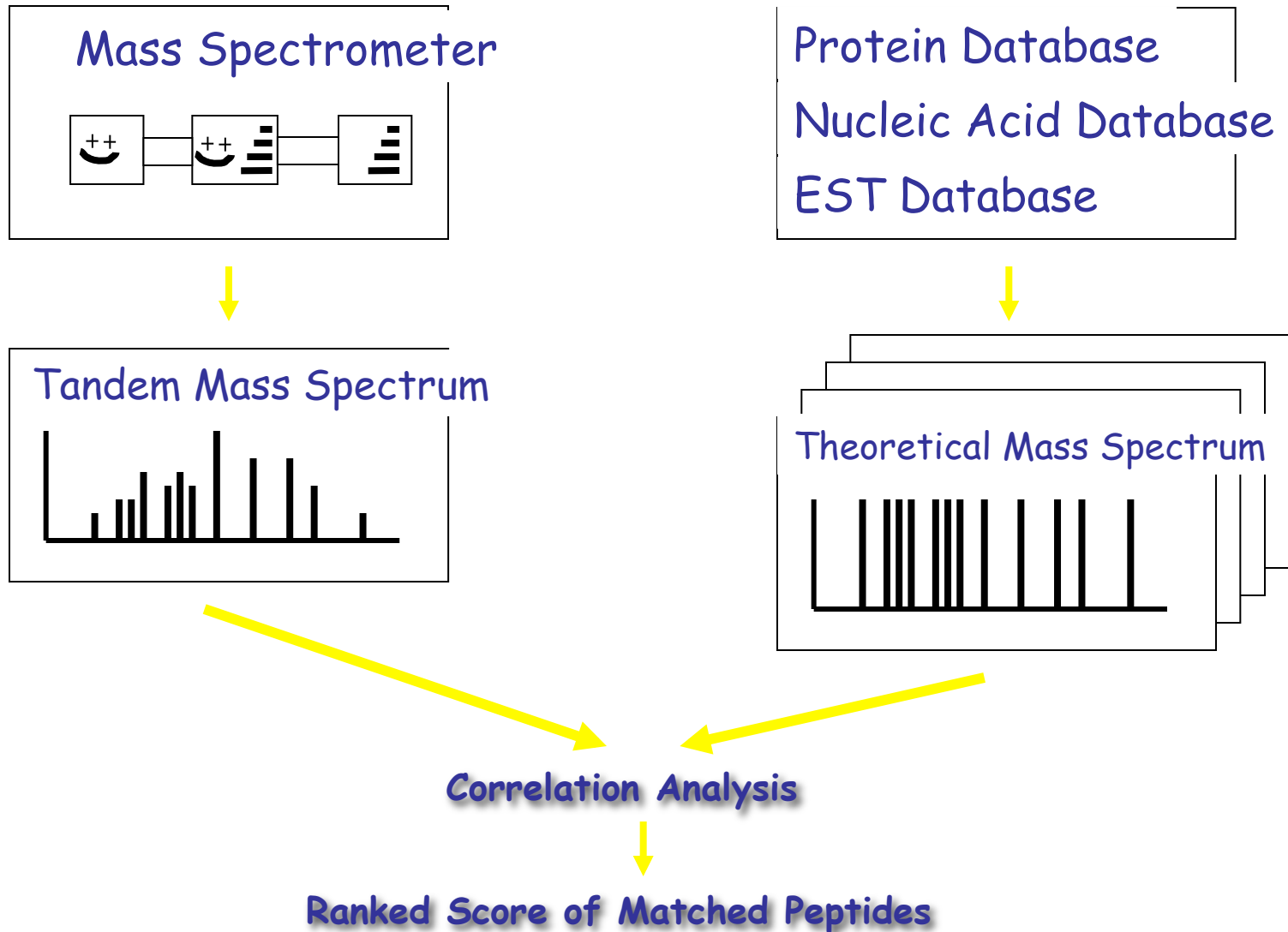
**S-P-A-F-D-S-I-M-A-E-T-L-K**  
(protonated mass 1410.6)

Mass <sup>+</sup>	b-ions	y-ions	Mass <sup>+</sup>
81.1	S	PAFDSIMAETLK	1323.6
185.2	SP	AFDSIMAETLK	1226.4
256.3	SPA	FDSIMAETLK	1155.4
403.5	SPAF	DSIMAETLK	1008.2
518.5	SPAFD	SIMAETLK	893.1
605.6	SPAFDS	IMAETLK	806.0
718.8	SPAFDSI	MAETLK	692.3
850.0	SPAFDSIM	AETLK	561.7
921.1	SPAFDSIMA	ETLK	490.6
1050.2	SPAFDSIMAE	TLK	361.5
1151.3	SPAFDSIMAET	LK	260.4
1264.4	SPAFDSIMAETL	K	147.2

b)



# Sequest Database Search



# Peptide database



ENNPCKLQYDYNNTNVTHGFGQEYPCETDIVERFSDTEGAQCDDKKIKDNSEGACAPYRRL  
HVCVRNLENINDYSKINNKHNLLEVCCLAKEYEGESITGRYPQHQETNPDTKSQCLTVLA  
RSFADIGDIIRGKDLYRGGNTKEKKRKKLEENLKTIFGHIYDELKNGKTNGEEELQKRY  
RGDKDNDFYQLREDWWDANRETVWKAITCNAQSYQYSQPCTGRGEIPYVTLKQCQIAGE  
VPPTYFDYVPOYLRFEEWAEDFCRKKKKKIPNVKTNCRQVQORGKEYCVRDGYNCVGTIR  
KQYIYRLD'TDC'TKCSLACKTFAEWIDNQEQFDKQKQYQNEISGGGRRQKRSTHSTKE  
YEGYEKHFNEELRNEGKDVRSLQLLSKEKICKERIQVGEEETANYGNFENESNTFSHTEY  
CDRCPLCGVDCSSDNCRKKPKKSCDEQITDKEYPPENTTKIPKLTAEKRKTGILKKYEFK  
CKNSDGNNGGQIKKWECHYEKNDKDDGNDINNCIQGDWKTSKNVYYPISYYSFFYGSII  
DMLNESIEWRERLKSCINDAKLGKCRKGCNPKCECYKRWVEKKKDEWDKIKEFFRKQKDL  
LKDIAGMDAGELLEFYLENI'FLEDMKNANGDPKVIKFKELGKENEVQDPLKTKKTID  
DFLEKELNEAKNCVEKNPDNECPKQKAPGDGAAPSDPPREDITHHDGEHSSDEDEEEEE  
EEQOPPAEGTEQGEEKSESKEVVEQOETPQKDTEKTVPTTTP'VTDVCDTVK'TALADTGS  
NAACSLKYVTCKNYQWRCIAPSGTITSGKDGATCVPPRTOELCLYYLKELESDTTQKQLLEA  
FIKTAQOETYLLEKMKEDIQNETASTLLEIIPCTDLIGEEIPEDFKRDIFYTFQDTRD  
LFLGRYIGNDELDKVNNTTAVFONCEHLPNGQKTRQRQEFWGTYYGKDIWKGMLCALQEA  
GGKKTLETETYNYSNWFNGHLTGTKLNEFASRPSFLRWMTEWGDQFCRERITQLOILKER  
CMVYQYNGDKGKDDKKEKCTEACTYYKEWLTNWQDNYKKQNRQRYTEVKGTS'PYKEDSDVK  
ESKYAHGYLRKILKNIIC'TSGTDIAYCNCMEGTSTTDSSNNDNIPESLKYPPIEIEEGCT  
CKDPSPGEV'PEKKVPEPKVLPKPPKLPKRQPKERDFP'TPALKNAMLSSTIMWSIGIGFA  
TF'TYFYLKKKTKSTIDLLRVINIPKSDYDIP'TKLSPNRYIPYTS'GKYRGKRYIYLEGDSG  
TDSGYTDHYS'DITSSSESEYEELDINDIYAPRAPKYKTLIEVVLEPSGNNTTASGNNTPS  
DTQNDIQNDGIPSSKITDNEWNTLKDEFISQYLQSEQPN'DVPNDYSSGDIPLNTQPN'TLY  
FDNPDEKPFITSIHDRDLYSGEEYSYVNMVNTNNDIPISGKNGTYSGIDLINDSLNSNN

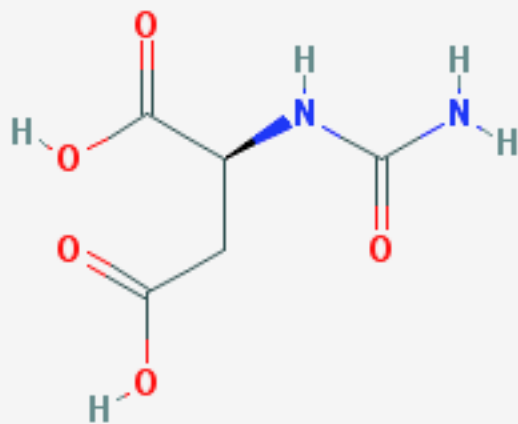


Note: ORFs in addition to predicted Genes must be searched

## Overview

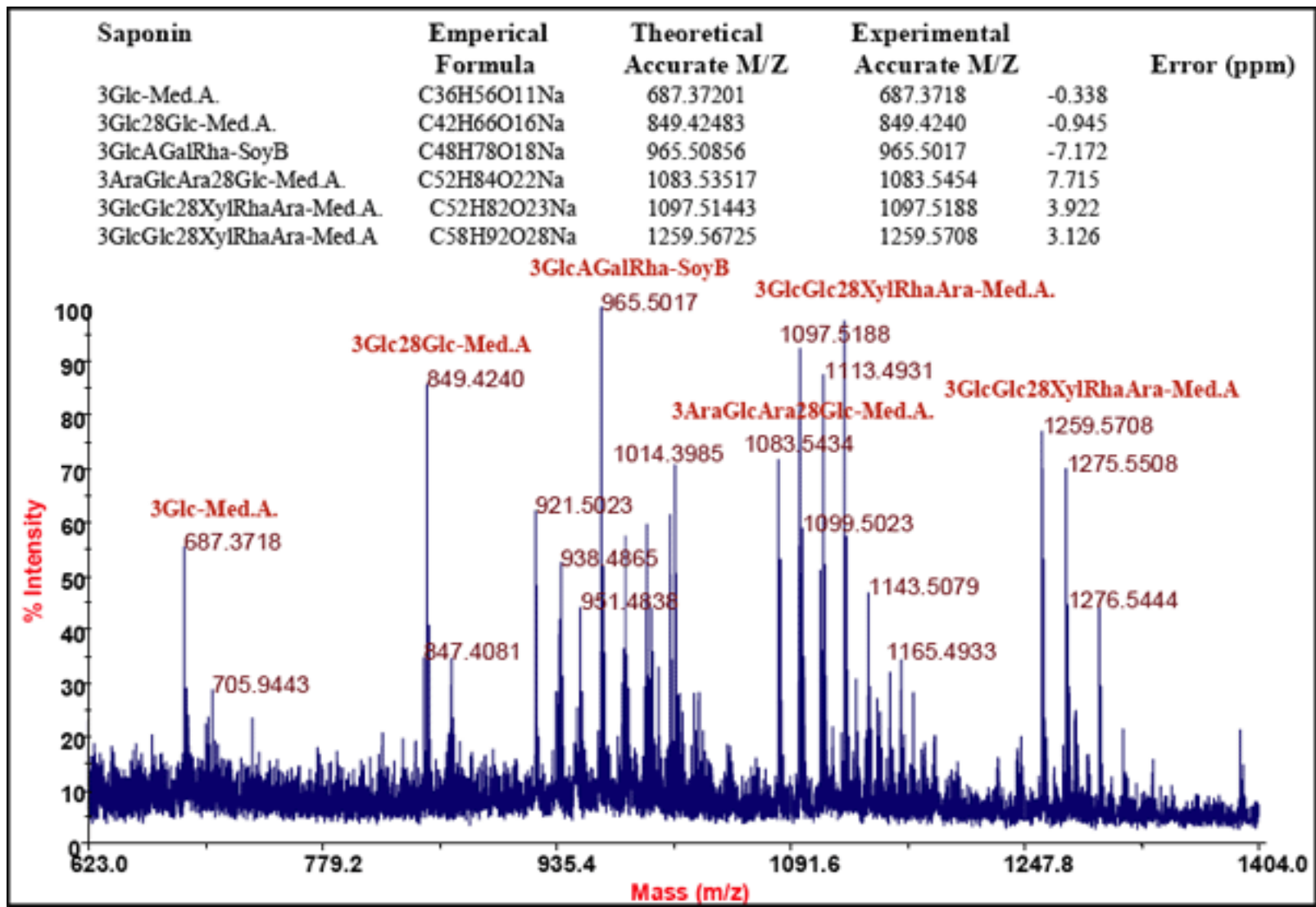
<b>PubChem Compound ID:</b>	<b>CID:93072</b>
<b>PubChem Substance ID(s):</b>	<b>3727</b>
<b>Synonyms:</b>	<b>N-Carbamoyl-L-aspartate</b>
<b>Molecular Weight:</b>	<b>176.12742</b>
<b>Molecular Formula:</b>	<b>C<sub>5</sub>H<sub>8</sub>N<sub>2</sub>O<sub>5</sub></b>

## 2D Structure



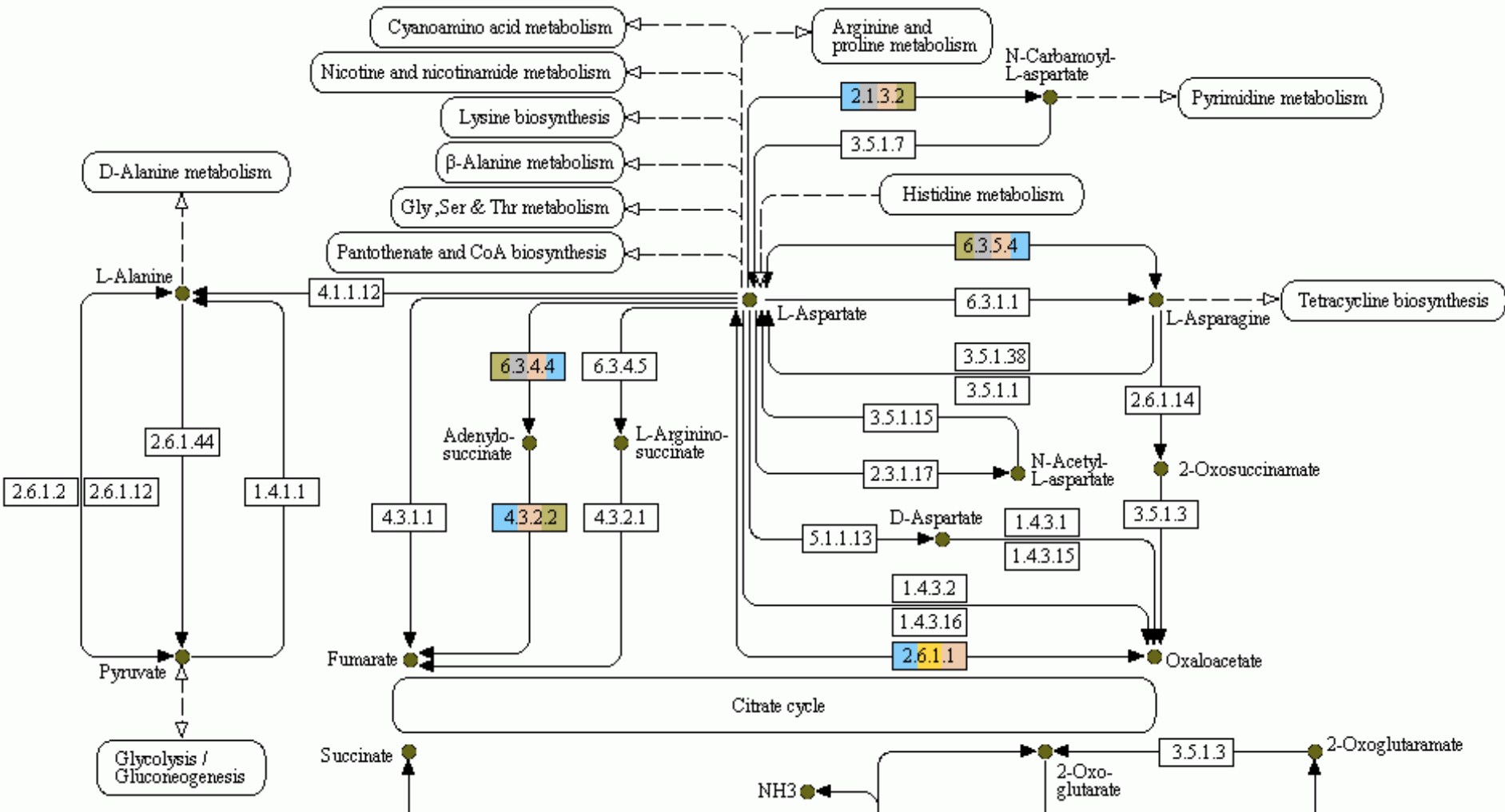
Mass Spectrometry can be used to measure metabolic and other chemical compounds

# Complex mixtures can be analyzed and interpreted

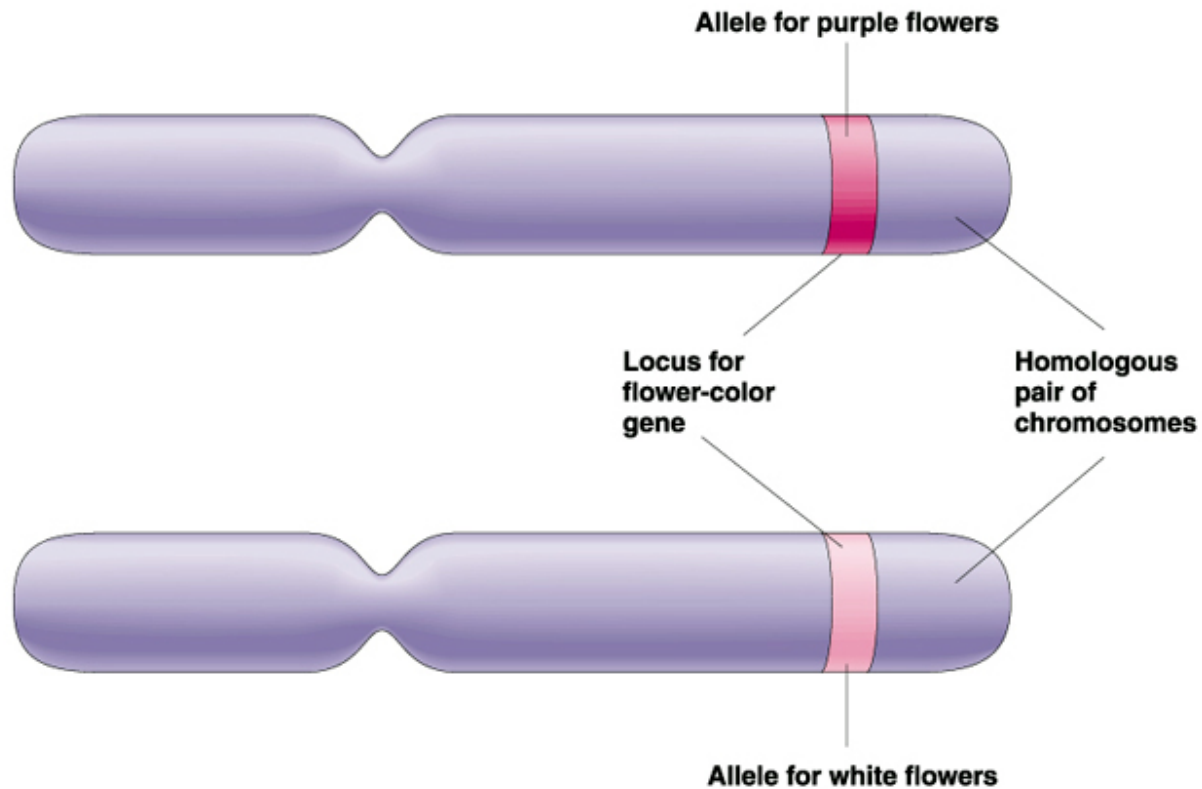


# Metabolites can be linked to metabolic pathways and enzymes

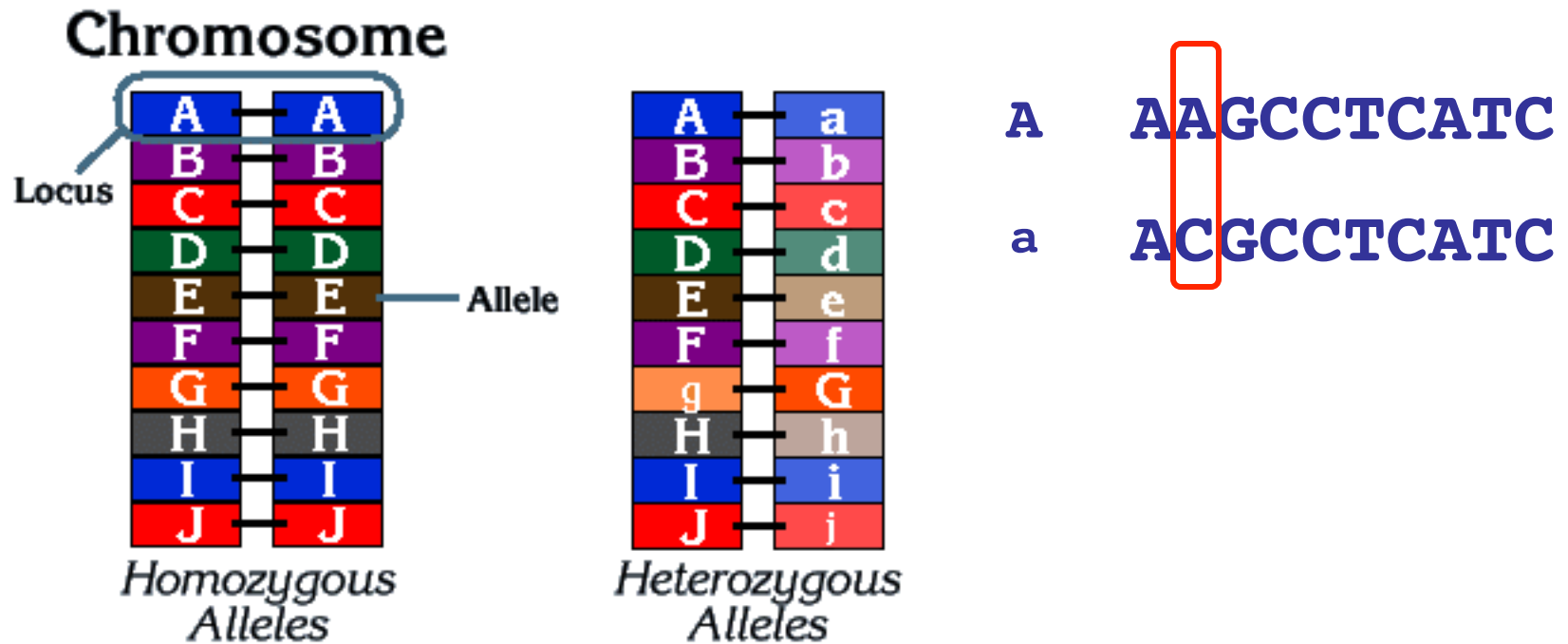
## ALANINE, ASPARTATE AND GLUTAMATE METABOLISM



# Homologous chromosomes (in a diploid)



# Loci, alleles and SNPs in a population



SNP = Single Nucleotide Polymorphism



# Alleles and Phenotype

- Some phenotypes are caused by a single locus in the genome and a single allele at that locus (e.g. some flower colors, or *Drosophila* eye color)
- Other phenotypes (Type-I diabetes, heart disease) are multi-locus or “complex” (i.e. many genes are involved, each potentially with many alleles)

# Population data

## Data

- Single Nucleotide Polymorphisms, SNPs
- Alleles
- Allele frequency
- Haplotypes

## Technology





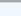
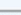















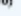








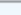
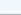


- Chip-Seq
- NGS

Allele Frequencies for **Exon 3 48 bp VNTR**

Locus [Dexamet receptor 1B](#)

[Information on Histograms](#)

Click on icon  for additional information.

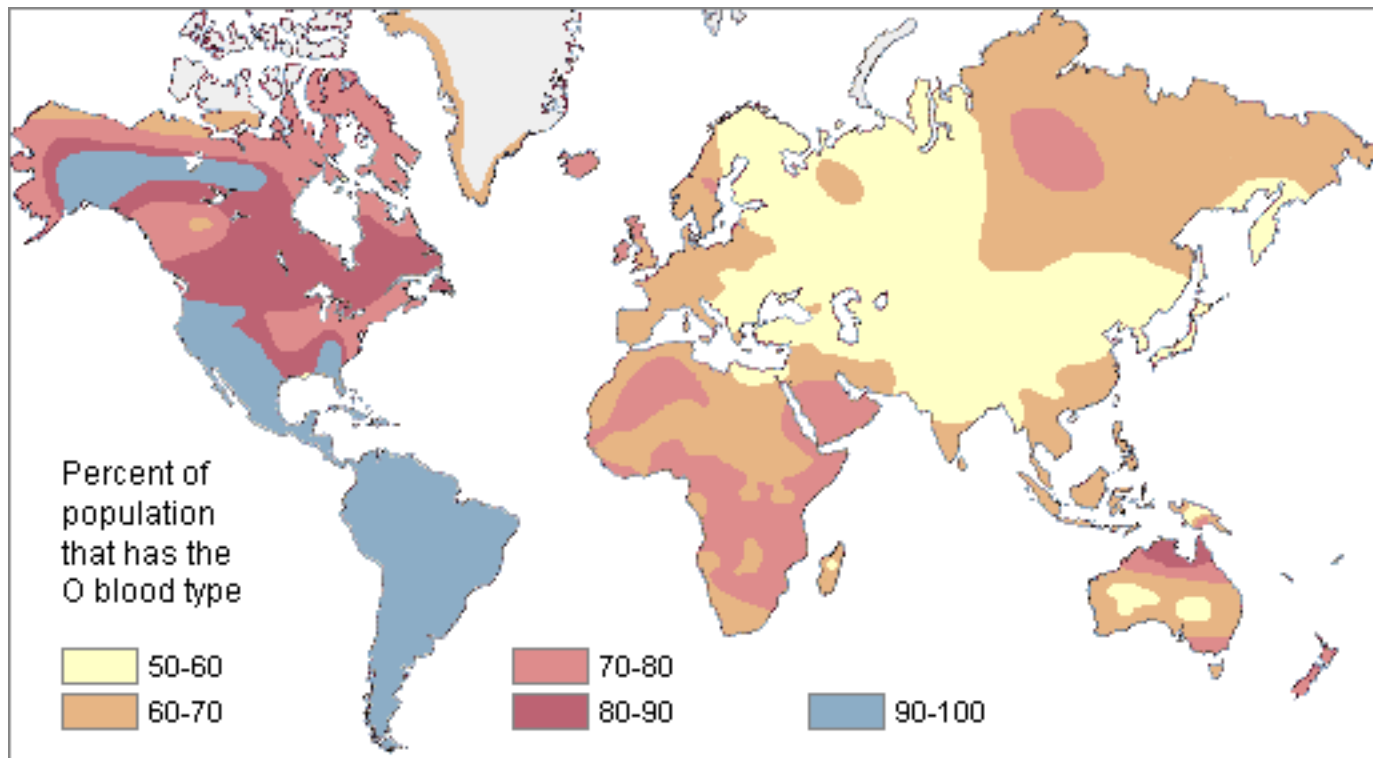
Geographic region	Population (Sample UID, Typed Sample Size (2N), entry date) Add Info	Allele Frequency Histogram
Africa	<a href="#">Baka(SA000005F)</a> , 134, 9/19/2000, 	
Africa	<a href="#">Mbuti(SA000006G)</a> , 72, 9/19/2000, 	
Africa	<a href="#">Jews, Ethiopian(SA000043N)</a> , 128, 9/19/2000, 	
Europe	<a href="#">Druze(SA000047L)</a> , 130, 9/19/2000, 	
Europe	<a href="#">Jews, Yemenite(SA000016H)</a> , 80, 9/19/2000, 	
Europe	<a href="#">Samaritans(SA000098R)</a> , 78, 9/19/2000, 	
Europe	<a href="#">Adyghej(SA000017I)</a> , 104, 9/19/2000, 	
Europe	<a href="#">Dane sj(SA000007H)</a> , 100, 9/19/2000, 	
Europe	<a href="#">Europeans, Mixed(SA000020C)</a> , 176, 9/19/2000, 	
Europe	<a href="#">Europeans, Mixed(SA001775U)</a> , 146, 6/21/2006, 	
Europe	<a href="#">Etrus(SA000018L)</a> , 66, 9/19/2000, 	
Asia	<a href="#">Brahmin(SA001831N)</a> , 156, 8/23/2006, 	
Asia	<a href="#">Keralite(SA001833P)</a> , 148, 8/23/2006, 	
Asia	<a href="#">Keralite(SA001834Q)</a> , 114, 8/23/2006, 	
Asia	<a href="#">Keralite(SA001835R)</a> , 130, 8/23/2006, 	
Asia	<a href="#">Marathas(SA001832G)</a> , 116, 8/23/2006, 	
Asia	<a href="#">Kachan(SA000040E)</a> , 36, 9/19/2000, 	
East Asia	<a href="#">Ami(SA000002C)</a> , 80, 9/19/2000, 	
East Asia	<a href="#">Atsuij(SA000021D)</a> , 84, 9/19/2000, 	
East Asia	<a href="#">Han(SA000000J)</a> , 96, 9/19/2000, 	
East Asia	<a href="#">Japanese(SA000010B)</a> , 100, 9/19/2000, 	
East Asia	<a href="#">Cambodians, Khmer(SA000023E)</a> , 50, 9/19/2000, 	
East Asia	<a href="#">Hakka(SA000003D)</a> , 32, 9/19/2000, 	
Oceania	<a href="#">Melanesian, Nasioi(SA000012D)</a> , 46, 9/19/2000, 	
Oceania	<a href="#">Micronesians(SA000063I)</a> , 58, 9/19/2000, 	
Siberia	<a href="#">Yakut(SA000011C)</a> , 92, 9/19/2000, 	
North America	<a href="#">Cheyenne(SA000023F)</a> , 96, 9/19/2000, 	
North America	<a href="#">Pima, Arizona(SA000025H)</a> , 94, 9/19/2000, 	
North America	<a href="#">Pima, Mexico(SA000026I)</a> , 104, 9/19/2000, 	
North America	<a href="#">Southern Amerindians(SA000024G)</a> , 86, 9/19/2000, 	
North America	<a href="#">Maya, Yucatan(SA000013E)</a> , 100, 9/19/2000, 	
South America	<a href="#">Gorillan(SA000148N)</a> , 58, 8/9/2005, 	
South America	<a href="#">Karitiana(SA000028K)</a> , 108, 9/19/2000, 	
South America	<a href="#">Surui(SA000014E)</a> , 90, 9/19/2000, 	

Alleles have frequencies in different populations



# Populations and alleles have geographic boundaries

A parasite isolate comes from a particular population, a particular location and will have a specific haplotype (e.g. representation of alleles) often characterized via SNPs



# Parasite Isolates

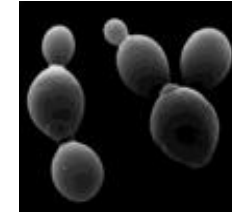
## Data

- Species, Strain,
- Isolate
- Location, Date
- SNP
- Sequence
- Allele
- phenotype

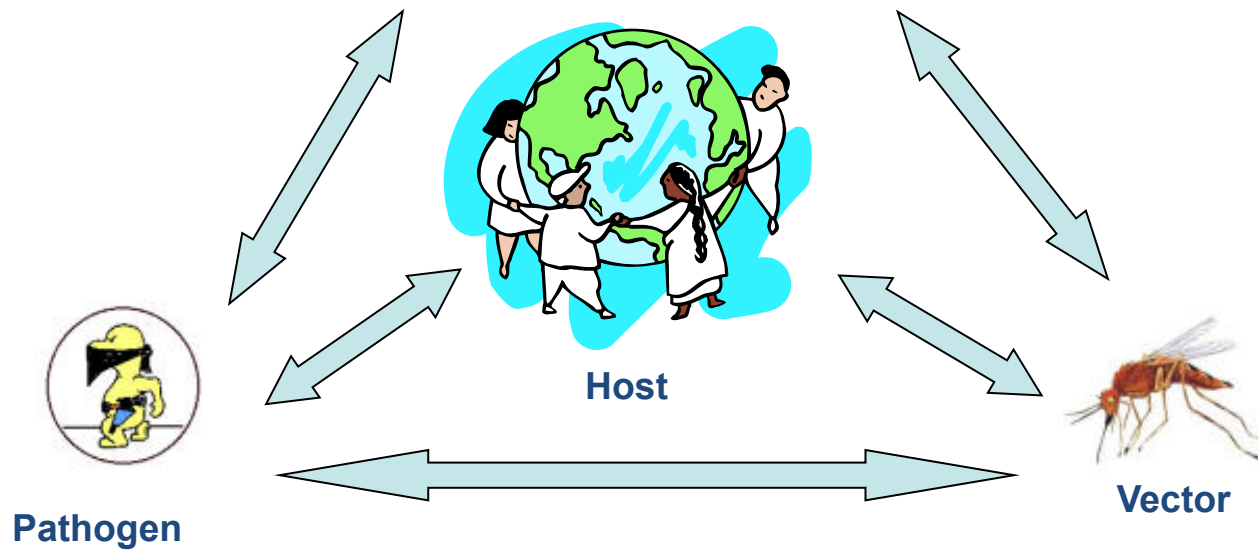
## Technology

- PCR-RFLP
- Microsatellites
- Sequencing
- SNP chip
- GPS

# Infectious Disease Paradigm



Experimental systems




# Important Information

- It might rain, especially Monday AM
- Air conditioned buildings may be cold, you may want a sweater
- You do not need your laptop for workshop sessions
- Buses are infrequent because we are in "Intersession"
- Guides (Brian and Omar) will leave hotel tomorrow lead folks to the workshop room via University bus

# UGA Campus Bus Routes Map

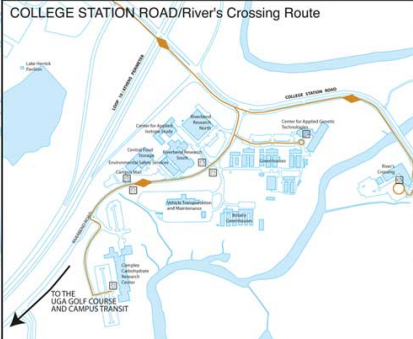
BUS ROUTES AND STOPS			
Route	Color	Route	Color
Ag Hill	AG	North/South	NS
East Campus Express	ECX	Orbit	O
East/West	EW	River's Crossing	RC
Family Housing	FH	Russell Hall	RH
Millidge Avenue	M	Weekender	WE
Bus Route travels in both directions	↔	Bus Route travels in direction of arrow	→
Bus Stops			



**Athens Transit**  
Routes and Numbers identified by **25**

For Further Information about these routes  
Call Athens Transit at 613-3430

**COLLEGE STATION ROAD/River's Crossing Route**

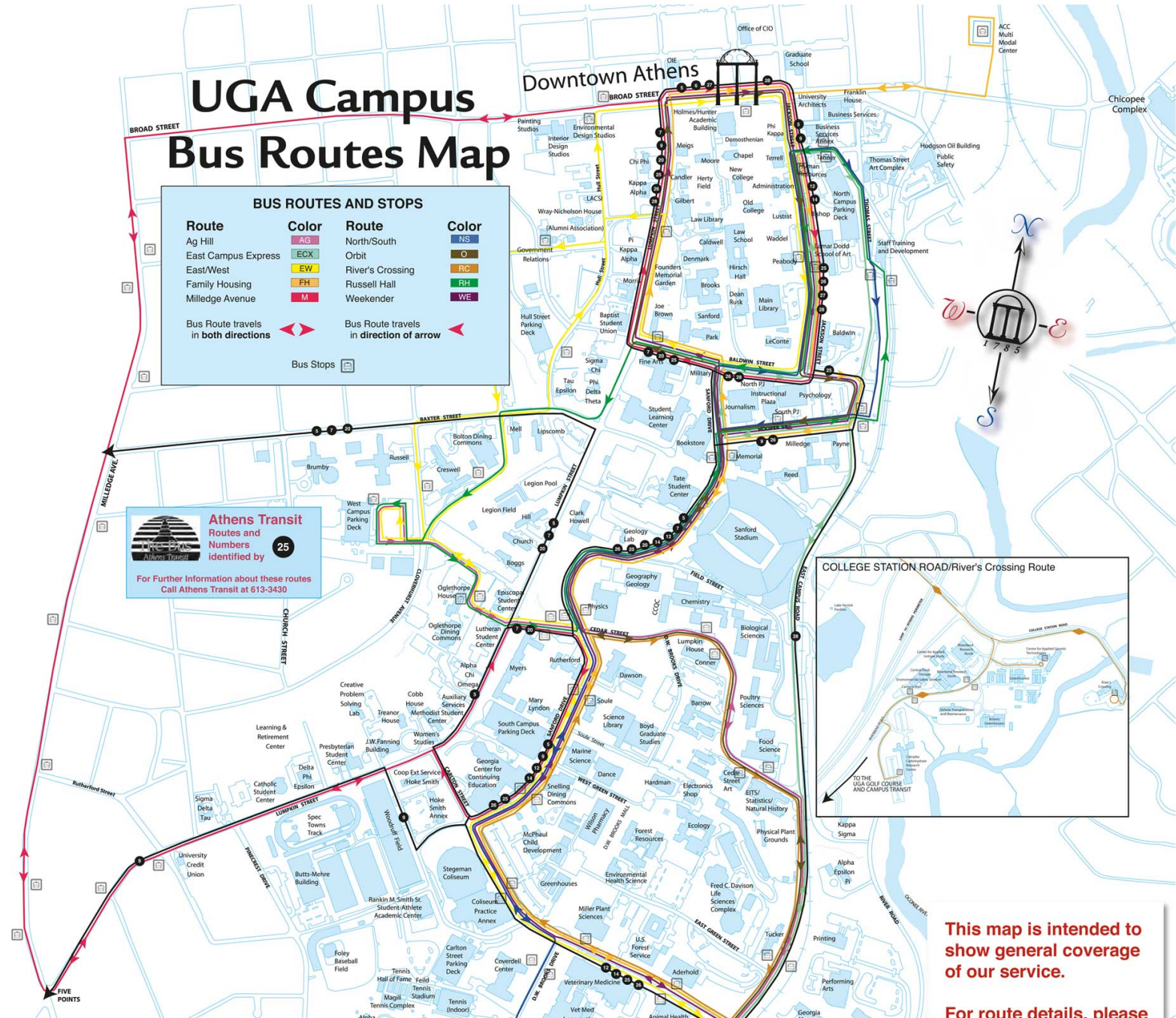


TO THE UGA GOLF COURSE AND CAMPUS TRANSIT

**General Information**

**This map is intended to show general coverage of our service.**

**For route details, please refer to the individual**

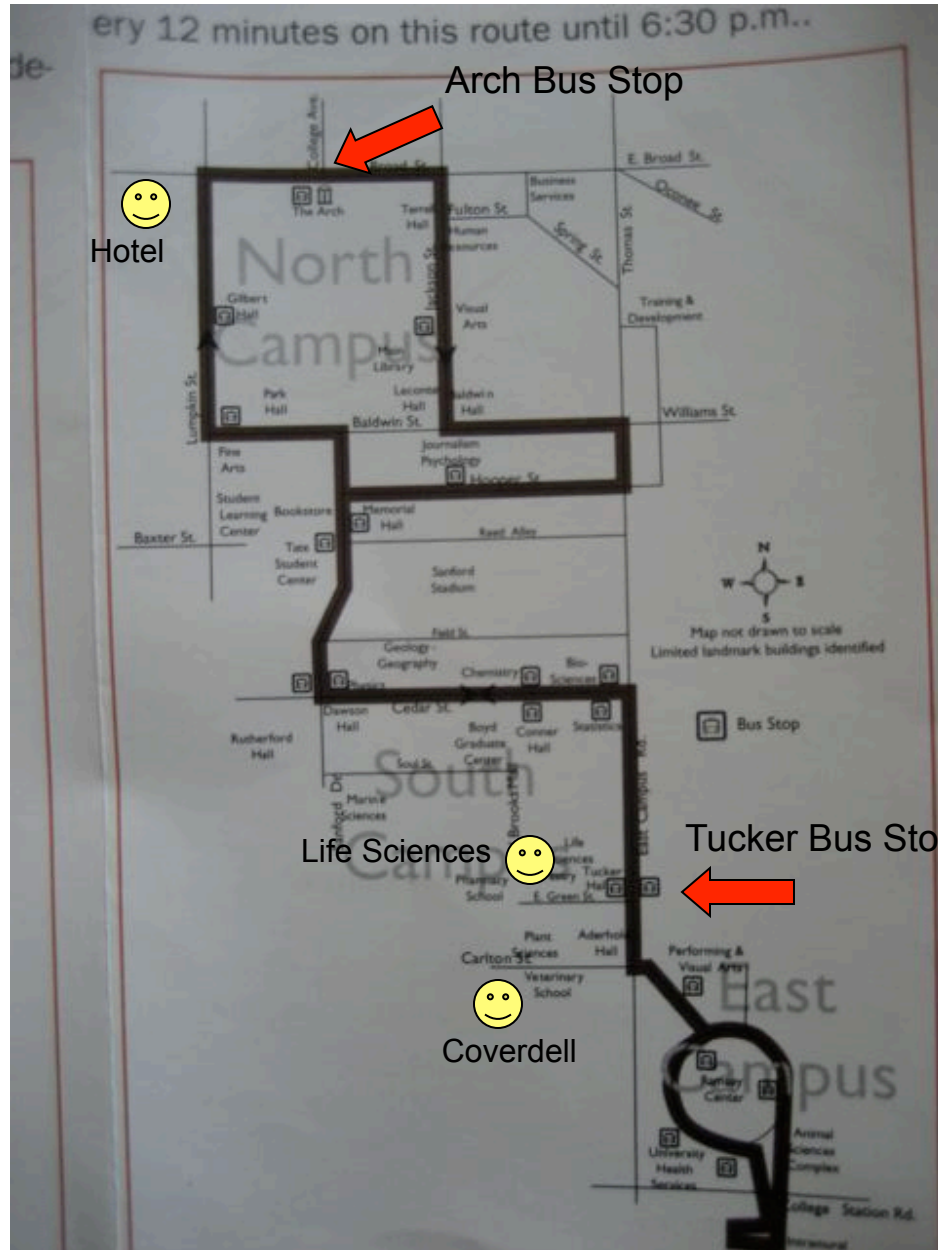




Orbit Bus Route  
**NOT Running**

~10 blocks  
1.1 miles

3 blocks, .6 Km



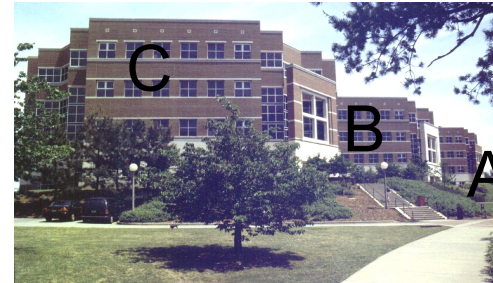
# Bad Bus



Get off at Tucker Hall



Return to Hotel



Life Sciences



# Good Bus



Coverdell

Take the Orbit bus from the Arch to Tucker Hall at East Campus and Green Street  
The Life Sciences building is on Greene street, up the hill on the Right Rm C128

# General Schedule

- Workshop all day Mon - Wed Schedule in your packets
- Dinner at the Coverdell Center on Monday (Informal!)
- Breakfast (2 choices)
  - Holiday Inn (coupons received at check-in)
  - Have nibbles with morning coffee break
- Coffee and lunch provided (M-W).  
Tuesday and Wednesday Dinner is on your own.

Questions?

# Now that you know us, We would like to know you

- Please tell us your name and a bit about your research.
- Also, if there is something specific that you came here to learn, please state it so that we can cover it if possible