



VEuPathDB

Eukaryotic Pathogen, Vector & Host Informatics Resources

Crash Course in Omics Terminology, Concepts & Data Types

Jessie C Kissinger
2024



@veupathdb

@jcklab

jkissing@uga.edu



UNIVERSITY OF
GEORGIA
Center for Tropical &
Emerging Global Diseases



Institute of Bioinformatics
UNIVERSITY OF GEORGIA



VEuPathDB

Eukaryotic Pathogen, Vector & Host Informatics Resources

758 genome sequences,
3083 datasets in Release 66!



PlasmoDB

Plasmodium Informatics Resources



ToxoDB

Toxoplasma Informatics Resources



CryptoDB

Cryptosporidium Informatics Resources



PiroplasmaDB

Piroplasma Informatics Resources



TrichDB

Trichomonas Informatics Resources



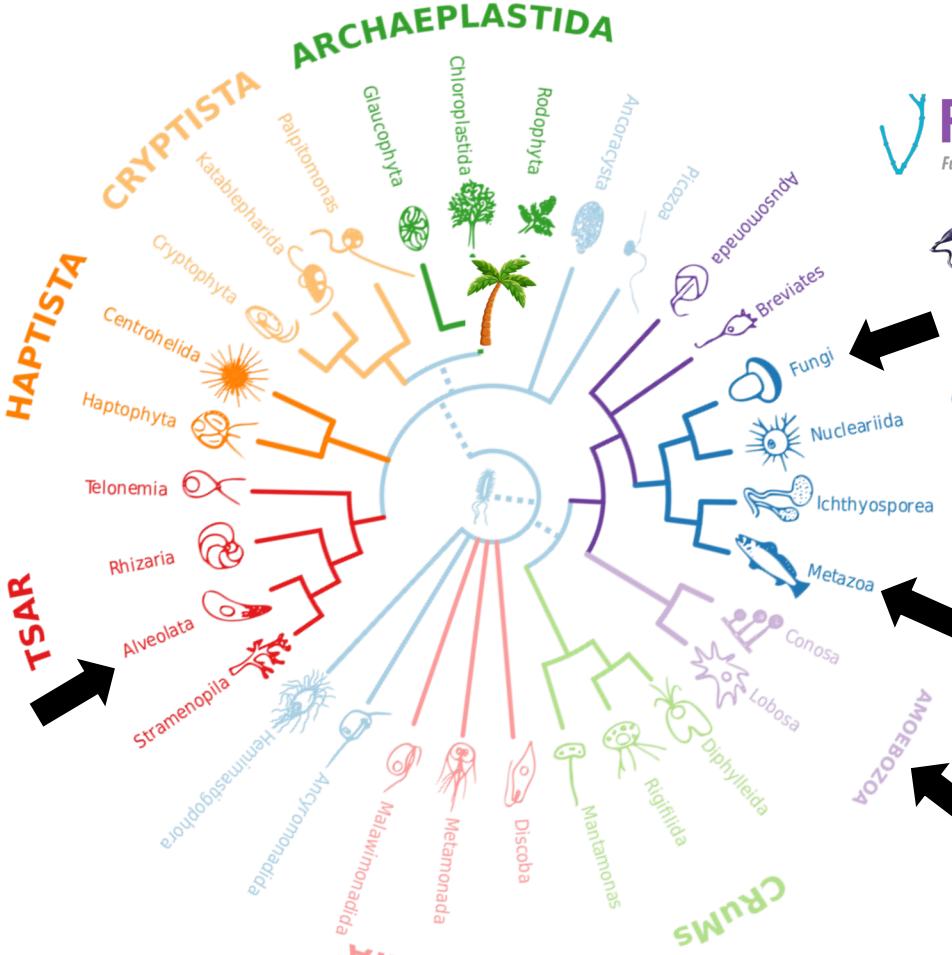
GiardiaDB

Giardia Informatics Resources



TriTrypDB

Kinetoplastid Informatics Resources



FungiDB

Fungal & Oomycete Informatics Resources



MicrosporidiaDB

Microsporidia Informatics Resources



HostDB

Pathogen Host Informatics Resources



VectorBase

Bioinformatics Resources for
Invertebrate Vectors of Human Pathogens



AmoebaDB

Amoeba Informatics Resources



MicrobiomeDB

A Microbiome Resource



ClinEpiDB

Clinical Epidemiology Resources



OrthoMCL DB

Ortholog groups of Protein sequences



KAYAK

Round-trip

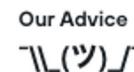
Atlanta ✈ + ↗

Tunis ✈ + ↗

Wed 11/29 – Wed 12/6 1 adult, Economy



Sign in



We're still gathering data for this route

Track prices Off

224 of 633 flights

Stops

- Nonstop
- 1 stop \$1,117
- 2+ stops \$894

Fee Assistant i

- Carry-on bag - 0 +
- Checked bag - 0 +

Book on KAYAK ⚡

Show offers instantly bookable on KAYAK.

Times

Take-off Landing

Take-off from ATL

Cheapest

\$894 • 18h 37m

Best i

\$1,121 • 13h 07m

Quickest

\$1,121 • 13h 07m

Other sort

priceline®

Land a great deal for less. Book your flight with confidence.

Go To Your Happy Price. Book now travel anytime.



10:10 am – 10:35 am⁺¹

ITA Airways

2 stops

IAD, FCO

18h 25m

ATL-TUN



\$904

Economy
Priceline

[View Deal](#)



11:25 am – 12:14 am⁺¹

ITA Airways

2 stops

FCO, JFK

18h 49m

TUN-ATL

Operated by Delta Air Lines

Ad

Best



3:40 pm – 9:50 am⁺¹

Delta

1 stop

CDG

12h 10m

ATL-TUN



\$1,121

Basic Economy
Delta

[View Deal](#)



5:30 am – 1:35 pm

Delta

1 stop

CDG

14h 05m

TUN-ATL

Operated by Air France

Main Cabin

\$1,301

Cheapest



10:10 am – 10:35 am⁺¹

ITA Airways

2 stops

IAD, FCO

18h 25m

ATL-TUN



11:25 am – 12:14 am⁺¹

ITA Airways

2 stops

FCO, JFK

18h 49m

TUN-ATL

\$894

Economy
ScholarTrip

[View Deal](#)

The Travel Site has Very Useful Data Filters!

OUR ADVICE
We're still gathering data for this route

Track Prices OFF

Stops

- Nonstop
- 1 stop \$1553
- 2+ stops \$2821

Times

Take-Off **Landing**

Take-Off from ATL
Sun 5:30 PM – 11:30 PM

Take-Off from TUN
Thu 8:00 AM – 8:00 PM

Airlines

- Alitalia \$3265
- Delta \$3260
- Frontier
- Qatar Airways
- Tunisair
- Multiple airlines ⓘ

Show 8 more airlines

Duration

Flight Leg
13h 45m – 41h 17m

Layover
0h 55m – 22h 55m

Price

\$1553 – \$15484

Cabin

- Economy \$1553
- Prem Econ \$4956
- Business \$10933
- Mixed \$6037

Flight Quality

- Show Wi-Fi Flights Only
- Show Hacker Fares¹
- Show Red-Eyes
- Show 65 Longer Flights

Alliance

- oneworld
- SkyTeam \$2821
- Star Alliance \$1779

Aircraft

- Narrow-Body Jet
- Wide-body jet

Layover Airports

- Algeria
- Algiers (ALG)
- Canada
- Toronto (YYZ)
- France
- Paris (CDG)
- Germany
- Frankfurt am Main (FRA)
- Stuttgart (STR)

Booking Sites

- Airlines Only
- Air France \$6863
- Alitalia \$3265
- Delta \$3260
- Expedia
- FlightHub \$1779

Show 6 more sites

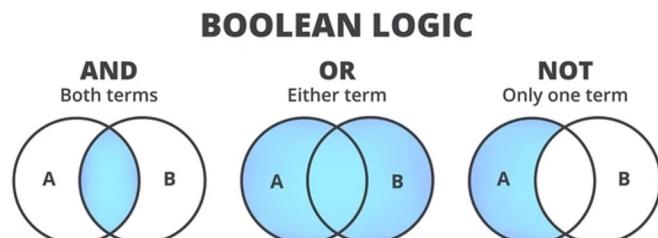
Filters vs Boolean operators

Filters

- Can only narrow down the original search
- They only return a subset of the original data
- Examples:
 - All genes on chromosome 4
 - All genes with "kinase in their name
 - All genes from *Trypanosoma cruzi*

Boolean operators (and, or & not)

- Intersect, union, subtract
- They can operate on two different searches!
- They can narrow down, or, expand the original search
- Examples:
 - All genes on Chr 4 that have kinase in their name
 - All genes on chr 4 or chr 8
 - All genes in *T. cruzi* that also have a signal peptide



The Biological Equivalent of Travel Search Engine with Filters and Boolean Logic

- Find all genes that....
 - That are near centromeres
 - That encode a predicted signal protein
 - That encode the amino acid motif CC..CC
- Which have evidence of expression ...
 - In developmental stage X
 - After treatment with drug Y
- That are phosphorylated in proteomic studies
- That show evidence of diversifying selection in population studies

Searching biological data is difficult because there are so many different technologies!

- Each technology e.g. genomics, transcriptomics, proteomics, metabolomics, etc.. has its own vocabulary that is more complicated than selecting a window or aisle seat.
- So,...to use the databases efficiently, you do not need to be a bioinformatician, rather you need to be an expert on the technologies related to the data you will mine so you can use the filters and Boolean operators well and interpret your results.
- Since nobody can keep up with all of the technologies and terminologies, and because we come from so many different backgrounds, we have created this crash course in omics

Most Genomic terminology in VEuPathDB refers to the following concepts:

Genome assembly: Reads, contigs, scaffolds, chromosomes, genome sequences, gaps, indels rearrangements, sequence

Genome annotation: Genes, sequence, coding and non-coding, intergenic regions, untranslated regions, introns, Promoters

Evolution: Sequence differences, SNPs, SNV, InDels, synonymous, non-synonymous, orthologs, paralogs, homology

Chromatin status: Epigenetics, Methylation, open chromatin, closed chromatin

Gene expression: Transcripts, splicing, alternative splicing, differential expression, expression levels (relative or absolute), transcript modifications. Analyses can bulk on a tissue or population of cells/organisms or can be single-cell

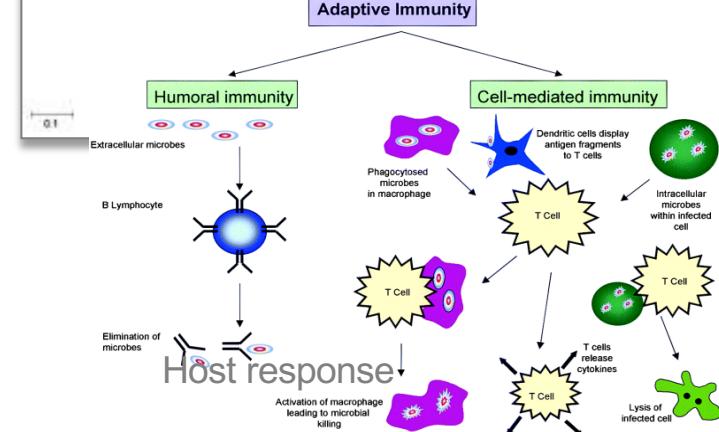
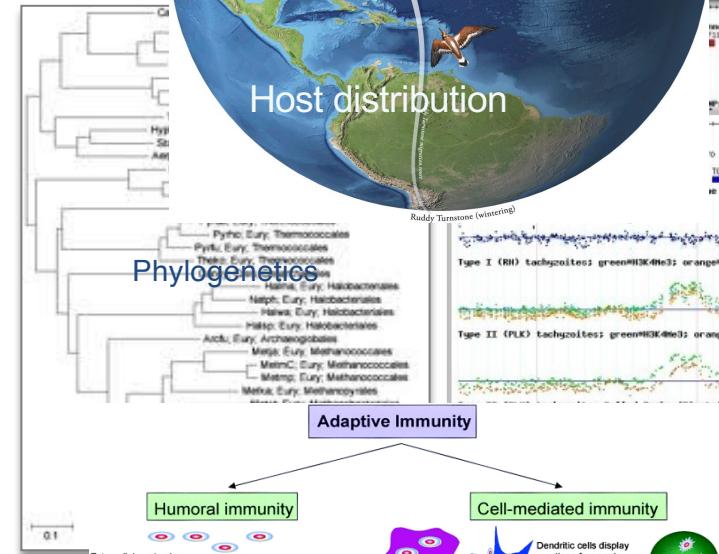
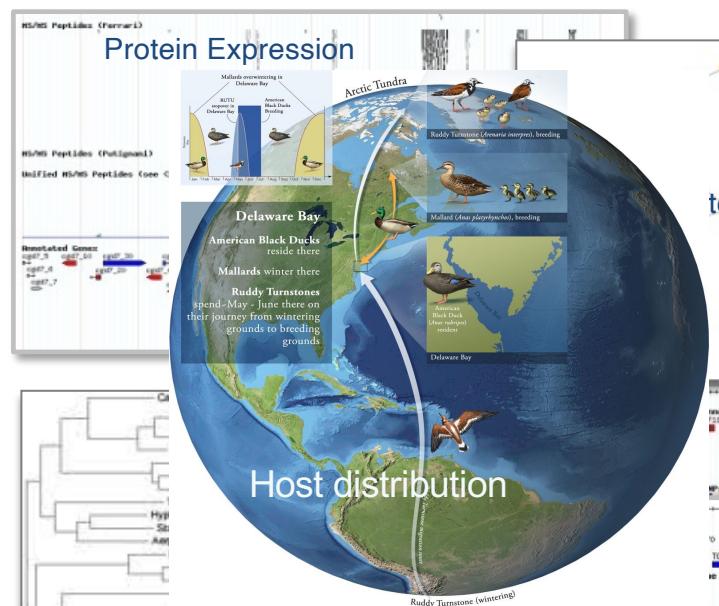
Proteins: Sequence, protein features (motifs, signal peptides, TM domains: chemical properties, chemical modifications (phosphorylation, glycosylation), expression, processing, localization

Metabolites: Chemical compounds, enzymes, pathways, flux

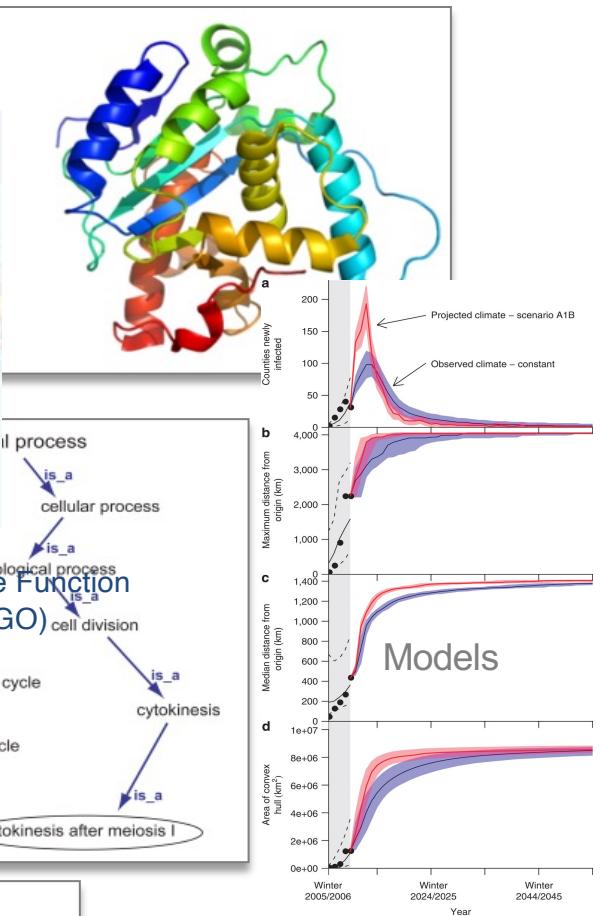
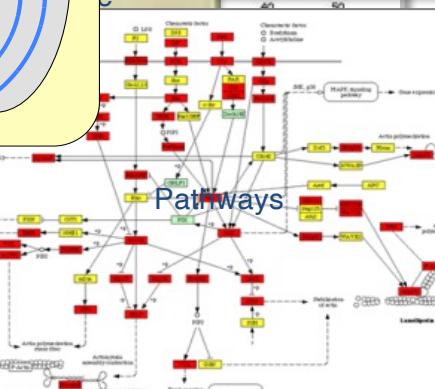
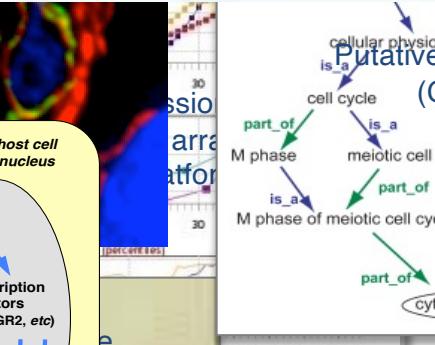
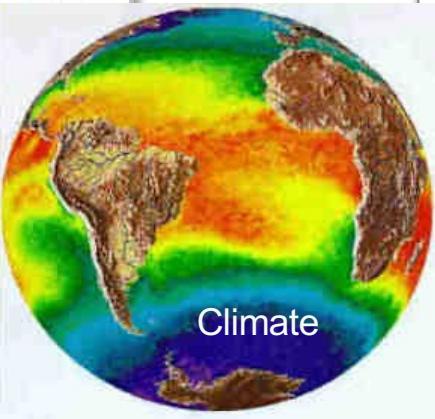
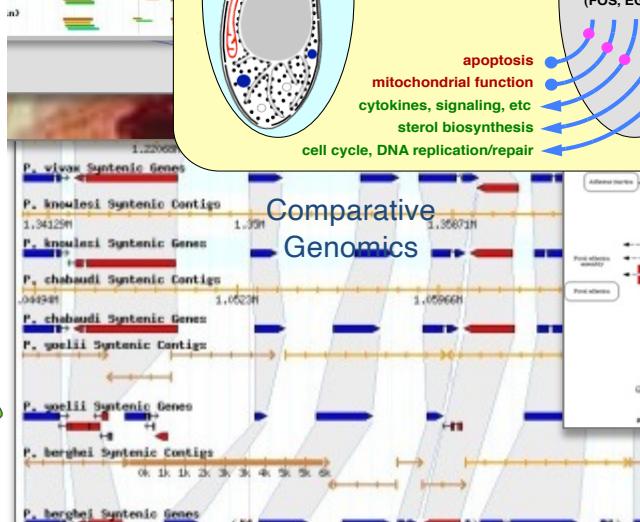
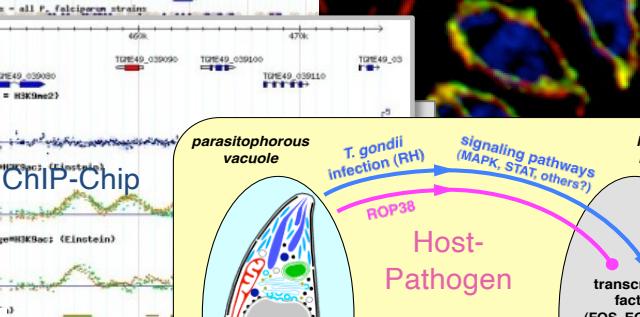
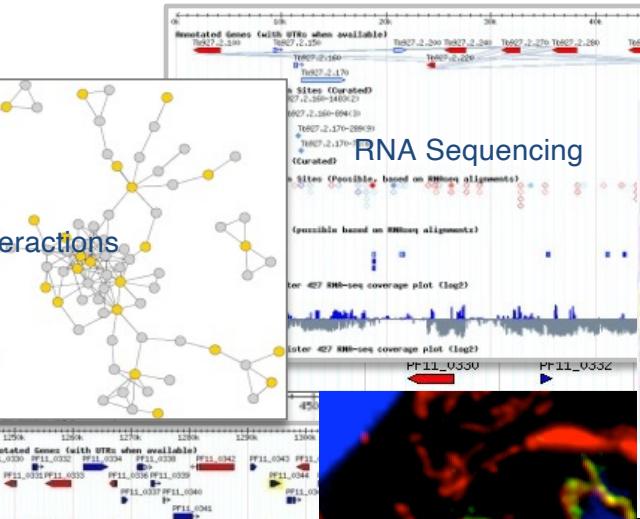
Host(s): Host response, immune responses, gene regulation responses, metabolic responses

Mutant analysis: Phenotypic response to gene knock-down or knock out, e.g. via CRISPR or other approach, or specific mutations

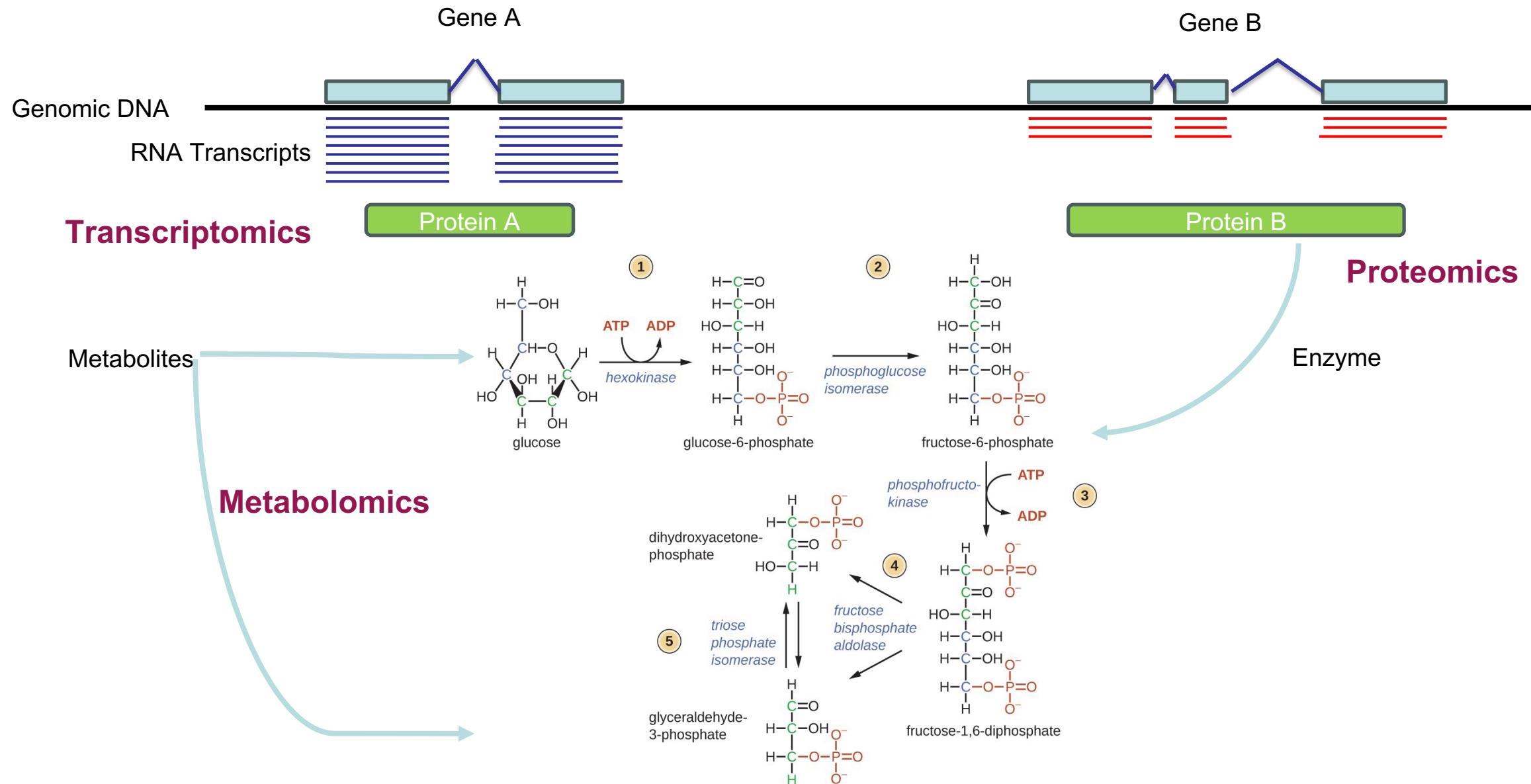
Metadata: Data about the data, e.g. the patient, source, environment or experimental condition



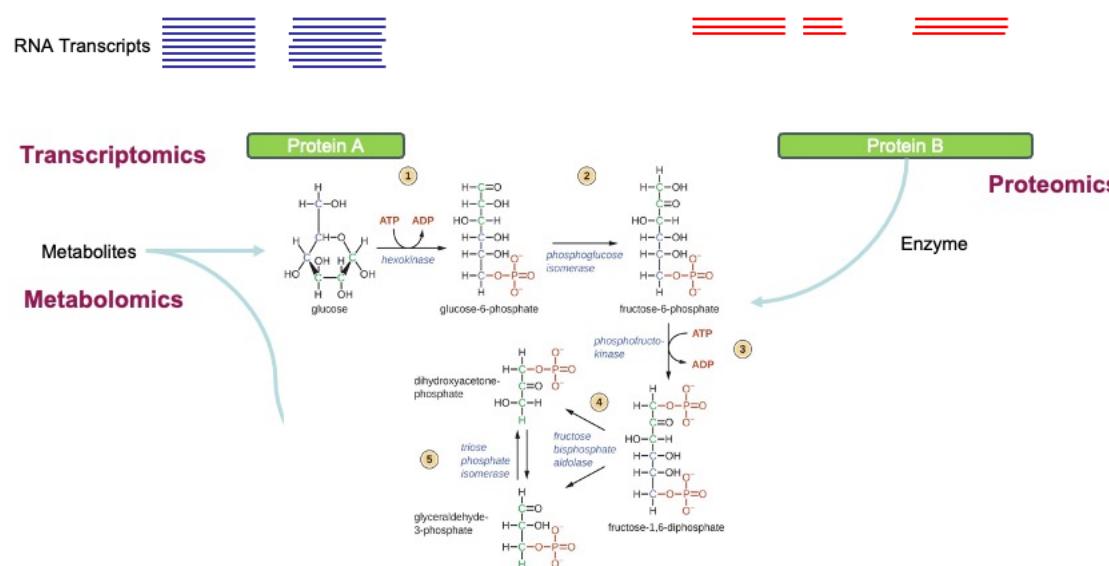
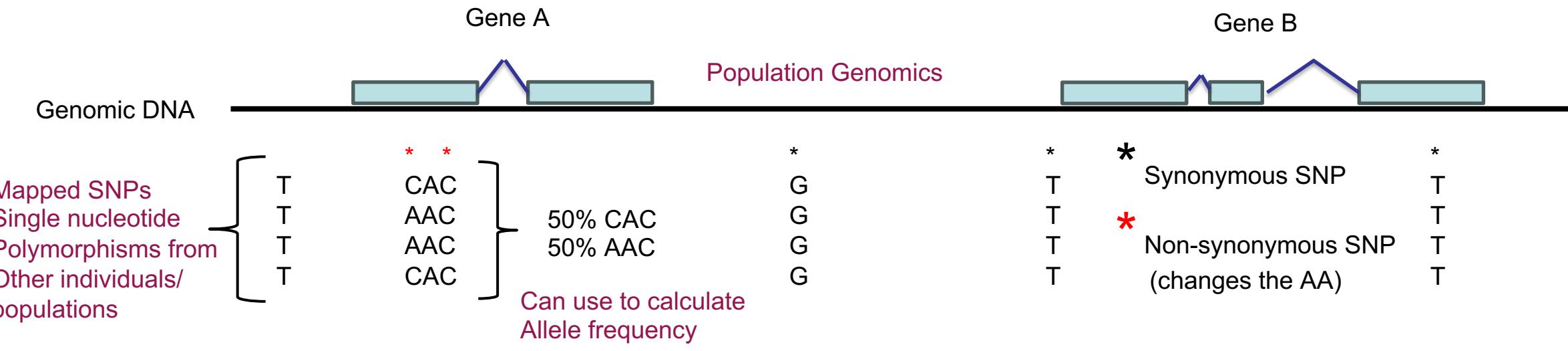
Modified from slide provided
by David Roos



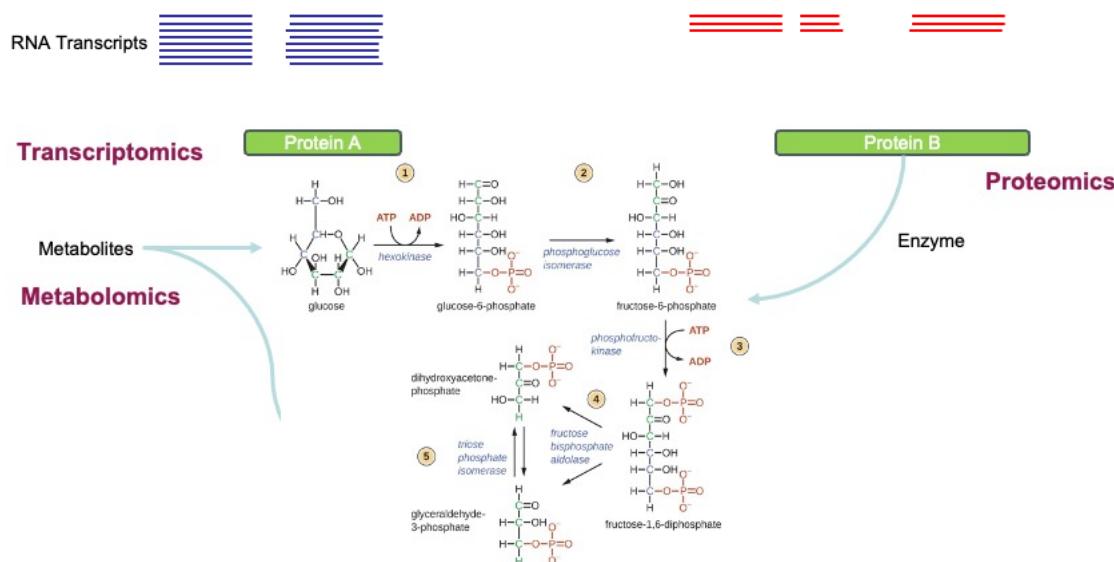
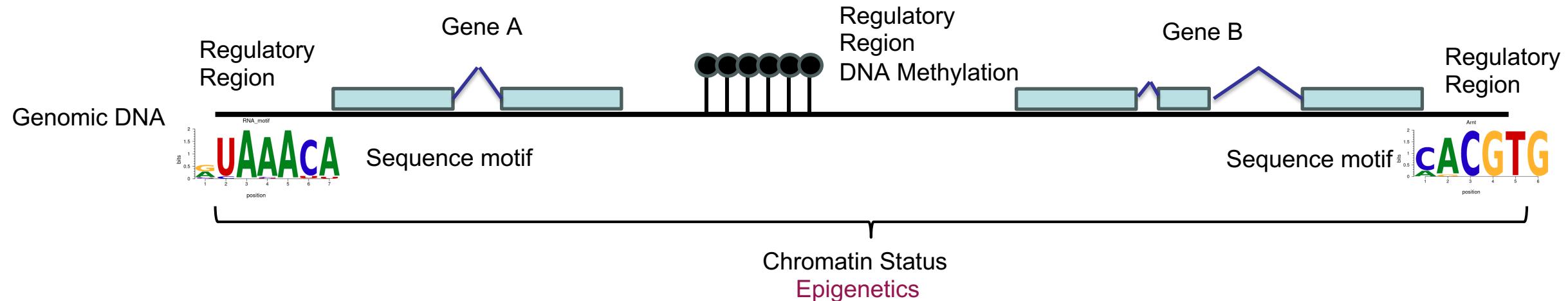
Conceptually, how are the data integrated?



Conceptually, how are the data integrated?



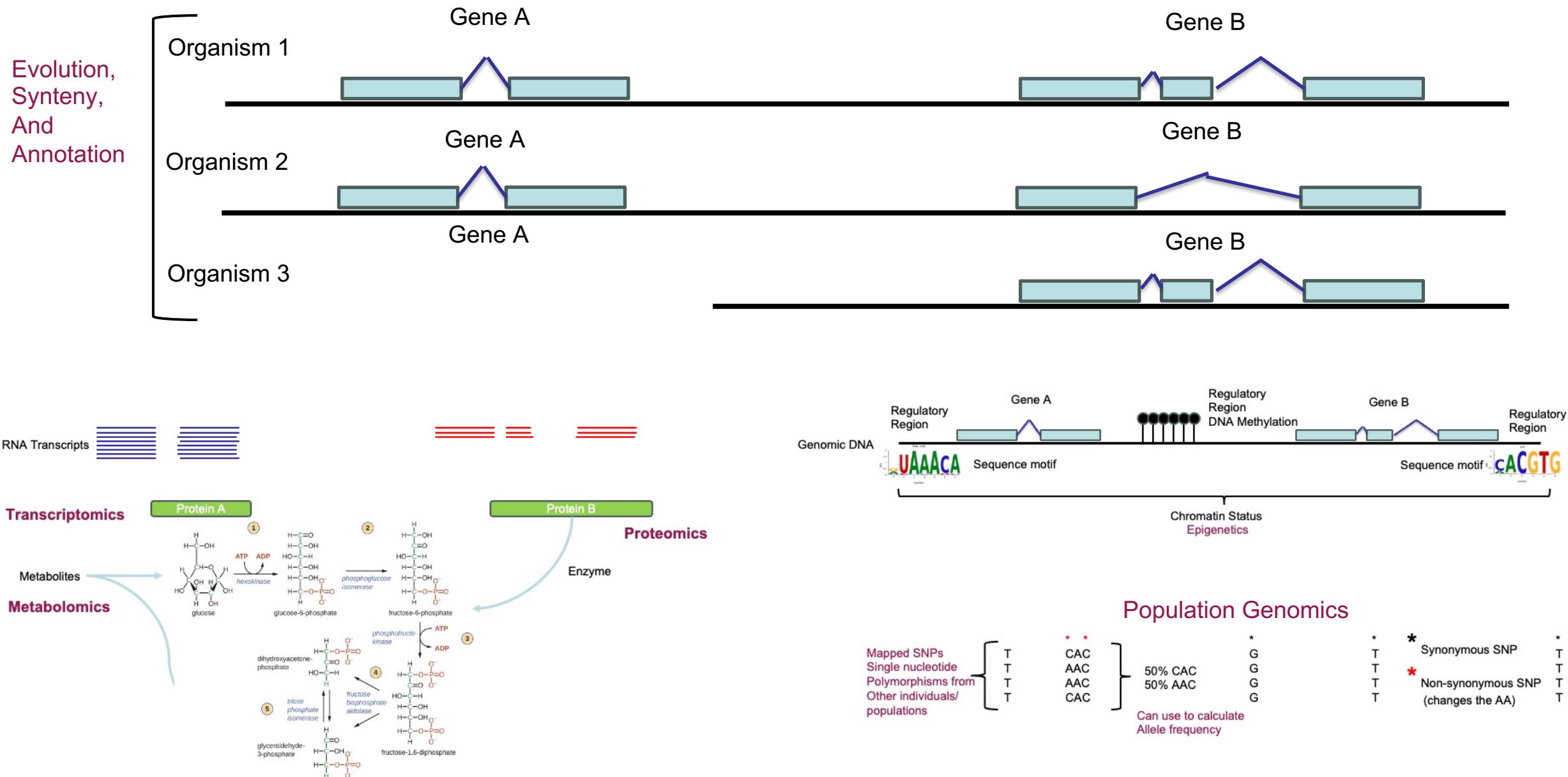
Conceptually, how are the data integrated?



Mapped SNPs Single nucleotide Polymorphisms from Other individuals/ populations	<table border="1"> <tr> <td>CAC</td><td>*</td></tr> <tr> <td>AAC</td><td>*</td></tr> <tr> <td>AAC</td><td>*</td></tr> <tr> <td>CAC</td><td>*</td></tr> </table>	CAC	*	AAC	*	AAC	*	CAC	*	50% CAC 50% AAC	Can use to calculate Allele frequency								
CAC	*																		
AAC	*																		
AAC	*																		
CAC	*																		
		<table border="1"> <tr> <td>T</td><td>*</td> </tr> <tr> <td>T</td><td>*</td> </tr> <tr> <td>T</td><td>*</td> </tr> <tr> <td>T</td><td>*</td> </tr> </table>	T	*	T	*	T	*	T	*	<table border="1"> <tr> <td>G</td><td>*</td> </tr> <tr> <td>G</td><td>*</td> </tr> <tr> <td>G</td><td>*</td> </tr> <tr> <td>G</td><td>*</td> </tr> </table>	G	*	G	*	G	*	G	*
T	*																		
T	*																		
T	*																		
T	*																		
G	*																		
G	*																		
G	*																		
G	*																		

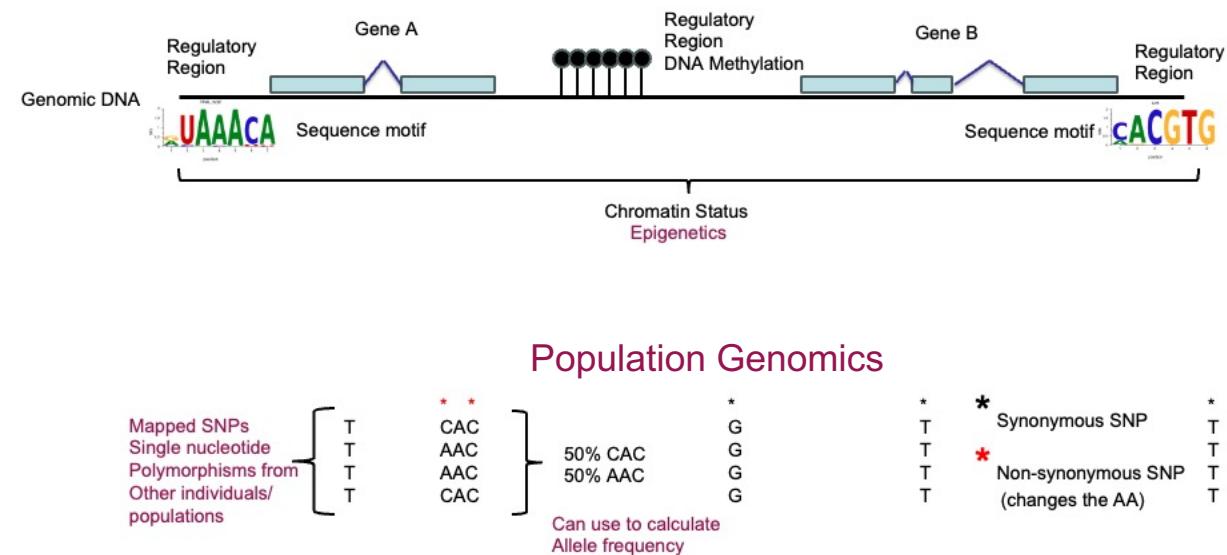
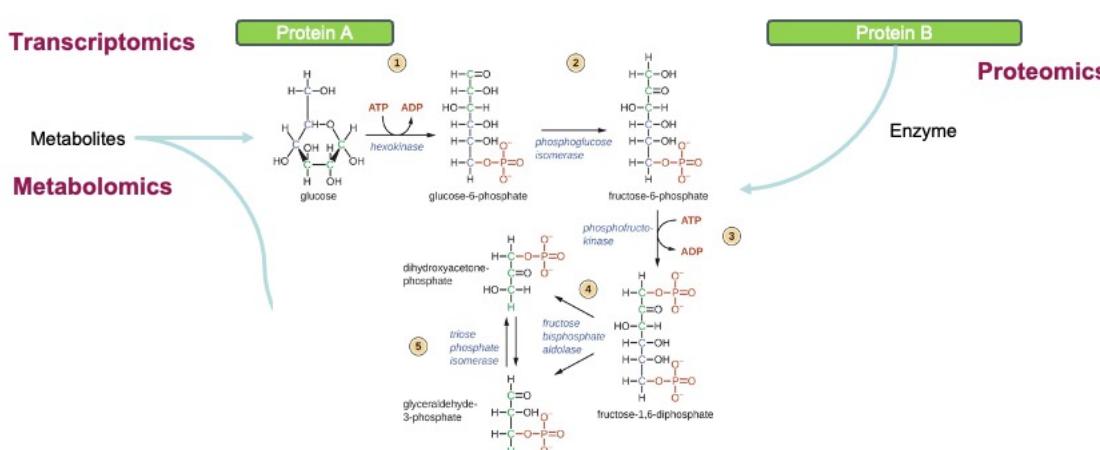
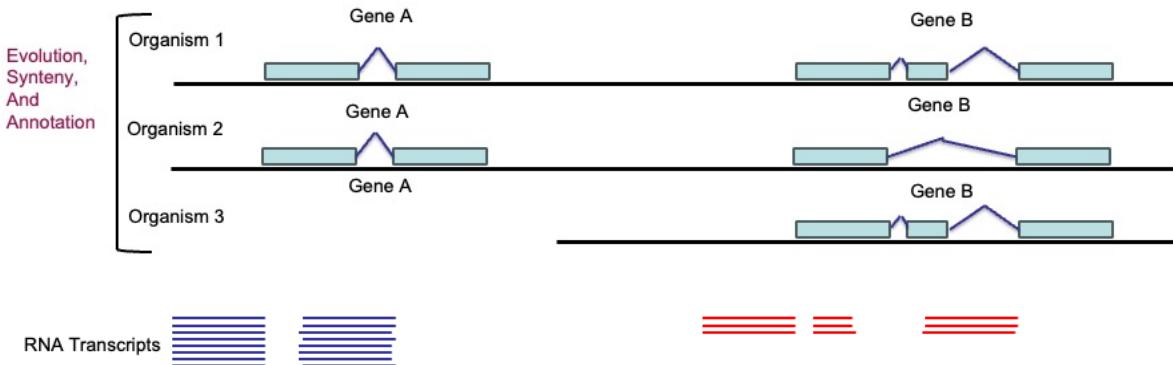
* Synonymous SNP
* Non-synonymous SNP (changes the AA)

Conceptually, how are the data integrated?



Conceptually, how are the data integrated?

Biologically! Most data can be mapped directly to a genome sequence or to a feature that is mapped to a genome sequence thus, we can use the genomic backbone of coordinates to relate genes, to transcripts and proteins and SNPs. We have also linked the meta data about the genomes, e.g. organisms, taxonomy and geographic location. Together the system is powerful.



So let's look at a few data types in more detail

Genome sequences & Assembly

FASTQ
format for
reads

Label
Sequence
Base = T, Q = A = 25
Q Scores (as ASCII charts)

```
@FORJUSP02AJWD1
CCGTCAATTCTATTAAAGTTAACCTTGCAGGCCGTACTCCCAGGCGGT
+
AAAAAAA::99@:::?:@::FFAAAAACCAA::::BB@@?A?
```

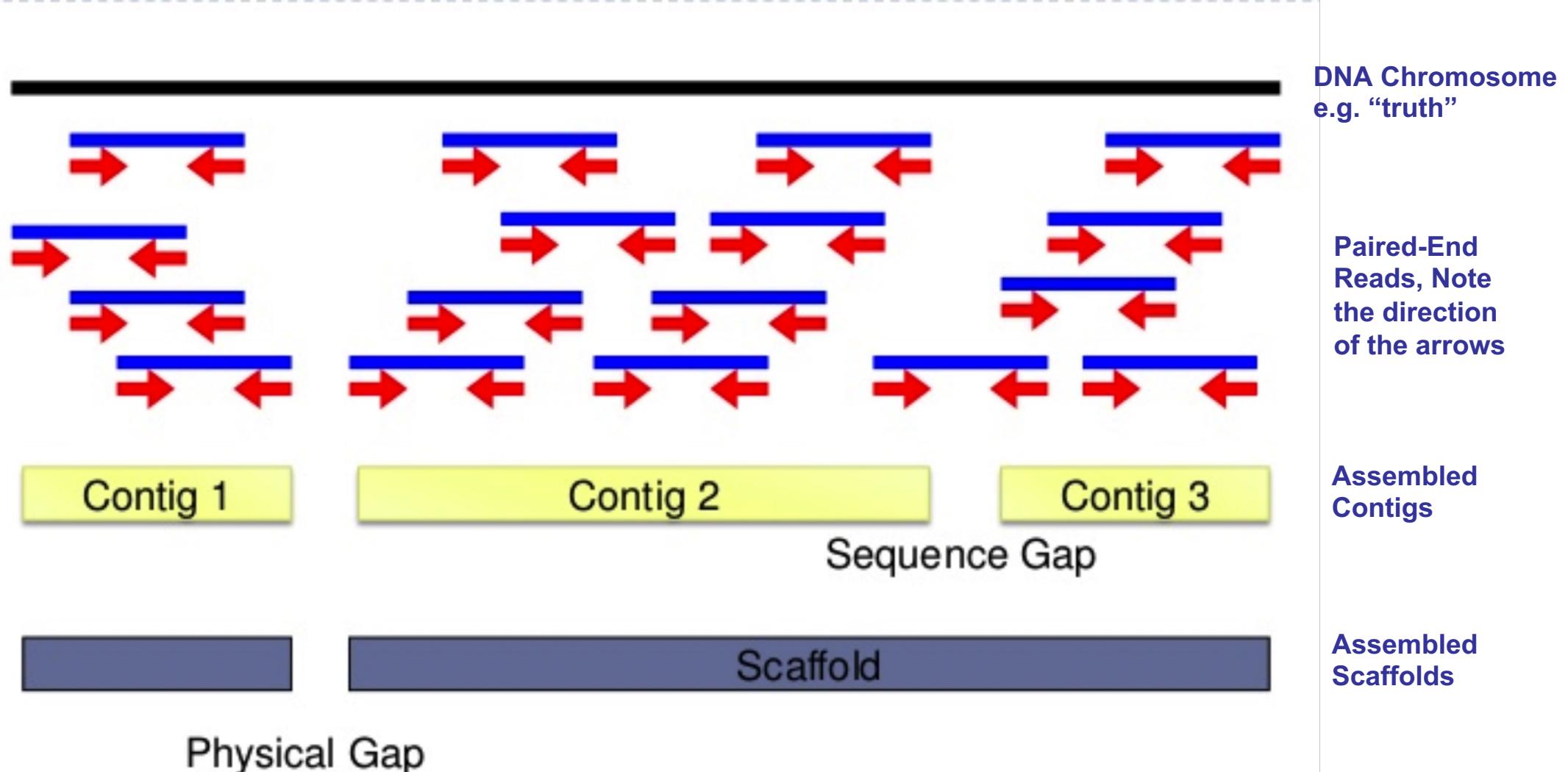
Figure 2 – Flowchart of an NGS workflow

$$Q = -10 \log_{10} P$$

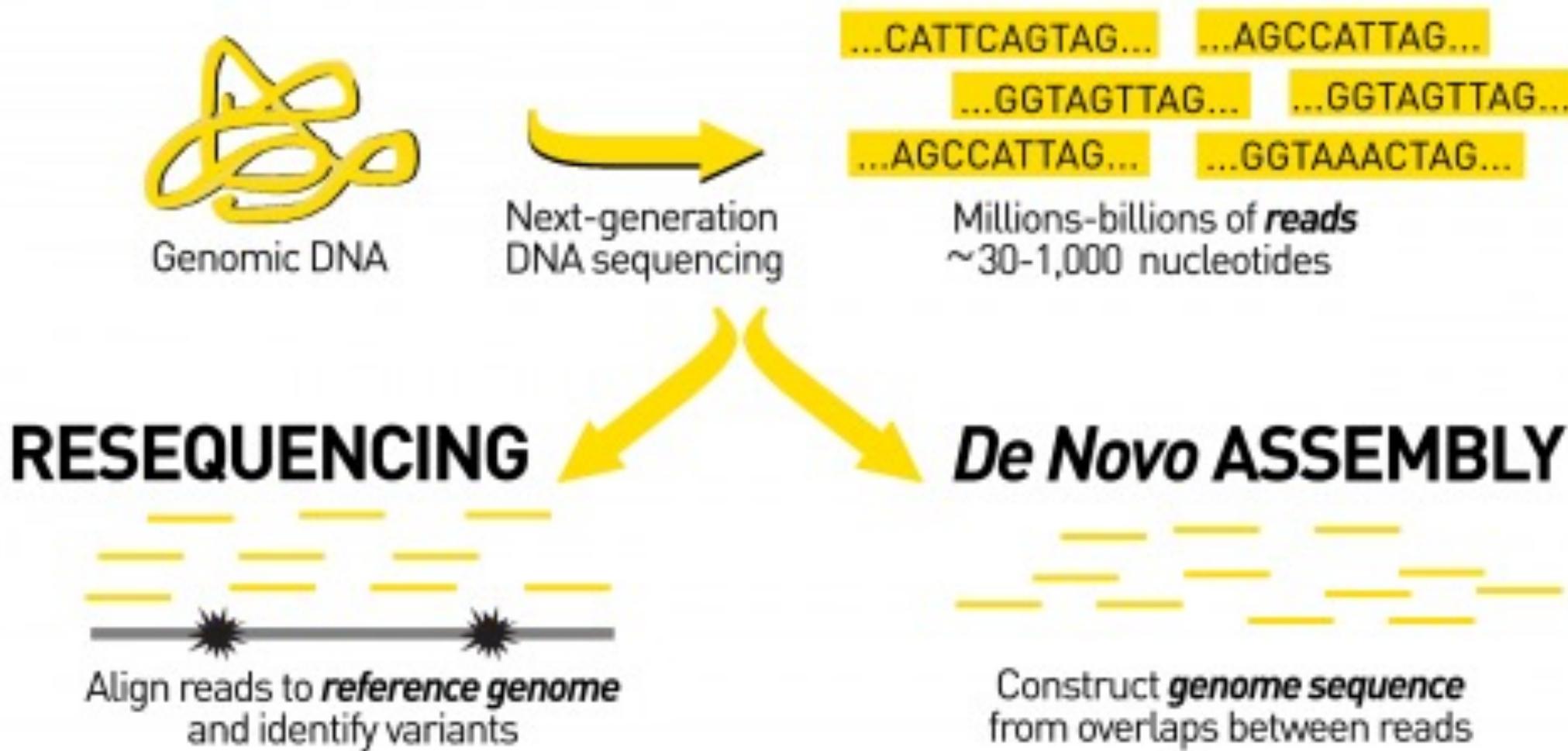
Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10,000	99.99%
50	1 in 100,000	99.999%
60	1 in 1,000,000	99.9999%

Figure 3 – Phred quality score chart

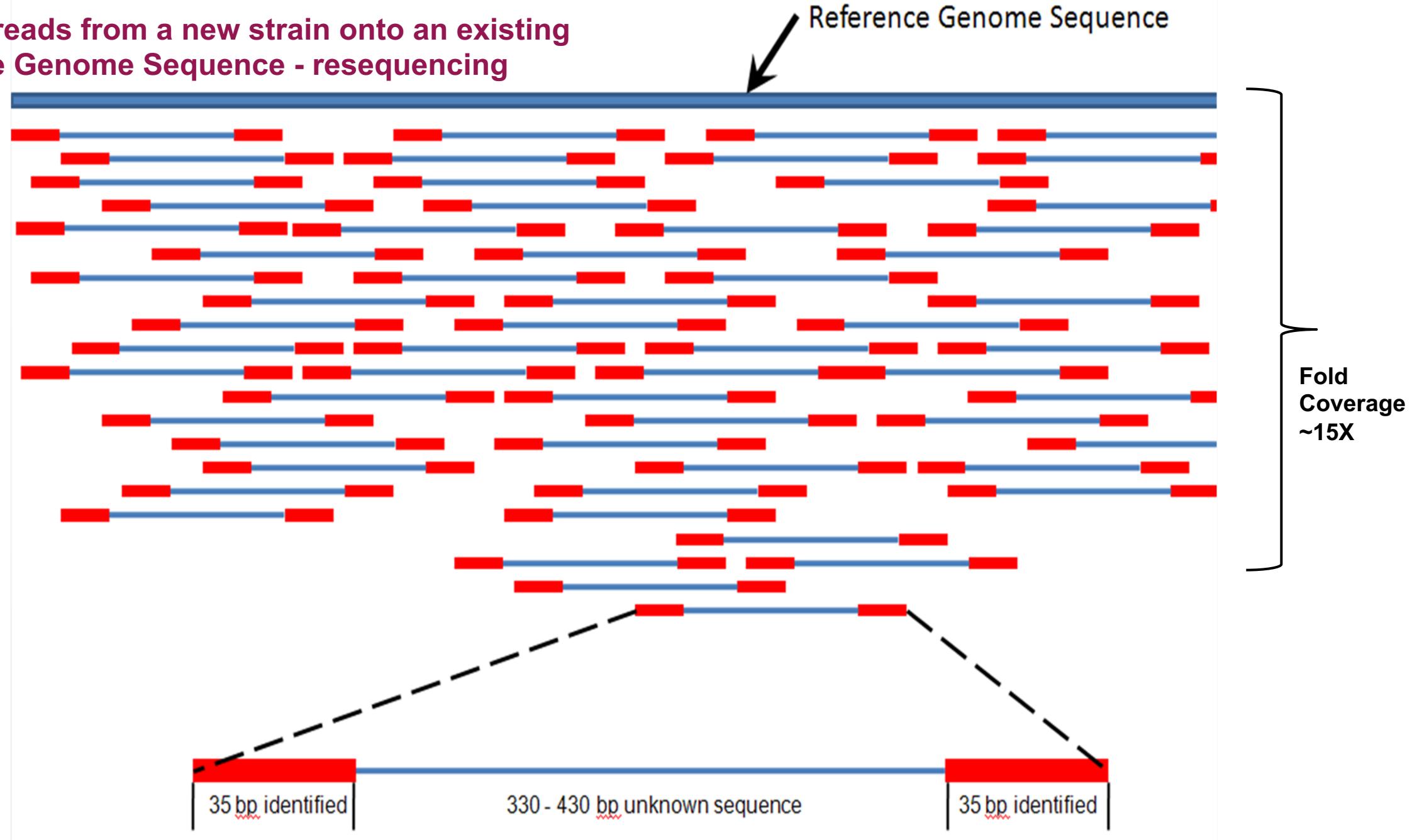
A *de novo* Short-Read Paired-End Genome Assembly



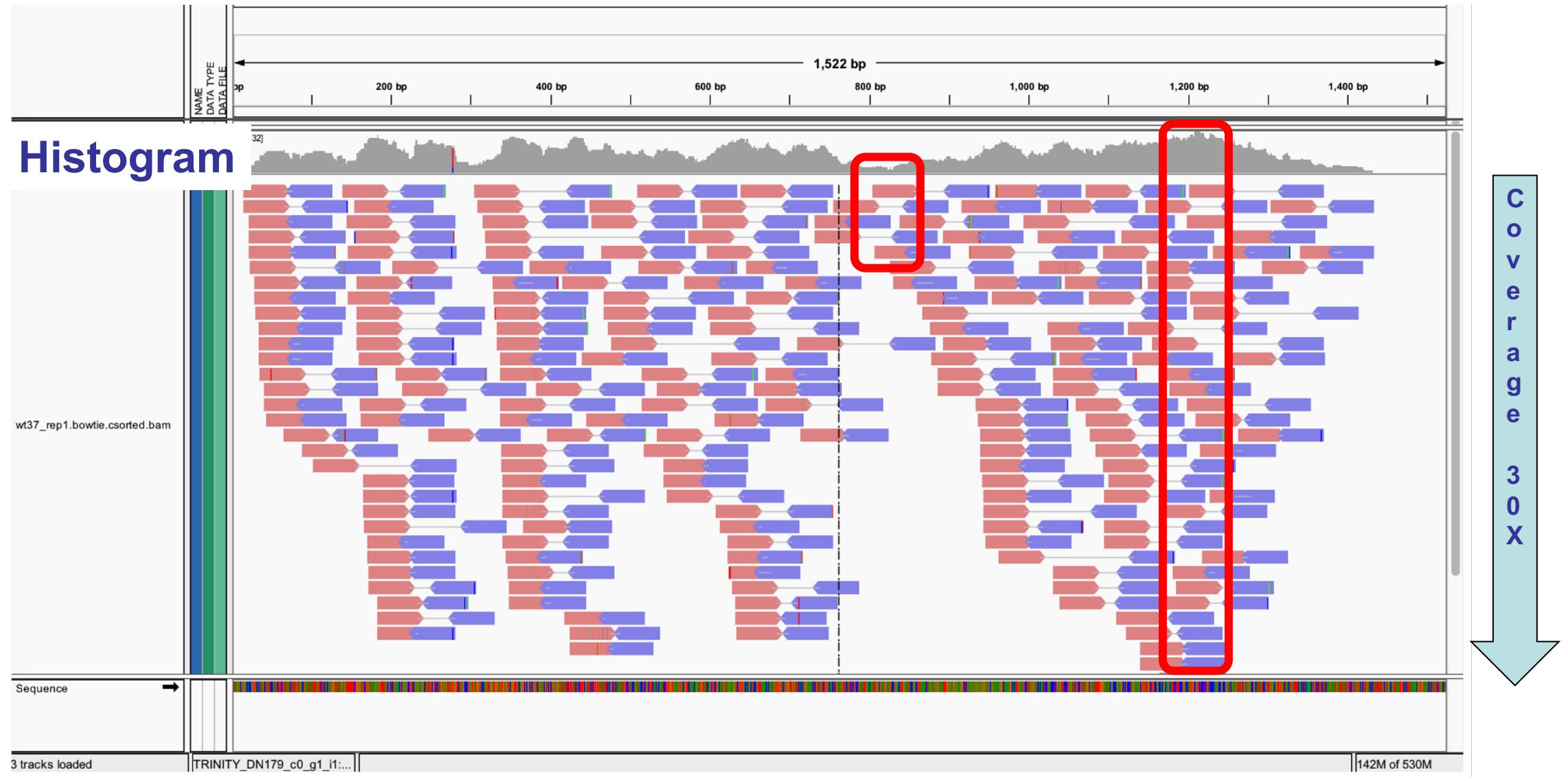
We use genome sequences in 2 ways



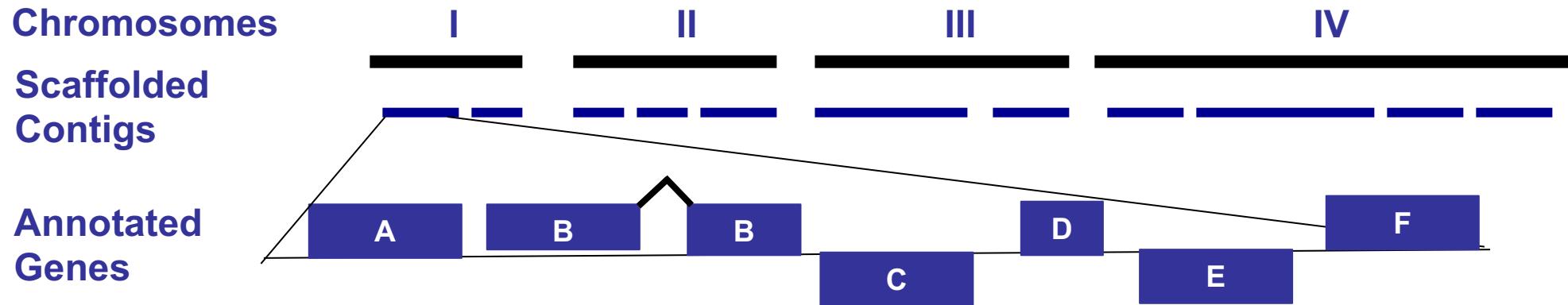
Mapping reads from a new strain onto an existing
Reference Genome Sequence - resequencing



RNA or DNA reads mapped to a reference genome sequence provide insight into “coverage”, the number of reads mapping to a specific region



30,000 ft View - Genome Annotation



The Genome Sequence

Genes can be located on either DNA strand Convention -
Gene location = non-template strand, i.e. the sequence of the
gene is the same as the mRNA (except U = T in DNA)

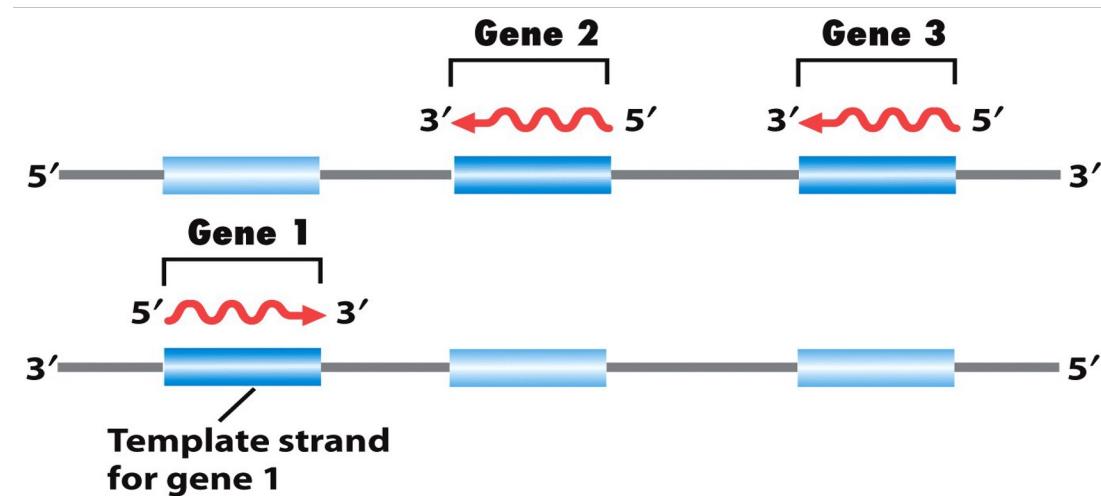


Figure 8-3
Introduction to Genetic Analysis, Ninth Edition
© 2008 W.H. Freeman and Company

Nontemplate

strand 5' — CTGCCATTGTCAGACATGTATAACCCGTACGTCTTCCCGAGCGAAAACGATCTGCGCTGC — 3' } DNA

Template

strand 3' — GACGGTAACAGTCTGTACATATGGGGCATGCAGAAGGGCTGCTTTGCTAGACGCGACG — 5' } DNA

5' — CUGCCAUUGUCAGACAUGUAUACCCGUACGUUUCCCGAGCGAAAACGAUCUGCGCUGC — 3' mRNA

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

Six Frame Translation Looking for Open Reading Frames, ORFs

1/1	31/11	61/21
M Y A L L I L Y Y 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 * T L V T S C M * G G L L Y L		
A C T L Y * F Y I I 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 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 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 – but they can be part of a gene

The “Coding Sequence” - CDS

AAGCTTCGCCAGGCTGTAAATCCCGTG AGTCGTCCTCACAAATCATCAAGCAGGTGTCCCTCAGGGAGACTGCCGACTGAGTTATGCTAATTCCTTCTACTTTGGCGTG
GTCACTGTAAACCATACTCGAACATTCATTCTCTAGCCCTACGAACAGGTAAAGAGCGTAGGGATGTCCTGGAGTAGTGCTTACTCGATAATATTCAAGTTGGACTAC
AGCGAGGCCTCGTTTGTCAACGCAATGCCGTAGACAGTTGAGAAATGTAACCGACAAACGCCGTTCATATGCTTCAAACTTAGTAGACCGTACTGTCTGA
AACTGGCGGTCAACAGGACACAGATAACGCCCTTGGCATCGGCATGTCGCTACAGAGGTCCGTATGTAGTGCCACAGCAGTGTCTTGTCTTGTCT
TTTACACGTATTAGCCCGCTGCGATTCTCGGAGCGCACCTGTTCAACACTAGAAAACGGAGTTTCTGTATCGAGAAAGCCACCCCTTCCAGAAGTTGAACGCTAGCA
TGTCAATTGATTTCACCCCCCGCTAGTTCTGTGAGACAACCTGTGTCCTCCATATGCGTACTTTCCGCAATTTTTCA
AGACTTTCAAGGAAAGACAGGCTCCGGAACGATCTCGTCCATGACTGGTAAATCCACGACACCGCAATGGCCCCCAGCACCTCTATCTCTCGTCCAGGGACTAACGTTG
TATGCGTCGCGCTTGTCTTCAAAAGAGAGCCATCCGTTCCCCCGCACATTCAACGCCGAGTGCGTTTTGTCTTTTGAGTGGTAGG
ACGCTTTTCATCGCGAACATCGTGGACATTAAAGTTCCATTCTCTTGTCACTTCAACCCGCCCGGAAGATCCGATTTGCTGCTGCTGTTGCG
CAGTCCCAGTAGCGTCTGTGCGCCGCGCCGTCTCTGTTGGTGGCGAGCCGCTACACCTGTTATCTGACTGCCGTGCGCGAAATGACGCCATTTTGGGAAAATCCGGG
AACTTCATTCATTAAAGTAGCGGAGGTTTCTCTCTGTTCTGGGTTTGATAACCGTGTTCGATGTAAGCACTTTCCGTCTCCCTCCGTC
TTTGTTGACATCGAGACCAGGTGCGAGATCTTCGTTGCGATCCGGAGACGCGTGCTGTAACCTTTTCACTTACCGACAGTGCGGAGCACTGCTGCTG
AGTCAGCAGGGACGGGTAGAAGTTTCGCTTTAGTAGTGCGTTCTGCTCTACGGGGCGTGTGCGTGGAGTCAGTGGTGTGCGATGAC
CCCCAAGAGGGGATCGGCATCAACACGGCTCCCGTGGCCCCACTTGGACCCACAGATTTCACAAACACTTTTCTCGTGTGACAAAAACGACGCCAGTGCC
TGAACGGGTGGCTTCCCAGGAAATTGCAAAGACGGCGACTCTGGACTTCCCTCTCCATCAGTGGCAAGAGATTCACGCCGTGTCATGGGACGGAAAACCTGGGAA
AGCATGCTCGAAAGTTAGACCCCTCGTGGACAGATTGAAACATCGTCGTTCTCTCTGTCAGCACACACAGTAGTCCGCCACCGCTGTTGAGACGTGTCATCT
CCAAGAGTGTGGACGCTGTTCCACGCTTCAAAATGTTTCCAAACATCCGTCGCTAGTAGACACACACAGAACGCGAACACGGCGAATCTGCTCATCGAGGAGCC
GGGGGGCACACAACATATCTCAACTCTCGAACGAACATATCCGGGGCCCGGAAGACGTTCCAGTCTCTCAATCCAAACCGGAAGCGAAACATTCTGCAAGTCACGA
TTGCCCGGTACCTCCATGTTAAGCAGTTCCATGAAACCTCCGGATATTACACACGACTGTGGATATGCAAGATGCAATAACTGAGACGCCAGTGCACACT
ATAGTTTCTCTGCCCCCTCCACATGGATATTTCAGACCTTCTCACATTGTTTGGCCCTACACCTCCGGTACCGCTTTTCTGCTGCTCTGCTGTTTATC
AGCAAAGAAGAACATTGCGCGGAGAAGCCCTAACGCTGAAGGCCAGCAGCGCGTCCAGTCTGTTGCTTCACTCCAGCAGCTCTGAGCCCTCTGGAGGAAGAGTACAA
GGATTCTGTCGACCAAGATTGTCGTTGGTATGTTGCTCTAAACTCTTGGACTCCATTCTGGTACAGAACGCTTACGAAACATGTTATACATGTTATACAGATGTTATG
GATAATATCTAGAGAAGATACAGGGAAAGACTGGCAAGGATGAAAGACATGCGACTTTAACGAGCAGAGGGCATPGCGAGAGGGACGCCGTTATGCTGTTGAGTGTG
GCTGTTGAATCTTACCTCGCCGTTGACTTGTCTGCAGCGCTTGTGCACTGAAACGTTGACTTCTTACCTTCCCAACGCCCTCTTATTCCTTCTCACTGCGAARCGC
CGCTCAGTGGGCCGTACCGAACACCCCTGGTTCTCGCTTCTGCTGCTGCTGCTGCTGCTTCTCACCTGTTG
GTGCGTCCAGACTATGTCGCTCTGCTTCTCCACCCCTCTCGGCTTGTGCTTCTGAGGAGCGGGACTGTACGAGGCAGCGCTGCTCTGGCTTCTCACCTGTTAC
ATCACGCGTGAGCCCGAGTTTCTCTGCGCTTCTCCGGAGATGACATTCTTCAAAACAAATCAACTGCTGCGCAGGCTGCAAGCTCTGCCGA
GTCCTGTGTTCTCTTCTGCTGCTGAGCTGGAAAGAGACAATGAGCGACGTTCTGCAAGACCTTCTCAGACAACGGGGTACCCCTACG
ACTTTTGTTCTCTGAGAAGAGAAGACTGACGACGCAACTGCGGAACCGGTAGATAAGAAAAACAACAAAGAGAAGGTGAAAC
ACGAAGAGAAGGAAAATGCGGAGAAACCGTGGATTACAAAGATATCAAGAGCAATGCTTGTGAGATTTTTTAATTCAGTAGAGACACCCGCCGTGCGAGGTGTG
TAGAAATAACTCGACCCCTGGAGACAGAGATGCCGGAGTACACCACTGTCGTTTTCTCTCTATGTTCTGACCGGTCTATCGTACTTAATTGGAGGAG
TCGTCCTCGAACAGCTTGGCTGCCATCGTGTTTGGCTTCTGAAAGCCAGAAGGCCGCTCCACAGTGAGGCCGATATACAGGGACGCCAACCCGGT
TTTCTGCCCTTGTGACTCTTGCAAGACAAACGCAATGAGCTCTTGACGTCACGCCAGAGACAACCTCCCGTGCACTGGGAGCTCCGCCAGGCCATTG
CCCCGGGTGGCGTGGAGGGACGAAGAGACGGGAAACTGATTCCGGCTTCCGCATGTTCACTTGTAGGCCATGAAAGAATTCCAGTAC
CTTGATCTCATGCGCACATTATTAACATGGAAAGGACAATGGTACCGAACGGGTAACGGCGACTGCGAGAAGAAAGGCCACACCGTTCTCTGCAATTCTGTCGCA
AGCCCTCTTGTGCTTCATCACCCTTGTCTATTCTCCGCCGCCCTTCTCTGCTTCAATTGCTGCTTCTGCTTCACTTCTCCCTGTTACCTCTG
TCATTCTGTTCTCTGCTCTATTAACTGTTCTACTCACAGTCTGCAATTGCGATAGACGAGCTTCCACGTCCTGCGTCTGCAATTGCTACCGC
CTCCCTCCACCGTGAATCGATTGTCGTTCCGCCGGTCTCTGCTGCTTCAACACGATGCCCTCTGCTGATGCTTCTGCTCTGCTTCTGCTTCTG
AGAGATATTCGCGCATCTTGTGCGGCCGTTCTCTGCTGCTTCAACACGATGCCCTCTGCTGATGCTTCTGCTCTGCTTCTGCTTCTGCTTCTG
CGTTGGTGTATCTCCAAATTCGGCTGCACTATGCGCTACTCGTACGACCTCTGCTGCTTCTGCTGCTTCTGCTGCTTCTGCTTCTGCTTCTG
GAATGCATATATTGACTTCAGACATTCTAATGTTGACAAACGATACAAATTGTTGTCGCTGCTGCTGACATGTCAGTATGTAAGAGTCGCTACTG
AGACTAACGACGACCAAGATTGTTATCTGCACTGCGCTGTCACCCGTTCTGAGTGCTGGAGTTCCGCAACCTTCTGAAATTCTGGGTTCTGTTATG

ATGCAGAAACCGGTGTCTGGCTGCGATGACCCCAAGGGGCATGGCATCAACAGGCCCTCCGGTGGCC
ACTTGACCAACAGATTCAAACACTTTCTCTGTCGACAAAAGACGCCGACTCTGGACTTCCCTCTCCATCAGTCGGCAAGAGATTCAAGGCCGTTGTCATG
TCCCAAGAAAATTCTCAAAGACGGGACTCTGGACTTCCCTCTCCATCAGTCGGCAAGAGATTCAAGGCCGTTGTCATG
GGACGAAAACCTGGAAAGCATGCCCTCGAAAGTTAGACCCCTCGTGGACAGATTGAACATCGTCGTTTCTCTCC
TCAAAAGAAGAACATGCCGGAGAACGCTCAAGCTGAAGGCCAGGCCGCTCCAGTCTGCTTCACTCTCC
AGCTCTAGCCCTTGAGGAAGAGTACAAGGATTCTGCGACCAGATTGGTGTGGAGGAGGCCACTGACGAG
GCCGCTGCTCTGGCGCTCCCTCTGACATCACGCCGTAGGGCGAGTTCCCTGCGACGTTTCTCC
CTGGCTTCCCGGGAGATGACATTCTCAAACAACTCAACTGCTGCGAGGCCGACTCTGGAGCTGTTCTGCTG
TCCCTTTGTCGGAGCTGGAGAAGAACATGACGGGACTATGACCCCTCTTCAATTCAAAGACCTCTCA
GACAACGGGTACCCCTACGACTTGTGTTCTGAGAAGAGAAGGAAACTGACGGGAGCCACTGCCGAGGGTCC
ACGCAATGAGCTGACGTCGACGGAGAGAACACTCCGTGACGGGAGCTGCGAGCTGCTGCTGCTG
TCCGGCTGTTGGCTGGAGTGGACAGGCAAGGAAACTGCCGTGAGGGTCCCTCTGCGAGGCC
TGGCCGGTGTGGCTGGAGTGGACAGGCAAGGAAACTGCCGTGAGGGTCCCTCTGCGAGGCC
GTTCACTTTAGAGGCCATGAAAGAATTCACTGATCTCATGGCGACATTATAACATGAAAGGACAATGGATG
ACCGAACGG

>Translation Frame 1

The Protein

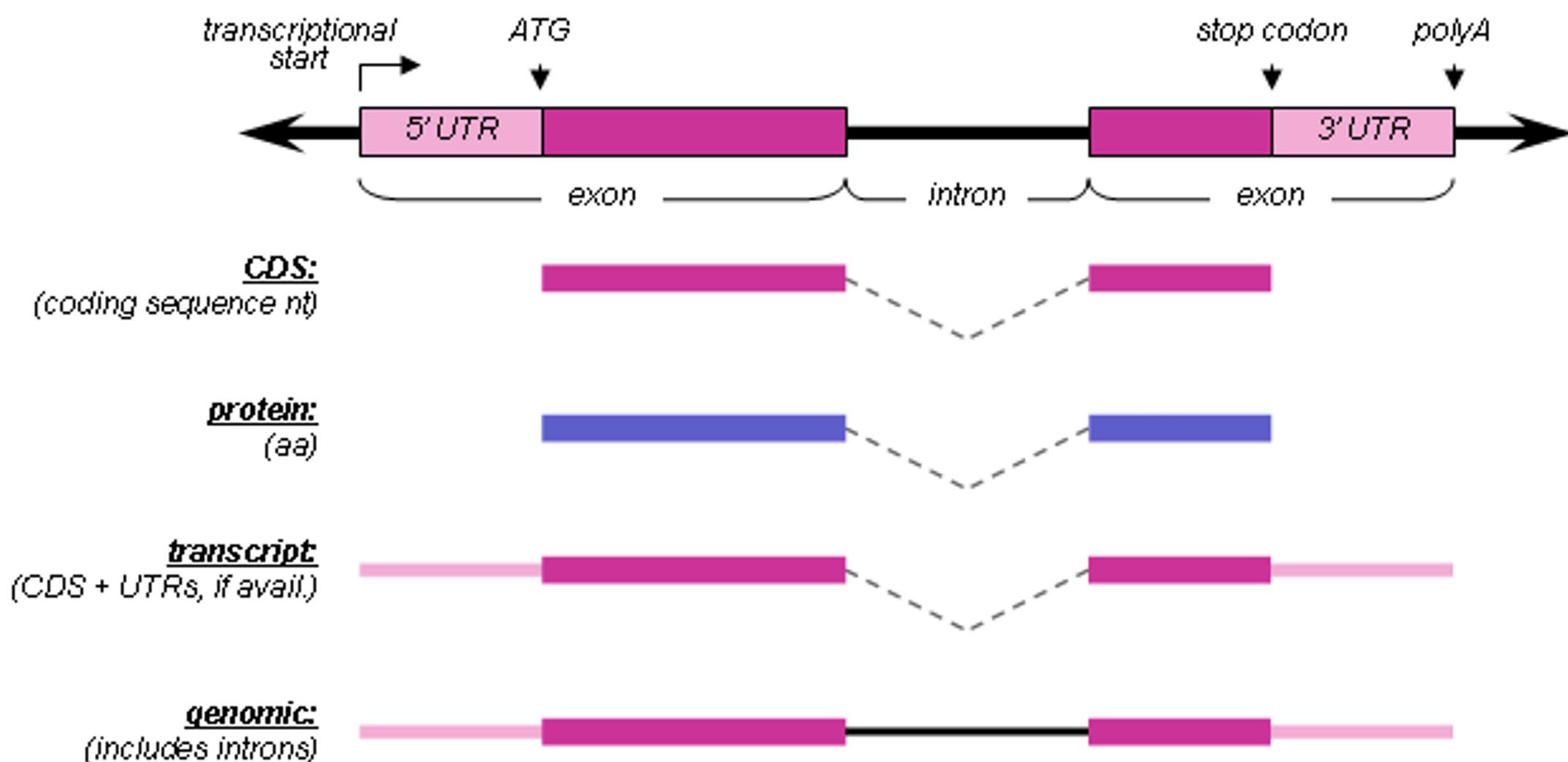
MQKPVCLVVAMTPKRGIGINNGLPWPHLTTDFKHFSRVTKTPPEASRLN
GWLPRKFAKTGDGLPSPSVGKRFNAVMGRKTWESEMPRKFRPLVDRLN
VVSSSLKEEDIAEKPQAEGQQRVRCASLPAALSLLEEEYKDSVDQIFV
VGGAGLYEAALSLGVASHLYITRVAREFPCDVFFPAFPGDDILSNKSTAA
QAAAPAESVFVPPFCPELGREKDNEATYRPIFISKTFSDNGVPYDFVLEK
RRKTDDAATAEPSNAMSLTSTRETTPVHGLQAPSSAAIAVPLAWMDEE
DRKKREQKELIRAVPHVFRGHEEFQYLDIADIINNGRTMDDRT

Green = UTRs

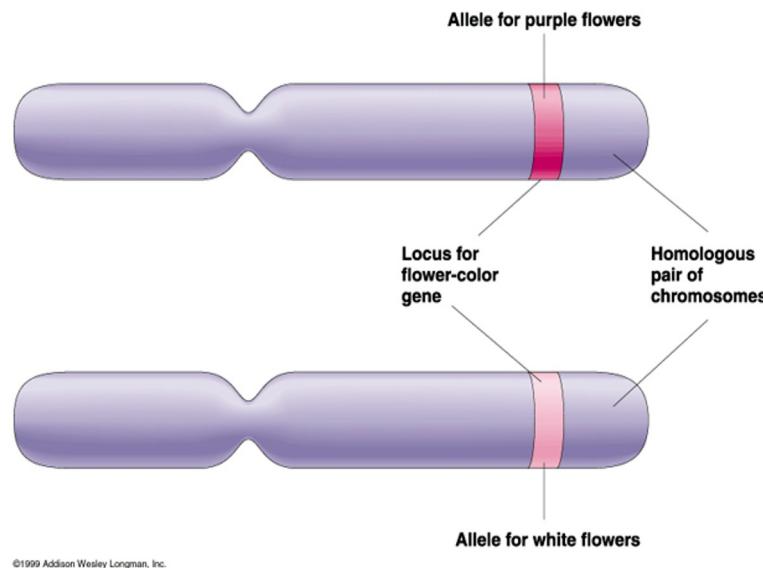
Red = CDS

Pink = Intron

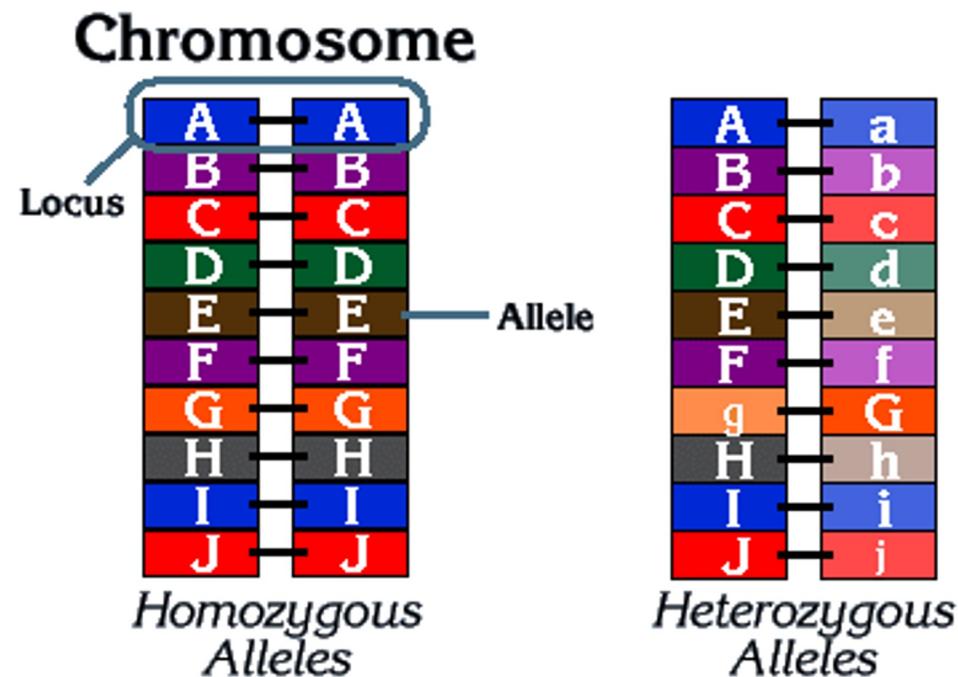
Terminology



Evolution Homologous chromosomes (in a diploid)



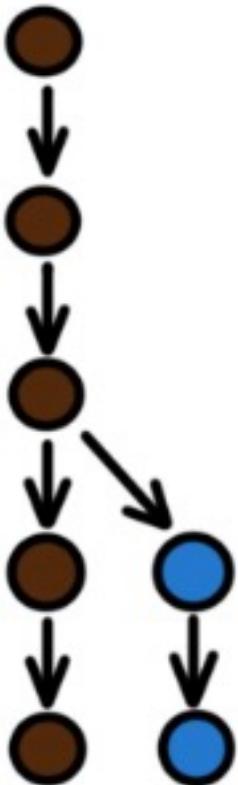
A AAGCCTCATC
a ACGCCTCATC



SNP = Single Nucleotide Polymorphism (a variant)

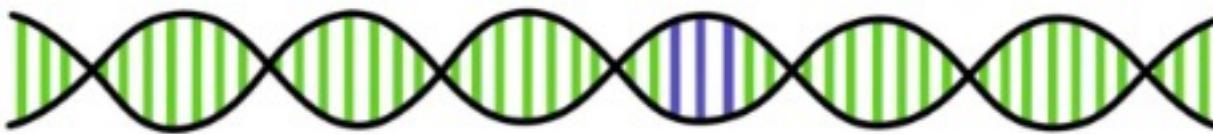
How Do Mutations Alter Allele Frequencies?

200,000 BC



At the start of the human race, brown was the only eye color people could have.

The allele frequency for the allele that caused brown eyes was 1.0 (100%).



A mutation in the main gene for eye color created a new phenotype: blue eyes!

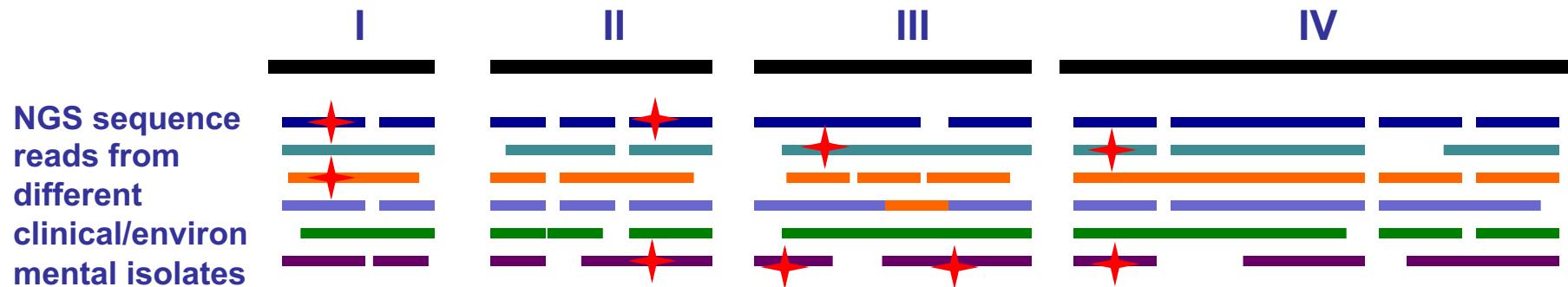
5,000 BC

Slowly, the allele frequency shifted.

Since then, many variations of eye color have emerged, though brown eyes remain the most common.

Image source: By Gabi Slizewska

30,000 ft View- NGS SNPs



★ = SNP
■ = SNV

reference	TGGTGATACT AAGCTGGGAA CTCCACTTCT	TTTTCTACTG CGGTGCTTCA
303.1	TGGTGATACT AAGCTGGGAA CTCCACTTCT	TTTTCTACTG CGGTGCTTA
309.1	TGATAATNCT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CAGTGCTTCA
RV_3600	TGGTGATACT AAACCTGGGAA CTCCACTTCT	TTTTCTACTG CGGTGCTTCA
RV_3606	TGATAATNCT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CAGTGCTTCA
RV_3610	TGATGATTCTT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CAGTGCTTCA
SenT119.09	TGGTGATACT AAACCTGGGAA CTCCACTTCT	TTTTCTACTG CGGTGCTTCA
SenT123.09	TGATRATTCTT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CAGTGCTTCA
SenT140.08	TGGTGATACT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CGGTGCTTCA
SenT142.09	TGGTGATACT AAACCTGGGAA CTCCACTTCC	TTTTCTACTG CAGTGCTTCA
SenT175.08	TGGTGATACT AAACCTGGGAA CTCCACTTCT	TTTTCTACTG CGGTGCTTA

Reference = A
6 isolate seq = A
2 isolate seq = T
2 isolate seq = N (no call)
% with base call = 80
Major allele = A
Major allele freq = 75% (6/8)

Reference = G
9 isolate seq = A
1 isolate seq = G
% with base call = 100
Major allele = A
Major allele freq = 90% (9/10)

Reference = C
8 isolate seq = C
2 isolate seq = T
% with base call = 100
Major allele = C
Major allele freq = 80% (8/10)

Alleles have frequencies in different populations

Population variation data

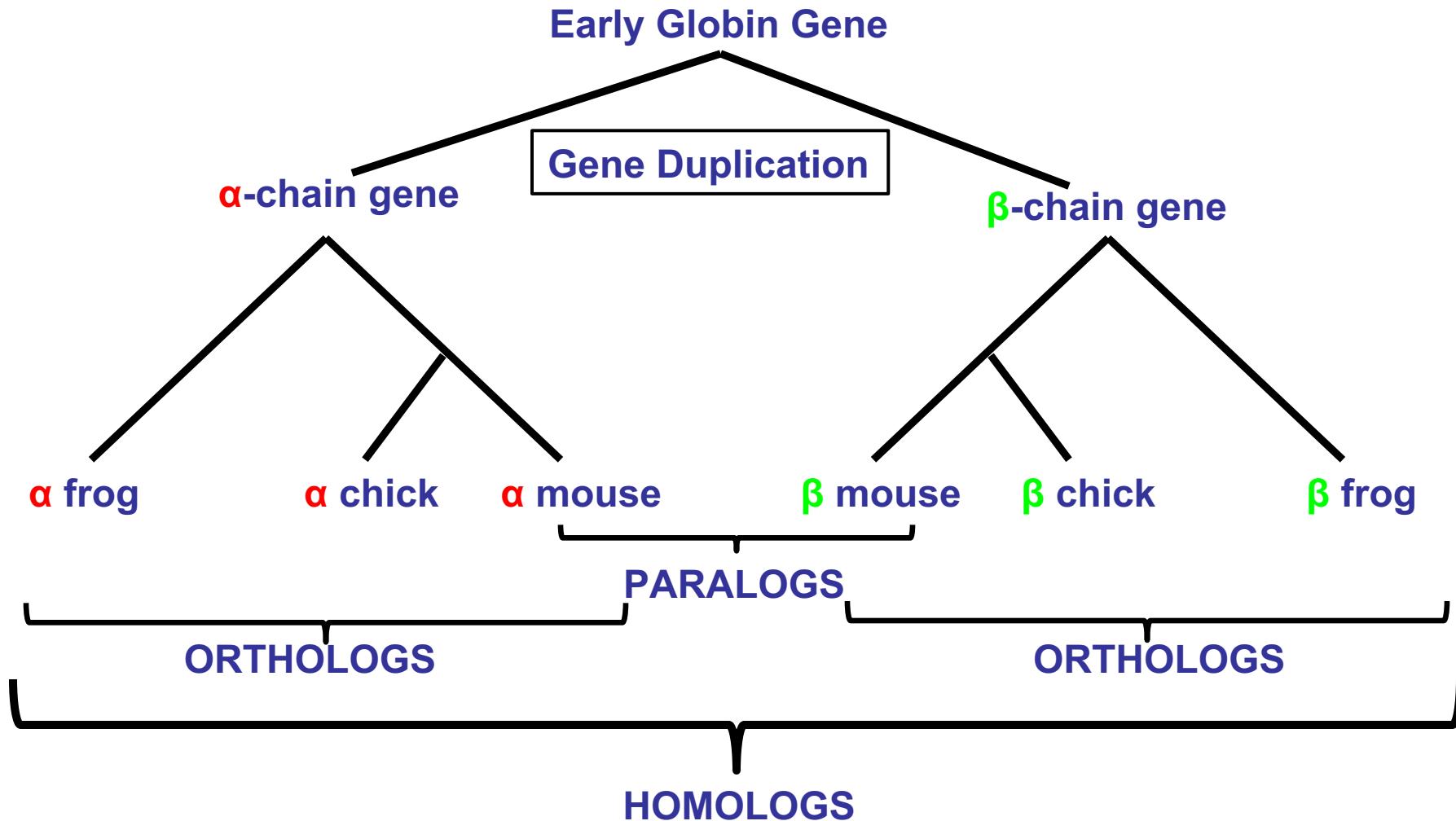
Data

- Single Nucleotide Polymorphisms, SNPs. SNVs
- Rearrangements
- Alleles
- Allele frequency
- Haplotypes (an organism's collection of variants)

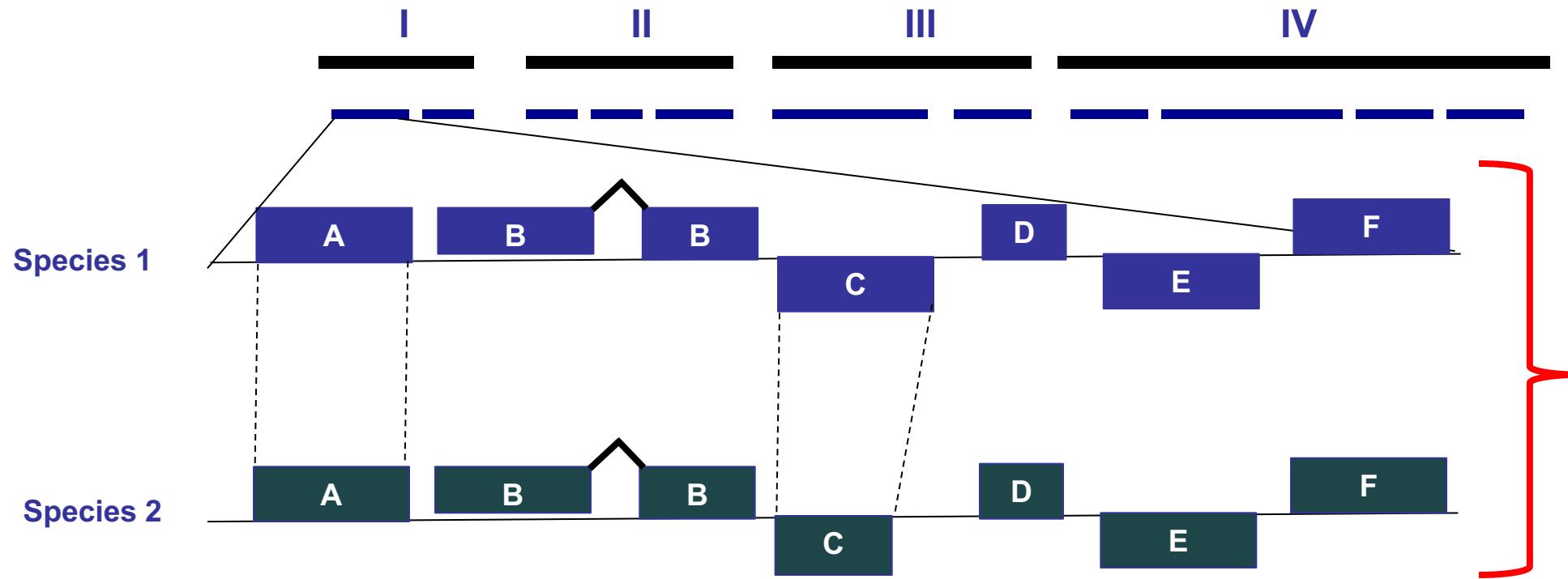
Technology

- Next Generation Sequencing, NGS
- Synteny (conserved positions on chromosomes)

Homology - a vocabulary for different types of evolutionary relationships

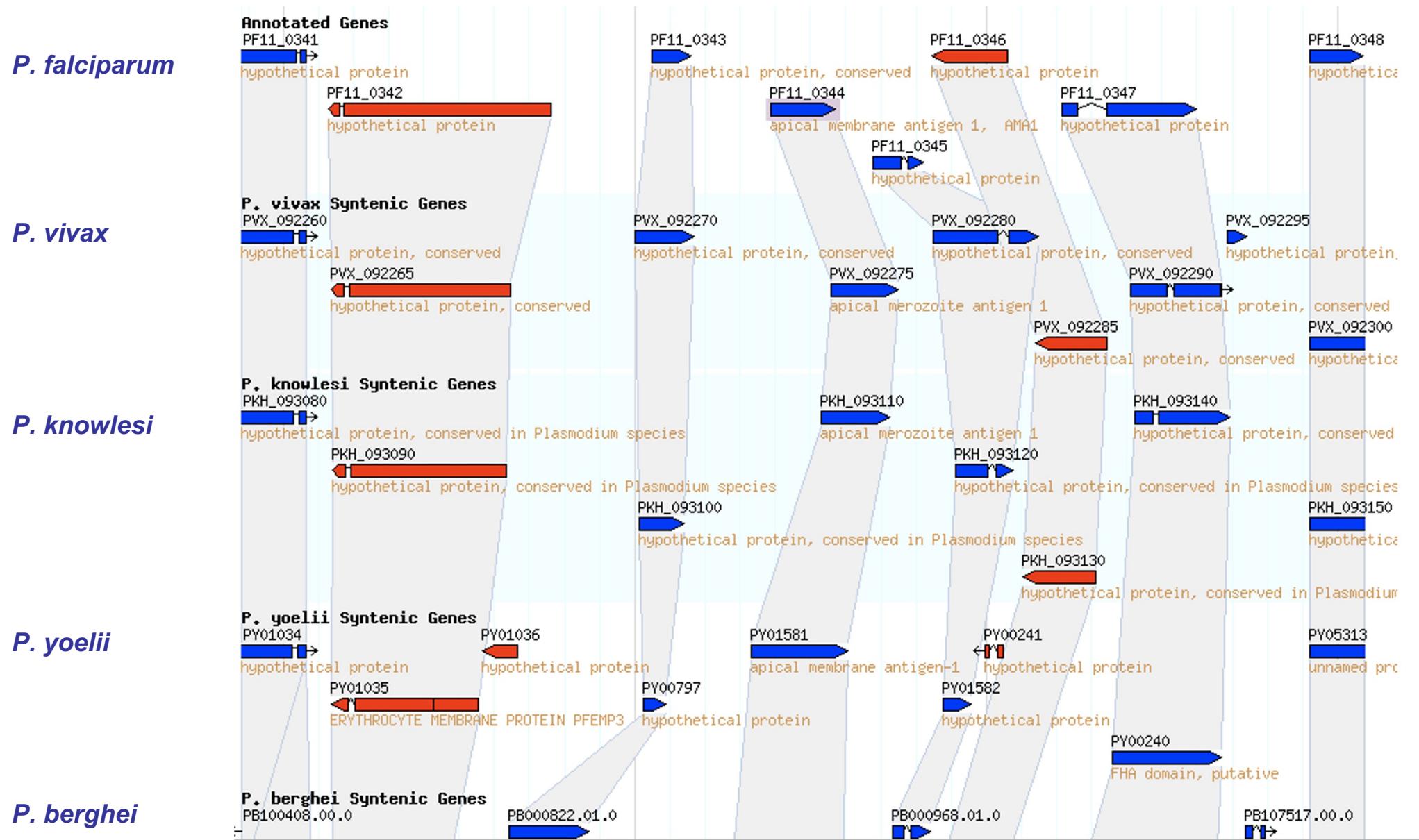


30,000 ft View - Synteny



Synteny = the majority of the same genes are present in the same order and orientation in another species. The chromosomal regions are evolutionarily related

Synteny among *Plasmodium* species



Synteny shows relationships in positioning: Ontologies show relationships in meaning

- The Gene Ontology - GO provides terms to link genes with similar functions and/or locations in the cell.
- An ontology was needed because the cultural traditions in different organisms led to different gene naming schemes that made it difficult to identify orthologous genes with the same function.

For Example:

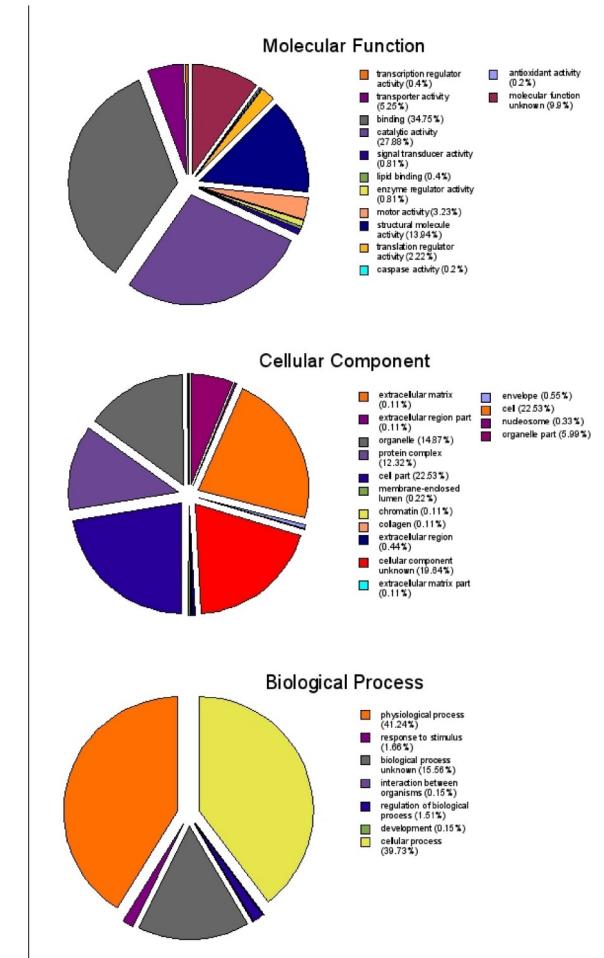
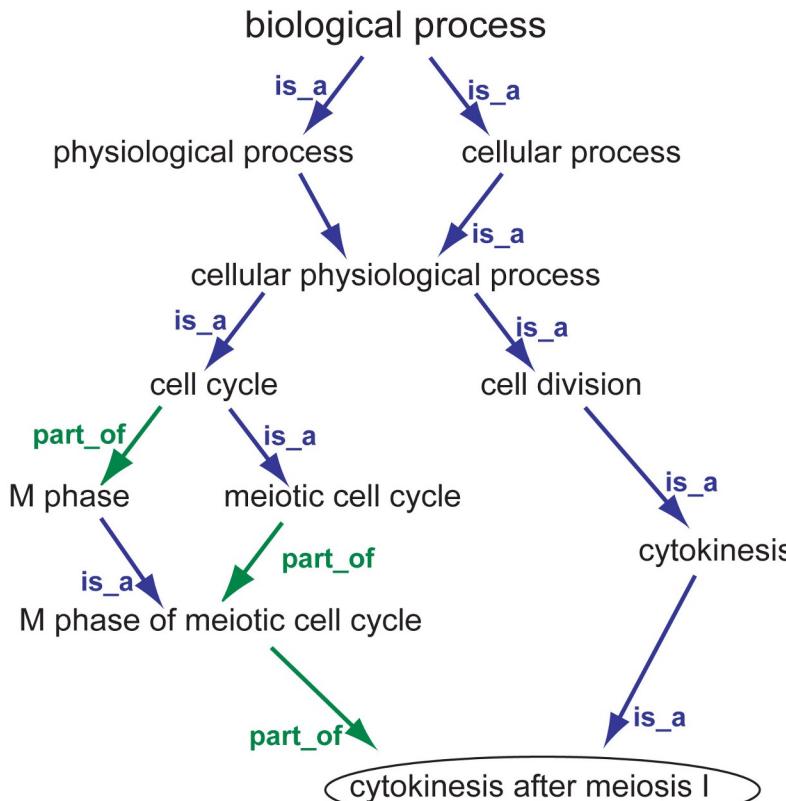
D. melanogaster gene CG3340 annotated as: "Kruppel" and *P. falciparum* gene PF3D7_1209300 annotated a "putative KROX1"

Both can be annotated with GO term:

GO:0003705 (RNA polymerase II distal enhancer sequence-specific DNA binding transcription factor activity)

Both proteins, functionally, are Zinc Fingers despite their different names

Note that the Gene Ontologies themselves contain only information about terms in the ontology and their relationships to other terms



Gene expression

Expression Profiles (RNA and Protein)

- 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, much like a photograph

RNA expression

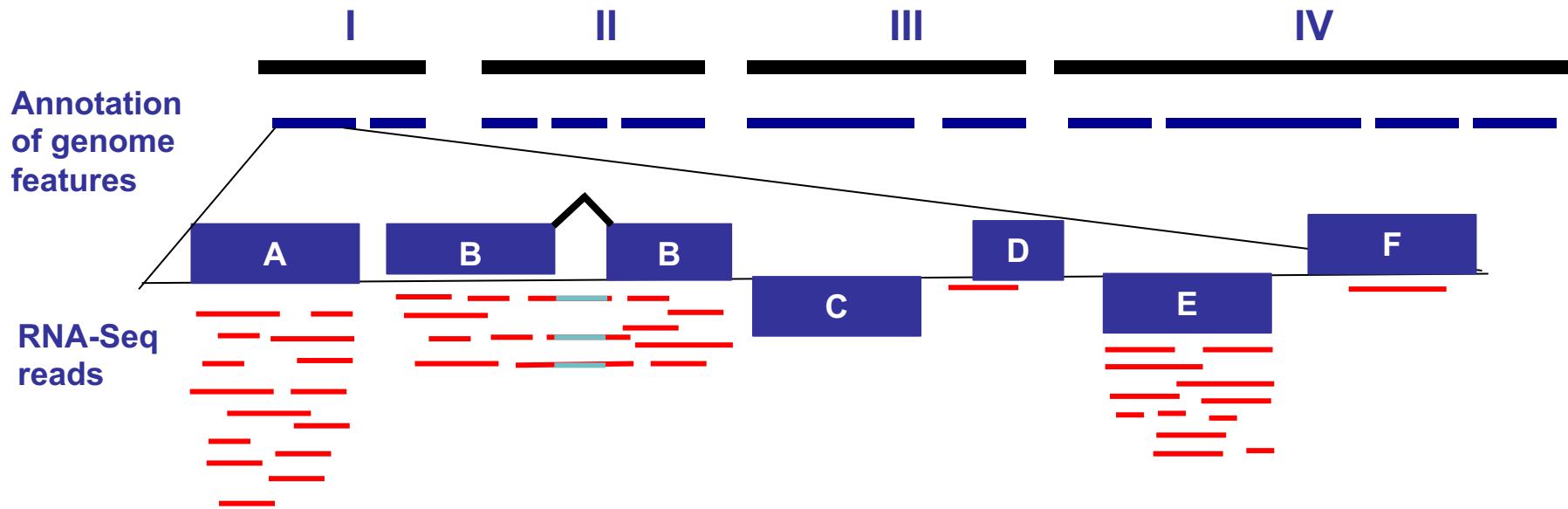
Bulk sequencing from many cells

- RNA-Seq (NGS)
 - Little sequence bias
 - Quantitative
 - Usually are strand-specific
- PacBio ISO-seq
 - Full-length transcripts from single molecules
- ONT Direct seq
 - Single-molecule, direct sequencing of RNA (or can sequence cDNA)
- All of these methods can be used to identify UTR's and exon splice junctions

Single-Cell Sequencing

- Examines the transcriptome inside each cell analyzed
- Often only detects a fraction of the transcripts within a cell
- Often analyzed with tSNE plots to categorize cells that have similar transcriptional profiles.
- Excellent for detecting cellular heterogeneity or differentiation

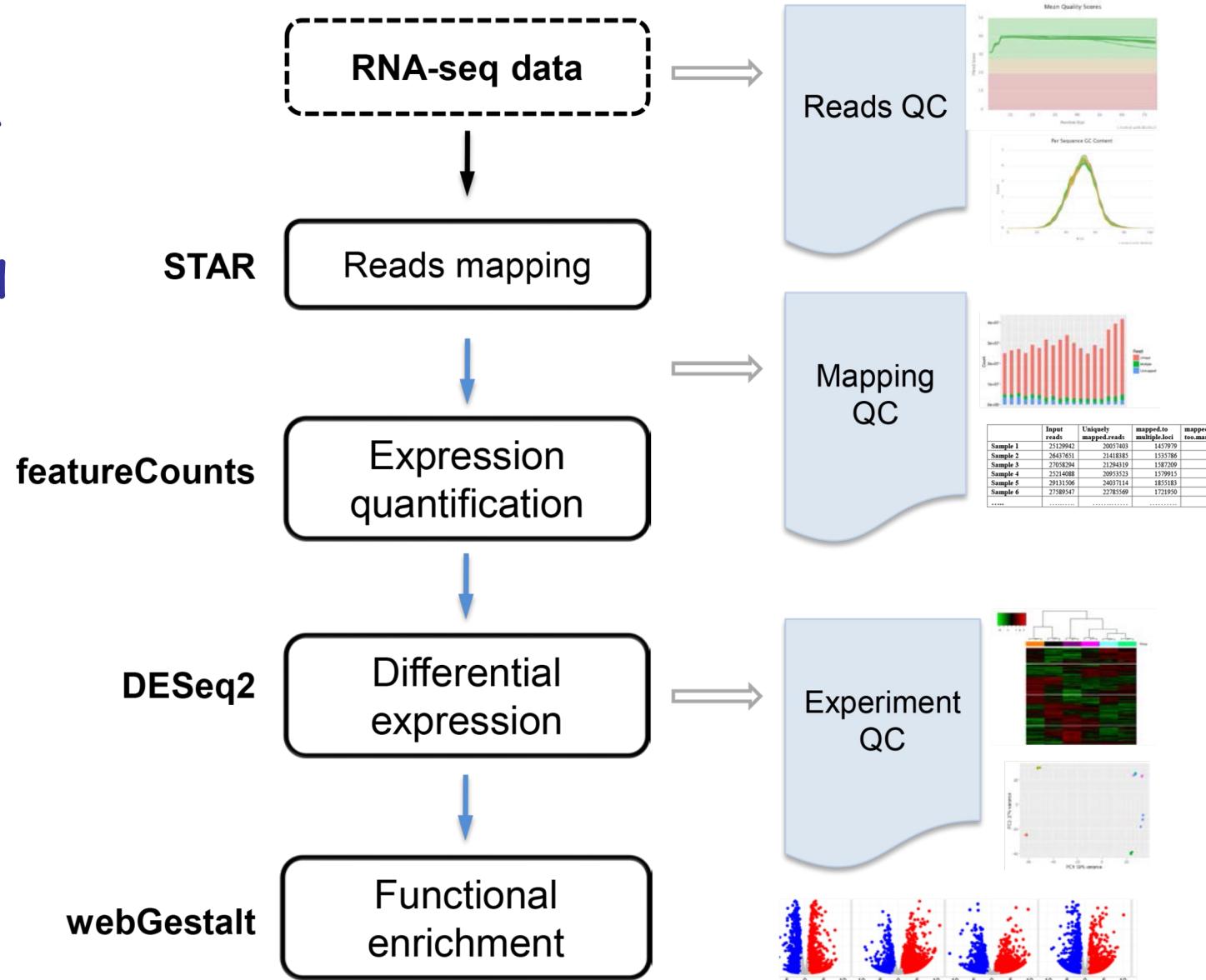
30,000 ft View - RNA-Seq



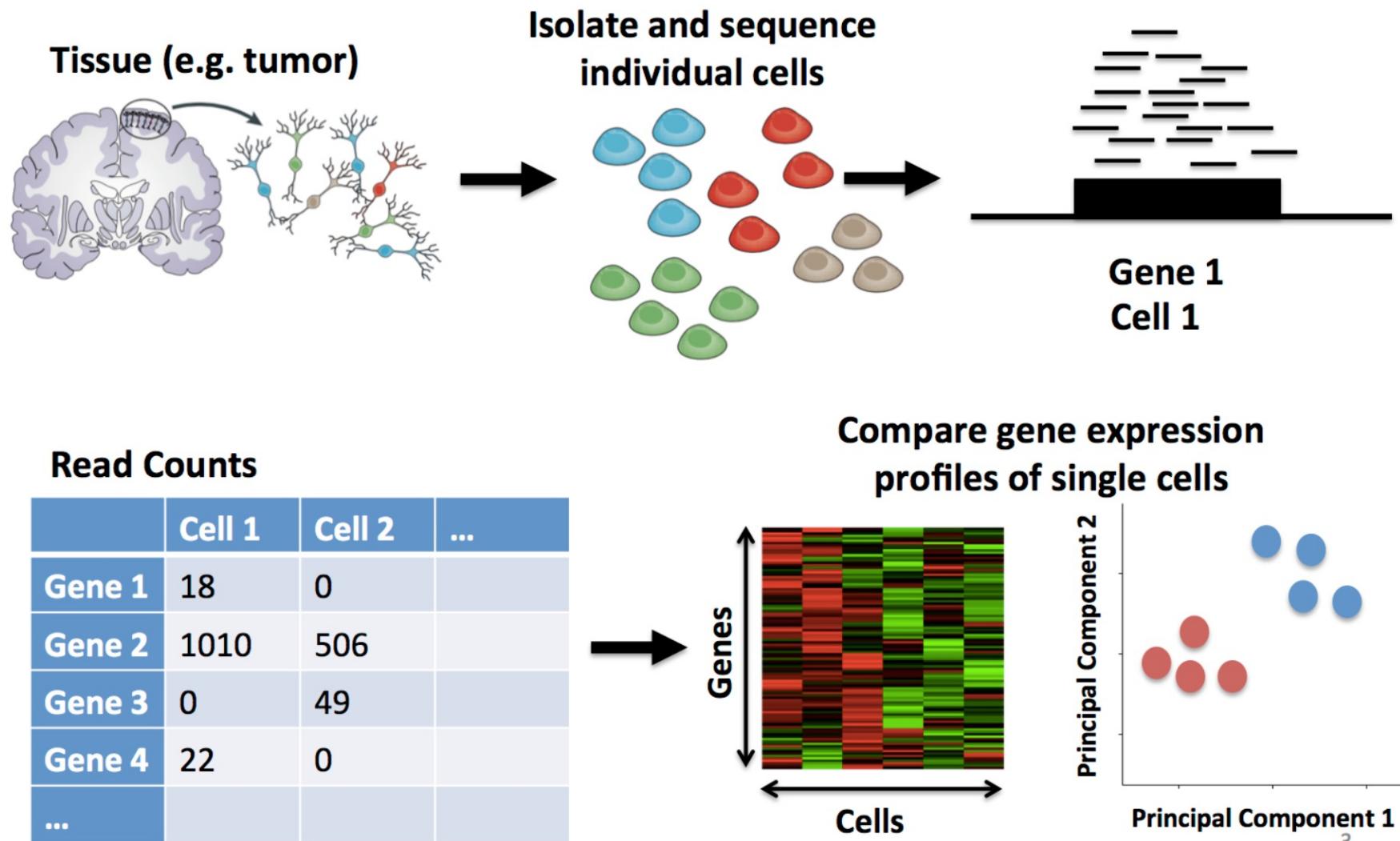
FPKM = Fragments per kilobase of exon per million fragments mapped (old calculation)

TPM = Transcripts per kilobase million (counts per length of transcript (kb) per million reads mapped)

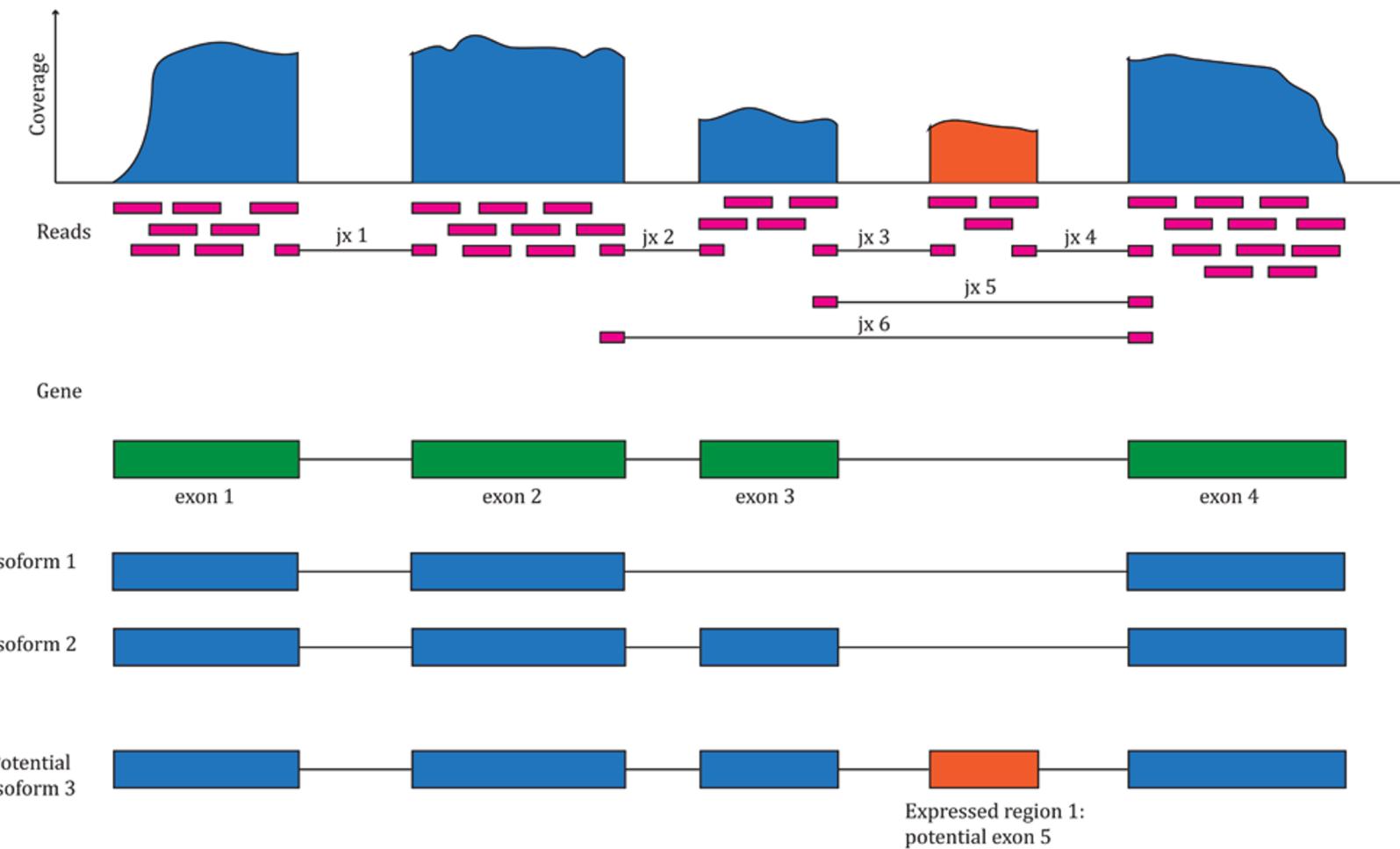
RNA-seq data
are very
powerful - you
can see how
expression
changes over
time or
different
conditions



Single-cell RNA-Seq (scRNA-Seq)



RNA-seq identifies splice junctions if present (remember context dependent)



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

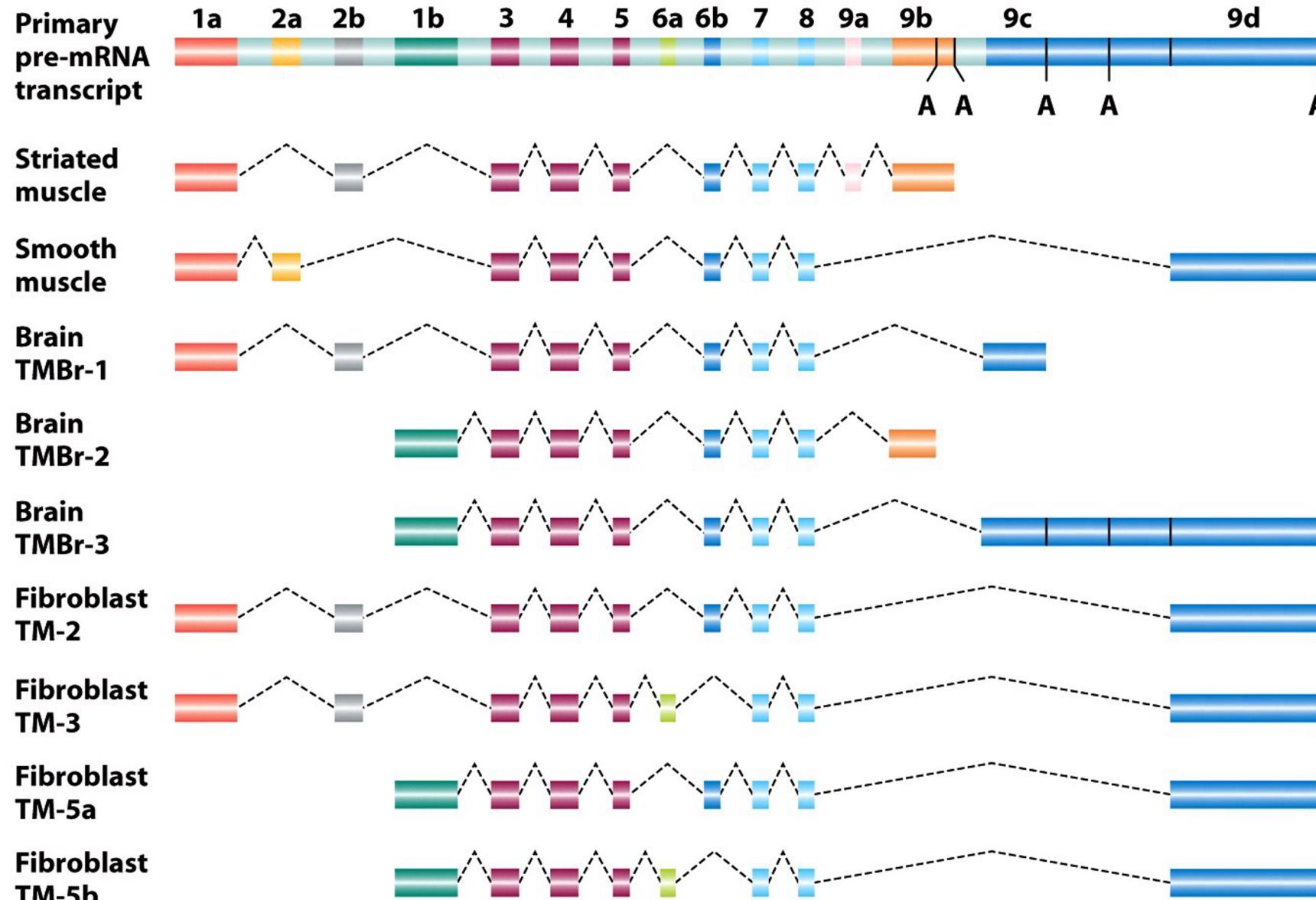
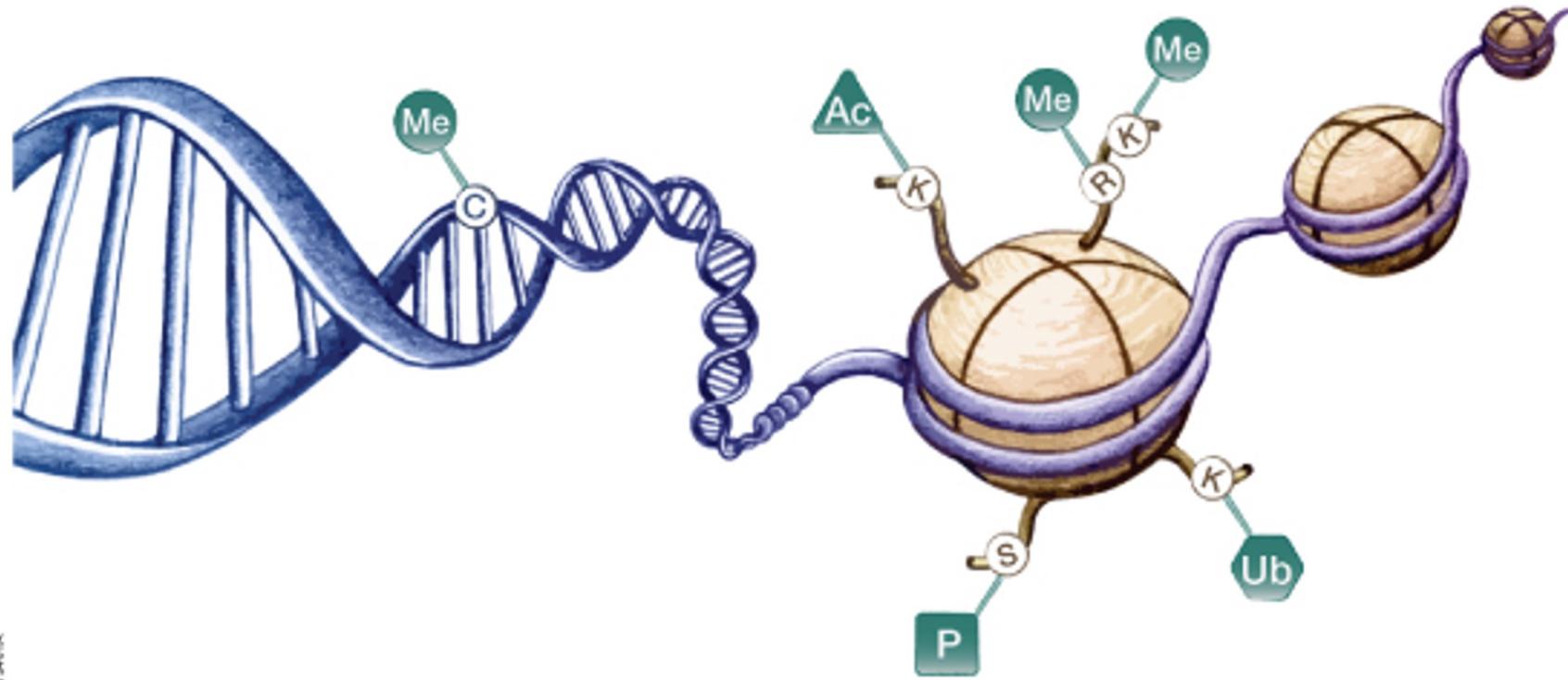


Figure 8-14

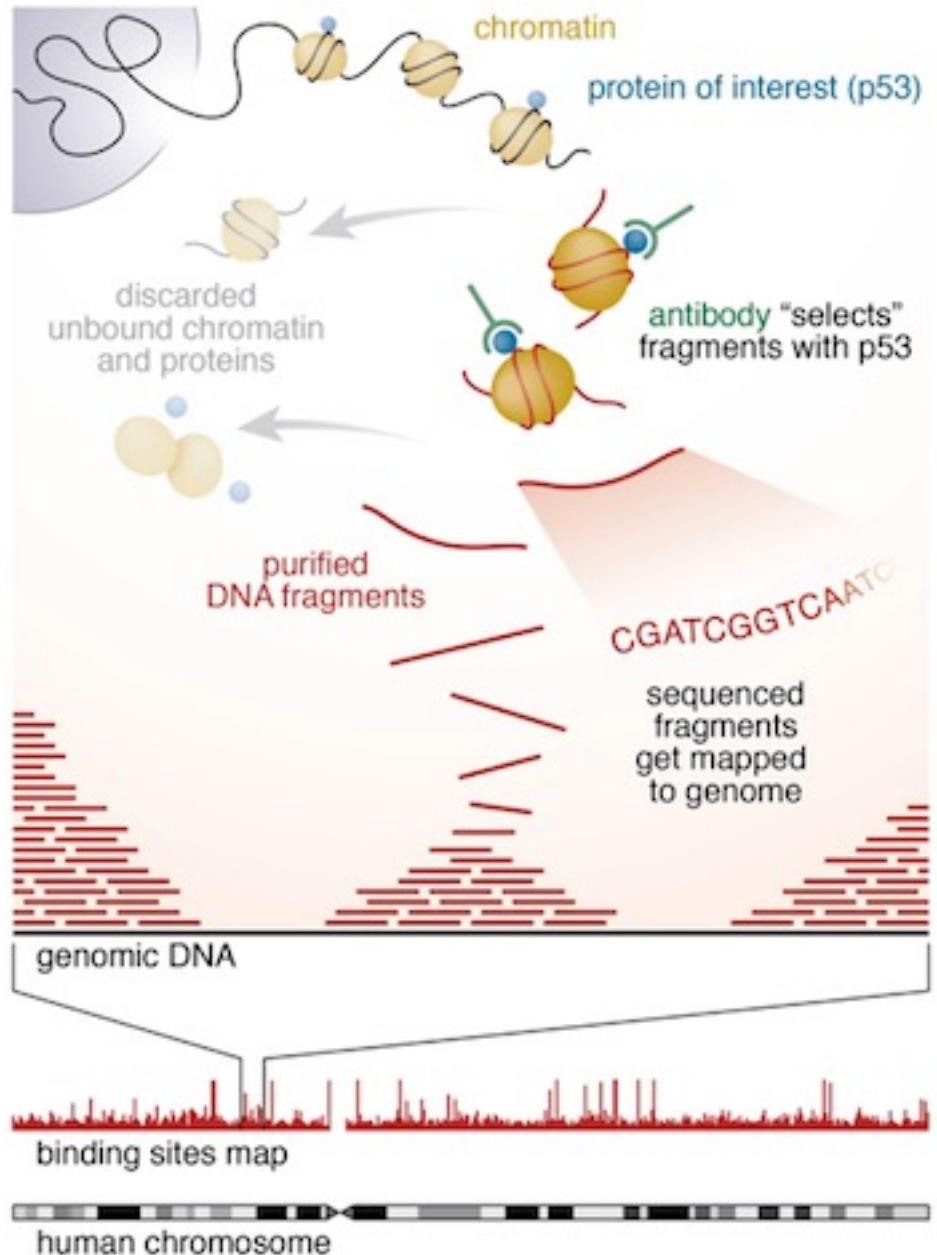
Introduction to Genetic Analysis, Ninth Edition

© 2008 W.H. Freeman and Company

Chromatin Status and Epigenetic Gene Regulation

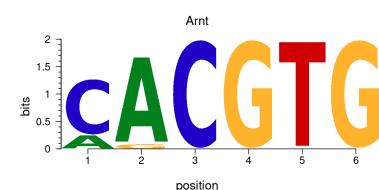


- DNA methylation at CpG islands – Bisulfite sequencing is a common assay
- H3K4me3 = transcriptionally active chromatin
- H3K27me3 = compact chromatin
- There are MANY other histone modifications
- ChIP-Seq (Chromatin ImmunoPrecipitation) is a common assay for histone markers



ChIP - Seq

- Chromatin Immuno Precipitation
- Can identify regulatory sequence motifs, e.g. transcription factor binding sites

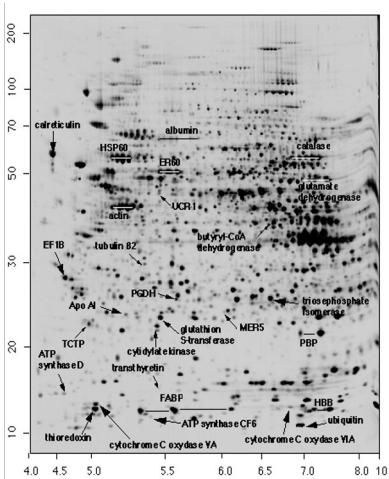


Protein Expression/Sequence

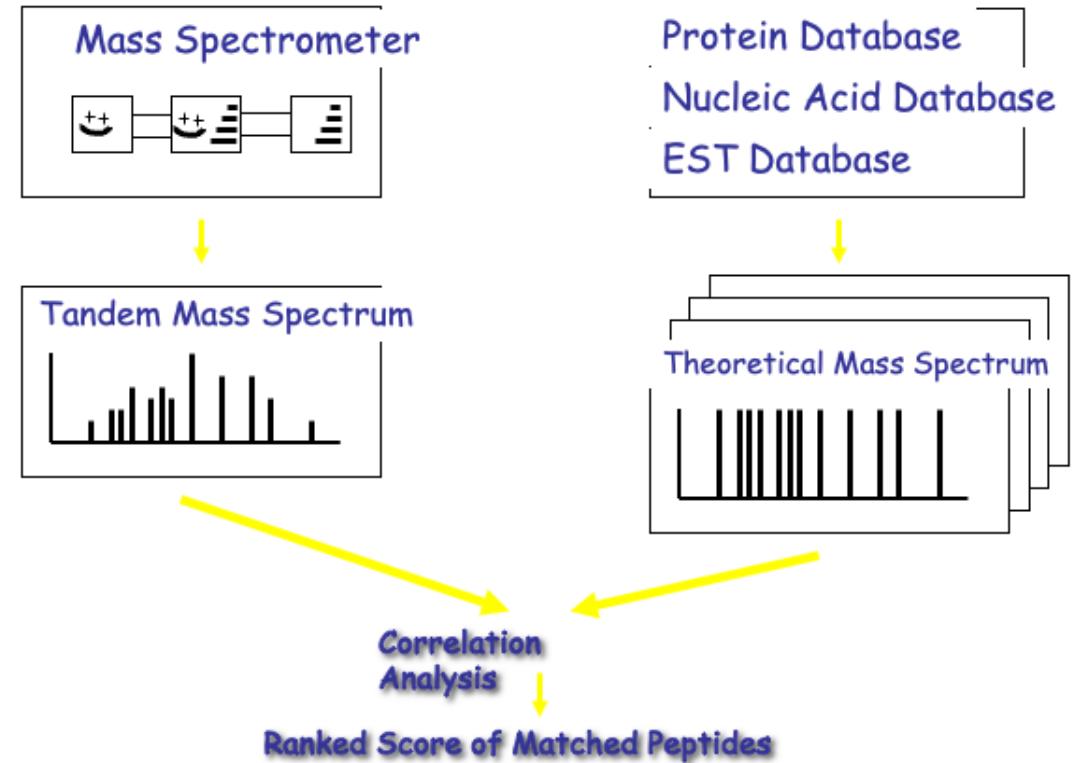
Technology

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

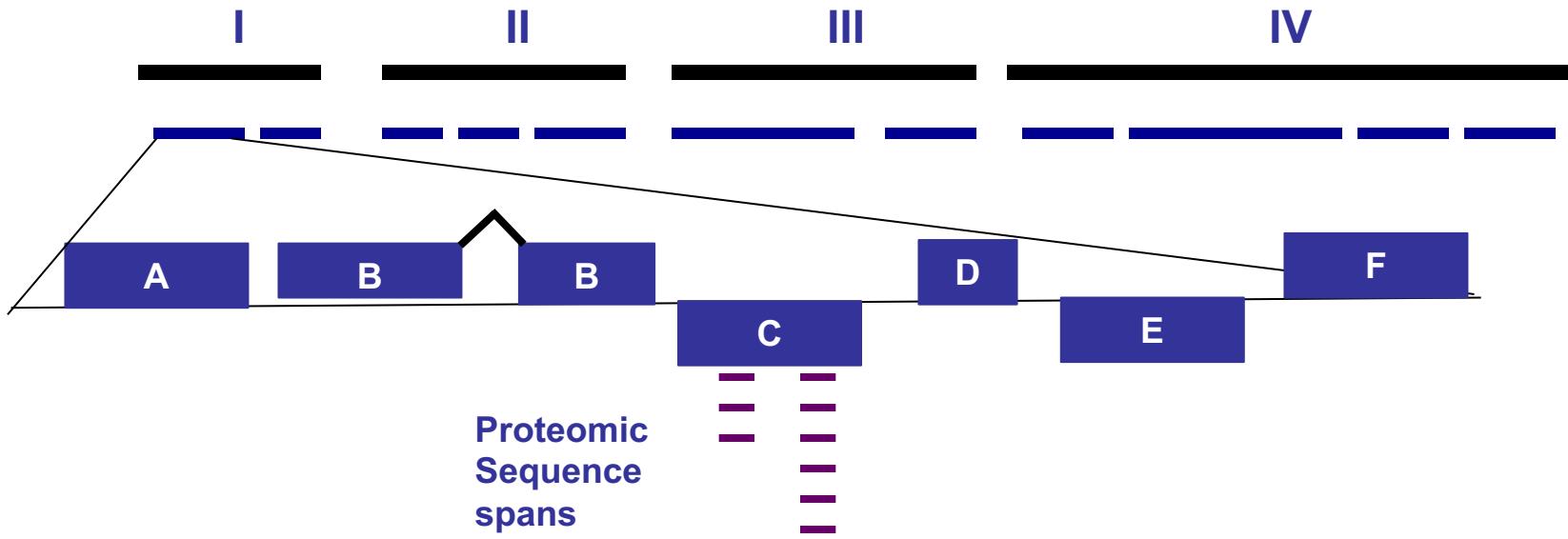
Typical 2 D gel



Tandem Mass Spectrometry analysis



30,000 ft View - Proteomics



When looking at protein mass-spec sequences it is common to only detect parts of proteins. Some regions are refractory to detection, so don't be alarmed.

Metabolomics measures all metabolites

Overview

PubChem Compound ID: [CID:93072](#)

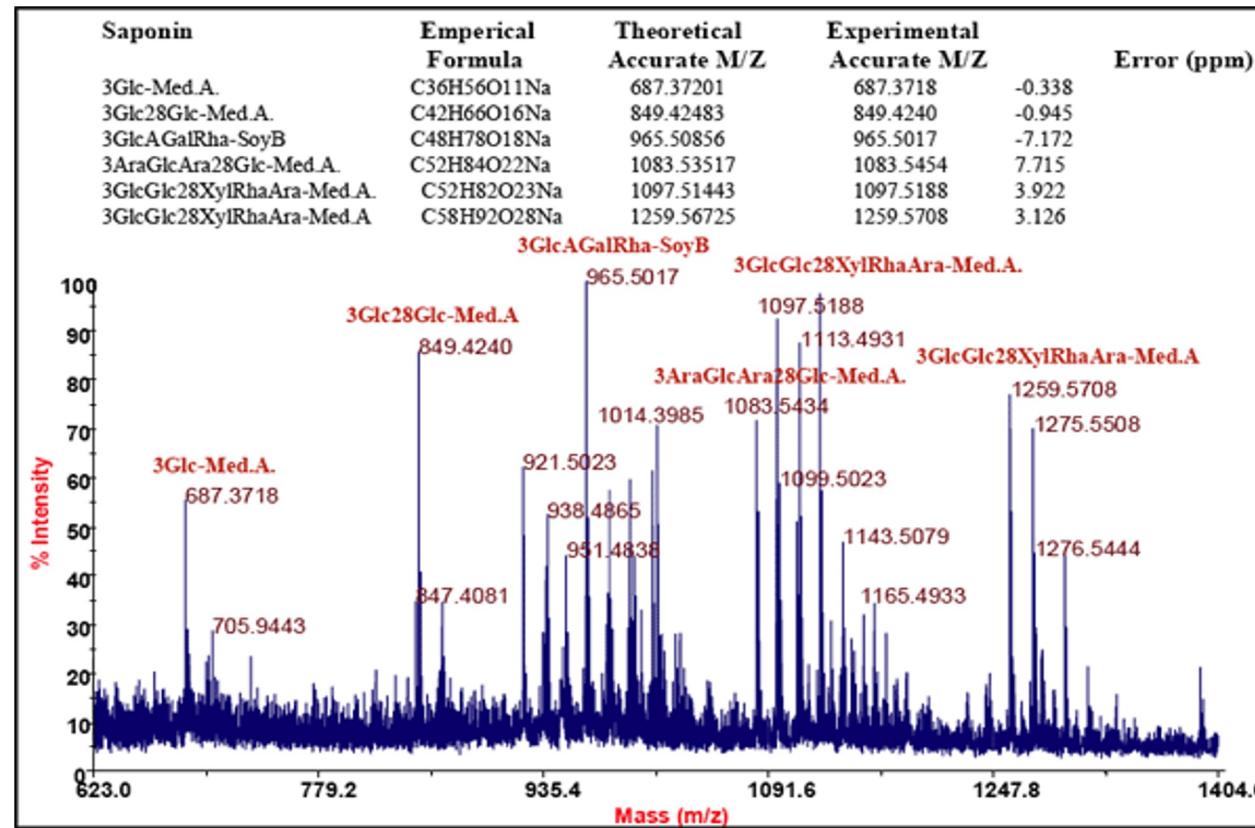
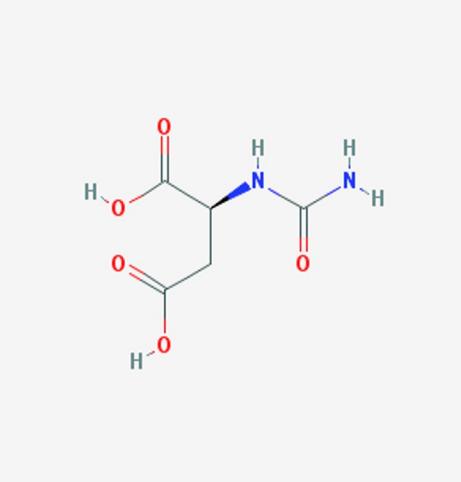
PubChem Substance ID(s): 3727

Synonyms: N-Carbamoyl-L-aspartate

Molecular Weight: 176.12742

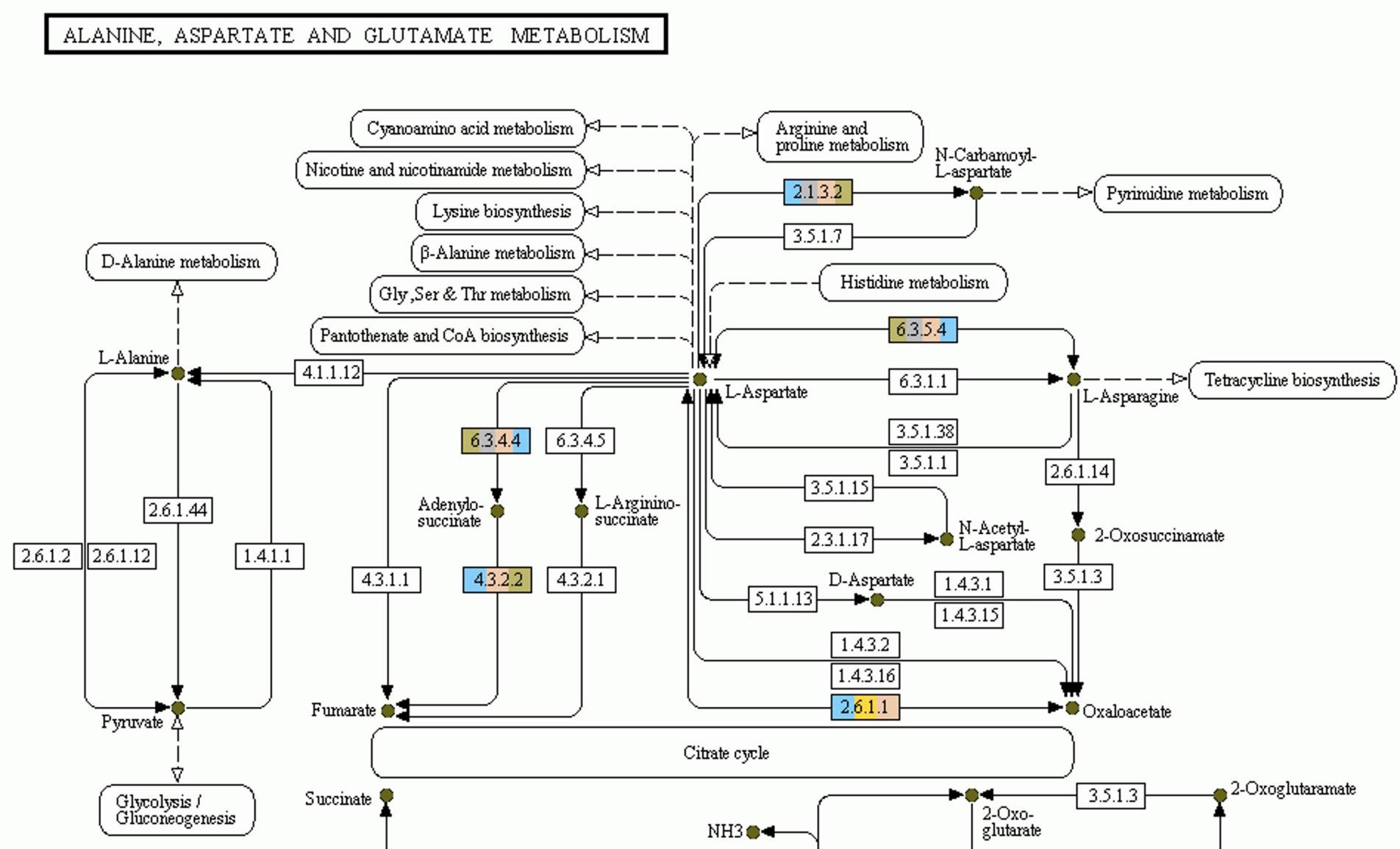
Molecular Formula: C₅H₈N₂O₅

2D Structure



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

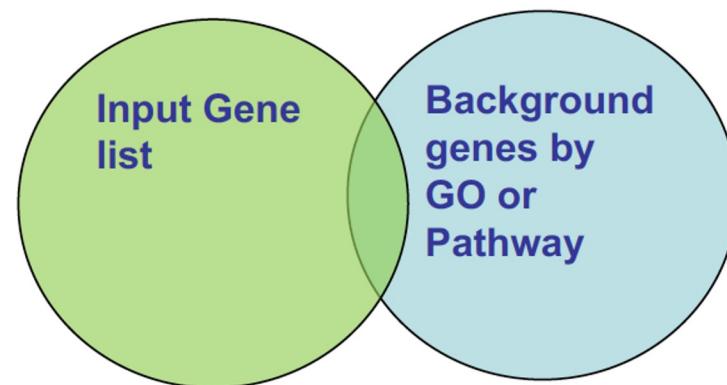
Metabolites can be linked to metabolic pathways and enzymes



Gene & Pathway Enrichment

Gene list:

Up/Down-regulated
based on some
experiment, e.g.
RNA-Seq



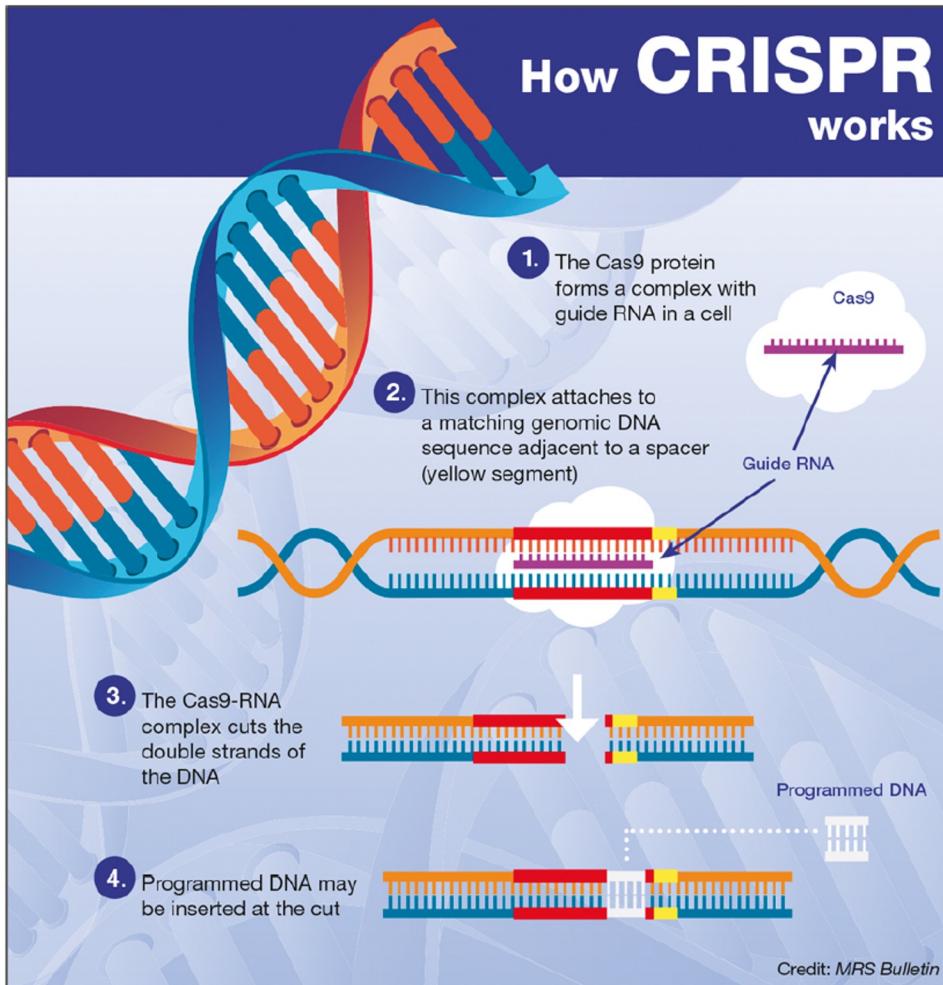
Background-Pathway information: All genes known to be involved in some process, e.g. glycolysis or cell signaling. ALL pathways are examined

Result: GO:ID or Pathway ID that is enriched

Statistics: Are more genes observed than expected (P-value)
Multiple hypothesis testing (Bonferroni, Benjamini-Hochberg)

Mutant analysis of one gene or all genes!

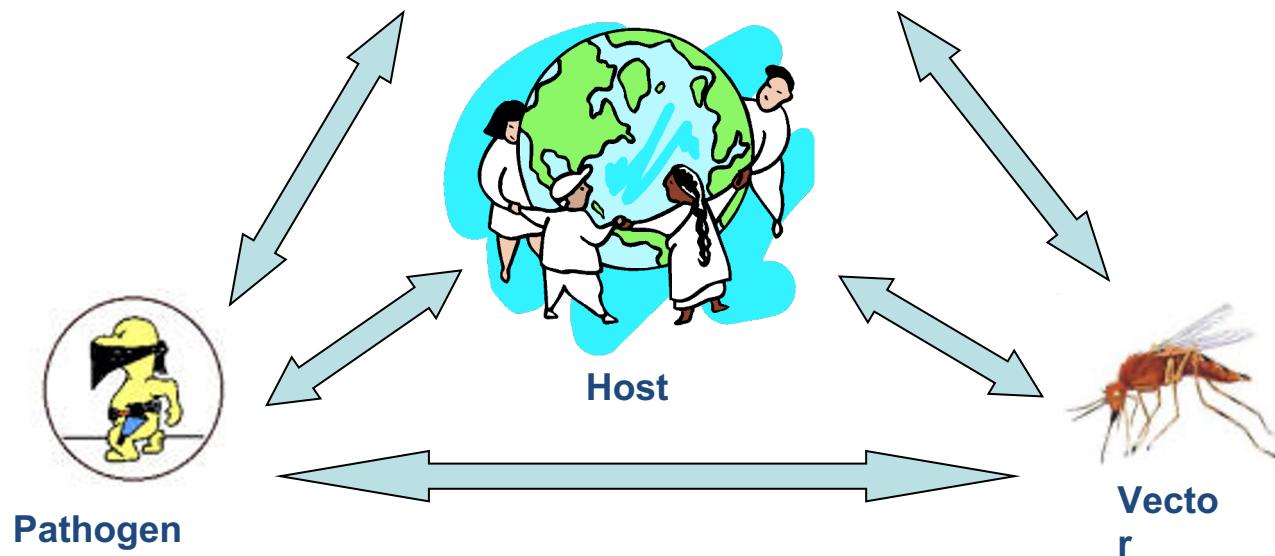
CRISPR-CAS



- Need to provide both the enzyme and the guide RNA to the cell
- Need to design the guide RNA to the gene of interest, ideally at multiple target locations per gene

Host(s)

Infectious Disease Paradigm of Host-Pathogen Interactions



Metadata - The next Frontier

- Data about the data are critical
- What makes a data set valuable? (The reason it was generated...but often this is missing)
- Introducing the "data set"
- How can you find the data set you need? Pull down Menu? A search of data set properties?
 - Person and technology that generated the data
 - Clinical outcome
 - Geographic location
 - Phenotype

Data sharing standards

OPEN ACCESS Freely available online

PLOS ONE

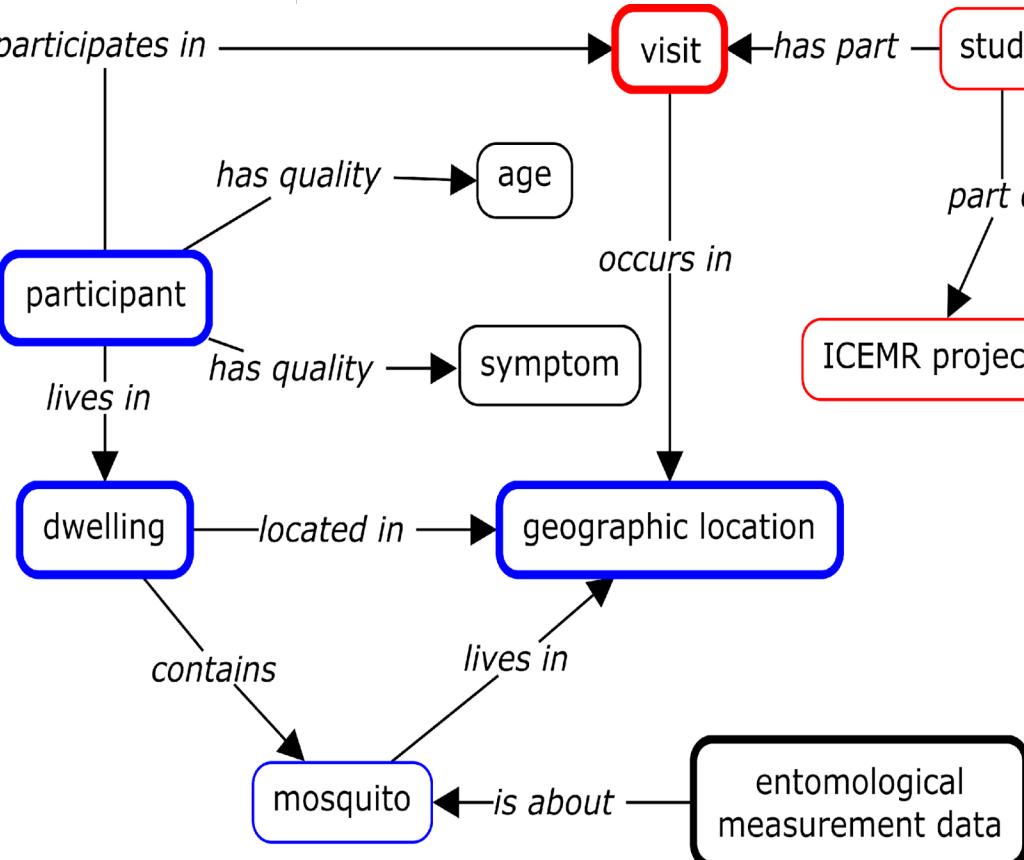
Standardized Metadata for Human Pathogen/Vector Genomic Sequences

Vivien G. Dugan^{1,2}, Scott J. Emrich³, Gloria I. Giraldo-Calderón³, Onyinye Okeke³, Brett E. Pickett¹, Lynn M. Schriml⁶, Timothy B. Stockwell¹, Christian Indresh Singh¹, Doyle V. Ward⁵, Alison Yao², Jie Zheng⁴, Tanya Barral¹, Vincent M. Bruno⁶, Elizabeth Caler^{1,12*}, Sinéad Chapman⁵, Frank H. Cross¹, Valentina Di Francesco², Scott Durkin¹, Mark Eppinger^{6,12*}, Michael J. Fidder¹, Florian Fricke⁶, Maria Giovanni², Matthew R. Henn^{5,12*}, Erin Hine⁶, Julian Mizrachi⁸, Jessica C. Kissinger⁹, Eun Mi Lee², Punam Mathur², Emma Murphy¹, Cheryl I. Murphy⁵, Garry Myers⁶, Daniel E. Neafsey⁵, Karen E. Nelson¹, David Rasko⁶, David S. Roos⁴, Lisa Sadzewicz⁶, Joana C. Silva⁶, Bruce A. Stevens¹¹, Rick L. Stevens¹¹, Luke Tallon⁶, Hervé Tettelin⁶, David Wentworth¹, Jennifer Wortman⁵, Yun Zhang¹, Richard H. Scheuermann^{1,12*}

1 J. Craig Venter Institute, Rockville, Maryland, and La Jolla, California, United States of America, 2 National Institutes of Health, Bethesda, Maryland, United States of America, 3 University of Notre Dame, Notre Dame, Indiana, United States of America, 4 University of Michigan, Ann Arbor, Michigan, United States of America, 5 Broad Institute, Cambridge, Massachusetts, United States of America, 6 Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, Maryland, United States of America, 7 Cyberinfrastructure Division, Virginia Bioinformatics Institute, Center for Biotechnology Information, National Library of Medicine, Bethesda, Maryland, United States of America, 8 University of Texas at Austin, Austin, Texas, United States of America, 9 University of Florida, Gainesville, Florida, United States of America, 10 Kelly Government Solutions, Rockville, Maryland, United States of America, 11 Argonne National Laboratory, Lemont, Illinois, United States of America, 12 Department of Pathology, University of California San Diego, San Diego, California, United States of America

Abstract

High throughput sequencing has accelerated the determination of genome sequences for many disease pathogens and dozens of their vectors. The scale and scope of these association studies to identify genetic determinants of pathogen virulence and transmission have led to the need for a standardized metadata schema.

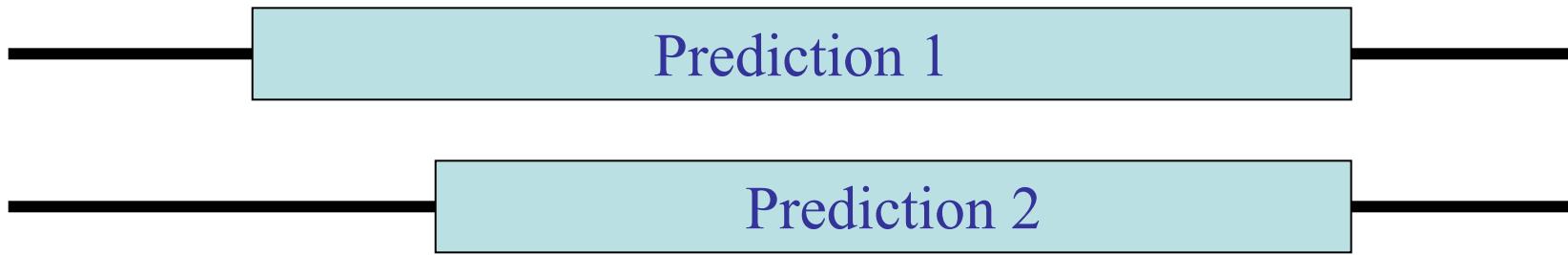


	Core Project	Core Sample	Project Specific	Pathogen Specific	Sequencing Assay
Investigation	■				
lost Characterization		■■■	■■■	■■■	■■■
Specimen Isolation			■■■	■■■	■■■
Pathogen Characterization			■■■	■■■	■■■
Specimen Processing			■■■	■■■	■■■
Pathogen Detection		■■■			
Pathogen Isolation			■■■	■■■	■■■
Sample Management			■■■	■■■	■■■
Data Transformation				■■■	■■■
Sample Shipment					
Sequencing Sample Preparation					
Sequencing Assay					

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).

Different algorithms often generate different results



We provide lots evidence so that you can decide or design an experiment to confirm!

Garbage in Garbage out!

- The algorithms will almost always return a result, it is up to you, the scientist to evaluate if it has made a mistake. Much of the data in the archival databases have errors. Not intentional errors but errors
- If you can't find the gene or answer it does NOT mean that it does not exist. It may be in a gap, or never have been annotated, or discovered after the annotation e.g. lncRNAs. Interpret carefully

Bioinformatics Resource Center Community Evolution

Browsing → Mining → Integrating → Facilitating



The End

- If you have questions, I and the other instructors will be around and we are happy to talk to you.
- These slides are available to you as a PDF on the workshop web site.



VEuPathDB

Eukaryotic Pathogen, Vector & Host Informatics Resources



Project Leadership:

David Roos – UPENN (joint-PI)

Mary Ann McDowell – Notre Dame (Joint-PI)

Andrew Jones – Liverpool

Jessie Kissinger – UGA

Sarah Dyer – EBI

Kathryn Crouch – Glasgow

George Christophides - Imperial



Our goal: enabling end users in the lab, field & clinic to make effective and appropriate use of large-scale datasets, expediting discovery research and translational application by making data FAIR

