



VEuPathDB

Eukaryotic Pathogen, Vector & Host Informatics Resources

Crash Course in Omics Terminology, Concepts & Data Types

Jessie C Kissinger

2023



@veupathdb
@jcklab
jkissing@uga.edu



UNIVERSITY OF
GEORGIA
Center for Tropical &
Emerging Global Diseases



Institute of Bioinformatics
UNIVERSITY OF GEORGIA

KAYAK

Hotels

Flights

Cars

Packages

Rentals

Cruises

More ▾



My Account

Atlanta (ATL)

↔ Seattle (SEA)

Tue 12/26 – Sun 12/31

4 travelers, Economy

 Our Advice: **Buy now**Prices predicted to rise
\$22 by **Nov 21** ⓘ**Track Prices** Receive emails with price changes and travel tips for this trip. OFF

Depart 1 day earlier

Depart 1 day later

Original Dates

Return 1 day earlier

Return 1 day later

\$447

\$631

\$561

\$648

\$717

391 out of 413

RESET

TOP FILTERS

MORE

Stops nonstop

\$631

 1 stop

\$628

 2+ stops

\$561

Times

Take-off Atlanta (ATL)

Tue 5:30a 10:30p



Take-off Seattle (SEA)

Sun 12:00a Mon 12:00a

PRICE

BEST + PRICE

▼

DURATION

▼

**Feel at home in Seattle**

Airbnb has over 3M homes. Find the one that's right for you.

View Deal

airbnb.com | Sponsored

BEST FLIGHTS ⓘspirit
Spirit Airlines

6:40 am

—□—□—

2:18 pm

10h 38m

**\$561**
KAYAKspirit
Spirit Airlines

6:00 am

—□—□—

8:07 pm

11h 07m

**View Deal** ▼

Share Watch

Alaska Airlines

6:17 am

nonstop

9:02 am

5h 45m

\$631
Alaska Airlines

The Travel Site has Very Useful Data Filters!

OUR ADVICE
ヽ(ಠ‿ಠ)
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

Alliance

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

Layover Airports

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

Flight Quality

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

Aircraft

- Narrow-Body Jet
- Wide-body jet

Booking Sites

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

Show 6 more sites

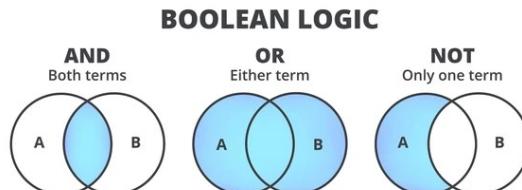
Filters vs Boolean operators

Filters - are very useful, but...

- 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



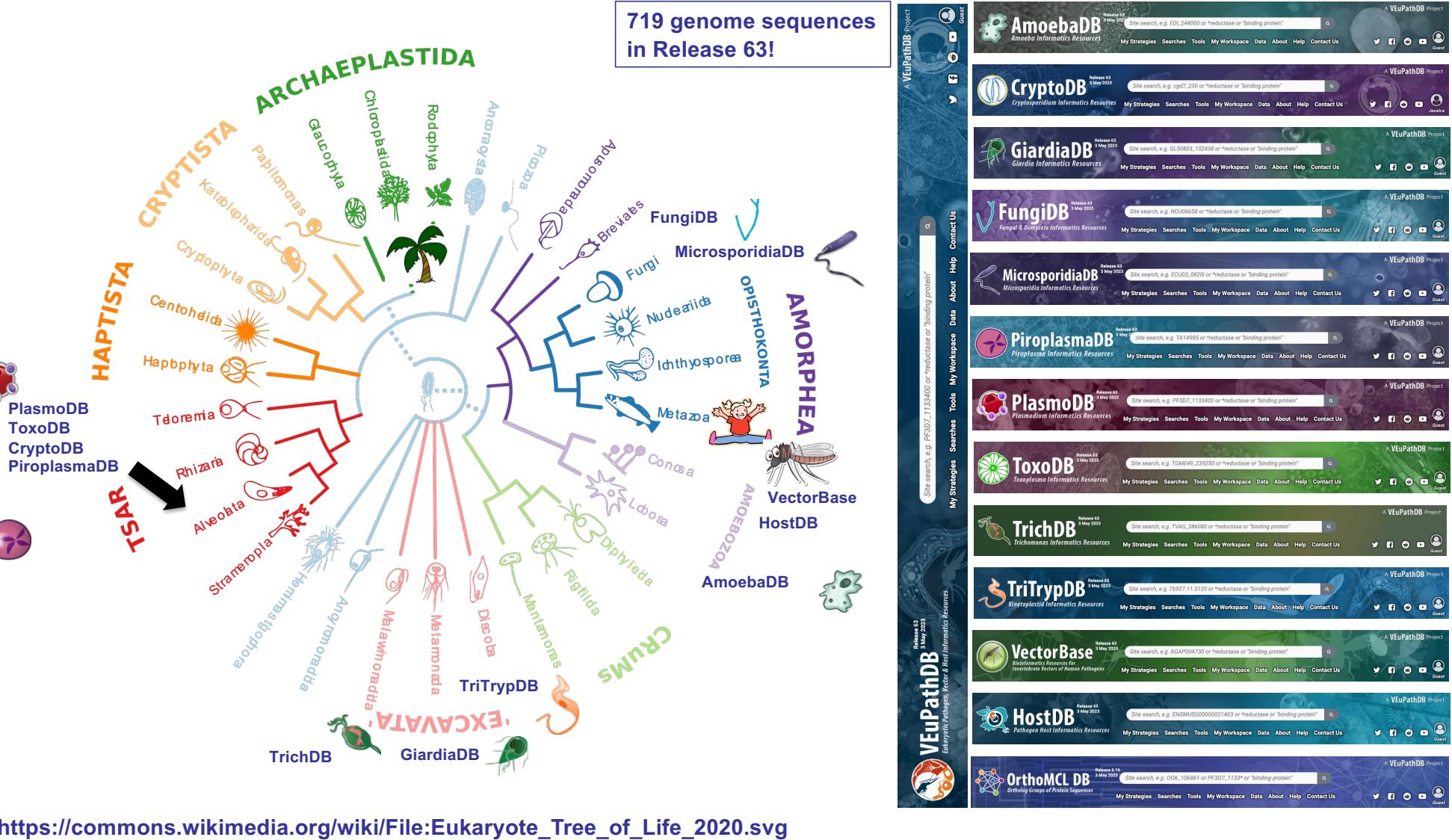
shutterstock.com · 2148907405

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
 - In tissue Y
- That is 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 biological 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

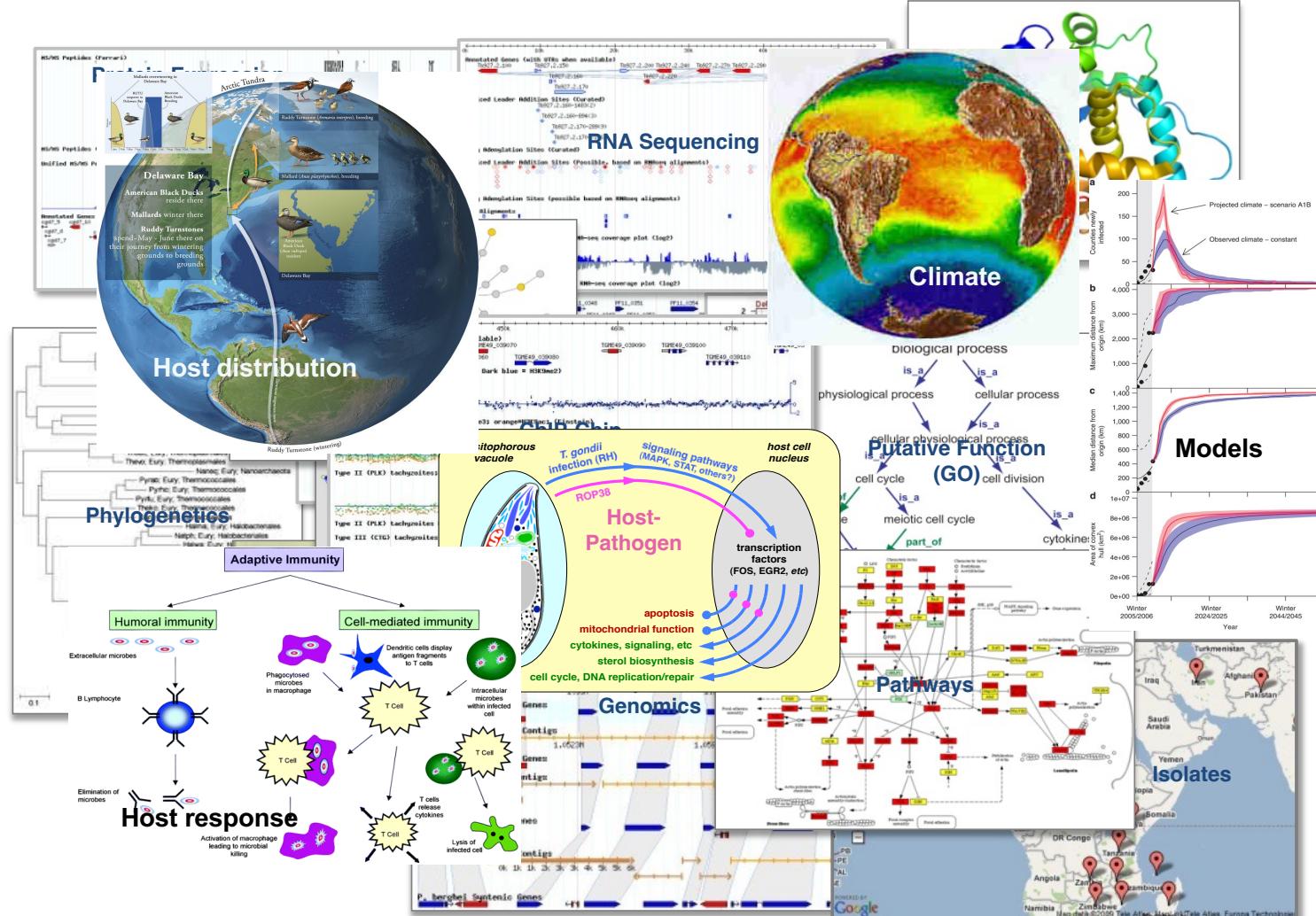
Proteins: sequence, protein features (motifs, signal peptides, TM domains: chemical properties, chemical modifications (phosphorylation, glycosylation), expression, processing, localization

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

Genome Assembly 30,000 ft View

FASTQ
format for
reads



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%

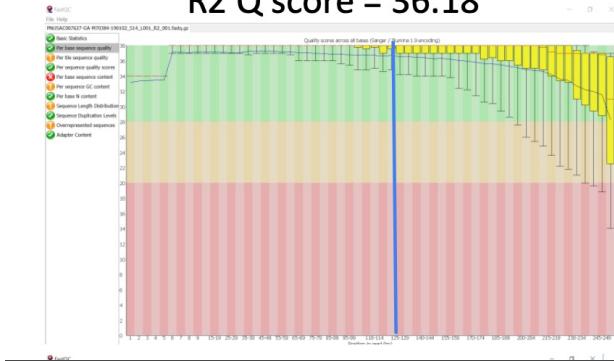
Figure 3 – Phred quality score chart

https://www.abmgood.com/marketing/knowledge_base/next_generation_sequencing_data_analysis.php

FastQC Analysis – Passing Q Scores

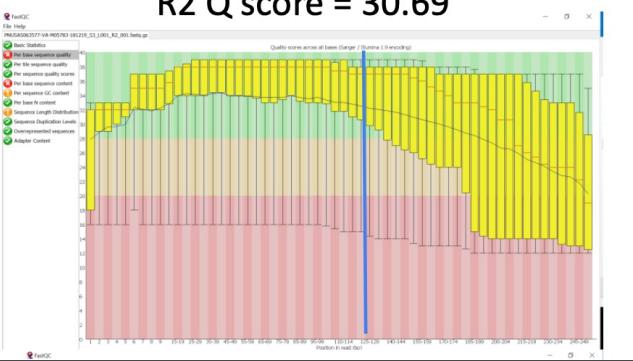
Per Base Sequence Quality & Per Sequence Quality Scores

R2 Q score = 36.18



- Horizontal red line: median Q score
- Horizontal blue line: mean Q score
- Yellow boxes: 50% of the reads
- Whiskers: 80% of the reads
- Vertical blue line: 125 bp of the read

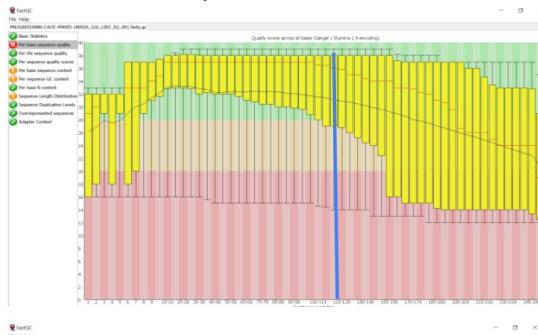
R2 Q score = 30.69



FastQC Analysis – Suboptimal Q Scores (pass with extra coverage)

Per Base Sequence Quality & Per Sequence Quality Scores

R2 Q score = 29.56

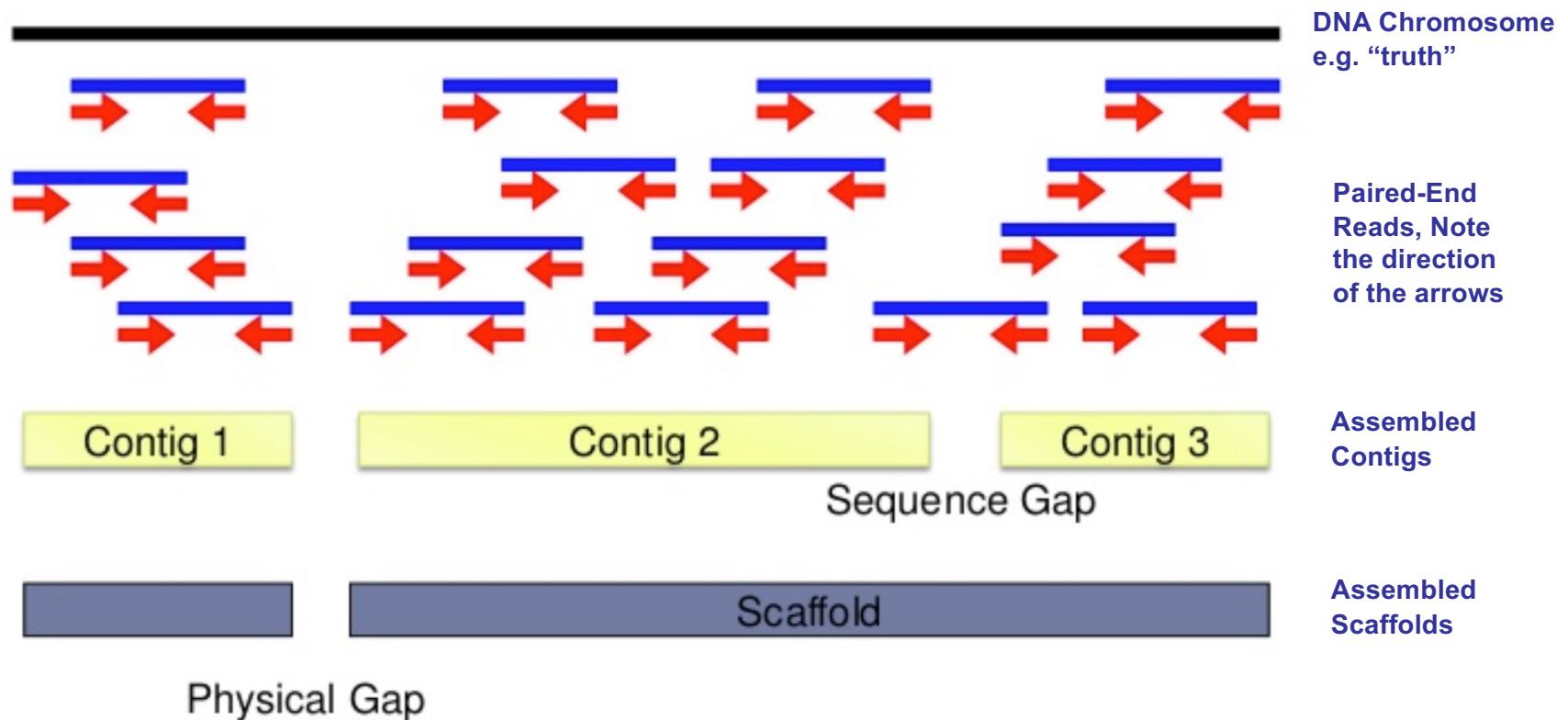


R2 Q score = 28.56



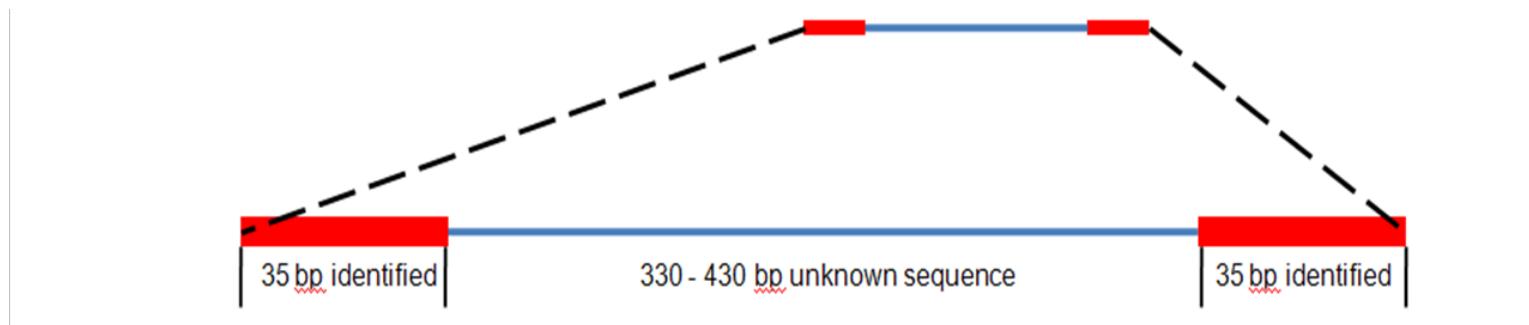
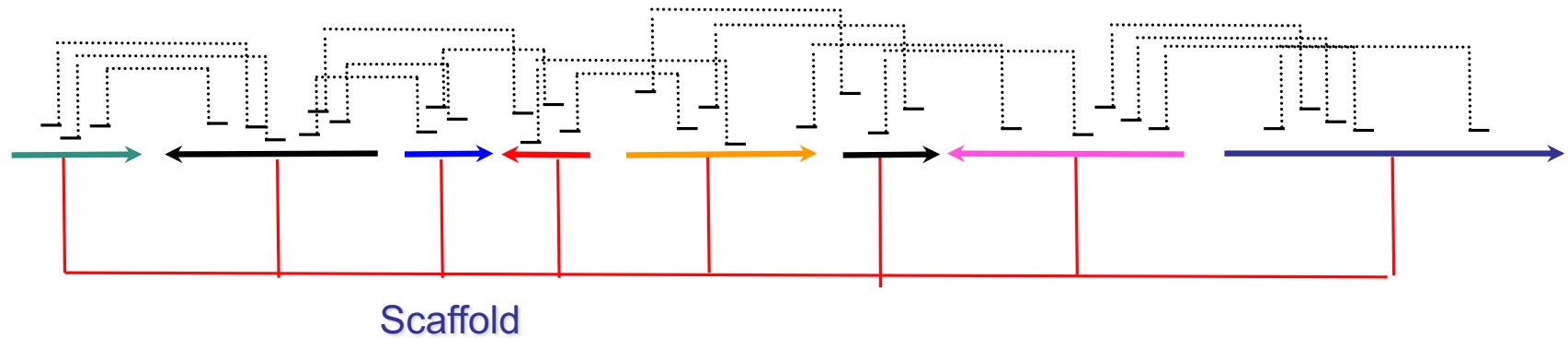
https://www.aphl.org/conferences/proceedings/Documents/2018/4_Eija%20Trees.pdf

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



<https://github.com/Ecological-and-Evolutionary-Genomics/eeg2016/wiki/Mar-21-Exercise-7----SPAdes-assembler>

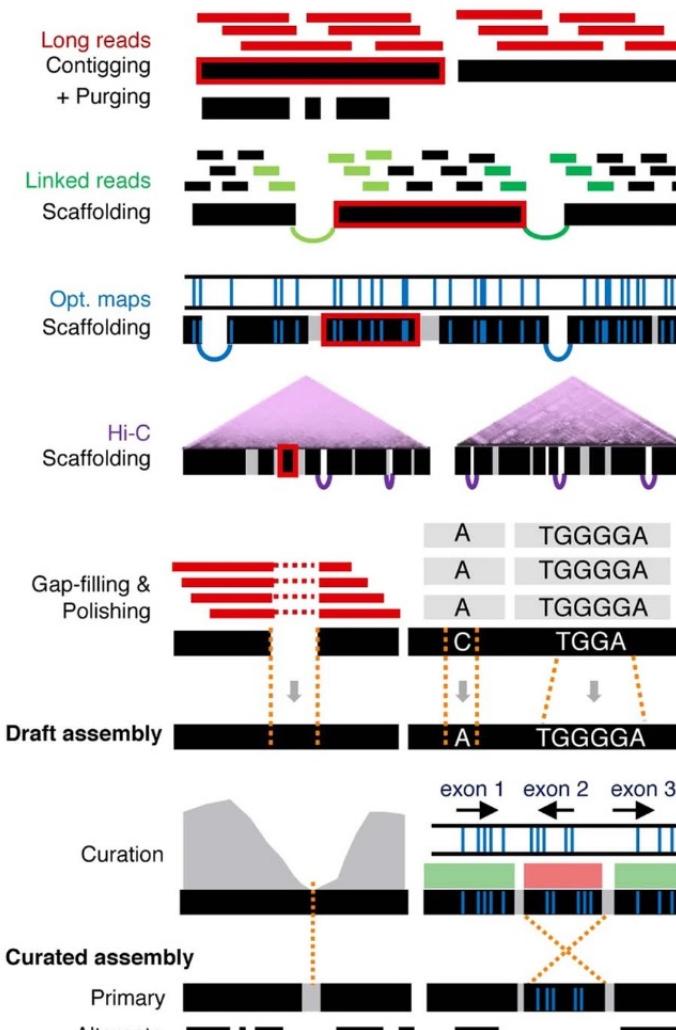
Paired-End Reads can Yield Order & Orientation of Contigs



<https://www.biostars.org/p/104218/>

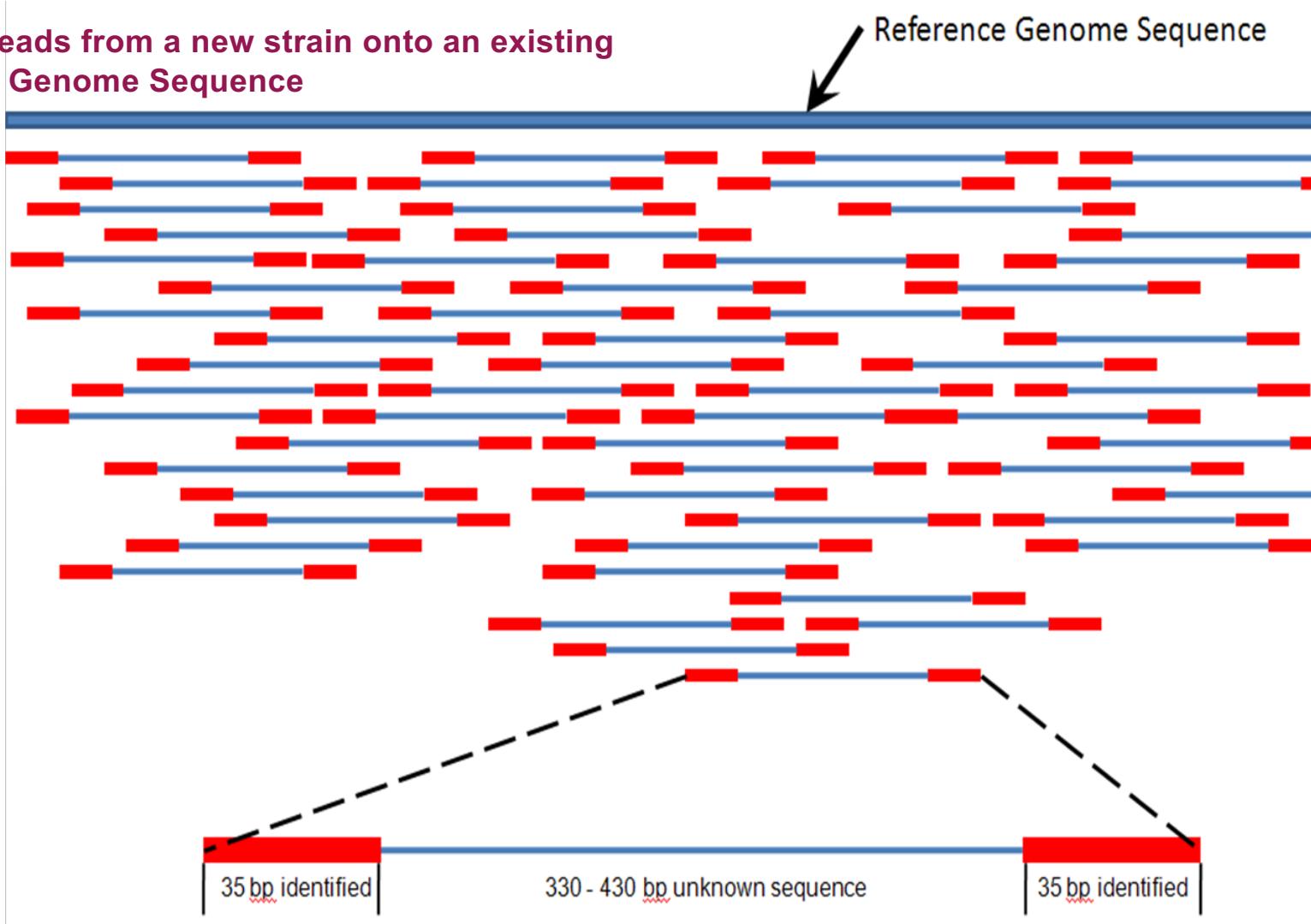
De novo assembly of a complex genome sequence from scratch requires many technologies:

- Deep long reads or Long reads and Illumina short-reads
- Some form of physical mapping, can be genetic or optical mapping for chromosome interactions captured with Hi-C
- All assemblies have gap and these need to be filled and/or corrected this phase is called polishing.
- Assemblies should be curated by a human to catch errors of mis-assembly (often apparent when read coverage is low as in the example



Rhie et al 2021 Nature

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



<https://www.biostars.org/p/104218/>

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

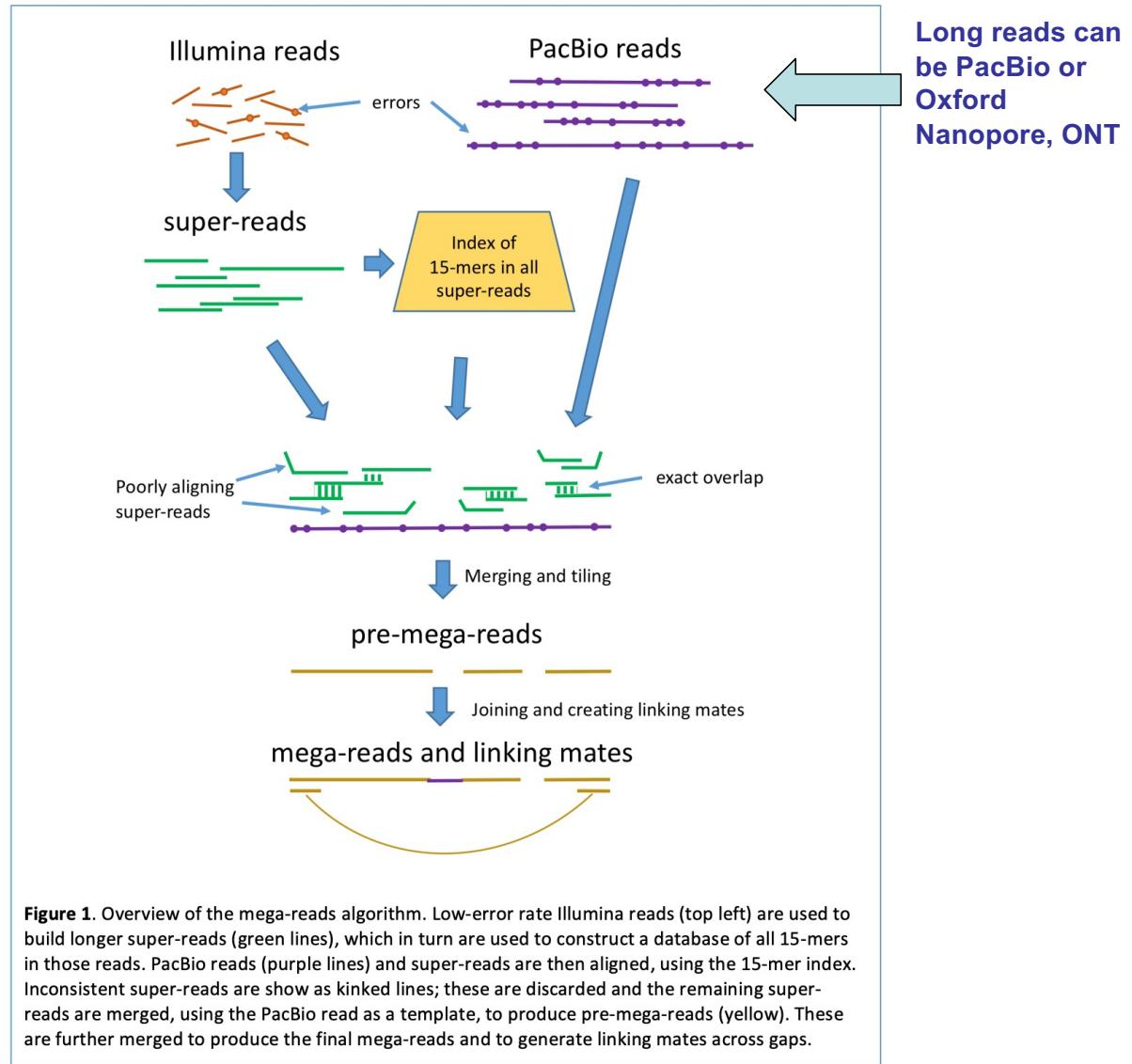


https://github.com/trinityrnaseq/NaplesWorkshop2016/wiki/Day_2

Hybrid Assembly = short + long reads from differing technologies

It is a **VERY** useful
approach for "correcting"
and completing telomere to
telomere (T2T) genome
assemblies

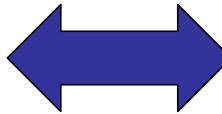
[https://genome.cshlp.org/content/early/2017/01/27/
gr.213405.116.full.pdf](https://genome.cshlp.org/content/early/2017/01/27/gr.213405.116.full.pdf)



Genomes: Important Considerations for Assembly and Interpretation

Biological

- Size
 - Mb
 - Gb
- Ploidy
 - Haploid
 - Diploid
 - Tetraploid
- Repeat content
 - Retrotransposons
 - Big gene families
 - AT content
- Clone vs population



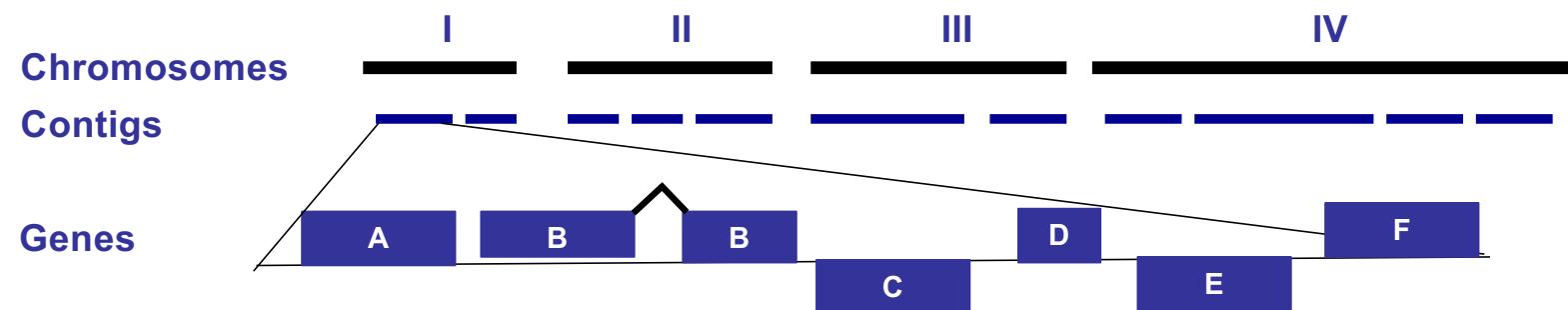
Technical

- Read length
 - Short
 - Long
- Coverage
 - 5X
 - 100X
- Read Quality
 - 20
 - 30
 - 40
 - Bias?

DNA sequencing technologies: 2006–2016 Elaine R Mardis NATURE PROTOCOLS | VOL.12 NO.2 | 2017 | 213 [Seq types chart download \(11 MB\)](#)

<https://www.dropbox.com/s/kfkkft5gilmxd68z/ForAllYouSeqMethods.pdf?dl=0> (PDF figures for next two slides)

30,000 ft View - Genome Annotation



The Genome Sequence

AAGCTTCGCCAGGCTGTAATCCGTGAGTCGTCCTCACAAATCATCAAGCAGGTGTCCTCAGGGAGACTGCCTGACTGAGTTATGCTAATTCCTTCTACTTTGGCGTGGTCACGTGTA
ACCATATCCGAATCATTCTCTAGCCCCTACGAACAGGTAAAGAGCGCTAGGGATGTCCTGGAGTAGTGTGCTTACTCGATAATAITCAGTTGGAGCTACCAAGCGGCGCTCGCTTGC
CACGCAATGCCCTGAGACAGTGCAGAAATGAATGGTAAACCGACAACCGTTCATATGCTTTCAAACCTAGTAGACGCTACTGTGCTGAAACTGCGGTACAGGCACCAGATAACGCC
CTTGGCATGGCATGTCCTACAGAGGTGGTATGTAGTCCCAGACTCTAAATCCGGCAGAGGCTGGCTTTGCTTACCCAGTATTAGCCCGCTGGATTCTCGGAGCGCAC
CTGTTCAACACTAGAAAACGGAGTTCTGATCGAGAACGCCACCTTCCAGAAGTTGAACGCTAGCATGTCTTCAGATTTCACCCCCCGTAGTTCTGTGTGTCATTGTTGTC
GAGACAACCTGTCGCCCGGTGCTGGCATATGGCTGAGCTTCAAGGAAGACAGCTCCGGAAACGATCTCCATGACTGGTAATCCACGACA
CCGCAATGCCCGCCACGGCCTTACATCTCTGCTGCCAGGACTAACCTGTTGCTGGCTTGTGCTTCAGGTTTCCAAAAAAAGAGACGCCATCCGTTCCCCCGCACATT
AACGCCGGAGTGGCGTTTGTCTTTTGAGTGTAGGAGCGCTTTCATCGCGAACTACGGCAACTAAGTCCATTCTCTTCGACAGCGAAACCTGATTCAAACCCCG
CCCGGGAAGATCCGATCTGCTGCTGCGAGTCCAGTAGCGTCTGCGCCGCGCTCTGTTGGTGGCGAGCGCTACACCTGTTATCTGACTGCCGTGCGGAAAATGACGC
CATTTTGGAAAATGGGAACCTCATTCTTAAAGATAGCGGAGCTTCTTCTGTTCTGTTCTGTTCTGCGGTTGATAACCGTGTGCTGATGTAAGCAGCTTCCGTCTC
TCCCTCCGTCTTGTGACATCGAGAACCGAGGTGTGAGATCCCTGGTGTGAGATCCGGTGTGCTGAGACCTTTACACACGGCAGTGGGAGCAGTGTCTG
AGTCAGCAGGGACGGTGAAGGTTGCTTAGTGTGCTTCTGCTGCTACGGGGCTGGTGTGCTGAGAAGTGCAGAAACCCGGTGTGCTGCGATGACCCCCAGAGG
GGCATCGGCATCAACACGGCCTCCGGCCACTGACCCAGATTCAACACACTTTCTGTTGACAAAAACGACGCCGAGAAGGCCAGTGGCTGAACGGGTGGCTCCAGG
AAATTGCAAAGACGGGACACTGGACTTCCCTCCATCGAGTGGCAAGAGATTAACCGCTGTCATGGGACGGAAAACCTGGGAAAGCATGCCCGAAAGTTAGACCCCTCGT
GACAGATGAAACATCGCTTCTCTCTGACACAGACTAGTGTGACCCACCGCTTGTGAGACGCTGTCATCTGAGGAGGGAGCGGGGGCACACAAACTCTCAACTCG
AACATCCCGTGTAGACACACCAAAAGACACGGGCAATCTGCTCATGGAGGGAGGGGGGGCACACAAACTCTCAACTCGAAGAACATATCCGGGCC
GAAGACGTGGAGTCTCAATCAACCCGAACGAAACATTCTCCATCAAGTCAAGGATCCGGCGTACCTCCATGTTGTAACGAGTTCATGAACACCTCCGATATTACACGACTG
TGGATATGAAATTATGCAAGATGCAATATACTGAGACGCCGATGCAACTATAGGTTCTGGCCCTCCATGGATATTCAGACCTCCTCACATTGGTTGCCGTACCTCCGT
TACGCTTTTCTGCTTCTCTGCTGTTATGCAAAGAAGACATGGCGGAGGAACCTCAAGCTGAAGGCCAGCAGGGCTCCGAGCTGTGCTTCACTCCAGC
AGCTCTCACGCTTCTGGAGGAAGAGATCTGCGACAGGATTCTGCTGGGGTATGTTCTGCTTCAACTCTGCTGAGAAACGTACTGTTCAATGTT
CATGTTATGATGATGTTGATAATCTAGAGAAAGATACGGGAAGACTGGCAAGGATGAAAGACATGCGCTTAAACGAGGAGGGCATTGGCAGAGGGACGCCGTTATGCT
GTGTGATGGCTGTGAATCTTACCTCGCCTTGACTTGCTGCAAGCGCTTGTGCACTGAACTGACTTCTGTTTACCTTCCCAACGCCCTCTATTCCCTCACTGCAAGCG
CGCTCAGTGGCCGTCACCGAACACCTTGGTCTTCGTTGAGCTGTCCTTCTCGCTTGCTCCGTGGCTCGGTTCTCTCTCTGTTGGTGGCTCCAG
ACTATGCTCCGTGTTCCCAACCTTCGCGCTTCTCAGGGAGGAGCGGGACTCTAGGAGCAGGCCGCTGTCCTGGCTGCTCCCTCACCGCTGTAACCCCGGA
GTTTCCGTGCGACGCTTCTCGCCGGAGGATGACATTCTTCAACAAACATCAACTGCTGCGAGCTGCTGAGCTGTTCTGCTTGTGCGGAGCT
CGGAAGAGAGAAAGGACAAATGAGCAGCTATGACCCATCTCATTCCAAGACCTCTCAGACAACGGGGTACCCCTACGACTTTGTTCTGAGAAGAGAAAGACTGACGCC
AGCCACTCGCGAACCGTAGAGGGAAACCGAACGCGTAGATAAGAAAAACACAAAGAGAAGGTGAAACACGAAGAGAAGGGAAAGCGGAGAAACCGTGGATTCAAAGATATCAA
GACCAATCTTGTGGAGATTTTTAAATTCACTGAGAGACACCCGGCTGCGAGGTGTGAGAAATACTGCCACCTCGGAGACAGATGCCGAGTACACCACTTGTGCTTCTTCC
TCCTATGTTCATGCGGCTGTAACGCTCTATGTAATGGAGGAGCTGCTCTGGCGAGCAGCTTGGCTGGCTGGCTGGCTTCTGAAAAGGCCAGAAGGCCCTCC
ACAGTGAAGGGCATATACAGGGACGCTACCGGAGGCCGTTCTGCTGACTCTTGAGAGAACGCAATGAGCTCTGAGCTCACAGGGAGAGCAACTCCGTCACGGGT
TGCAGGCTCTTCTCGCCGAGCATTGCCCCGGTGTGGCGTAGGAGCAAGAACGCCGAAAAGCGAGCAAAAGGAACGTTGGCCGTTCCGATGTTCACTTAGAG
GCCATGAAGAATTCCAGTACCTTGATCTATTGCCGACATTAAACATGAAAGGACAATGGATGACCGAACGGTAACGGCAGTGCAGGAAAAGCCACACCGTTCTCTGTTGAT
TCTGTCCTCAAGGCCCTCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTG
TCATTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTG
CGTGAATCGATTGCGTTGCGGTTCTGGGTAGAAAAGGCTGCGCAGTATTCTGAAATAACCCCTGCCATTGTTAAAGGGCAAGGAACAAAGAGATATTGCGGCTCATCT
TTTGTGCGGCGCGTTCTCGTGTCTCACACCGATGCCCTCTGCTGATGTTCTGCTCTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTG
TATGCGCTACTGCTGGATCAGGCCCTTCCACCTTCTCAGGAAAGGGTAGGGTAAGGGCTGCTTCTGAGATGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTG
CAACCAACCTACAAATTGTTGTCGCTGCGTAGATGTCAGATGTCAGTAAAGACTGAGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTGCTTCTG
TTCTGAGTGTCTGGAGTTCTGGCAACCTTCTTGAATTCTGGGTTGTTTATGCGGCCACTGGTTGCTGAGAGACAGATGCAAGGTGGGTGATGCGCT
GCTGAGAGAAACTCCGGCAAGGGCAGATAAAAGGAGAGTGGAAATCATTGAACAGTGTGCGGTGCTGTTGCTGAGGGCTCTCGAAGAGTGTGCTGTTGCTTCTGCTTCTG
CGAAGCAACCATCTTCTGAGAAGGGCTGAAAGGCAAGTACGTTGTAACCTCTGCTCTGCCAGCTGAGATGTCATGCCAGGGCTGGTTCTGCTTCTGCTTCTGCTTCTG
TTACCATGAGTCAACACTCATGTTGCGTGTCTACATGTTTCTAGAACGTCGGTTGCTGCTGCGACCGGGCGAGGTGATGTAACCTCTGCGCTGCAAGAGTTGATCCT

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)

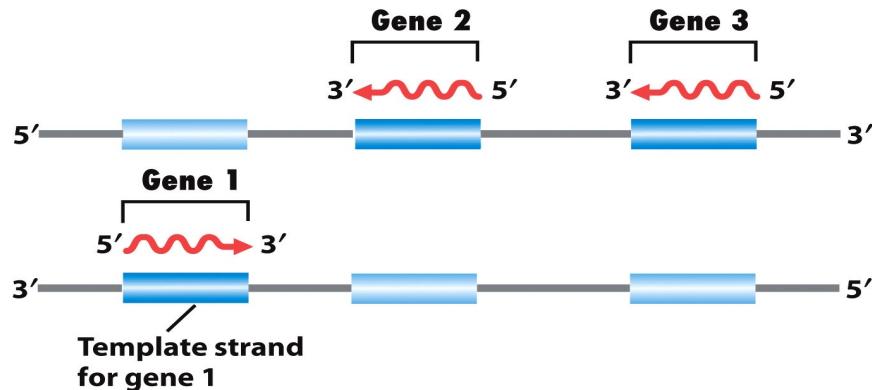


Figure 8-3
Introduction to Genetic Analysis, Ninth Edition
 © 2008 W.H. Freeman and Company

Nontemplate
 strand 5' — **CTGCCATTGTCAGACATGTATAACCCGTACGTCTTCCCAGAGCGAAAACGATCTGCCTGC** — 3' } DNA
 Template
 strand 3' — **GACGGTAACAGTCTGTACATATGGGGCATGCAGAAGGGCTCGTTGCTAGACGCGACG** — 5' }
 5' — **CUGCCAUUUGUCAGACAUGUAUACCCGUACGUCUUCCCGAGCGAAAACGAUCUGCGCUGC** — 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 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 GTG 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

AAGCTTCCGCAAGCTGTAATCCCGTGAGTCGTCTCACAAATCATCAAGCAGGTGTCTCAGGGAGACTGCCTGACTGAGTTATGCTAATTCTTCTACTTGGCGTG
GTCACTGTAAACCATATCGAACATCTAGCCCTACGAACAGGTAAAGGCCCTAGGGATGTGGAGTAGTGTCTACTCGATAATATTCACTGGGACTACC
AGCGAGGGCCTCCTTCTCACCCAATCGTAGACACTGGACAATGAATGGTAACCGACAAACCGCTTCATATGCCTTCAAACTTAGACCGCTACTGCTG
AACTGGCGGTACAGGCCACAGATAACGCCCTTGGCATCGCATCGTACAGAGGTCTGTATGTAGTGCACAGACTCTAAATCCGGCAGAGCTGGTCTTGT
TTACCAACGTATTAGCCCGCTGGGATTTCTCGGAGCGCACCTGTCAACACTAGAAAACCGGAGTTCCCTGATCGAGAAGGCCACCTTCCAGAAGTGTG
TGTCACTCGATTTCACCCCCCGCTAGTTCTGTGTCATTCTGTGAGACAACTCTGCCCAGCCCGGTGCTGTTCCATATGCCTGACTTCCCGCAATTNTTC
AGACTTTCAAGGAAAGACAGCTCCGAACGATCTGTCATGACTGGTAATCCACGACACCGCAATGGCCCCAGCACCTCTATCTCGTGCAGGGACTAACCTTG
TATGCGTCTGCGTCTGTCTTGCATTGCTTCAAAAGAGAGGCCATCGTCCCCCGCACATTCAACGCCCGAGTCGGTTTTGTCTTTTGAGTGGTAGG
ACCTTTTCACTGCCGAACACTGGGACATTAAGTTCATTCTTTCGACACGCCAAACCTTGCATTCAAACCCGCCCGAGAGATCCGATTTGCTGCTGCTG
CACTCCCACTGAGCTGCGCCGCGCTCTGTGGTGGGCCAGCCGCTACACCTGTATCTGACTGCCGTGCGGCAAATAGGCCATTTCGGGAAATTCGGG
AACTTCATTAAAGATGCCGAGTTCTCTTCTCTCGGTTGATAACCGCTTGTGATTCAGACTTCTCGTCTCCGATTTACACCGCAGTGGAGACTGCTG
TTTGTGACATCGAGACCAAGGGTGTGAGATCTTGTGATCCGGAGACCGCTGTCTGTAGAACCTTTTACACACGGCAGTGGAGACTGCTG
AGTCAGCAGGGAGGGTGAAGTTGCTTTAGTAGTGCCTTCTGCTTACCGGGCGTTGCTGAGATGCAAAACCGGTGTCGTCGCGATGAC
CCCCAAGAGGGCATGCCATCAACAAACGCCCTCCGTGGCCCACTTGCACACTTCCCTCTCCATCGCGCAAGAGATTCAACCCGGTTGCTATGGAGGAAAATCGGG
TGAACGGGTGGCTTCCCAGAAAATTGCAAAGAGGGCACTTGGACACTTCCCTCTCCATCGCGCAAGAGATTCAACCCGGTTGCTATGGAGGAAAATCGGG
AGCATGCCCTGAAAGTTAGACCCCTGTCGGACAGATTCAACATCGTCTTCCCTTCCCTGTGAGCACACAGACTGCCCACACGCTGTTGAGACCTGTCATCT
CCAAGAGTGTGGACGCTTCCACGTTCAATGTTCCACATCCGTGCTAGTAGACACACCAACAAAAGCACACGGCAATCTGCTCATCGAGGGAGGCC
GGGGGGCACACAACTTCTCACTCGAACGACATTCGGGGCCGCCAGACGCTCCAGCTCTCAATCCAACCCGGAAACGCAACATTCTCGCATCAAGTCAGA
TTGCGCCGGTACCTCCATGTAAAGCAGTTCCATGAAACCTCCGATATTACACAGCTGTGAGATTAATGCGATATATACTGAGACGCCGTCAG
ATAGGTTTCTGCCCTCCATGATTTCAGACCTTCTCTCAATTGGTTTCCCGTACACCTCCGTTACCGTCTTCTCTGCTTCTCTCGTCTGTGTTATC
ACCAAAGAAGAACATTCGGGGAGAACCTCAAGCTGAAGCCAGCAGCGCTCCGAGTCTGCTTCACTCCAGCAGCTCTGAGGAGAACAGTACAA
GGATTCTGCAACAGATTGGTGTGGTAGTTGTGCTCTAAACTCCCTGGAAACTCCATTCTGTCAGAAACGACTGAAACTGATGATATACAGATGTATG
GATAATATCTGAGAACATACAGGGAAAGACTGCCAGGATGAAAGAACATGCCAGCTTAAACGAGCAGAGGGCATTGGCGAGAGGGACGCCGTTATGCTGTG
GCTGTGAATCTTACCTCGCGTTGACTGTCGAGCGCTTGTCACTGAGCTGACTCTTCTACCTGGCTTCACTCCCAACGGCTTCTATTCCCTCACTCG
CGCTCAGTGGCCGTCACCGAACACCTTGGTTCTTCGCTGAGCTGCTCTCTGAGGCTGCTGGCTGGCTCTCTCTTCTG
GTGCGTCCAGACTATGTCGCTGTTCCCCACCTCTCGGCTTGTGCTTCAAGGAGGAGGGACTGTACGAGGCAGCGCTGTCTGGGCTTCTCACCTGTAC
ATCACCGCTGTACCCCGAGTTCTCGCTGCTTCTCCCTGCTGGGAGATGACATTCTTCAAAACAACTGCTGCGCAGGGCTCTCG
GTCGTGCTGTTCTCCCTTGTGAGCTCGGAAGAGAGAACGACATTCGAGGACACTCTCATTTCAGGACCTCTCAGACAACGGGTACCTCTCG
ACTTGTGCTTCTGAGAACAGAACGAGACTGAGCAGCAGCAGGCACTCGGGAAACGGTAGATAAGAAAACACAAAGAGAACGGTAAAC
ACGAAGAGAACGGAAATCGGGAGAACCGTGATTACAAAGATATCAAGAGCAATGCTTGTGGAGATTTTTTAATTCACTGAGAGACACCCGGTGCAGGGTGT
TAGAAATACTGGCACCTGGAGACAGAGATGCCGAGTACACCACTGTGTTTCTCTATGTTGATGACGGGTGCTGAACCTCTATCGTACTTAATGGAGGAG
TCGTCCTCCAGACGACTTGGCTGGCATCCGGTGTGTTGGCTTGTGAAAGGCCAGAGGGCTCCACAGTGGGCGATATAAGGGACGCCCTACCGGAGCCG
TTTCGCTTCTGAGACTTGGAGACAGAACGCACTGACCTCTGAGCTGCCAGGGAGACAACTCCGCTGACGGGTTGCGCTCCAGCTTCTCGGCCGAC
CCCCGGTGTGGCGTGGATGGAGCAAGAGACCGGAAAAACCGCAGAACAGGAACCTGATTTGCTTCACTTGTGAGGCGATGAAGAAATCCAGTAC
CTTGATCTCATTGCCGACATTAAACAATGGAGAACATGGATGACCGAACGGTAACGGCAGTGCAGAAAAGCCACACCGTTTCTCTGATTTGCTGCC
AGCCCTTCTGCTCATCCACCTTCTGCTATCTCCGCCCTTCTCTGCTCCAGTCTCCAGTCTTCCACCTTCTCTGTTACCTCTG
TCATTGCTTCTGCTCTATTAACTGTGTTCTACTCAGACTGCTGCACTCCGGCATAGACGAGCTTCCACGCTTGTGCTGACAGACTGCTATTG
CTCCCTCCACCGTGAATCGGATTGCGGTTCTGGGTCAGAAAAGGGCTCGCCAGTATTCTGATAATAACCTCTGCCATTGTAAGAGGGAGAACAA
AGAGATATTGCGGCATTTGTGCGGCCGTTCTCGTCTCACCGATGCCCTCTGCTGATGCTTCTGCTCTCTGCTCTTCTCTGTTAGG
CGTGGGTGTCATCTCAAATTGGCTGCACTATCGCTACTCGCTGGATCAGGCCCTTCCACCAAAACCGTGTGTTCTGAAAGGGTAAGGGCTCTTCA
GAATGCAATATTGACTTCAGACATTAACTGTTGACACCAACGTAACATTGTTGTCGCTGCGTGTGCTGACATGTAAGTATGTAAGAGCTGCTACTGT
AGACTAACGCACGAACCAAGATTGTTATCTGCATGCGCTGTGACCCGTTCTGAGTGTGCTGGAGTTCCGACACCTTCCCTGAAATTCTGGGTTCTGTTTATGC

>Translation Frame 1 The Protein

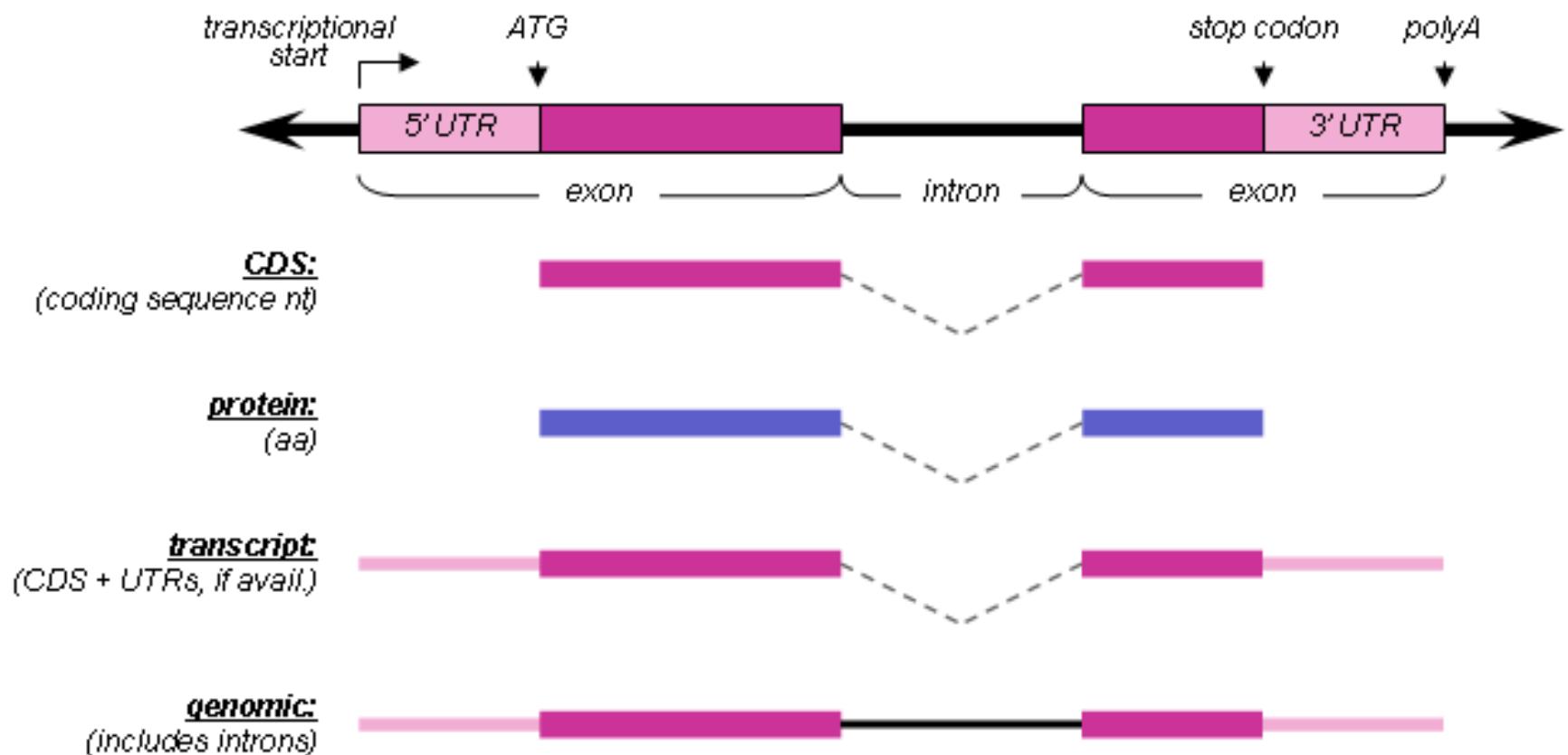
MQKPVCLVVAMTPKRIGINNGLPWPHTDFKHSRVTKTPEEASRLN
GWLPRKFAKTGDSLPSVGKRFNAVMGRKTWESMPRKFRPLVDRLNI
VVSSSLKEEDIAAEKPQAEGQQRVRVCASLPALSLEEEYKDSVDQIFV
VGGAGLYEAALSLGVASHLYITRVAREFPCDVFFPAFPGDDILSNKSTAA
QAAAPAESVFPFCPELGREKDNEATYRPIFISKTFSDNGVPYDFVLEK
RRKTDDAATEPSNAMSLTSTRETPVHGLQAPSSAAAIPVLAWMDEE
DRKKREQKELIRAVPHVFRGHEEFQYLDIADIINNGRTMDDRT

Green = UTRs

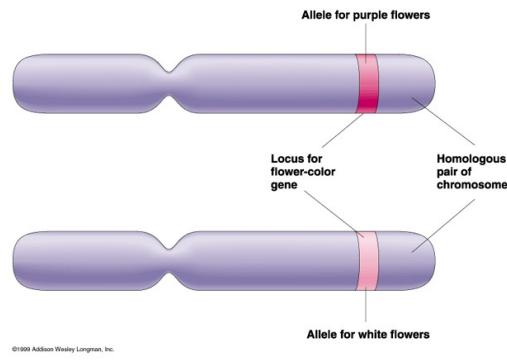
Red = CDS

Pink = Intron

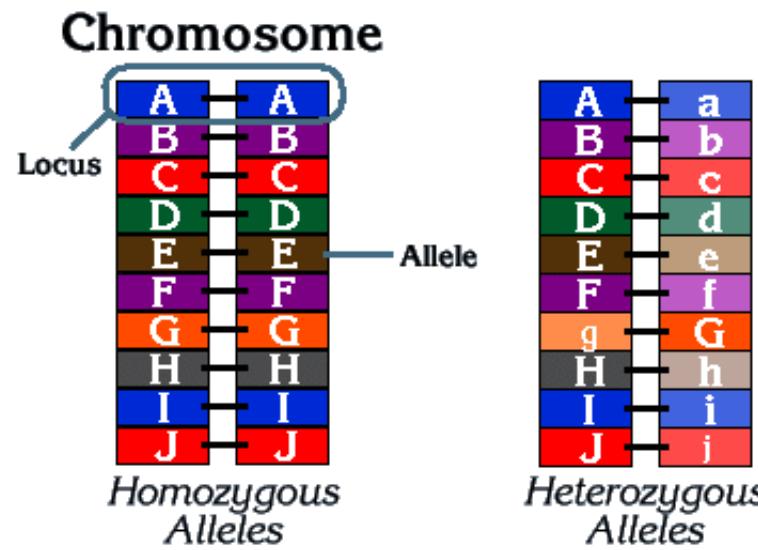
Terminology



Evolution Homologous chromosomes (in a diploid)



A AAGCCTCATC
a ACGCCTCATC

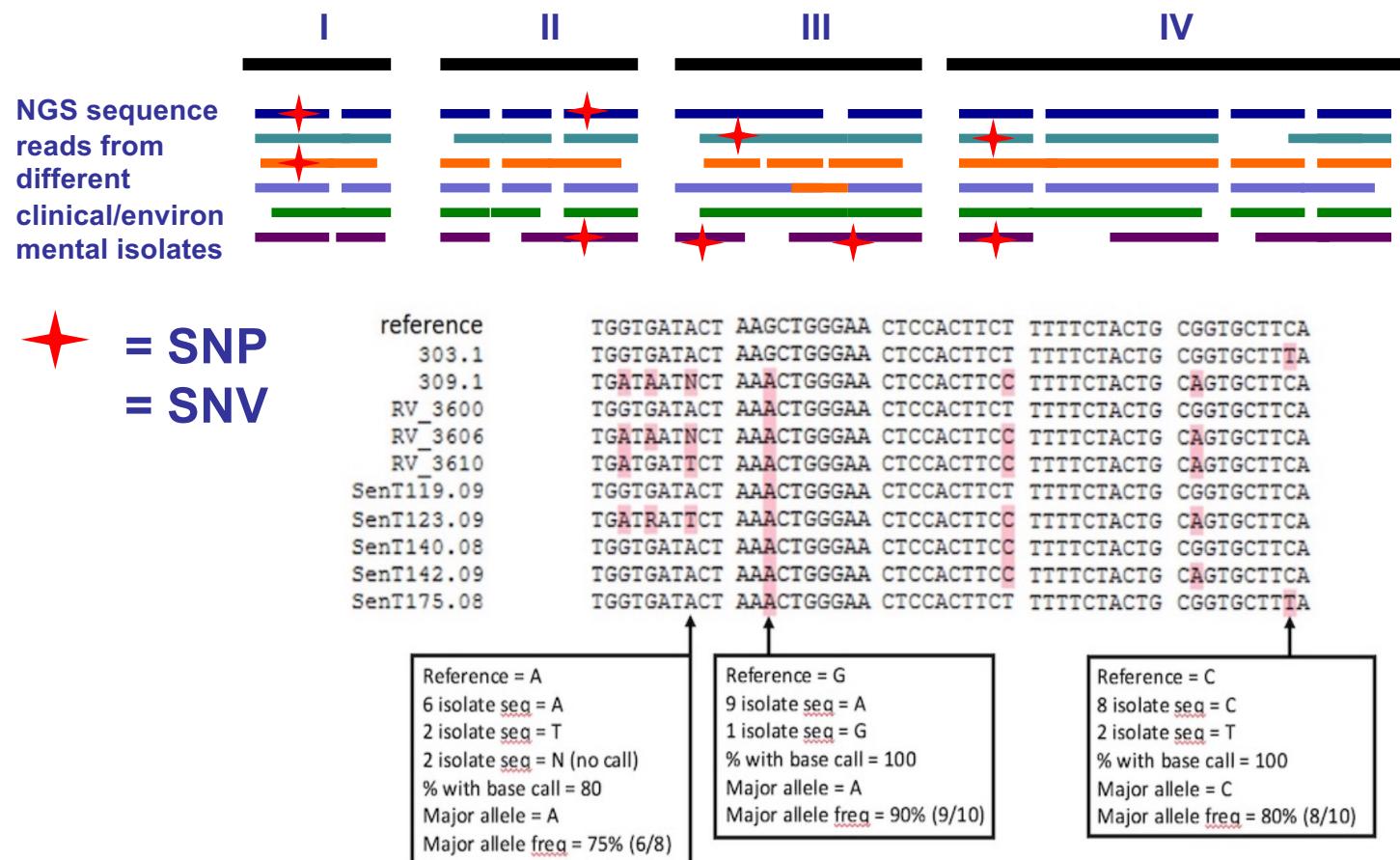


SNP =Single Nucleotide Polymorphism (a variant)

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)

30,000 ft View- NGS SNPs



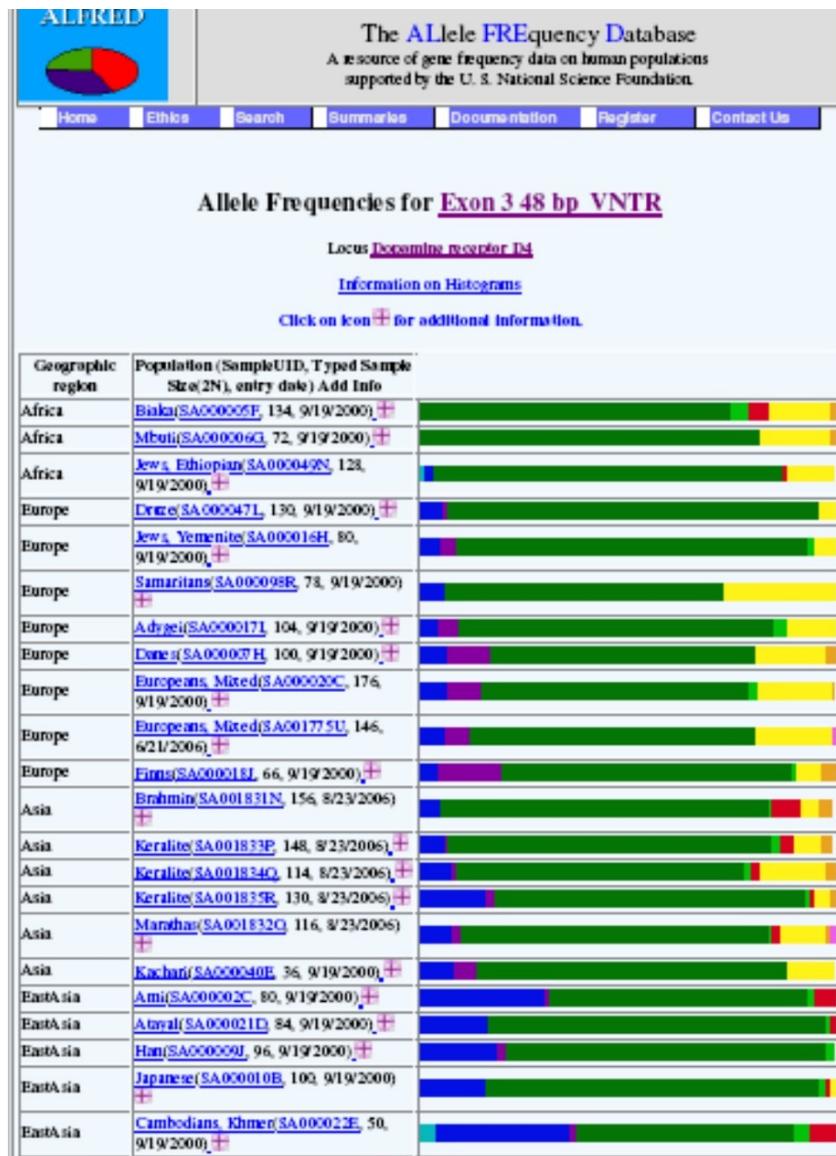
Population variation data

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)

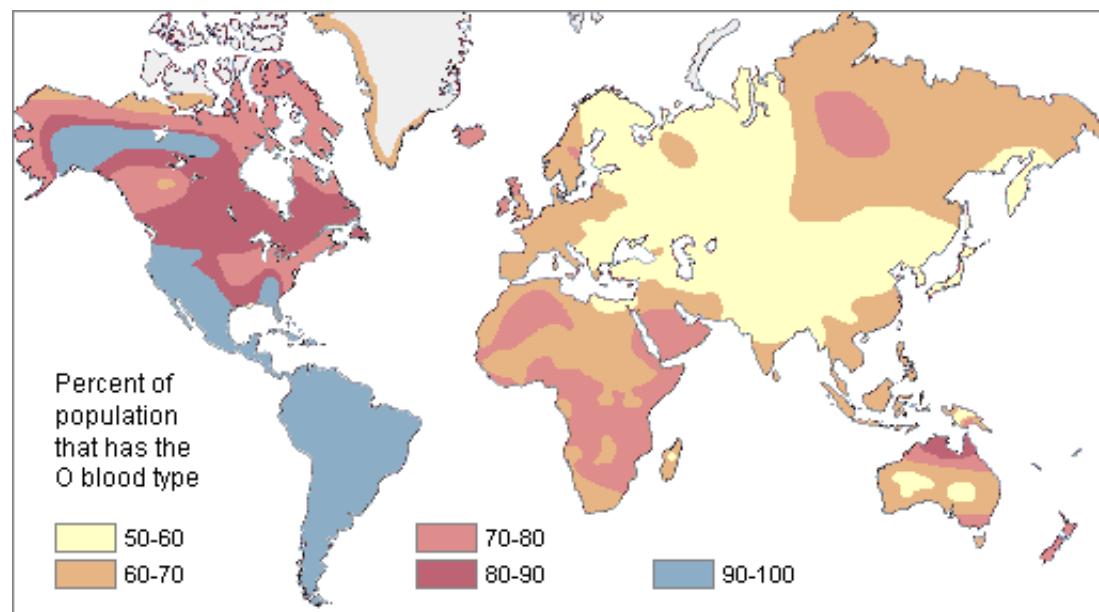


Alleles have frequencies in different populations

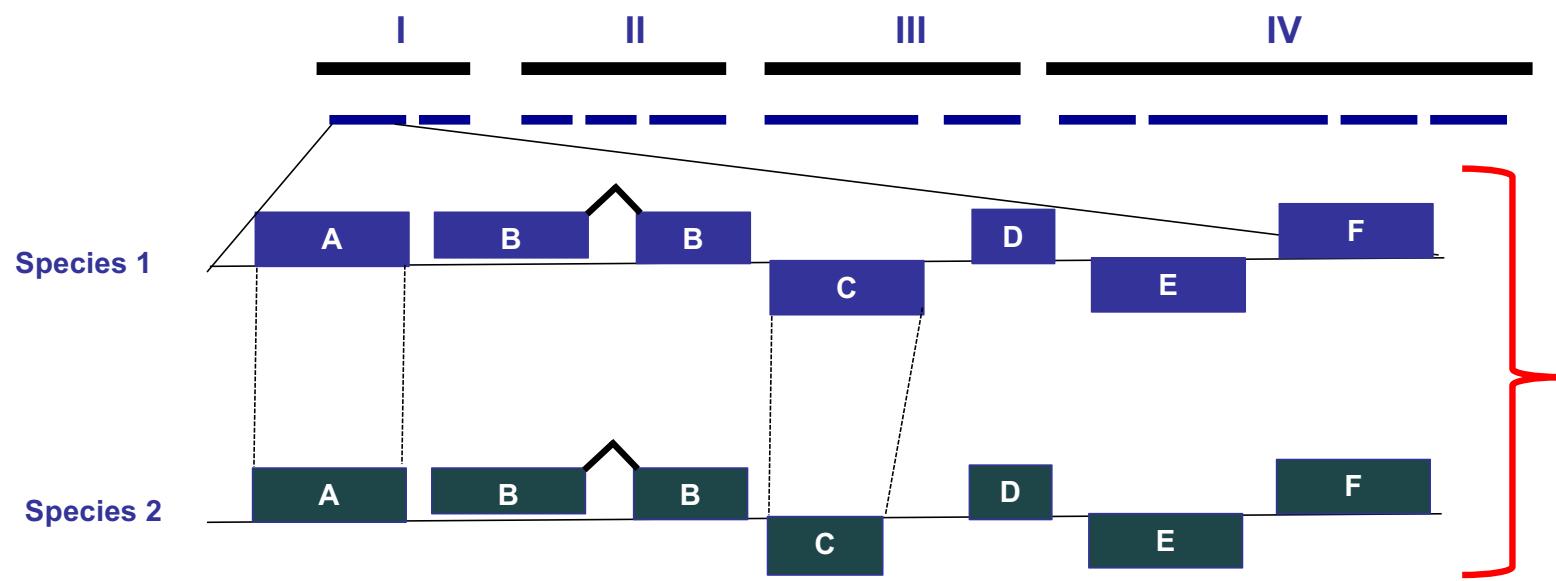


Populations and alleles can have geographic boundaries

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

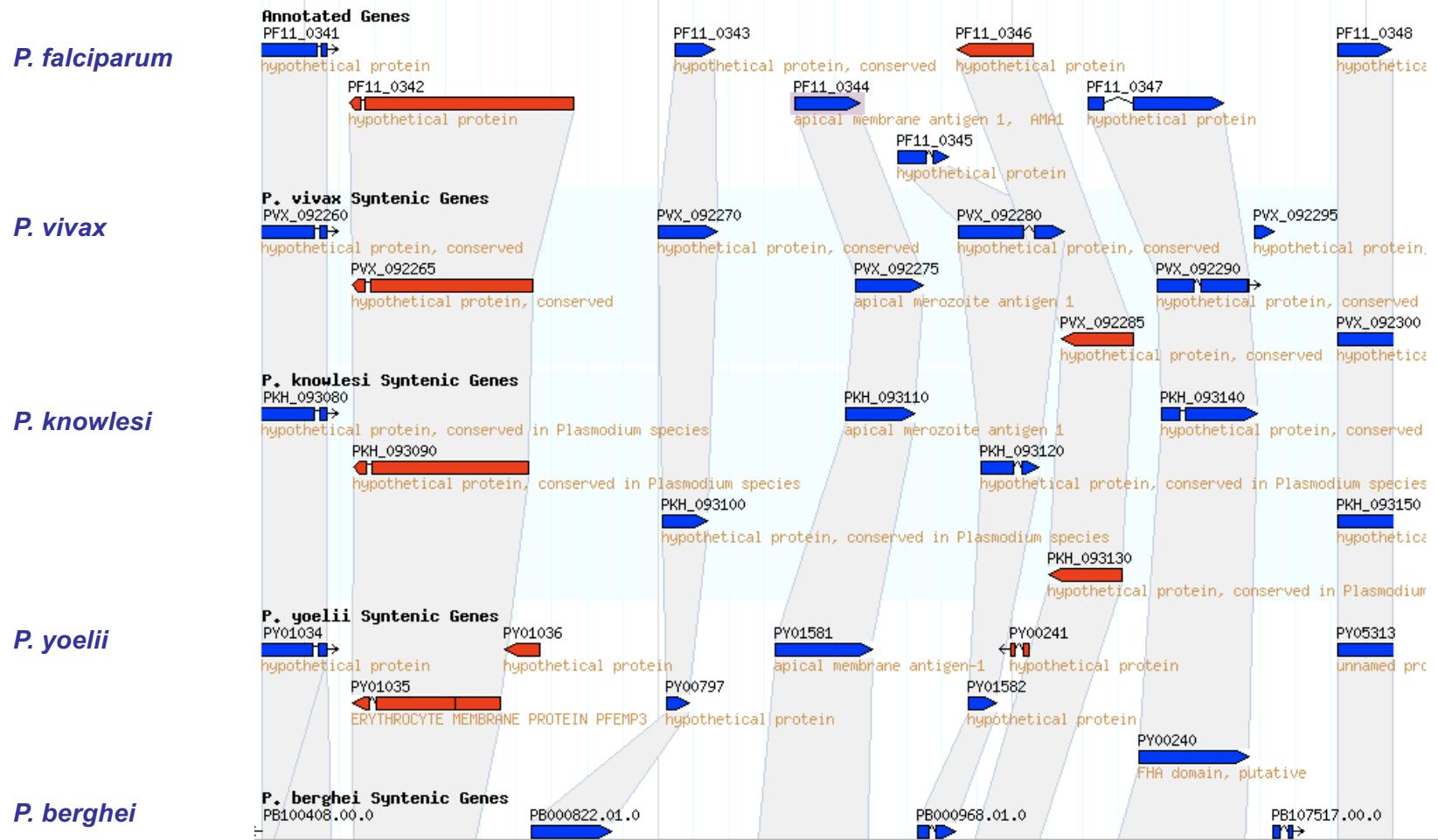


30,000 ft View - Synteny

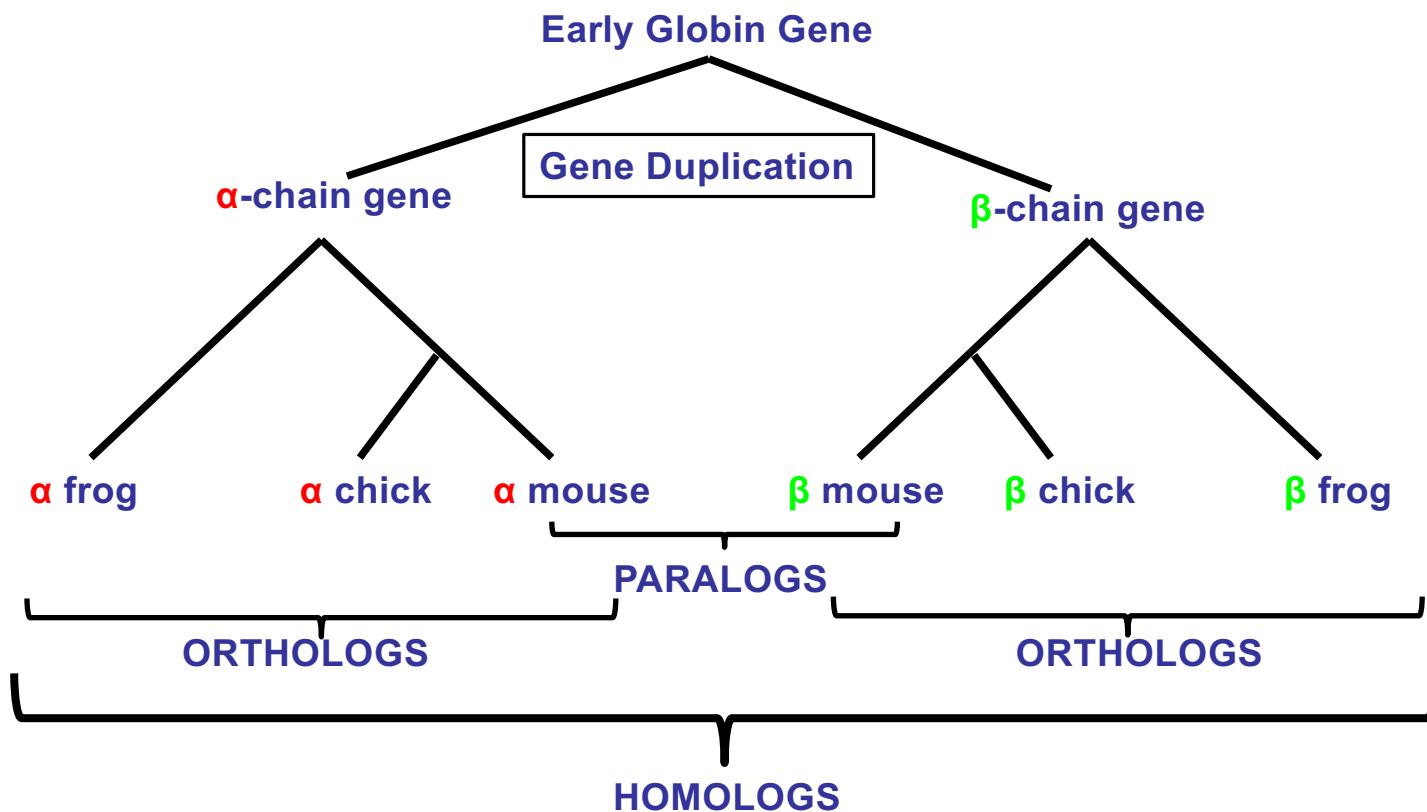


Synteny = the majority of the same genes are present in the same order and orientation in another species. The chromosomal regions are evolutionarily related

Synteny among *Plasmodium* species



Homology - a vocabulary for different types of evolutionary relationships



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:

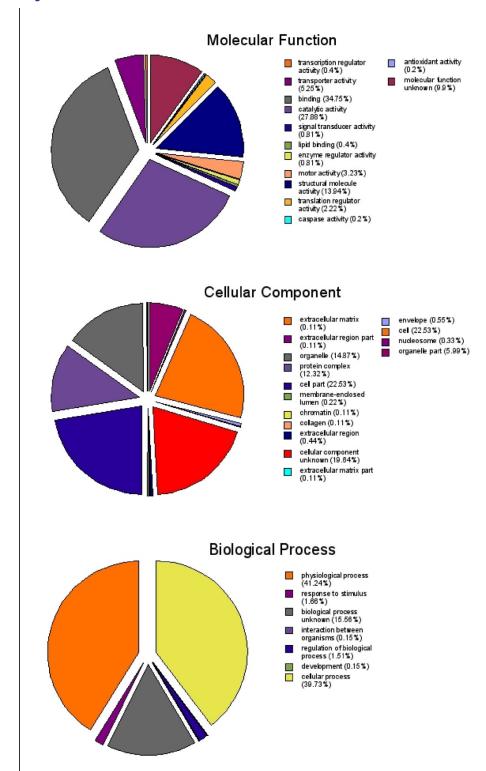
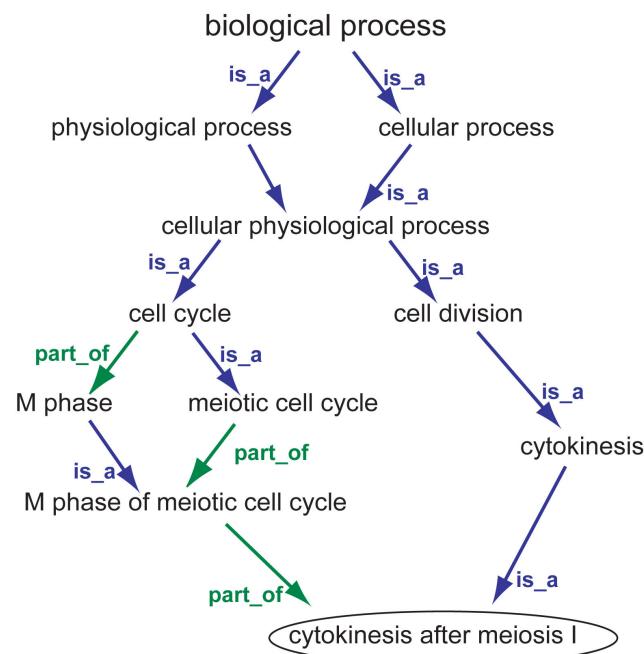
D. melanogaster gene CG3340 annotated as: "Kruppel" and *P. falciparum* gene PF3D7_1209300 annotated a "putative KROX1"

Both can be annotated with GO term:

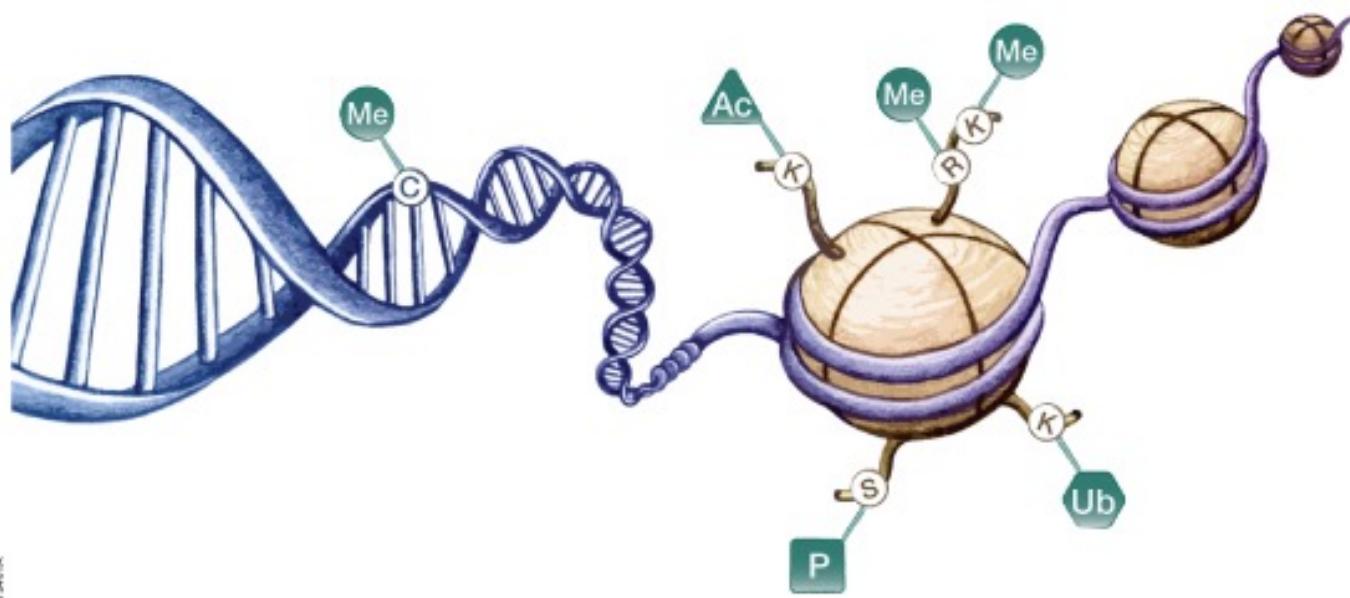
GO:0003705 (RNA polymerase II distal enhancer sequence-specific DNA binding transcription factor activity)

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



Chromatin Status and Epigenetic Gene Regulation



113477A

- 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

<https://www.promega.com/resources/guides/nucleic-acid-analysis/introduction-to-epigenetics/>

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

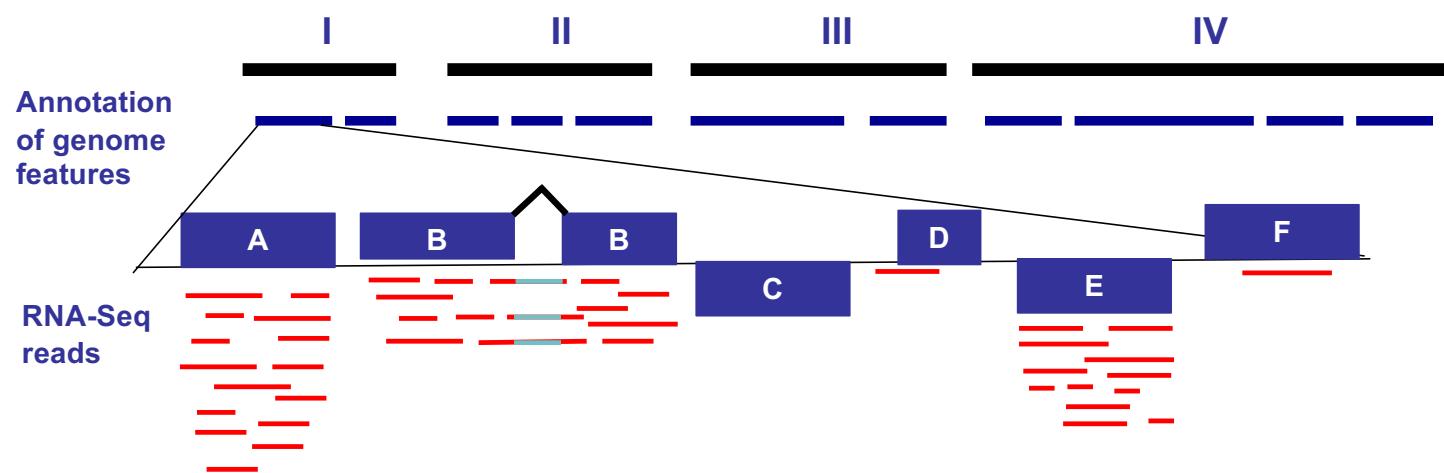
Bulk sequencing from many cells

- 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

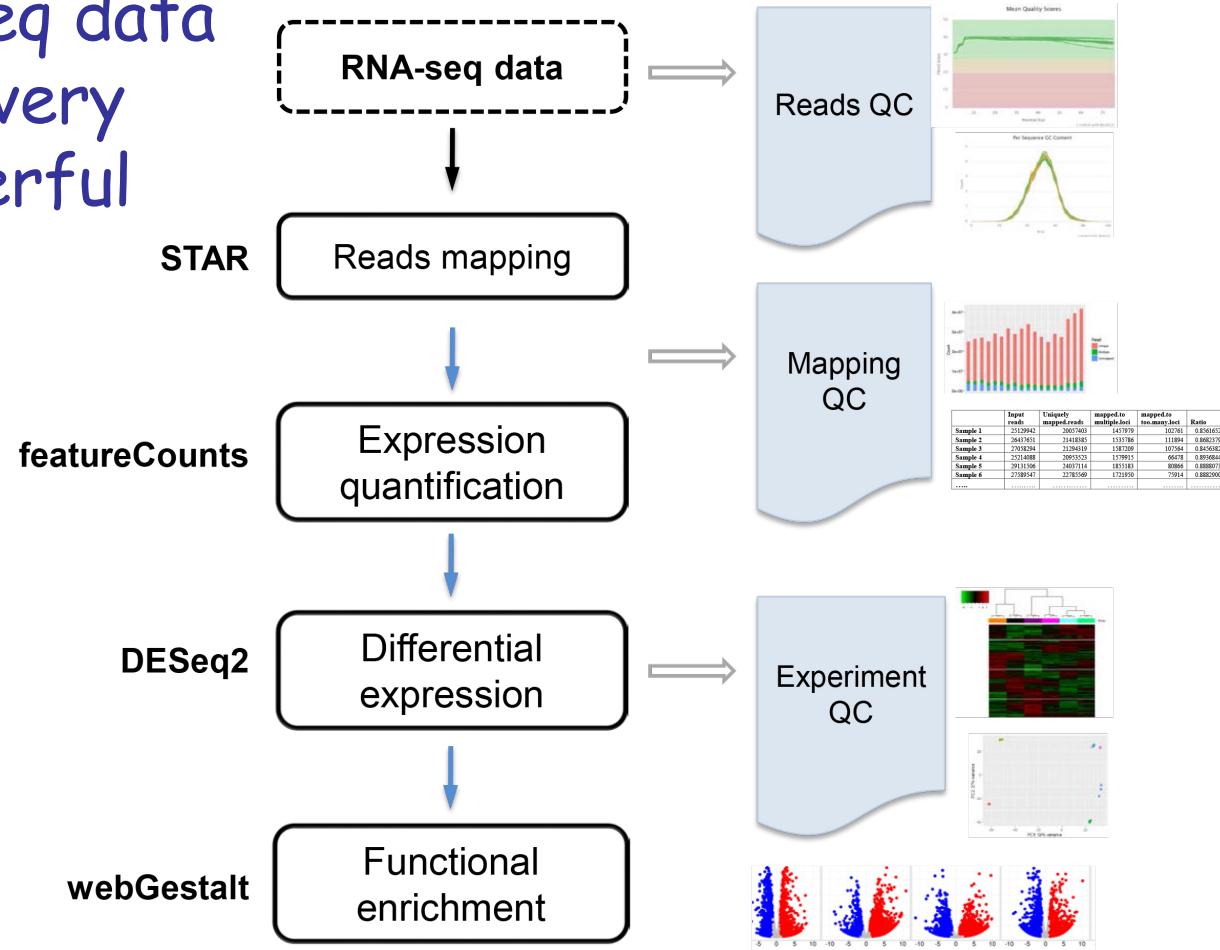
- Examines the transcriptome inside each cell analyzed
- Excellent for detecting cellular heterogeneity or differentiation
- Often only detects a fraction of the transcripts within a cell
- Often analyzed with tSNE plots to categorize cells that have similar transcriptional profiles.

30,000 ft View - RNA-Seq



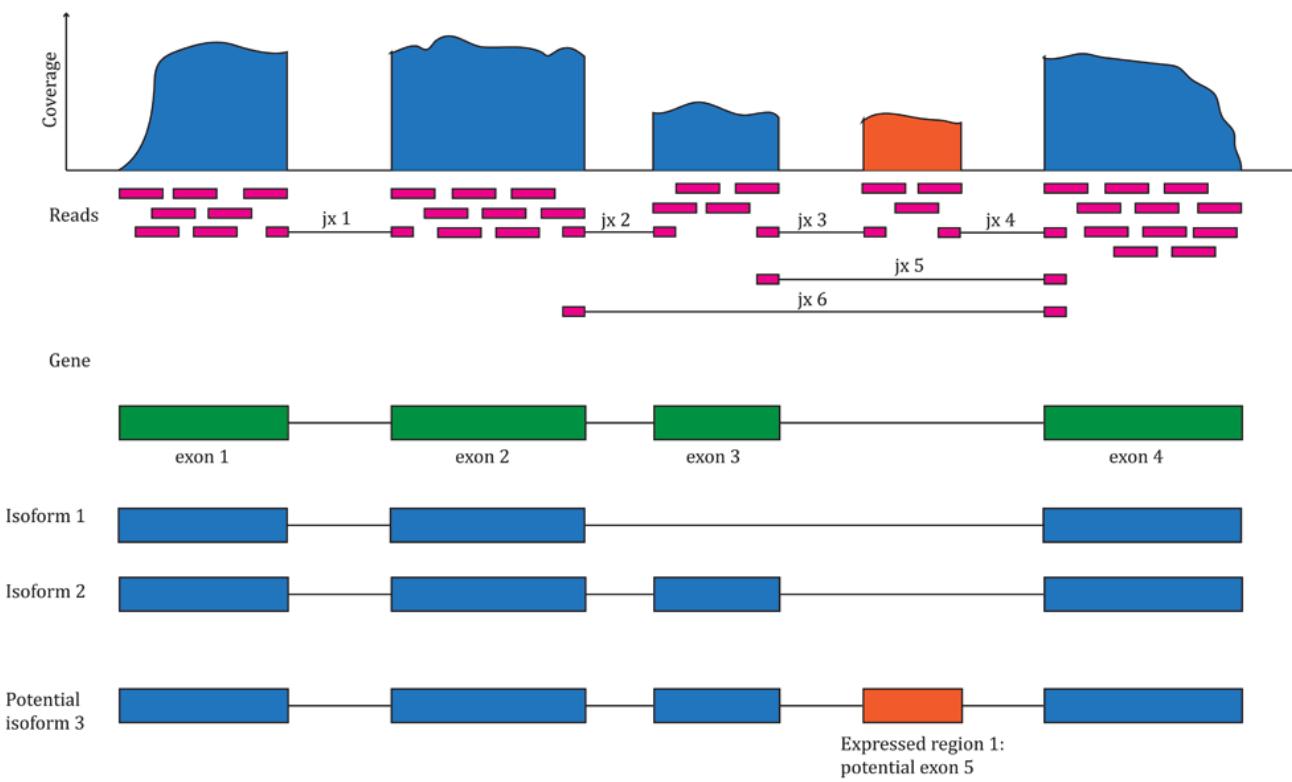
FPKM = Fragments per kilobase of exon per million fragments mapped

RNA-seq data are very powerful



<http://bioinfo.vanderbilt.edu/vanguard/services-rnaseq.html>

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



<https://bioconductor.org/packages/devel/workflows/vignettes/recountWorkflow/inst/doc/recount-workflow.html>

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

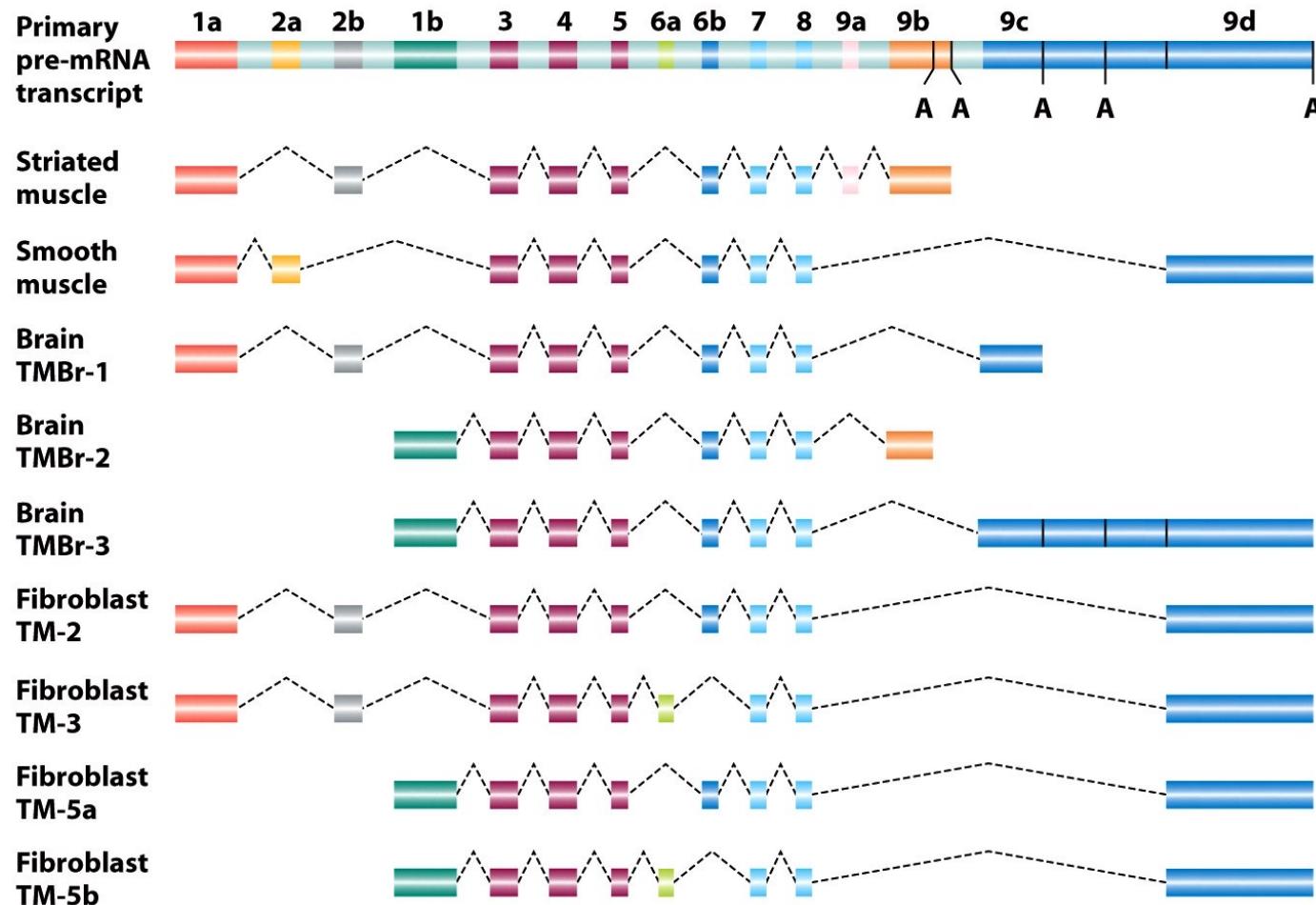
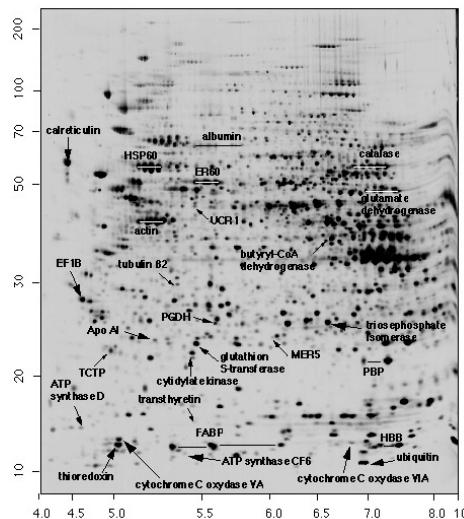


Figure 8-14
Introduction to Genetic Analysis, Ninth Edition
© 2008 W.H. Freeman and Company

Protein Expression/Sequence

Data

- MW-Isoelectric point
- MW
- Sequence/spans

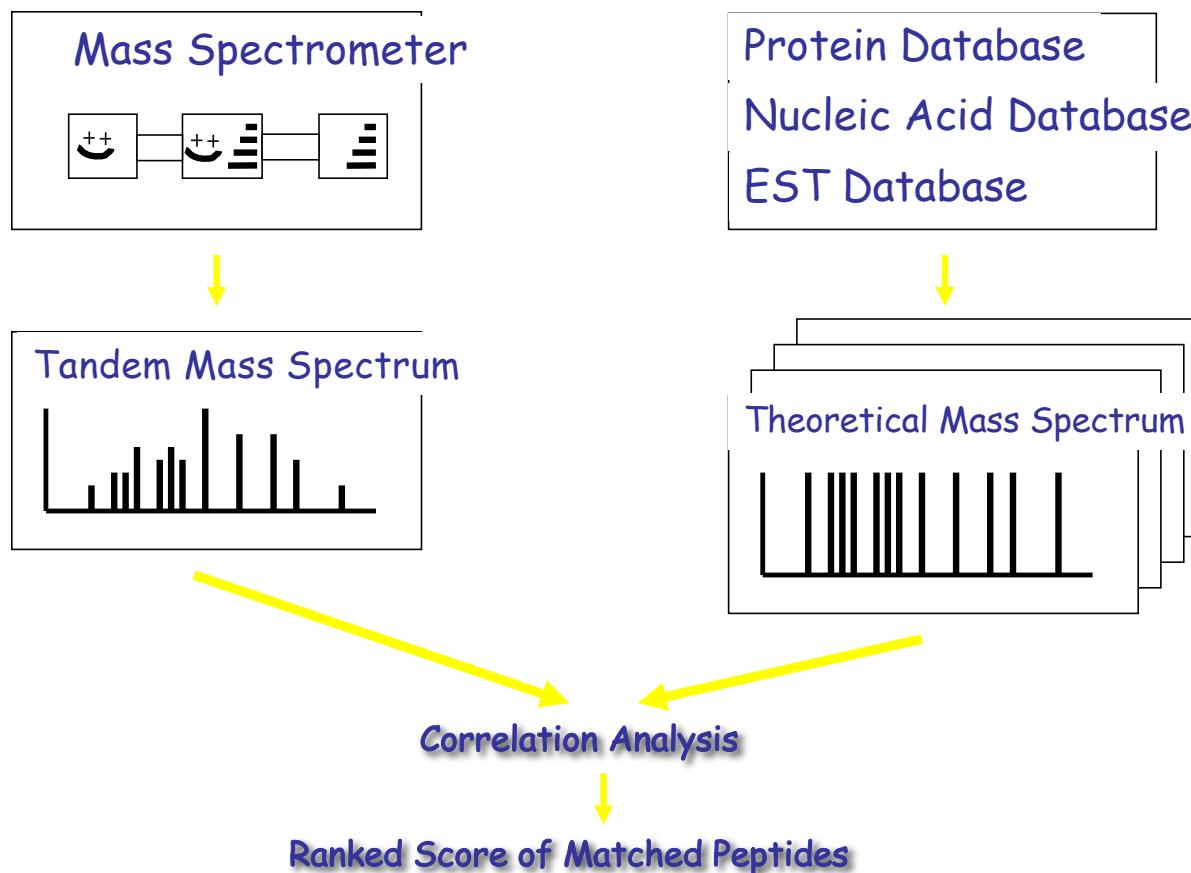


Technology

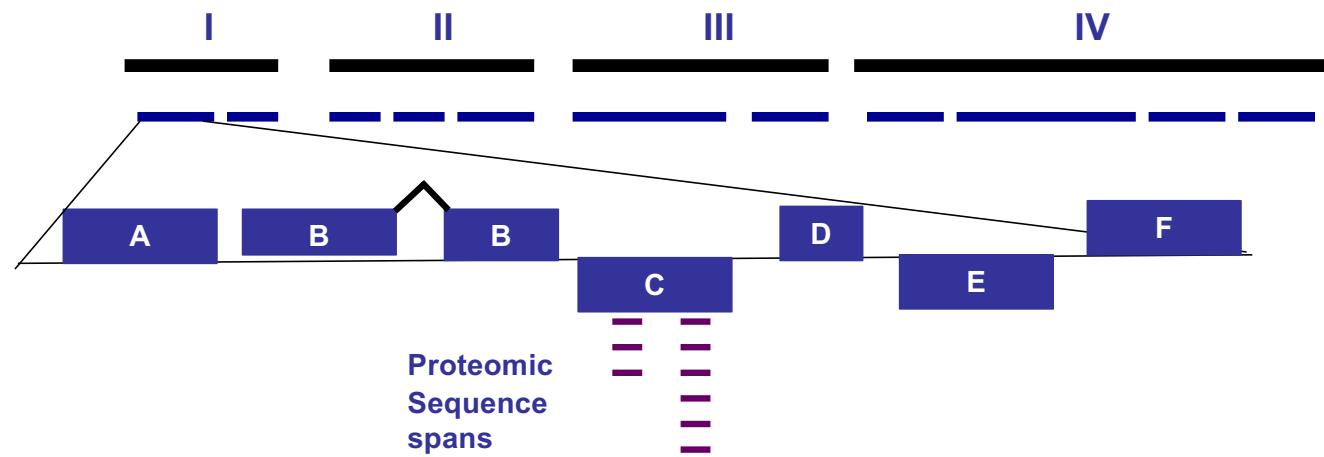
- 2D gel electrophoresis
- Mass spectrometry
- Tandem MS (MS-MS, LC MS-MS etc)

Typical 2 D gel

Sequest Database Search



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.

Overview

PubChem Compound ID: CID:93072

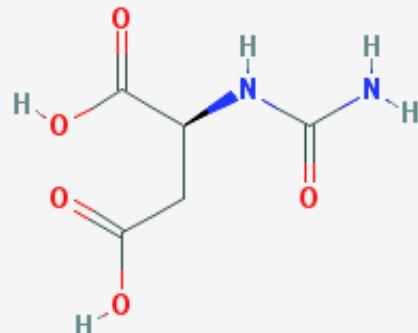
PubChem Substance ID(s): 3727

Synonyms: N-Carbamoyl-L-aspartate

Molecular Weight: 176.12742

Molecular Formula: C₅H₈N₂O₅

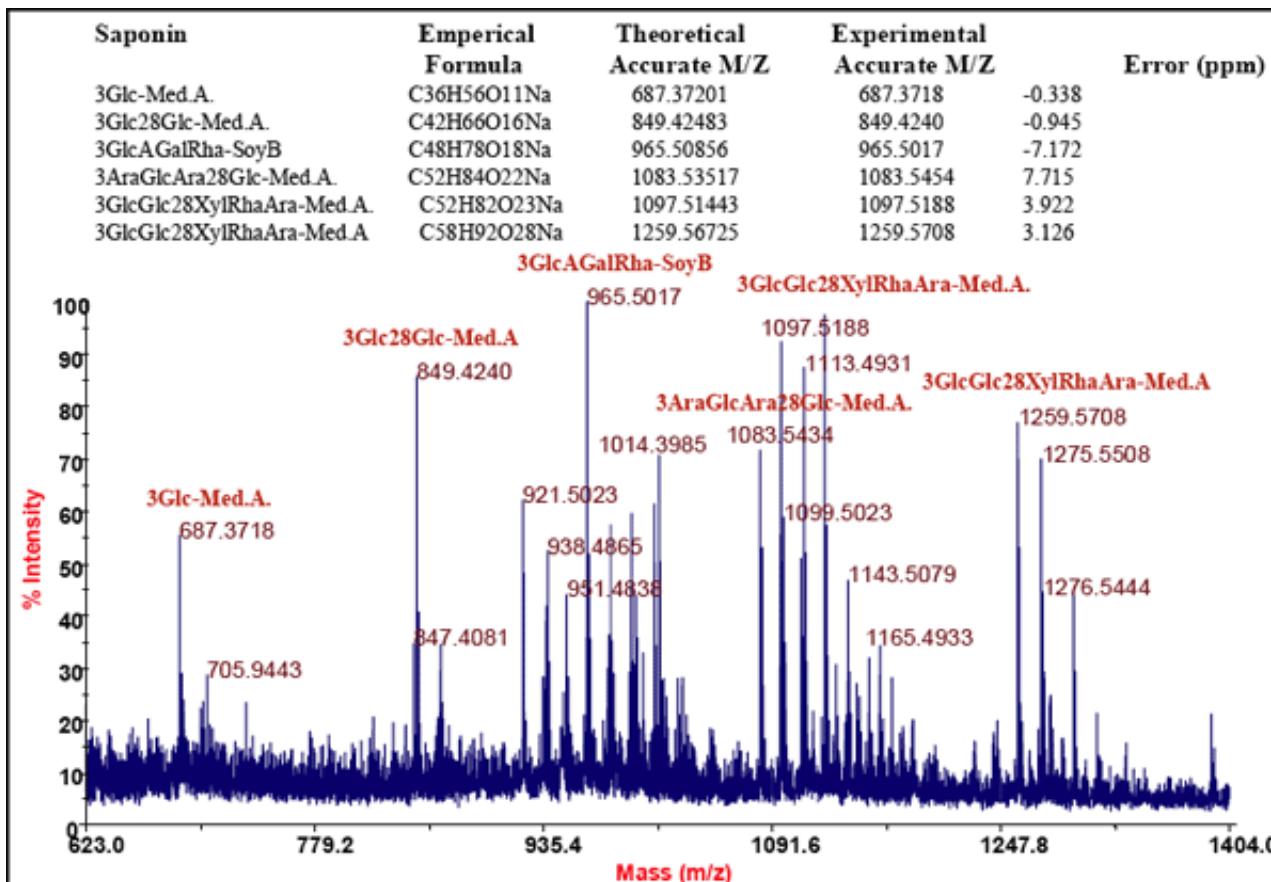
2D Structure



Metabolites

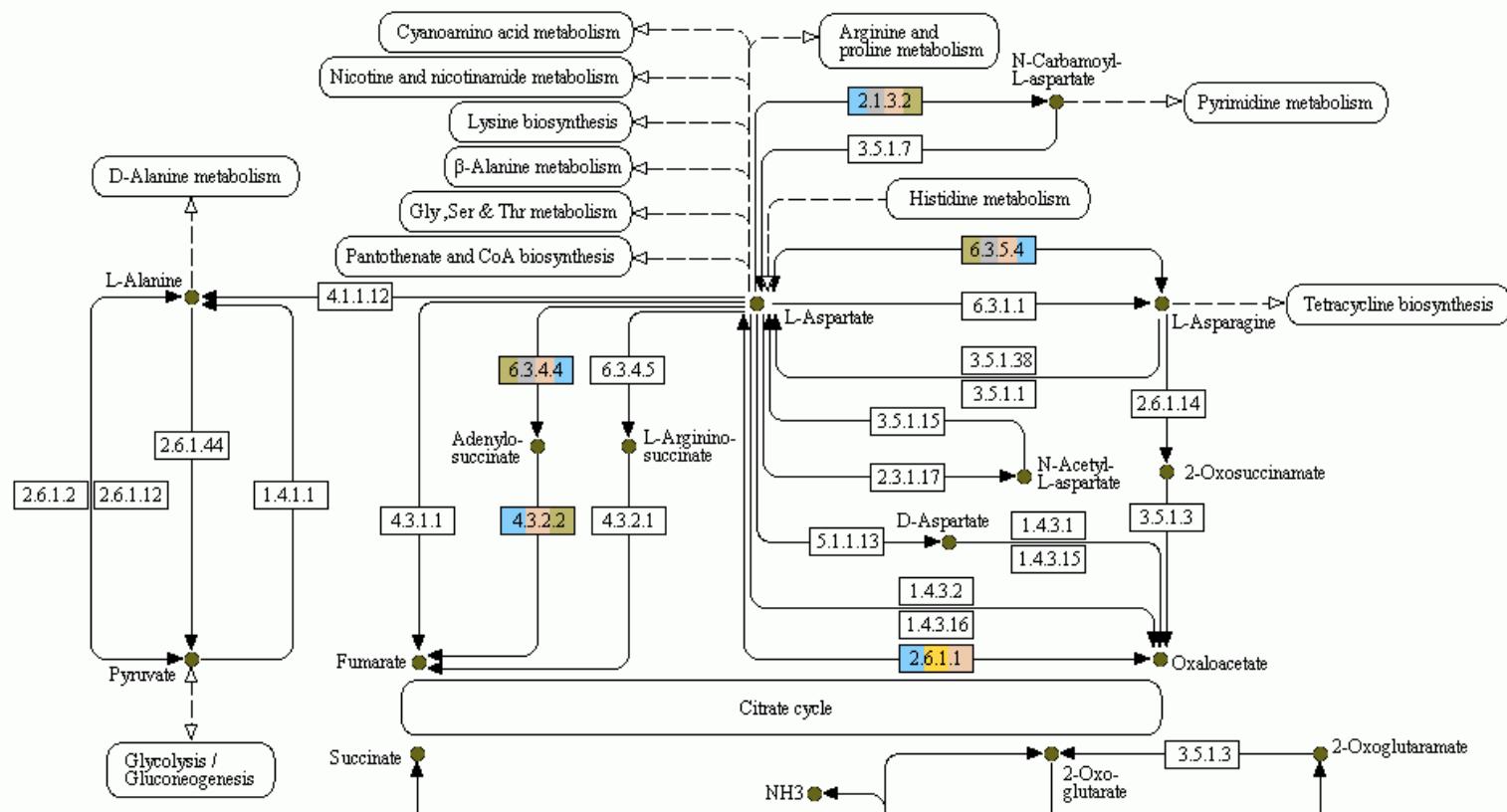
Mass Spectrometry
can be used to
measure metabolic and
other chemical
compounds

Complex mixtures can be analyzed and interpreted



Metabolites can be linked to metabolic pathways and enzymes

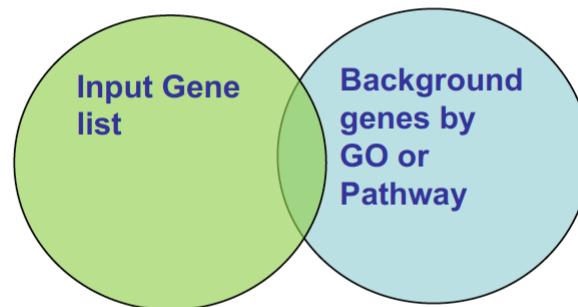
ALANINE, ASPARTATE AND GLUTAMATE METABOLISM



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

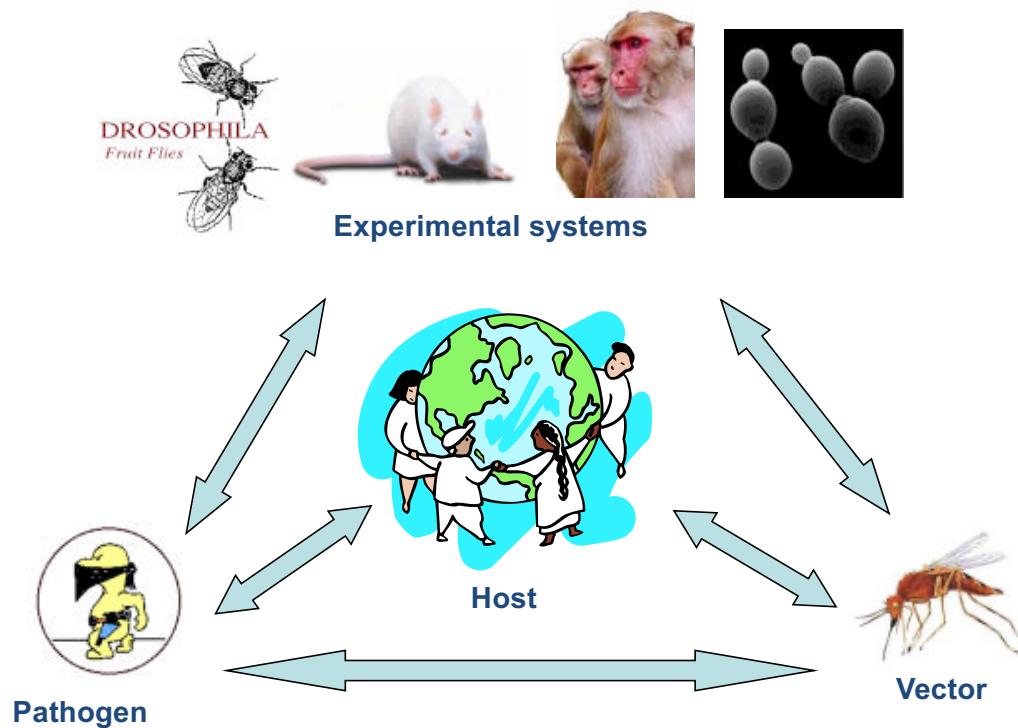
Result: GO:ID or Pathway ID that is enriched

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

Atul Butte Review: <http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002375>

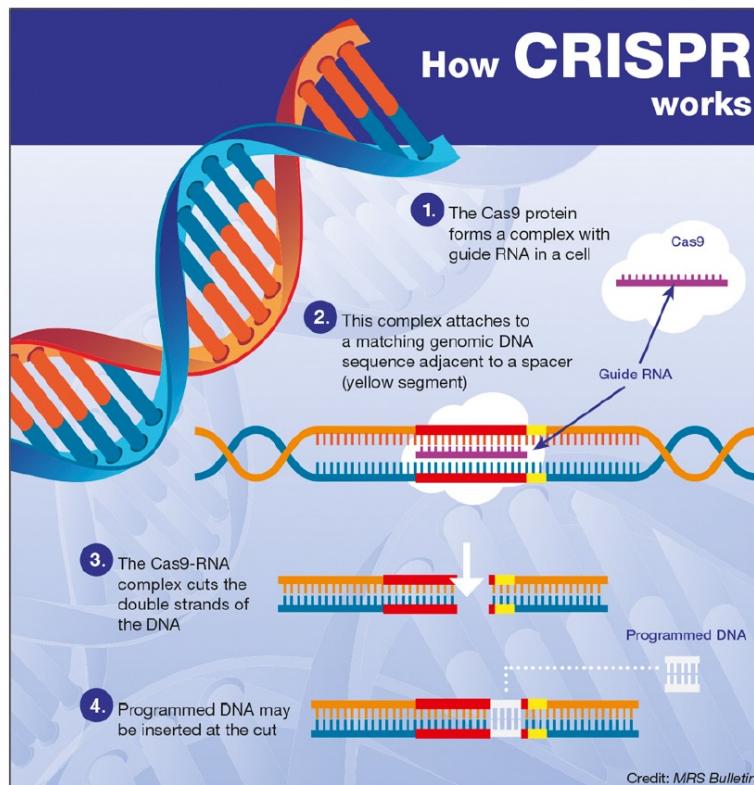
Host(s)

Infectious Disease Paradigm of Host-Pathogen Interactions



Mutant analysis

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

Ball et al., MRS Bulletin November 2016

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

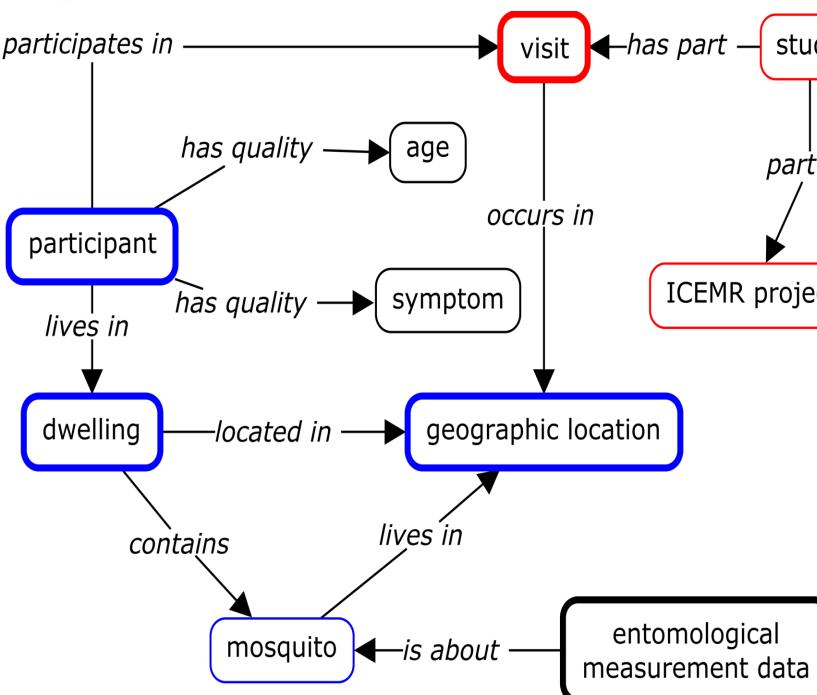
Standardized Metadata for Human Pathogen/Vector Genomic Sequences

Vivien G. Dugan^{1,2}, Scott J. Emrich³, Gloria I. Giraldo-Calderón³, On...
 Brett E. Pickett⁴, Lynn M. Schriml⁵, Timothy B. Stockwell¹, Christian Indresh Singh¹, Doyle V. Ward⁵, Alison Yao², Jie Zheng⁴, Tanya Barr Vincent M. Bruno⁶, Elizabeth Cader^{10a}, Sinead Chapman⁵, Frank H. Co...
 Valentina Di Francesco², Scott Durkin¹, Mark Eppinger^{6,ob}, Michael I. Florian Fricke⁶, Maria Giovanni⁷, Matthew R. Henn^{5,nc}, Erin Hine⁶, Ju Mizrahi⁸, Jessica C. Kissinger⁹, Eun Mi Lee², Punam Mathur², Emma Cheryl I. Murphy⁵, Garry Myers⁶, Daniel E. Neafsey⁵, Karen E. Nelson¹, David Rasko⁶, David S. Roos⁴, Lisa Sadzewicz⁶, Joana C. Silva⁴, Bruce L. Stevens¹¹, Luke Tallon⁶, Hervé Tettelin⁶, David Wentworth¹, Jennifer Wortman⁵, Yun Zhang¹, Richard H. Scheuermann^{1,12a}

1. J. Craig Venter Institute, Rockville, Maryland, and La Jolla, California, United States of America, 2National Institute of Allergy and Infectious Diseases, Bethesda, Maryland, United States of America, 3University of Notre Dame, Notre Dame, Indiana, United States of America, 4University of Michigan, Ann Arbor, Michigan, United States of America, 5Broad Institute, Cambridge, Massachusetts, United States of America, 6Institute for Genomics, Proteomics and Bioinformatics, University of Maryland, Baltimore, Maryland, United States of America, 7Cyberinfrastructure Division, Virginia Bioinformatics Institute, Blacksburg, Virginia, United States of America, 8Center for Biotechnology Information, National Library of Medicine, Bethesda, Maryland, United States of America, 9National Institute of Allergy and Infectious Diseases, Bethesda, Maryland, United States of America, 10Kelly Government Solutions, Rockville, Maryland, United States of America, 11Argonne National Laboratory, Lemont, Illinois, United States of America, 12Department 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 requires a standardized approach to share data across projects.

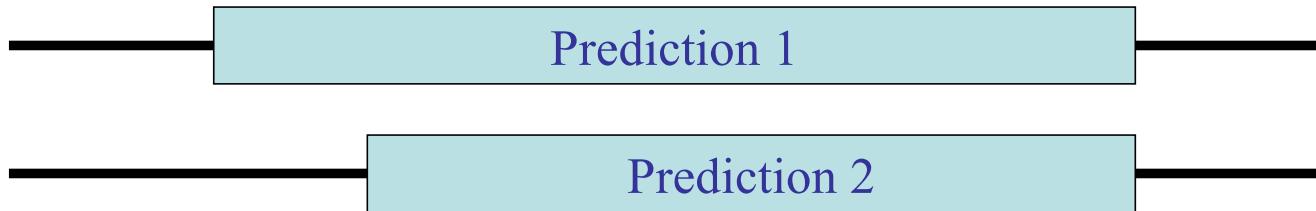


	Core Project	Core Sample	Project Specific	Pathogen Specific	Sequencing Assay
Investigation	█				
lost Characterization		█			
Specimen Isolation			█		
Pathogen Characterization				█	
Specimen Processing					█
Pathogen Detection			█		
Pathogen Isolation				█	
Sample Management					█
Data Transformation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

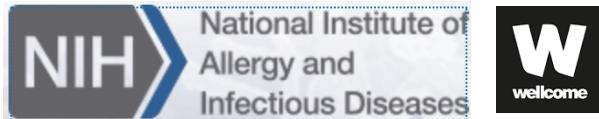
*Bioinformatics Resource Center
Community Evolution*

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.



Project Leadership:
David Roos – UPENN
Mary Ann McDowell – Notre Dame
Andrew Jones – Liverpool
Jessie Kissinger – UGA
Sarah Dyer – EBI
Kathryn Crough – Glasgow
George Christophides - Imperial

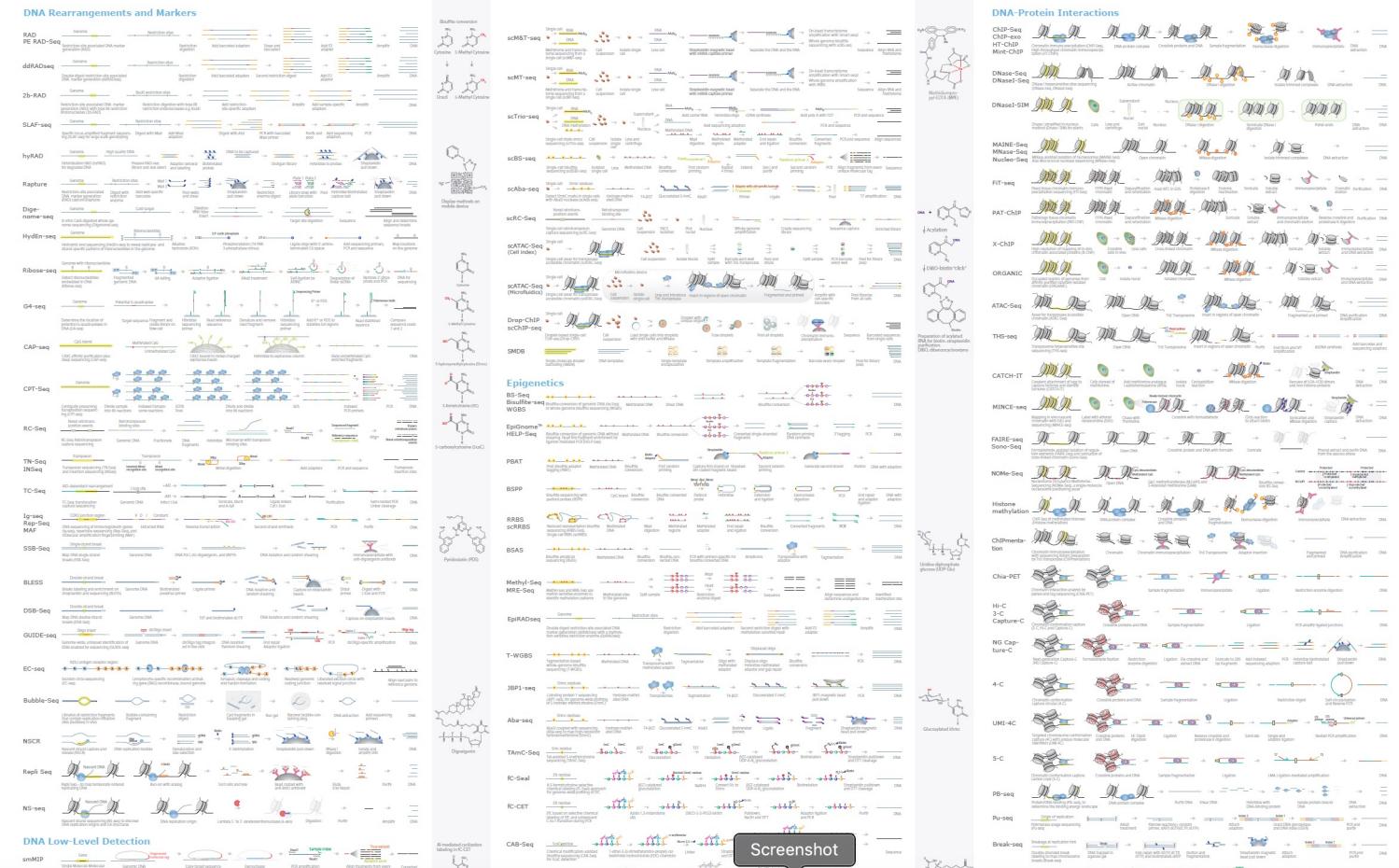


Our goal: enabling end users in the lab, field & clinic to make effective and appropriate use of large-scale datasets, expediting discovery research and translational application by making data FAIR



DNA

For all you seq...



RNA

