# 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

1

1

---

**KAYAK** — Round-trip · Atlanta × + ⇄ Tunis × + · Wed 11/29 – Wed 12/6 · 1 adult, Economy 🔍 · ♥ · ⊕ Sign in · 🇺🇸 $

**Our Advice**
¯\\_(ツ)_/¯
We're still gathering data for this route

Track prices — Off

**Cheapest** $894 · 18h 37m
**Best** ⓘ $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.

**224** of 633 flights

**Stops**
- Nonstop
- ✓ 1 stop — $1,117
- ✓ 2+ stops — $894

**Fee Assistant** ⓘ
- Carry-on bag — – 0 +
- Checked bag — – 0 +

**Book on KAYAK** ⚡
Show offers instantly bookable on KAYAK.

**Times**
Take-off | Landing
Take-off from ATL

ITA Airways — 10:10 am – 10:35 am⁺¹ — 2 stops — IAD, FCO — 18h 25m — ATL-TUN — 🧳0 🛍1 — **$904** — Economy — Priceline
ITA Airways — 11:25 am – 12:14 am⁺¹ — 2 stops — FCO, JFK — 18h 49m — TUN-ATL — **View Deal** — Ad
Operated by Delta Air Lines

**Best**
Delta — 3:40 pm – 9:50 am⁺¹ — 1 stop — CDG — 12h 10m — ATL-TUN — 🧳1 🛍0 — **$1,121** — Basic Economy — Delta
Delta — 5:30 am – 1:35 pm — 1 stop — CDG — 14h 05m — TUN-ATL — **View Deal** — Main Cabin $1,301
Operated by Air France

**Cheapest**
ITA Airways — 10:10 am – 10:35 am⁺¹ — 2 stops — IAD, FCO — 18h 25m — ATL-TUN — 🧳1 🛍0 — **$894** — Economy — ScholarTrip
ITA Airways — 11:25 am – 12:14 am⁺¹ — 2 stops — FCO, JFK — 18h 49m — TUN-ATL — **View Deal**

2

# The Travel Site has Very Useful Data Filters!

---

# 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



BOOLEAN LOGIC

AND — Both terms   OR — Either term   NOT — Only one term

shutterstock.com · 2548907405

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

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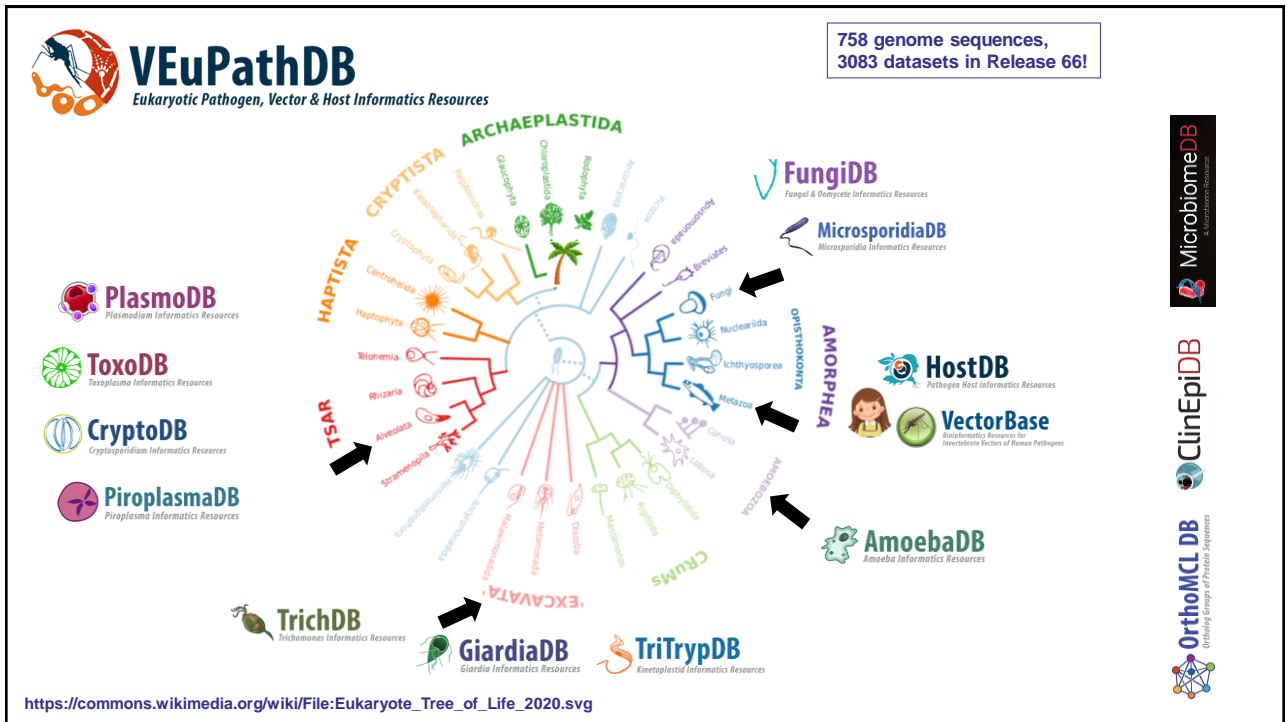  - That are near centromeres
  - That encode a predicted signal protein
  - That encode the amino acid motif CC..CC
- Which have evidence of expression …
  - In developmental stage X
  - After treatment with drug Y
- That are phosphorylated in proteomic studies
- That show evidence of diversifying selection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

5

## 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

6

**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



Label
@FORJUSP02AJWD1

Sequence
CCGTCAATTCATTTAAGTTTTAACCTTGCGGCCGTACTCCCAGGCGGT
+
AAAAAAAAAAAA::99@::::??@::FFAAAAACCAA::::BB@@?A?

Base = T, Q = A = 25

Q Scores (as ASCII charts)

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                    R2 Q score = 30.69

- 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

## 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

11

---

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



DNA Chromosome e.g. "truth"

Paired-End Reads, Note the direction of the arrows

Contig 1        Contig 2        Contig 3

Assembled Contigs

Sequence Gap

Scaffold

Assembled Scaffolds

Physical Gap

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

12

# Paired-End Reads can Yield Order & Orientation of Contigs

**Paired-end reads**

**Contigs**

Scaffold

35 bp identified  330 - 430 bp unknown sequence  35 bp identified

13

---

**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 gaps 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**

Long reads
Contiging
+ Purging

Linked reads
Scaffolding

Opt. maps
Scaffolding

Hi-C
Scaffolding

Gap-filling &
Polishing

Draft assembly

A    TGGGGA
A    TGGGGA
A    TGGGGA
C    TGGA
A:   TGGGGA

Curation

exon 1  exon 2  exon 3

Curated assembly
Primary
Alternate

Rhie et al 2021 Nature

14

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

Reference Genome Sequence



| 35 bp identified | 330 - 430 bp unknown sequence | 35 bp identified |

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

15

---

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

**Histogram**



Coverage 30X

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

16

## Slide 17

### Hybrid Assembly = short + long reads from differing technologies

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

Illumina reads          PacBio reads

Long reads can be PacBio or Oxford Nanopore, ONT

errors

super-reads

Index of 15-mers in all super-reads

Poorly aligning super-reads          exact overlap

Merging and tiling

pre-mega-reads

Joining and creating linking mates

mega-reads and linking mates

**Figure 1**. Overview of the mega-reads algorithm. Low-error rate Illumina reads (top left) are used to build longer super-reads (green lines), which in turn are used to construct a database of all 15-mers in those reads. PacBio reads (purple lines) and super-reads are then aligned, using the 15-mer index. Inconsistent super-reads are show as kinked lines; these are discarded and the remaining super-reads are merged, using the PacBio read as a template, to produce pre-mega-reads (yellow). These are further merged to produce the final mega-reads and to generate linking mates across gaps.

**17**

## Slide 18

# 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?

**18**

# 30,000 ft View – Genome Annotation

**Chromosomes**    I    II    III    IV

**Scaffolded Contigs**

**Annotated Genes**

A   B   B   C   D   E   F

# The Genome Sequence

```
AAGCTTCGCCAGGCTGTAAATCCCGTGAGTCGTCCTCACAAATCATCAAGCAGGTGTCCTCAGGGAGACTGCCTGACTGAGTTATGCTAATTCCTTTCTACTTTGGCGTGGTCACGTGTA
ACCATATCCGAATCATTTCTCTAGCCCTACGAACAGGTAAGAGCGCTAGGGATGTCCGTGGAGTAGTGTGCTTACTCGATAATATTCAGTTGGGACTACCAGCGAGGCGCTCGCTTTGCT
CACGCAATGCCTGAGACAGTTGCAGAATGAATGGTAACCGACAAACGCGTTCATATGCGTTTTCAAACTTAGTAGACGCGTACTGTCTGAAACTGGCGGTCACAGGCACCAGATAACGCC
CTTGGCATCGGCATGTCTCGTACAGAGGTCCGTATGTAGTGCCACGACTTCTAAATCCGGCGACAGGCTGGTCTTTTGTCTTACCACGTATTAGCCCGCGTGCGATTTCTCGGAGCGCAC
CTGTTCAACACTAGAAAACGGAGTTTCCTGATCGAGAAGCCACCACCTTTCCAGAAGTTGAACGCTAGCATGTCATTCGATTTTCACCCCCGCGTAGTTCCTGTGTGTCATTCGTTGTC
GAGACAACTCTGTCCCGCCCGGTGCTGTTCCATATGCGTGACTTTCCCGCAATTTTTTCAGACTTTCAGGAAAGACAGGCTCCGGAACGATCTCGTCCATGACTGGTAAATCCACGACA
CCGCAATGGCCCCCAGCACCTCTATCTCTCGTGCCAGGGGACTAACGTTGTATGCGTCTGCGTCTTGTCTTTTTGCATTCGCTTTCCAAAAAAGAGAGCCATCCGTTCCCCCGCACATTC
AACGCCGCGAGTGCGGTTTTTGTCTTTTTGAGTGGTAGGACGCTTTTCATGCGCGAACTACGTGGACATTAAGTTCCATTCTCTTTTTCGACAGCACGAAACCTTGCATTCAAACCCGC
CCGCGGAAGATCCGATCTTGCTGCTGTTCGCAGTCCCAGTAGCGTCCTGTCGGCCGCGCCGTCTCTGTTGGTGGGCAGCCGCTACACCTGTTATCTGACTGCCGTGCGCGAAAATGACGC
CATTTTTGGGAAAATCGGGGAACTTCATTCTTTAAAAGTATGCGGAGGTTTCCTTTTTCTTCTGTTCGTTTCTTTTTCTCGGGTTTGATAACCGTGTTCGATGTAAGCACTTTCCGTCTC
TCCTTCCGTGCTTTGTTTCGACATCGAGACCAGGTGTGCAGATCCTTCGCTTGTCGATCCGGAGACGCGTGTCTCGTAGAACCTTTTCATTTTACCACACGGCAGTGCGGAGCACTGCTCTG
AGTGCAGCAGGGACGGGTGAAGTTTCGCTTTAGTAGTGCGTTTCTGCTCTACGGGGCGTTGTCGTGTCTGGGAAGATGCAGAAACCGGTGTGTCTGGTCGTCGCGATGACCCCCAAGAGG
GGCATCGGCATCAACAACGGCCTCCCGTGGCCCCACTTGACCACAGATTTCAAACACTTTTCTCGTGTGACAAAAACGACGCCCGAAGAAGCCAGTCGGCCTGAACGGGTGGCTTCCCAGG
AAATTTGCAAAGACGGGCGACTCTGGACTTCCCTCTCCATCAGTCGGCAAGAGATTCAACGCCGTTGTCATGGGACGGAAAACCTGGGAAAGCATGCCTCGAAAGTTTAGACCCCTCGTG
GACAGATTGAACATCGTCGTTTCCTCTTCCCTGTGAGCACACACTAGTAGTCGCCACACGCTGTTTGAGACGTGTCAATCTCCAAGAGTGTGGACGCTGTTCCACGTCTTCAAATGTTTCC
CAACATCCGTCGTCTAGTAGACACACCAACAAAAAGCACACGGCGAATCTGCTCATCGGAGGGAGGAGCCGGGGGGCACACAACTATCCTCAACTCTCGAACGAACATATCCGGGGCCGC
GAAGACGTCCAGTCTCTCAAATCCAAATCCAAACCGGAACGCAAACATTTCTGCATCAAGTCACGATTGCGCCGGTACCTCCATGTGTAAGCAGTTCCATGAAACCTCCGATATTACACACGACTG
TGGATATGAATTATATGCAGATGCATATATACTGAGACGCCGATGCAACTATAGGTTTCCTGGCCCTCCATGGATATTTCAGACCTTCCTCTCACATTTGGTTTGCCCGTACACCTCCGT
TACGCTTTTTTTCTGGCTTTCTTCTTCGTCTCTGTTTATCAGCAAAGAAGAAGACATTGCGGCGGAGAAGCCTCAAGCTGAAGGCCAGCAGCGCGTCCGAGTCTGTGCTTCACTCCCAGC
AGCTCTCAGCCTTCTGGAGGAAGAGTACAAGGATTCTGTCGACCAGATTTTTGTCGTGGGTATGTTGTCCTAAACTCCTTGGAACTCCATTCTTGGTCAGAAACGTACTGAAACTGTATA
CATGTATATACAGATGTATGGATAATATCTAGAGAAGATACAGGGAAGACTGGCAAGGATGAAAGACATGCAAGTTTAACGAAGCAGAGGGCATTGGCGAGAGGGACGCCCGTTATGCT
GTGTGATGTGGCTGTGAATCTTACCTCGCCGTTTGACTTGCTGCAGCGCTTTGTCCACTTGAACGTGACTTCTTGTTTCTACCTTCCCCAACGCCTTCTATTCCCTTCACTGCGAAAGCG
CGCTCAGTGGGCCGTCACCGAACACCCTTGGTTCTTTCGTTCAGCTGTTGTCCTCTTTCTCGTGGCGTCGTGGCTCGGTGGCTCGGCTTCTCTCTCTTTCCTGTTGGTGCGTCCAG
ACTATGTCGCCTGTTTCCCCACCCTTCTCGGCTTGTGCTTTCAGGAGGAGCGGGACTGTACGAGGCAGCGCTGTCTCTGGGCGTTGCCTCTCACCTGTACATCACGCGTGTAGCCCGCGA
GTTTCCGTGCGACGTTTTCTTCCCTGCGTTCCCCGGAGATGACATTCTTTCAAACAAATCAACTGCTGCGGCAGGCTGCAGCTCCTGCCGAGTCTGTGTTCGTTCCCTTTTGTCCGGAGCT
CGGAAGAGAGAAGGACAATGAAGCGACGTATCGACCCATCTTCATTTCCAAGACCTTCTCAGACAACGGGGTACCCTACGACTTTGTGGTTCTCGAGAAGAGAAGGAAGACTGACGACGC
AGCCACTGCGGAACCGGTAAGAGGCAACCGAAGCGCGTAGATAAGAAAAACAACAAAGAGAAGGTGAAACACGAAGAGAAGGGAAAATGCGGAGAAACCGTGGATTTACAAAGATATCAA
GAGCAATGCTTTGTGGAGATTTTTTTTAATTCAGTAGAGACACCCGCCGTGCGAGGTGTGTAGAAATAACTGCGACCCTGGAGACAGAGATGCCGCGAGTACACCACTTGTCGTTTTTCC
TCCTATGTTCATGACGGGTGCTGAACGTCTATCGTACTTAATTGGAGGAGTCGTCTCCGAAGCAGCTTTGGCTGGCCATCCGTGTGTTTGCCTTGTTCCTGAAAAGCCAGAAGGCGCTCC
ACAGTGAGGCGATATACAGGGACGCCTACCGGAGCCCCGTTTTCTGCCTTTGTCGACTCTTGCAGAGCAACGCAATGAGCTCCTTGACGTCCACGAGGGAGACAACTCCCGTGCACGGGT
TGCAGGCTCCTTCTTCGGCCGCAGCCATTGCCCCGGTGTTGGCGTGGATGGACGAAGAAGACCGGAAAAAACGCGAGCAAAAGGAACTGATTCGGGCCGTTCCGCATGTTCACTTTAGAG
GCCATGAAGAATTCCAGTACCTTGATCTCATTGCCGACATTATTAACAATGGAAGGACAATGGATGACCGAACGGGTAACGGCGACTGCGAGAAAAAGCCACACCGTTTTCTCCTGTGAT
TCTGTCCGCAAGCCCTCTTTTGCTTCATCCACCCTTTGCTATTCTCCGCCGCCTTCCTTTTCTGCTCCATGTTCAATTCGTTCGCTTCTTCAGTCTTTCCATCTTCCCCTGTTACCTCTG
TCATTCGTTTTCTTGCCTCTATTTAACTGTGTTCTACTCACAGTCTGCATTCCGCGATAGACGAGCTTCCACGTCTTGCGTCTCGACAAGCAACTGTCATTTGTACGCGCCTCCCTCCAC
CGTGAATCGGATTGTCGGTTCGGTTCCTGGGTCAGAAAAGGCCTGCGCCAGTATTCTGAATAATACCCTTCGCCATTGTAAAGAGGCGAAGGAACAAAGAGATATTTCGGCGCATCT
TTTGTGCGGCGCGTTTCCTCGTGCTTCACACCGATGCCCTTCTGTGCATGTCTTCTGCTCCTCGTCCTTCTCTCTTTTTCCCTGTTTAGGCGTTGGTGTCATCTCCAAATTCGGCTGCAC
TATGCGCTACTCGCTGGATCAGGCCTTTCCACTTCTCACCACAAAGCGTGTGTTCTGGAAAGGGTAAGGGCGTCTTCAGTGAATGCATATATTTGACTTCAGACATTCTTAACTGTTTGA
CAACCAACGTACAAATTTGTTTGTCCGTGTGCGTGTTCGACATGTCAAGTATGTGAAGAGTCGCTACTGTAGACTAACGCACGAACCAGATTTGTTTATCTGCATGCGCTGTGCACCCGT
TTCTGAGTGTCTGGAGTTTCCGCAACCTTCCTTTGAATTTCTGGGTTCGTTTTTTTTTTATGCGCGCACTGGTTTGCATGTGGCCTGAGAGAGCACAGATCGAAGGTGGGGTGATGTGGCGTC
GCTGCAGAGAAACTCCGGCGAAGGCGACAGATAAAGGAGAGTGGAAATCATTGAACAGTGTCGGTCGTCGTCTGTTGGTTTCGCAGGGTCCTCGAAGAGTTGCTGTGGTTCATTCGCGGCGACA
CGAACGCAAACCATCTTTCTGAGAAGGGCGTGAAGGCAAGTCTACGTTGTACCTCTTGTCTCTGCCGAAGCTCAGATGTCTCCACGGCGTTGGTTTCTTTTCGTTTTTTGCTTTCGTGGCA
TTACCATCGAGTCACCACTCATAGTTGCGTGTGTCTACATGTTTTCTAGAACGTCCGTTGTGTTGCCTCGTGGCGACCGGCGCGAGTGTATGTACCCTGCGCTGTCAGAAGTTGATCCTT
```

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)



Gene 2    Gene 3

Gene 1

Template strand
for gene 1

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

Nontemplate
strand 5' — CTGCCATTGTCAGACATGTATACCCCGTACGTCTTCCCGAGCGAAAACGATCTGCGCTGC — 3'  }  DNA
Template
strand 3' — GACGGTAACAGTCTGTACATATGGGGCATGCAGAAGGGCTCGCTTTTGCTAGACGCGACG — 5'

5' — CUGCCAUUGUCAGACAUGUAUACCCCGUACGUCUUCCCGAGCGAAAACGAUCUGCGCUGC — 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



**ORFs ≠ Genes – but they can be part of a gene**

**The "Coding Sequence" - CDS**

**Green = UTRs**   **Red = CDS**   **Pink = Intron**

>Translation Frame 1   **The Protein**

MQKPVCLVVAMTPKRGIGINNGLPWPHLTTDFKHFSRVTKTTPEEASRLN

GWLPRKFAKTGDSGLPSPSVGKRFNAVVMGRKTWESMPRKFRPLVDRLNI

VVSSSLKEEDIAAEKPQAEGQQRVRVCASLPAALSLLEEEYKDSVDQIFV

VGGAGLYEAALSLGVASHLYITRVAREFPCDVFFPAFPGDDILSNKSTAA

QAAAPAESVFVPFCPELGREKDNEATYRPIFISKTFSDNGVPYDFVVLEK

RRKTDDAATAEPSNAMSSLTSTRETTPVHGLQAPSSAAAIAPVLAWMDEE

DRKKREQKELIRAVPHVHFRGHEEFQYLDLIADIINNGRTMDDRT

---

# Terminology

**transcriptional start**  **ATG**  **stop codon**  **polyA**

5' UTR   3' UTR

exon   intron   exon

**CDS:**
(coding sequence nt)

**protein:**
(aa)

**transcript:**
(CDS + UTRs, if avail.)

**genomic:**
(includes introns)

# Evolution  Homologous chromsomes (in a diploid)



Chromosome

Locus

Allele

Homozygous Alleles

Heterozygous Alleles

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



Slide 27



Alleles have frequencies in different populations

Slide 28

# 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

Percent of population that has the O blood type

- 50-60
- 60-70
- 70-80
- 80-90
- 90-100

# 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

**Early Globin Gene**

Gene Duplication

**α-chain gene**   **β-chain gene**

**α frog**  **α chick**  **α mouse**   **β mouse**  **β chick**  **β frog**

PARALOGS

ORTHOLOGS          ORTHOLOGS

HOMOLOGS

31

# 30,000 ft View - Synteny

I    II    III    IV

| Species 1 | A | B | B | | D | | F |
| Species 2 | A | B | B | | D | | F |

C    E

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

32

## Synteny among *Plasmodium* species



33

## 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.

34

# 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 <u>Zinc Fingers</u> 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
- 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



**Annotation of genome features**

**RNA-Seq reads**

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

39

# RNA-seq data are very powerful



http://bioinfo.vanderbilt.edu/vangard/services-rnaseq.html

40

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

41

---

**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

42

# 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

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

43

---

# Protein Expression/Sequence

**Data**
- MW-Isoelectric point
- MW
- Sequence/spans

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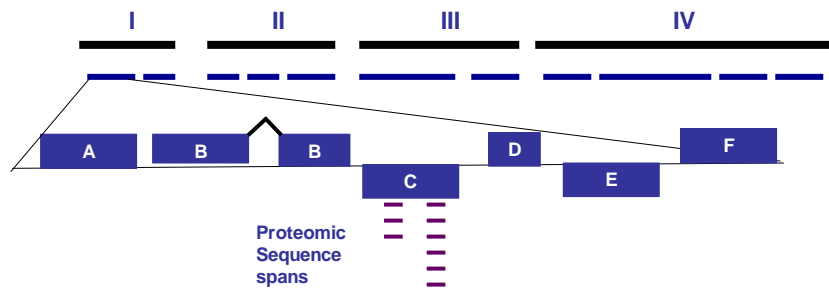- Tandem MS (MS-MS, LC MS-MS etc)



**Typical 2 D gel**

44

# Sequest Database Search



Mass Spectrometer

Protein Database
Nucleic Acid Database
EST Database

Tandem Mass Spectrum

Theoretical Mass Spectrum

*Correlation Analysis*

*Ranked Score of Matched Peptides*

45

# 30,000 ft View - Proteomics



I     II     III     IV

A    B    B    D    F

C

E

**Proteomic Sequence spans**

**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.**

46

# Metabolites

## Overview

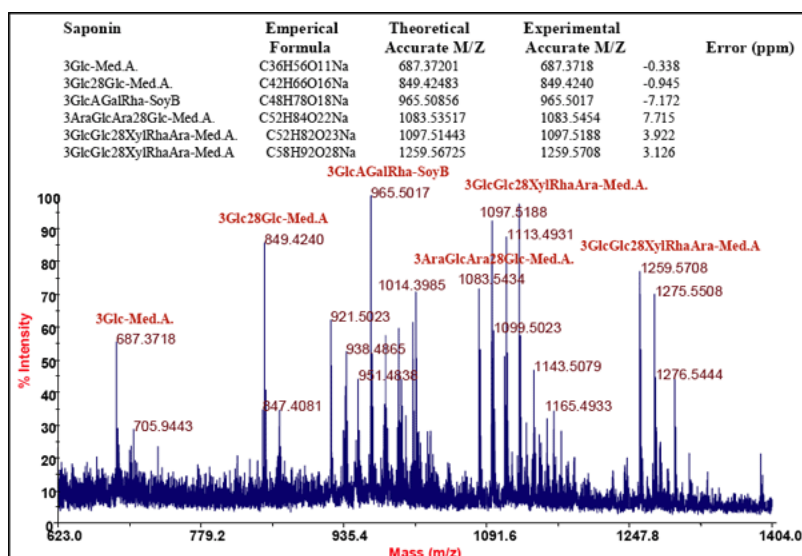| | |
|---|---|
| **PubChem Compound ID:** | CID:93072 |
| **PubChem Substance ID(s):** | 3727 |
| **Synonyms:** | N-Carbamoyl-L-aspartate |
| **Molecular Weight:** | 176.12742 |
| **Molecular Formula:** | $C_5H_8N_2O_5$ |

### 2D Structure



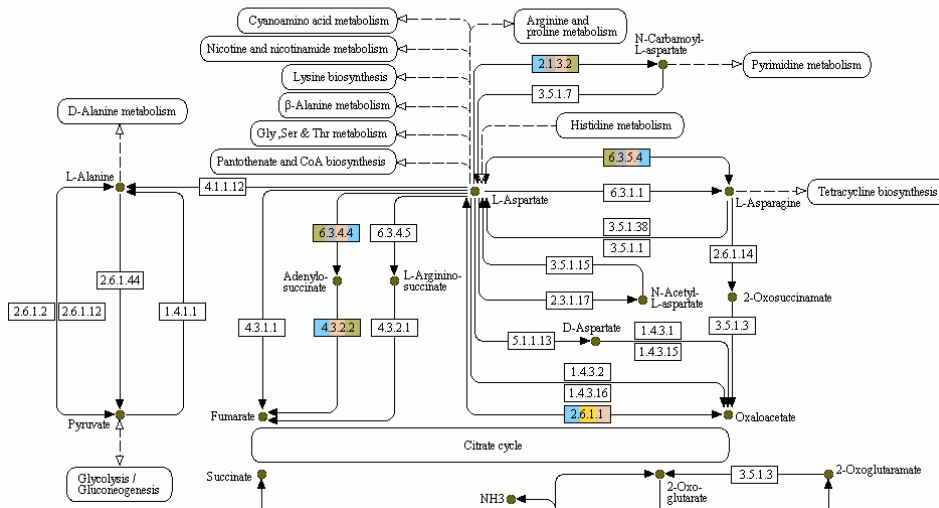**Mass Spectrometry can be used to measure metabolic and other chemical compounds**

47

# Complex mixtures can be analyzed and interpreted

| Saponin | Emperical Formula | Theoretical Accurate M/Z | Experimental Accurate M/Z | Error (ppm) |
|---|---|---|---|---|
| 3Glc-Med.A. | C36H56O11Na | 687.37201 | 687.3718 | -0.338 |
| 3Glc28Glc-Med.A. | C42H66O16Na | 849.42483 | 849.4240 | -0.945 |
| 3GlcAGalRha-SoyB | C48H78O18Na | 965.50856 | 965.5017 | -7.172 |
| 3AraGlcAra28Glc-Med.A. | C52H84O22Na | 1083.53517 | 1083.5454 | 7.715 |
| 3GlcGlc28XylRhaAra-Med.A. | C52H82O23Na | 1097.51443 | 1097.5188 | 3.922 |
| 3GlcGlc28XylRhaAra-Med.A | C58H92O28Na | 1259.56725 | 1259.5708 | 3.126 |



48

## Metabolites can be linked to metabolic pathways and enzymes
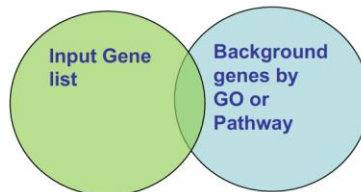
# Gene & Pathway Enrichment

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

Input Gene list

Background genes by GO or Pathway

Result:  GO:ID or Pathway ID that is enriched
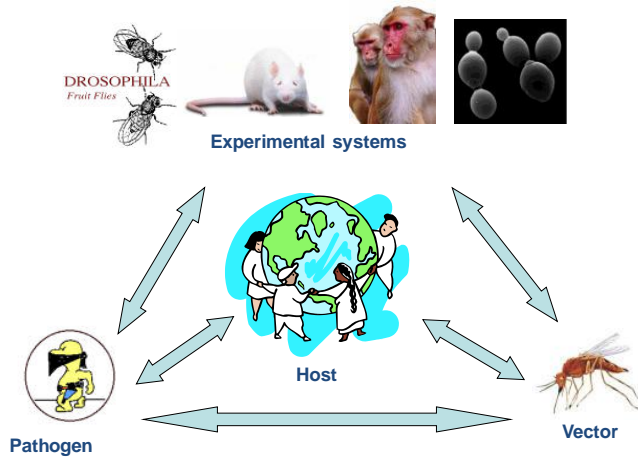
Statistics:  Are more genes observed than expected (P-value)
Multiple hypothesis testing (Bonferroni, Benjamini-Hochberg)

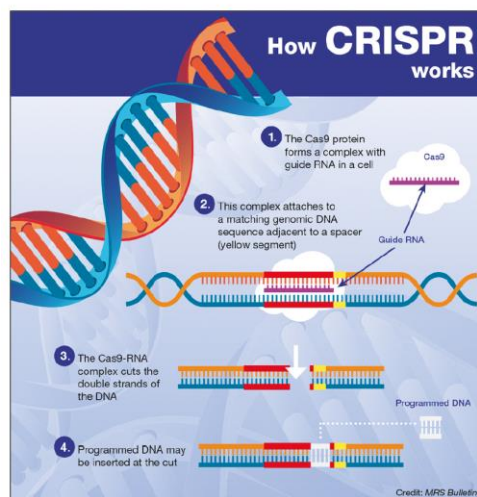Atul Butte Review:  http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002375

# Host(s)

## Infectious Disease Paradigm of Host-Pathogen Interactions

**Experimental systems**

DROSOPHILA
Fruit Flies

**Host**

**Pathogen**

**Vector**

51

# Mutant analysis

## CRISPR-CAS

**How CRISPR works**

1. The Cas9 protein forms a complex with guide RNA in a cell

Cas9

2. This complex attaches to a matching genomic DNA sequence adjacent to a spacer (yellow segment)

Guide RNA

3. The Cas9-RNA complex cuts the double strands of the DNA

Programmed DNA

4. Programmed DNA may be inserted at the cut

Credit: MRS Bulletin

**Ball et al., MRS Bulletin November 2016**

- 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
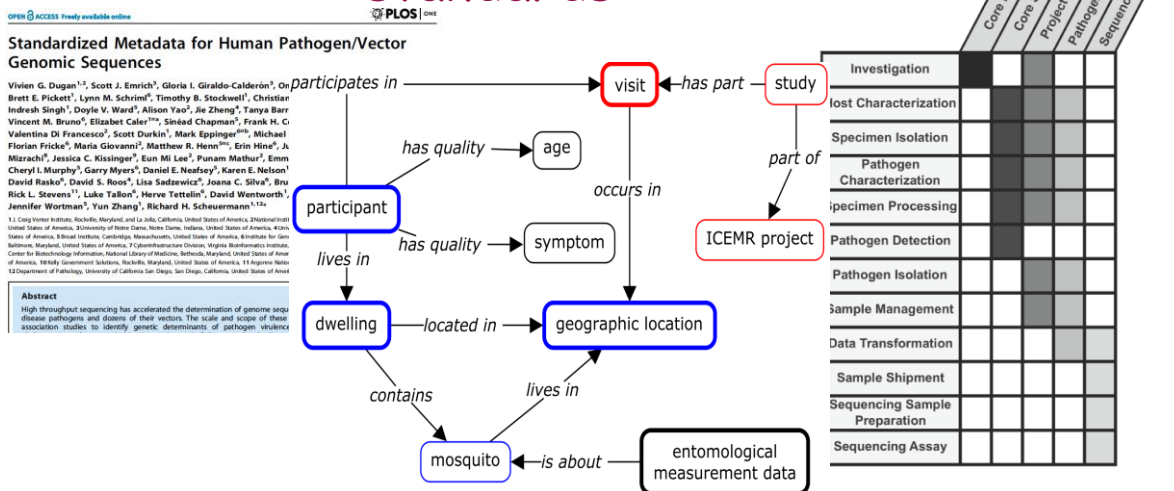
52

# Metadata – The next Frontier

- Data about the data are critical
- What makes a <u>data set</u> 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
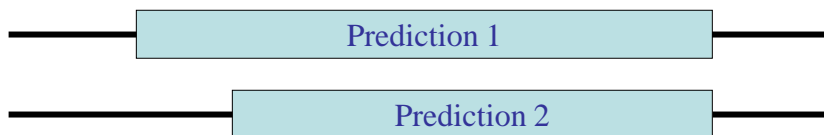  - Phenotype

53

# Data sharing standards



54

# 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.

55

# Different algorithms often generate different results

| Prediction 1 |
| Prediction 2 |

**We provide lots evidence so that you can decide or design an experiment to confirm!**
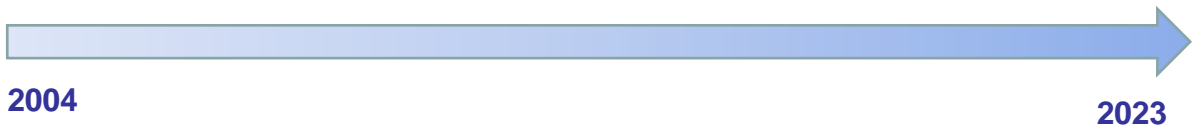
56

# 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 annotate, or discovered after the annotation e.g. lncRNAs. Interpret carefully

57

---

# *Bioinformatics Resource Center Community Evolution*

Browsing ⟹ Mining ⟹ Integrating ⟹ Facilitating

**2004**

**2023**

58

58

# The End

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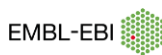- These slides are available to you as a PDF on the workshop web site.

59

**VEuPathDB**
*Eukaryotic Pathogen, Vector & Host Informatics Resources*

**NIH** National Institute of Allergy and Infectious Diseases

**W** wellcome

**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

**Penn** UNIVERSITY OF PENNSYLVANIA

UNIVERSITY OF NOTRE DAME

UNIVERSITY OF LIVERPOOL

UNIVERSITY OF GEORGIA

EMBL-EBI

University of Glasgow

Imperial College London

Thank you to the data providers, participants and community for their feedback

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

60

60