

Intro to Bioinformatics

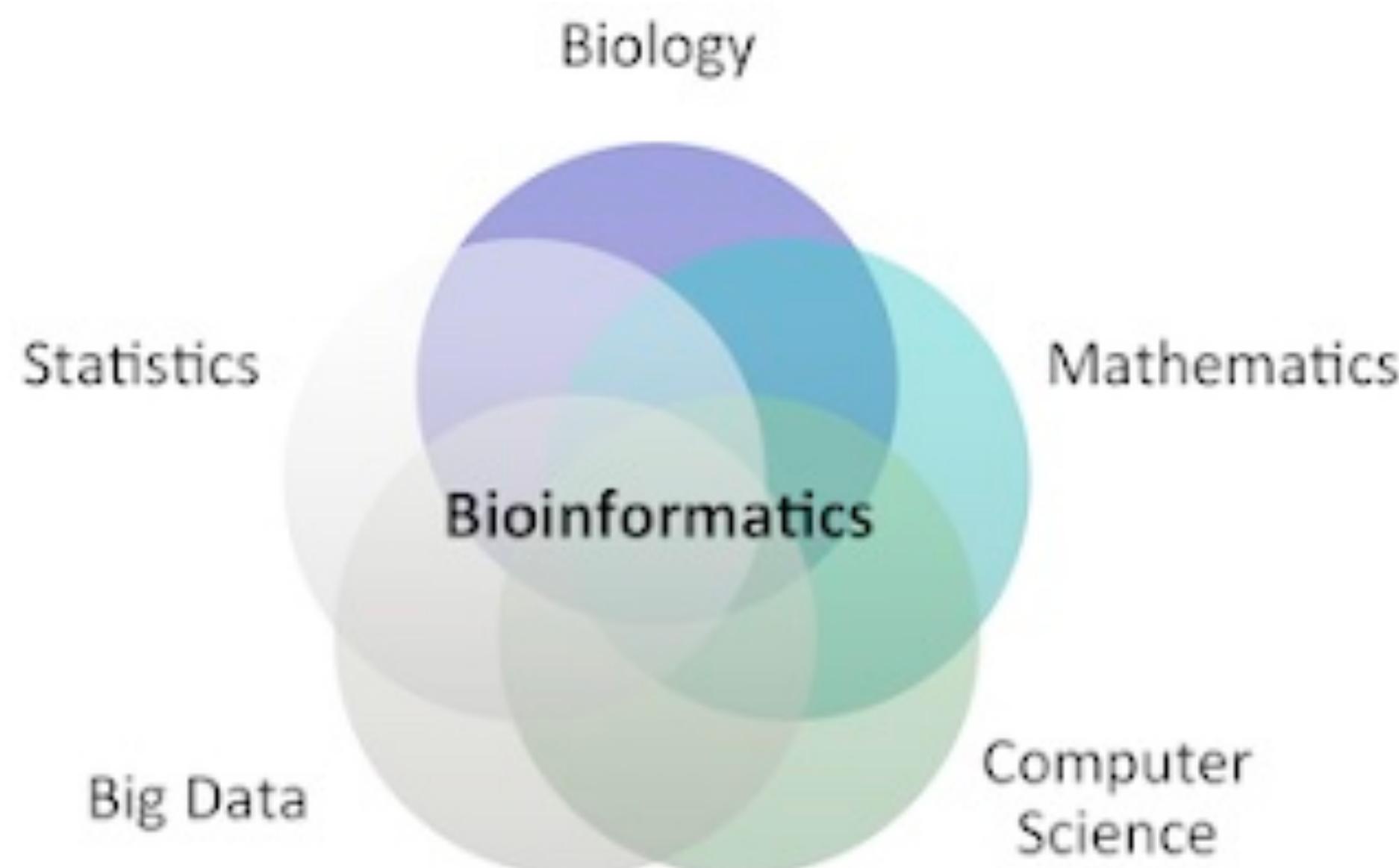
Susan Strickler



What we will cover:

1. Background
2. 'Omics
3. Databases
4. Operating systems used in bioinformatics
5. Programming languages used in bioinformatics
6. Data Best Practices
7. Computational infrastructure
8. Careers

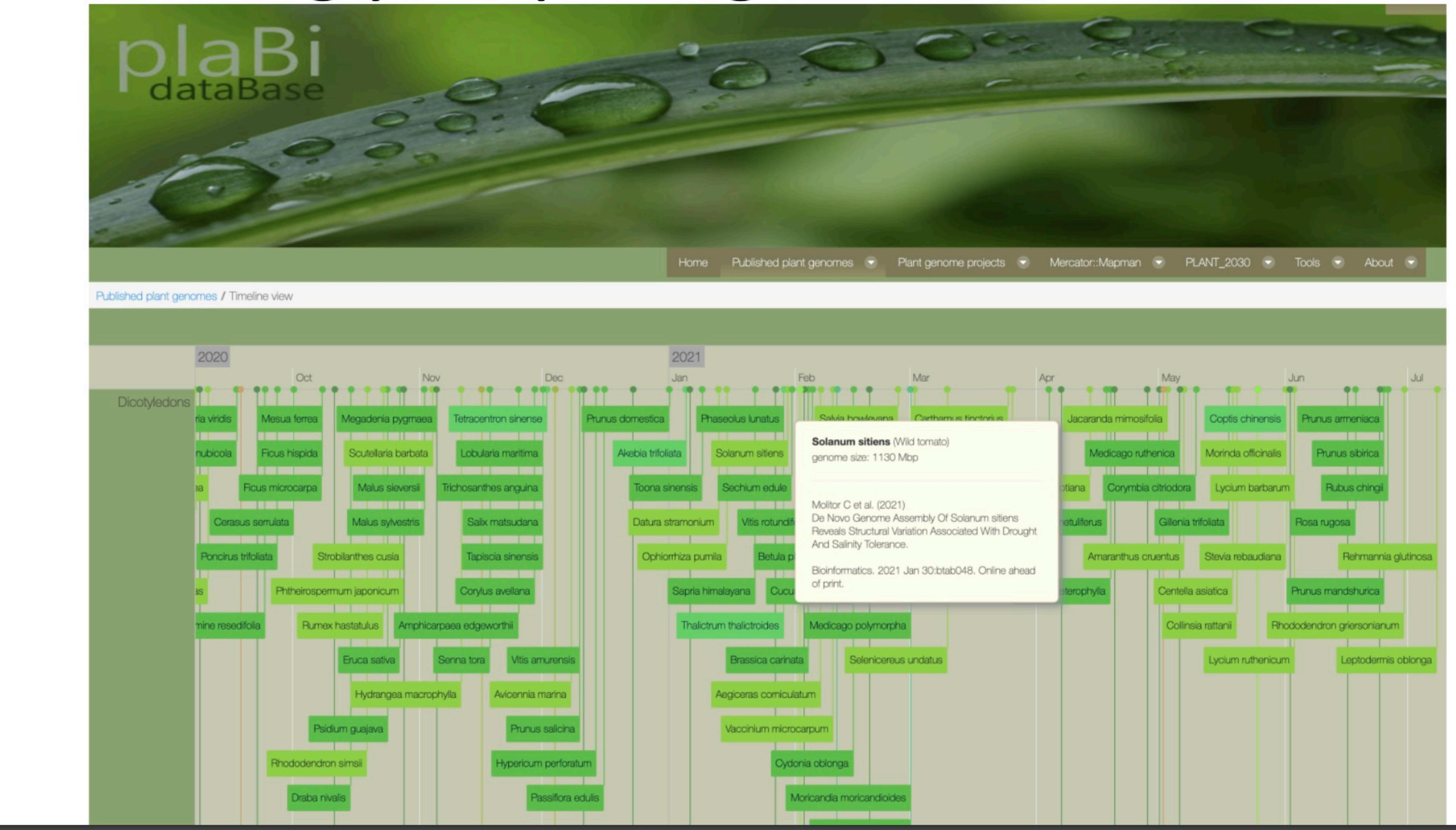
1. Background: What is bioinformatics?



Bioinformatics can...

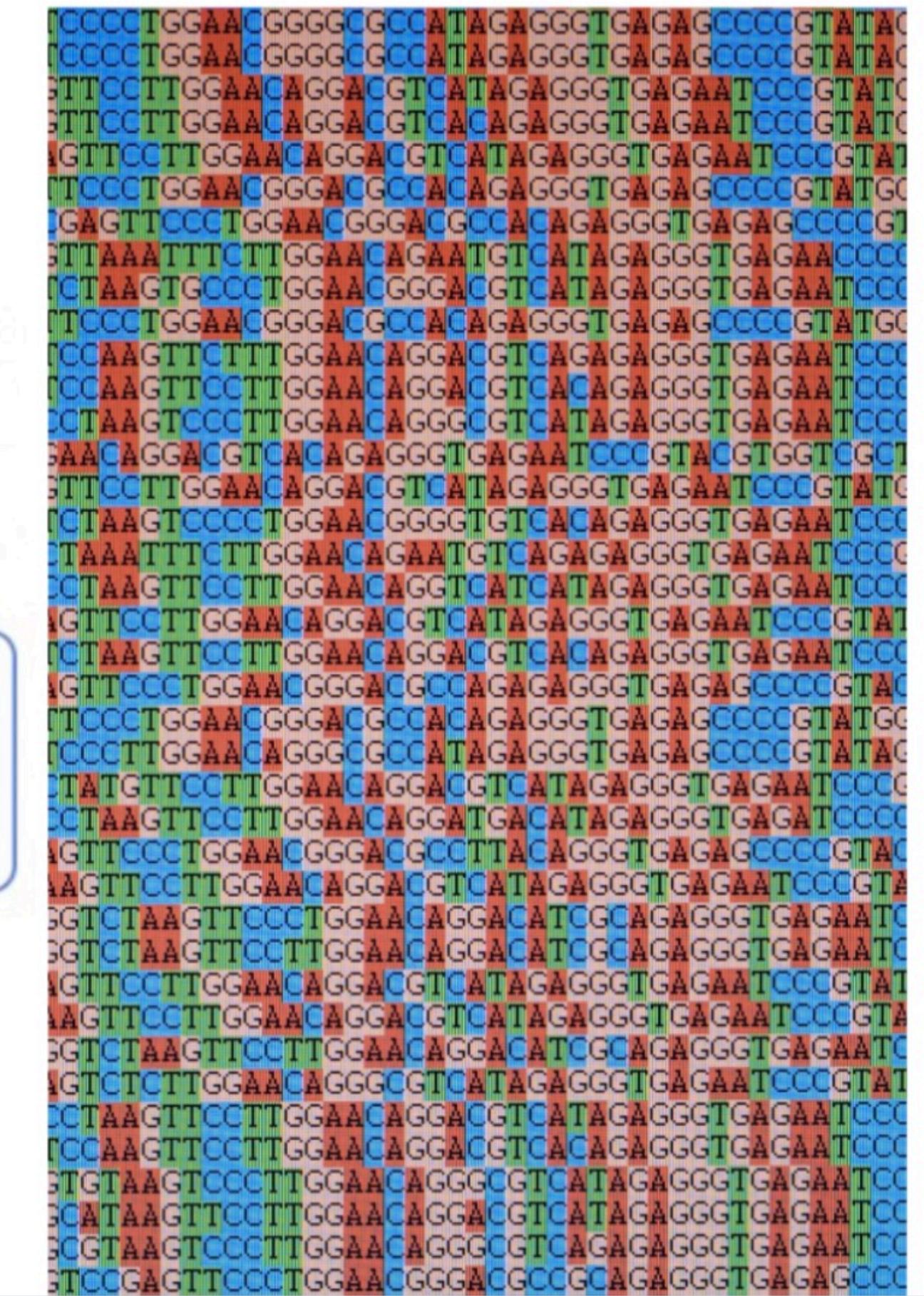
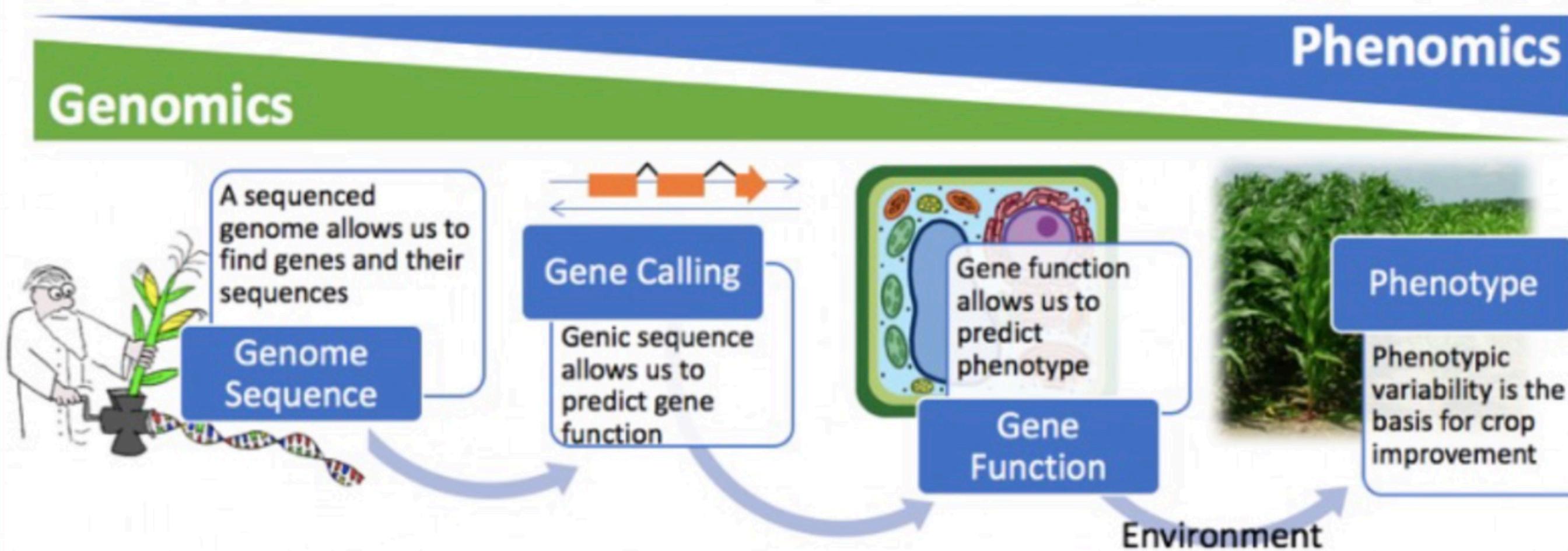
- Identify similar sequences
- Provide a putative function for a sequence
- Assemble sequences (genomes, transcriptomes)
- Annotate genomes
- Identify differentially expressed genes
- Build networks of genes or metabolites
- Determine phylogenetic relationships
- ‘Omics: genomics, transcriptomics, proteomics, metabolomics.....
- Mine literature for biological information
- Uncover differences between two genomes
- Calculate how a protein folds
- Find trends in large datasets
- Speed up your research
- Enable you to ask new questions

Genomics

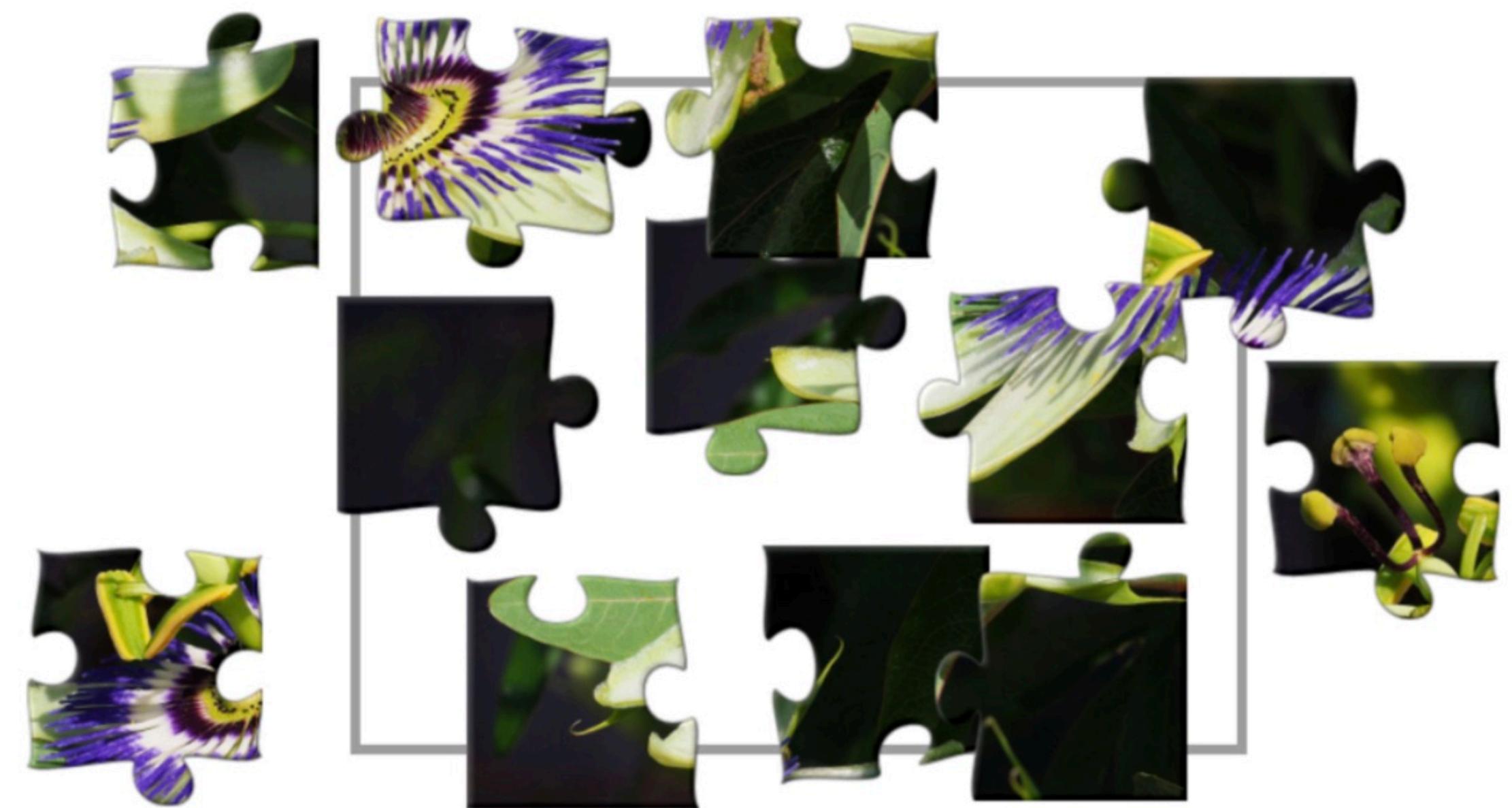
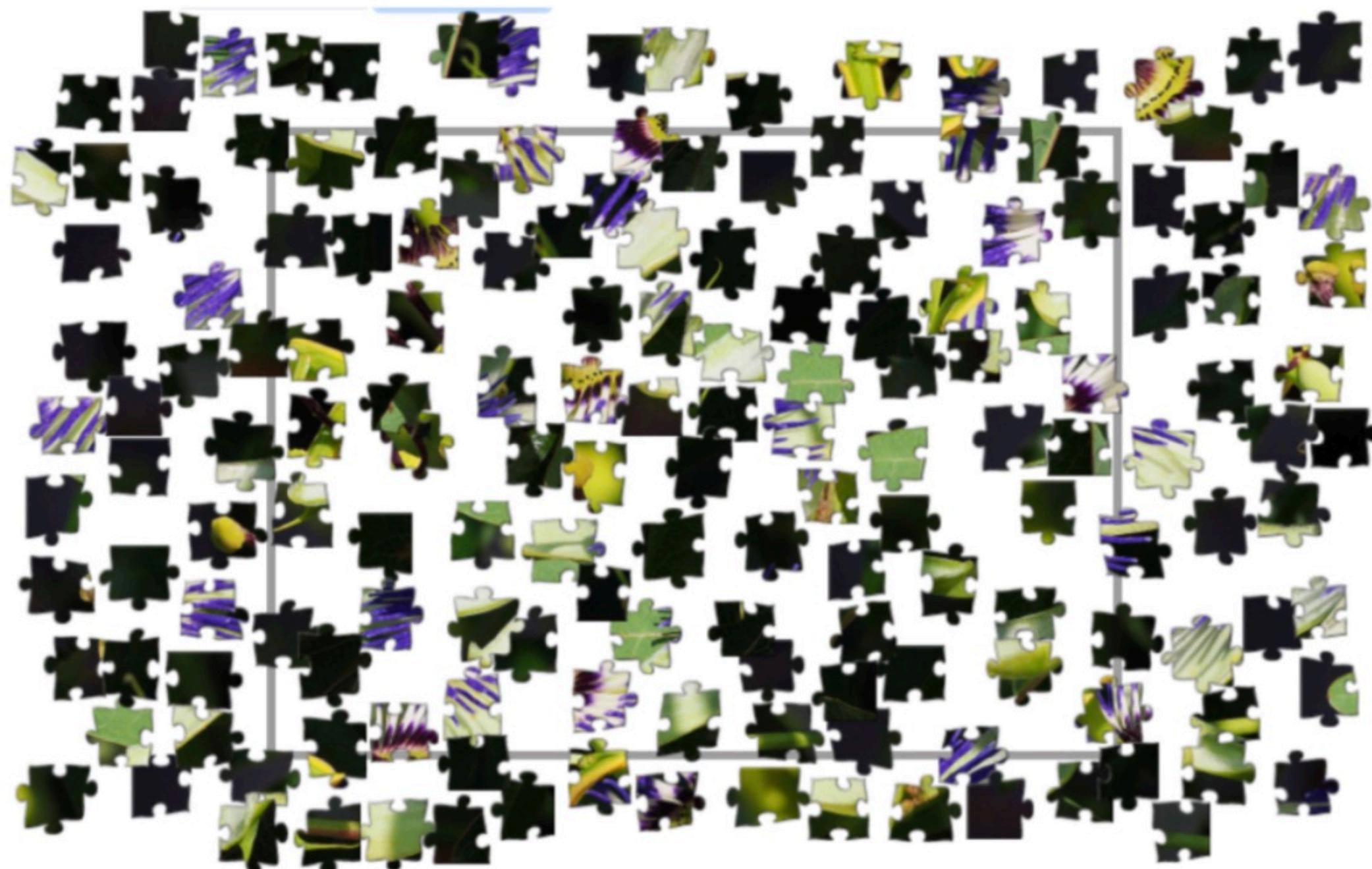


- The study of genomes

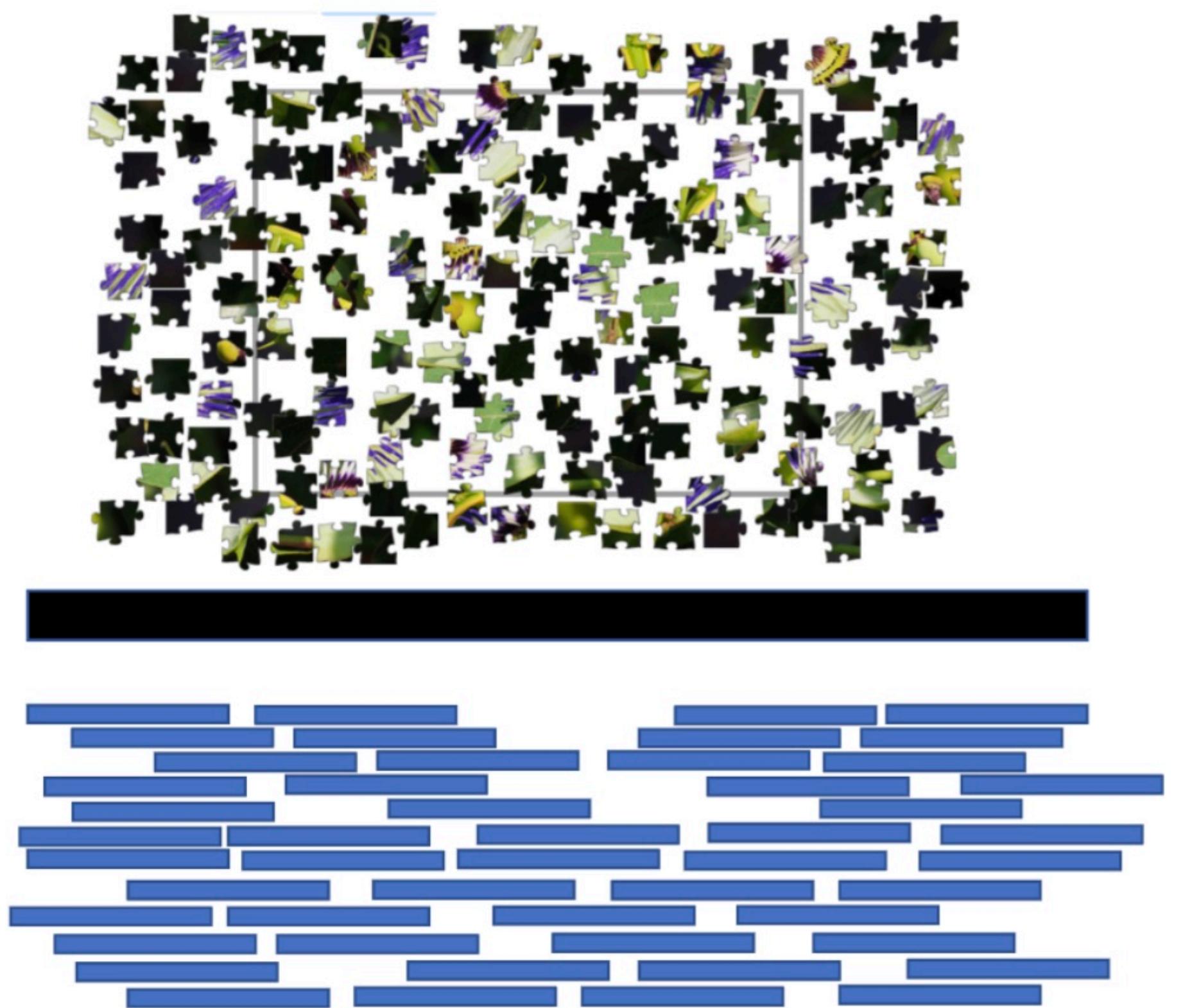
Genome as a puzzle



Which puzzle is easier to put together?

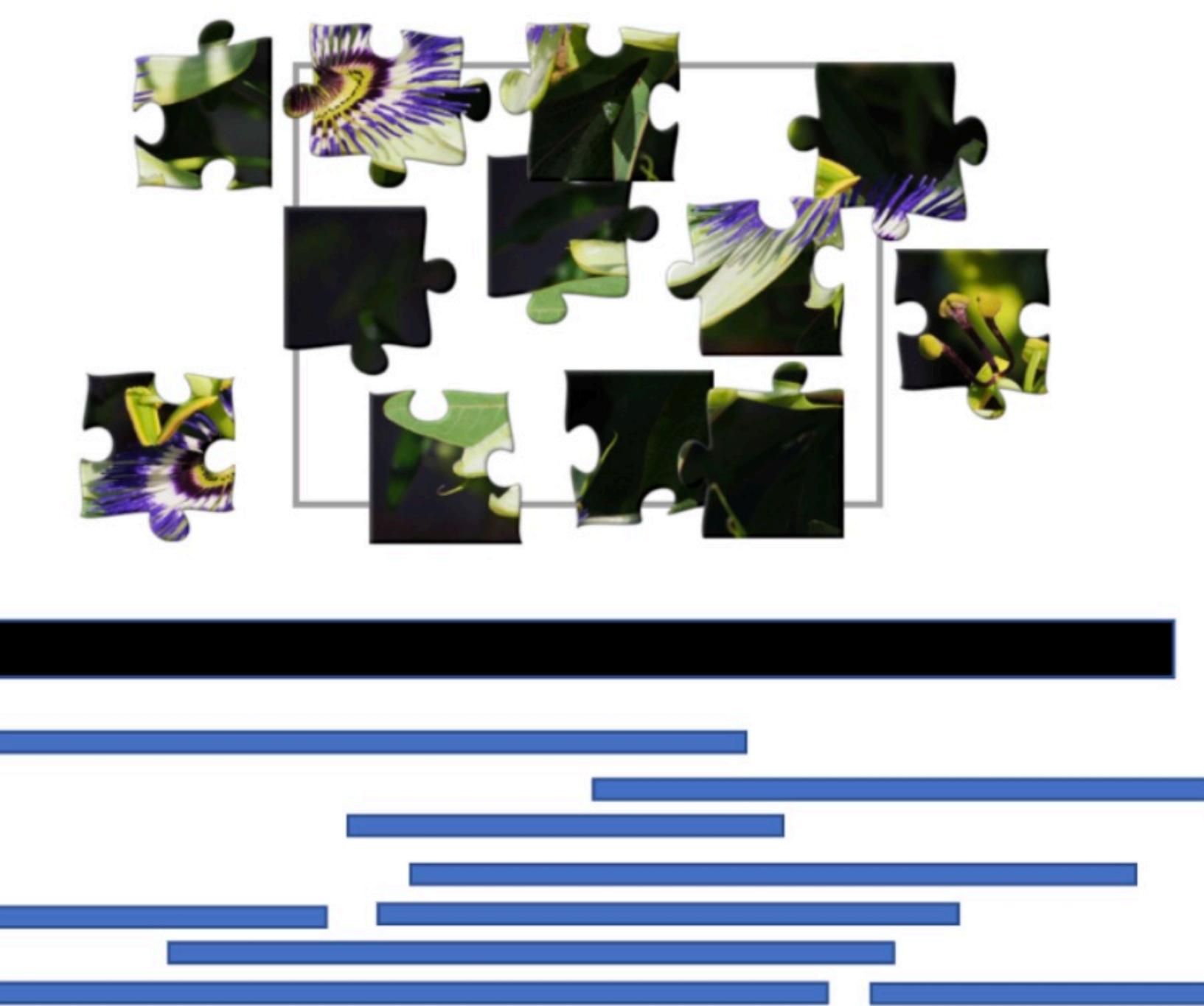


In the world of genomes



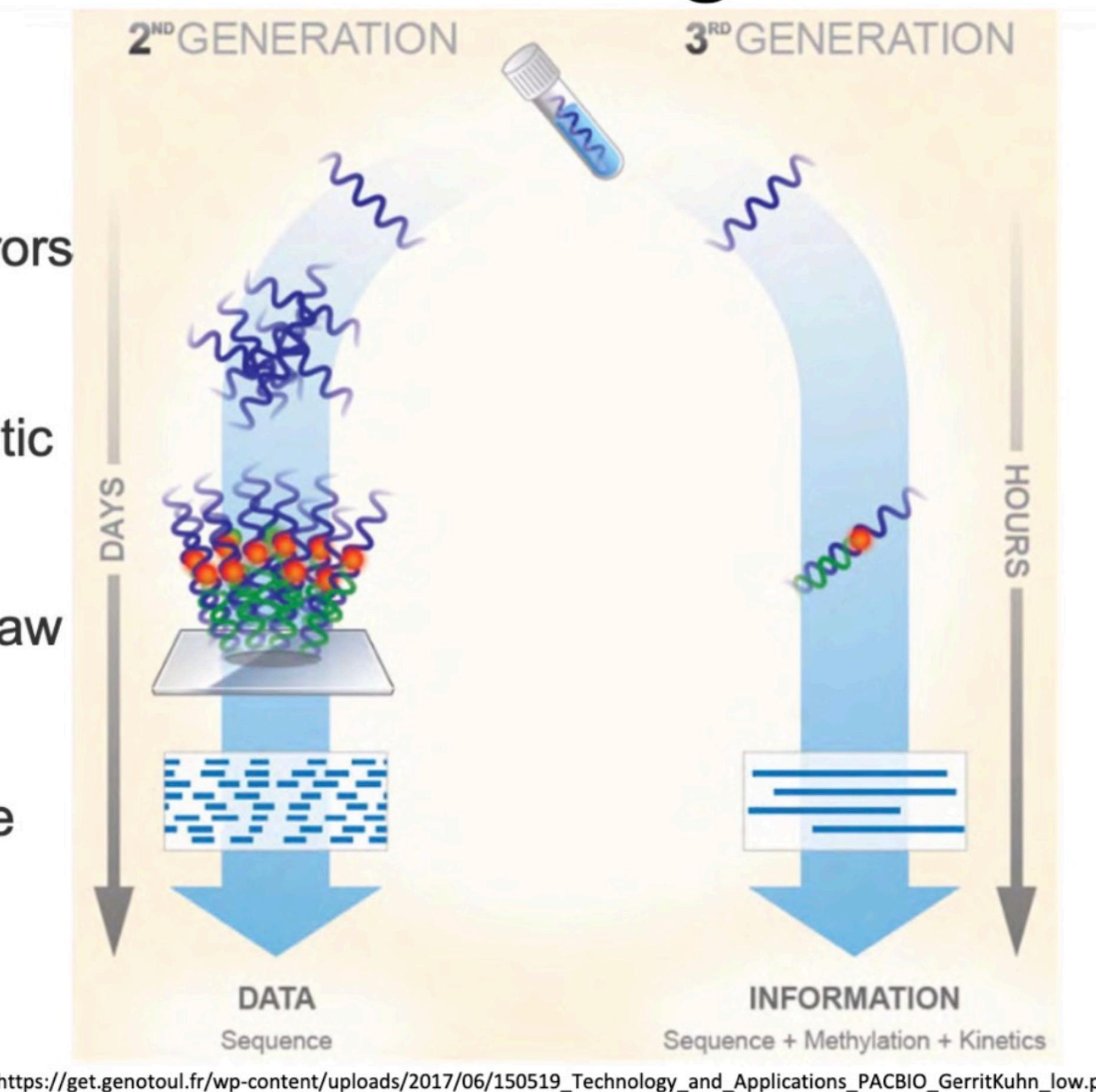
Reference
Genome

Sequencing
Reads



Short vs Long-reads

- Short reads
- Amplification errors and bias
- Several enzymatic steps
- Multi-molecule raw accuracy
- Errors tend to be systematic
- More coverage required



- Long reads
- No required amplification
- Simple sample prep
- Single molecule raw accuracy
- Errors tend to be random (vs. systematic)
- Less coverage required

https://get.genotoul.fr/wp-content/uploads/2017/06/150519_Technology_and_Applications_PACBIO_GerritKuhn_low.pdf



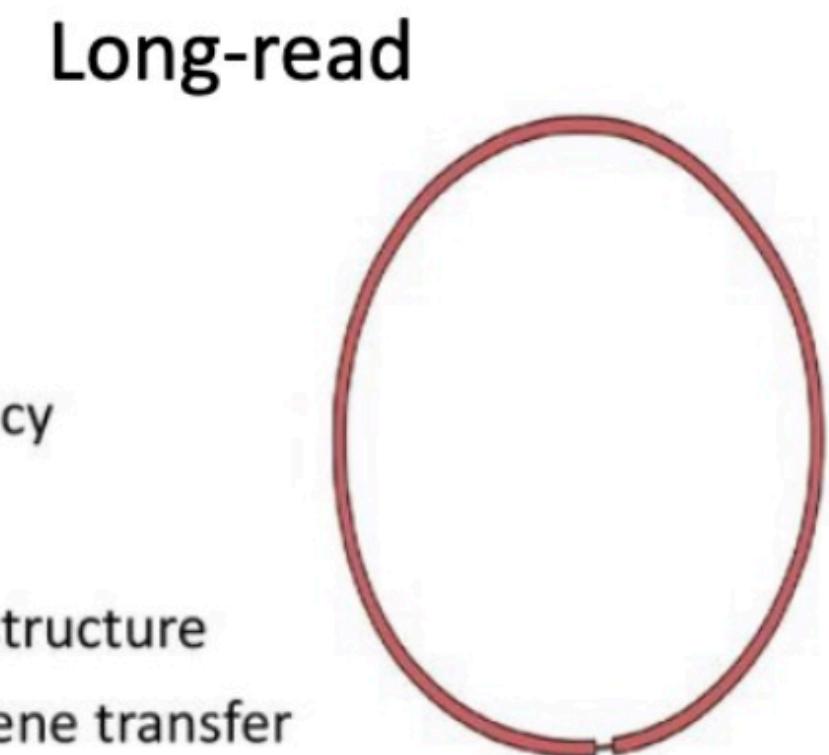
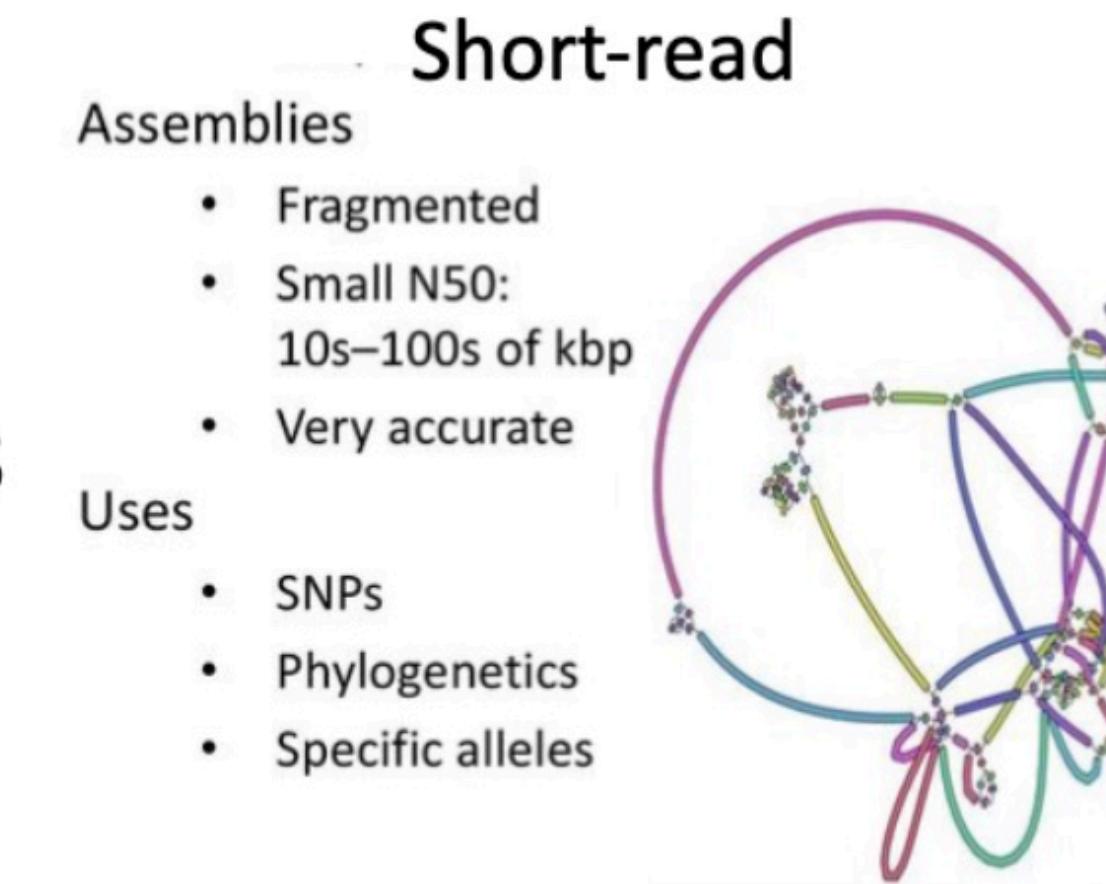
Flowcell

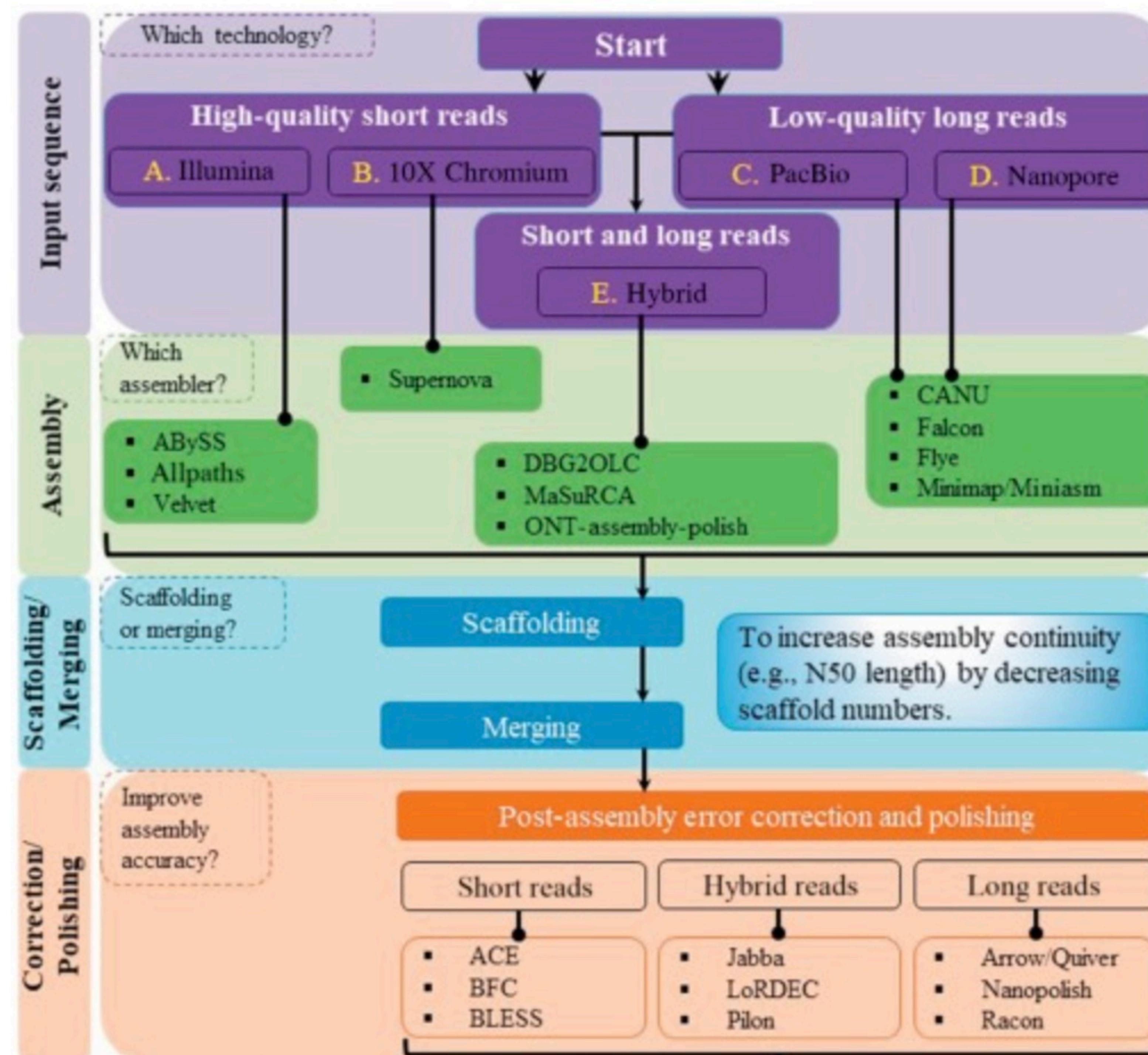


MinION device

de novo assembly

- Illumina only
 - High quality reads with fewer errors
- Hybrid option
 - Nanopore or PacBio + Illumina
 - Either raw or error corrected long-reads
- Long-read only
 - Raw typically works better
 - Need to polish after with Illumina data to fix errors
- Recommendation is 50x coverage short-reads and 50x long-reads
- So how much data do I need?



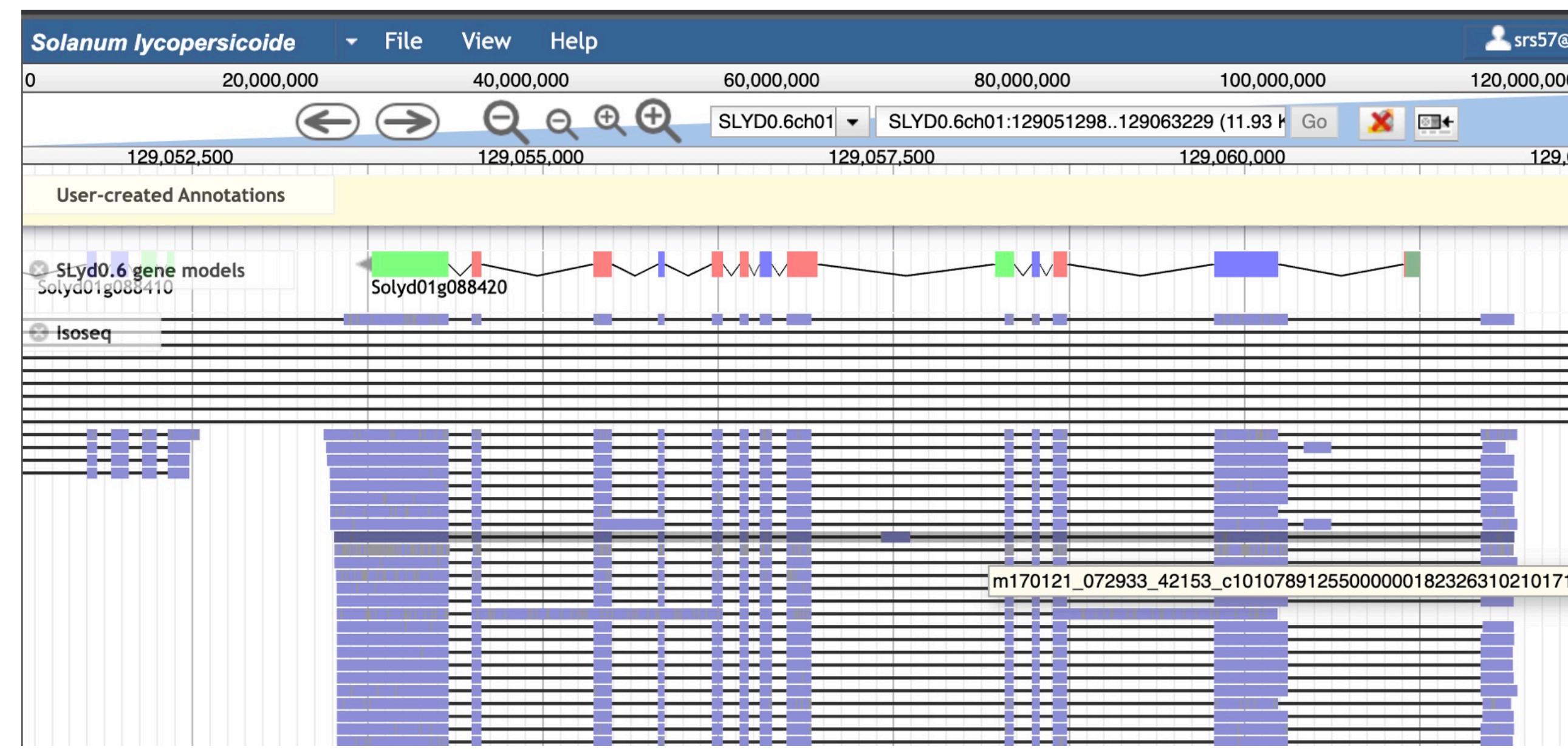


For more on genome assembly

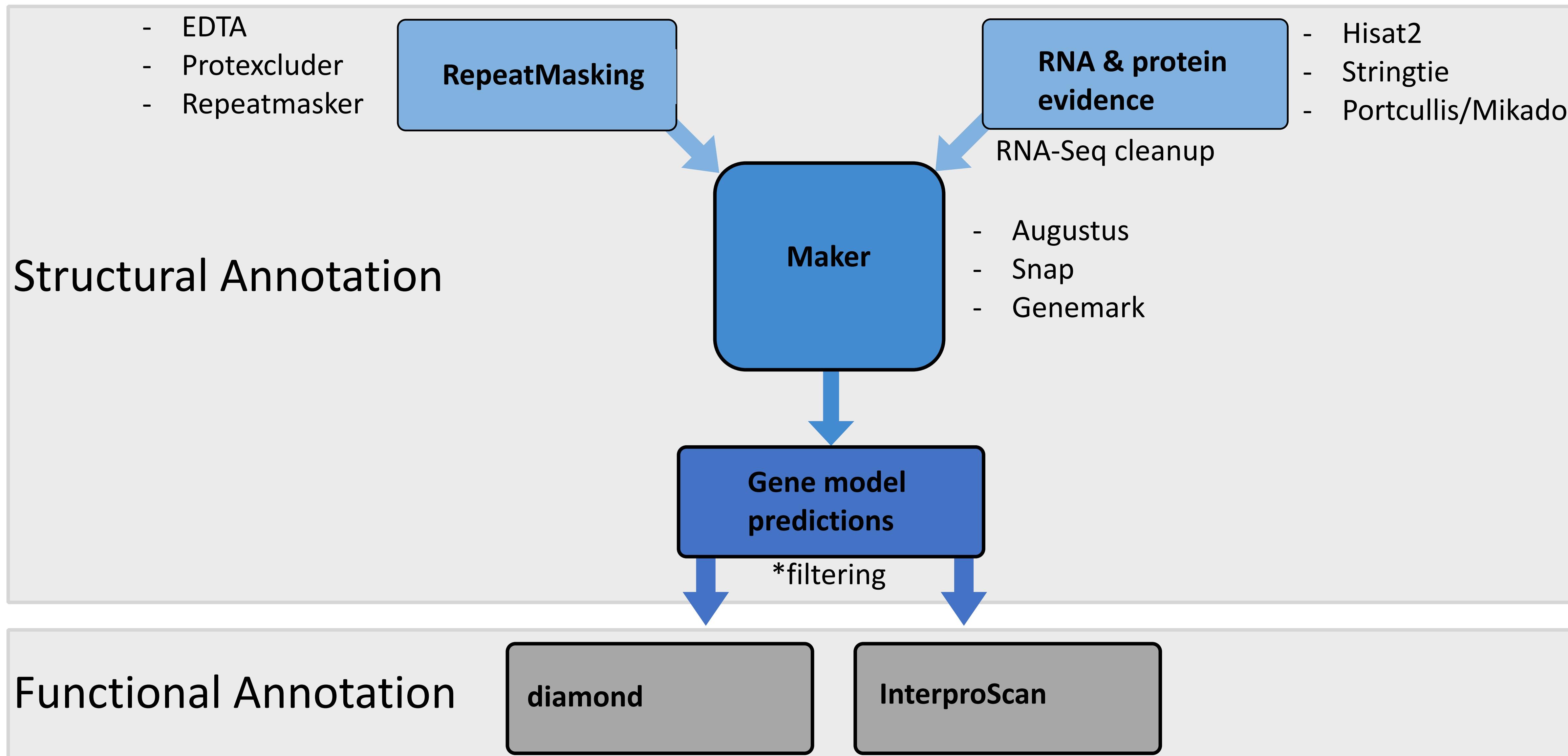
- https://github.com/bcbc-group/Botany2022NMGWorkshop/blob/main/4.GenomeAssembly/Botany2022_assembly_workshop_slides.pdf

Goals of genome annotation

- Predict, categorize, and mask repetitive elements
- Determine gene structures as accurately as possible
- Predict possible functions of predicted genes
- Associate GO terms, domains, etc for downstream analyses



Annotation pipeline



More info on annotation

- <https://github.com/bcbc-group/Botany2022NMGWorkshop/tree/main/5.Annotation>
- Can launch on CyVerse VICE

Variant Detection

Why Call SNPs?

[Home](#) > Current Issue > vol. 108 no. 17 > Jesse A. Poland, 6893–6898, doi: 10.1073/pnas.1010894108



Genome-wide nested association mapping of quantitative resistance to northern leaf blight in maize

Jesse A. Poland^{a,1}, Peter J. Bradbury^{a,b}, Edward S. Buckler^{a,b}, and Rebecca J. Nelson^{a,c,2}

UNIT 7.18 Next-Gen Sequencing-Based Mapping and Identification of Ethyl Methanesulfonate-Induced Mutations in *Arabidopsis thaliana*

Xue-Cheng Zhang¹, Yves Millet², Frederick M. Ausubel¹, Mark Borowsky¹

Published Online: 1 OCT 2014

DOI: 10.1002/0471142727.mb0718s108



Current Protocols in
Molecular Biology

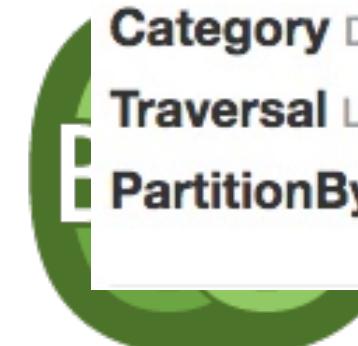
ASEReadCounter

Calculate read counts per allele for allele-specific expression analysis

Category Diagnostics and Quality Control Tools

Traversal LocusWalker

PartitionBy LOCUS



SNPhylo: a pipeline to construct a phylogenetic tree from huge SNP data

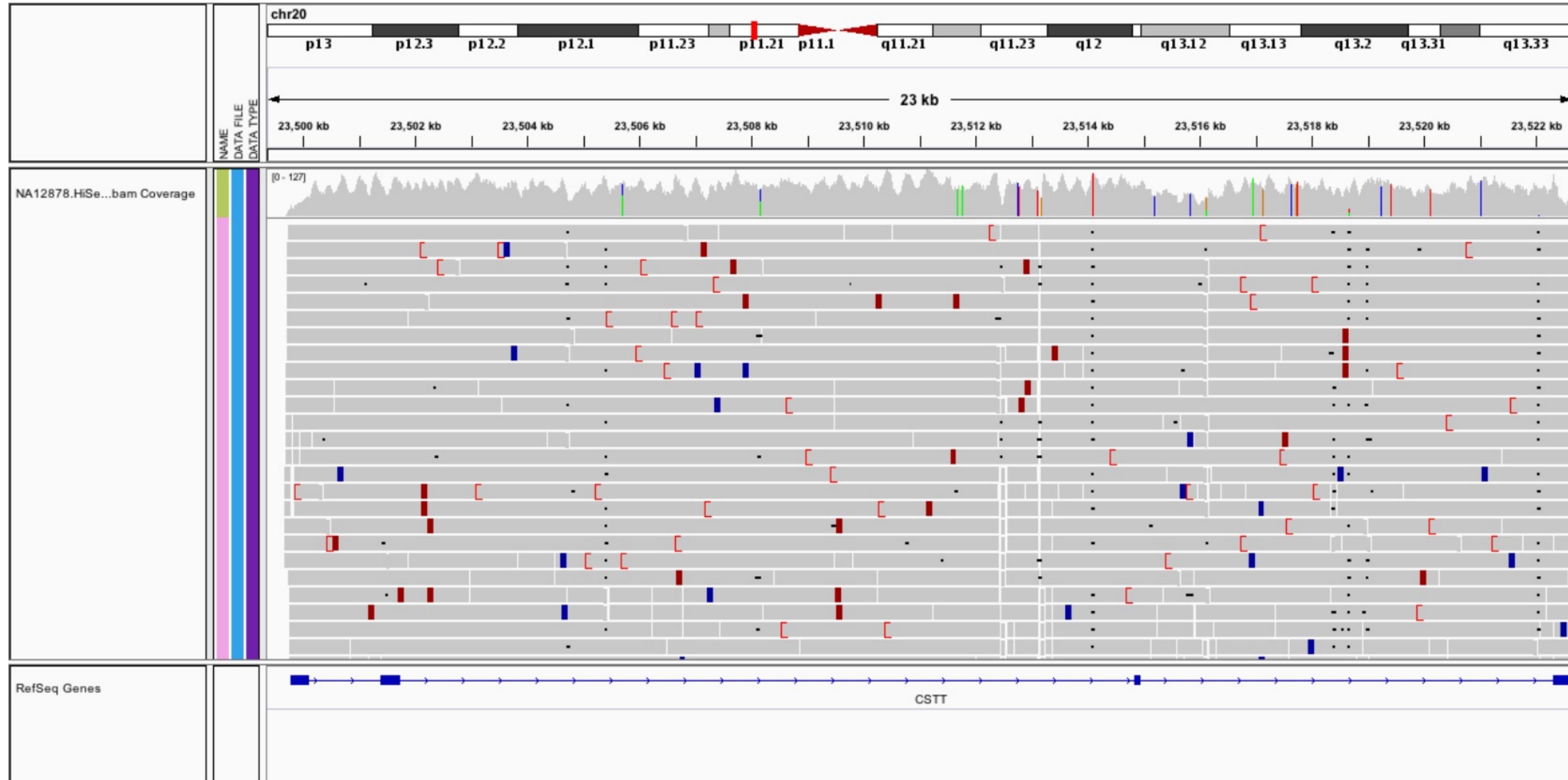
Tae-Ho Lee, Hui Guo, Xiyin Wang, Changsoo Kim and Andrew H Paterson

BMC Genomics 2014 15:162 | DOI: 10.1186/1471-2164-15-162 | © Lee et al.; licensee BioMed Central Ltd. 2014

Received: 25 September 2013 | Accepted: 18 February 2014 | Published: 26 February 2014



Which mismatches are real mutations and which are noise/error?



https://www.broadinstitute.org/gatk/events/3391/GATKw1310-BP-0A-Intro_to_NGS.pdf

Table 1 Algorithms and short descriptions of the seven variant calling tools

From: [Evaluation of variant calling tools for large plant genome re-sequencing](#)

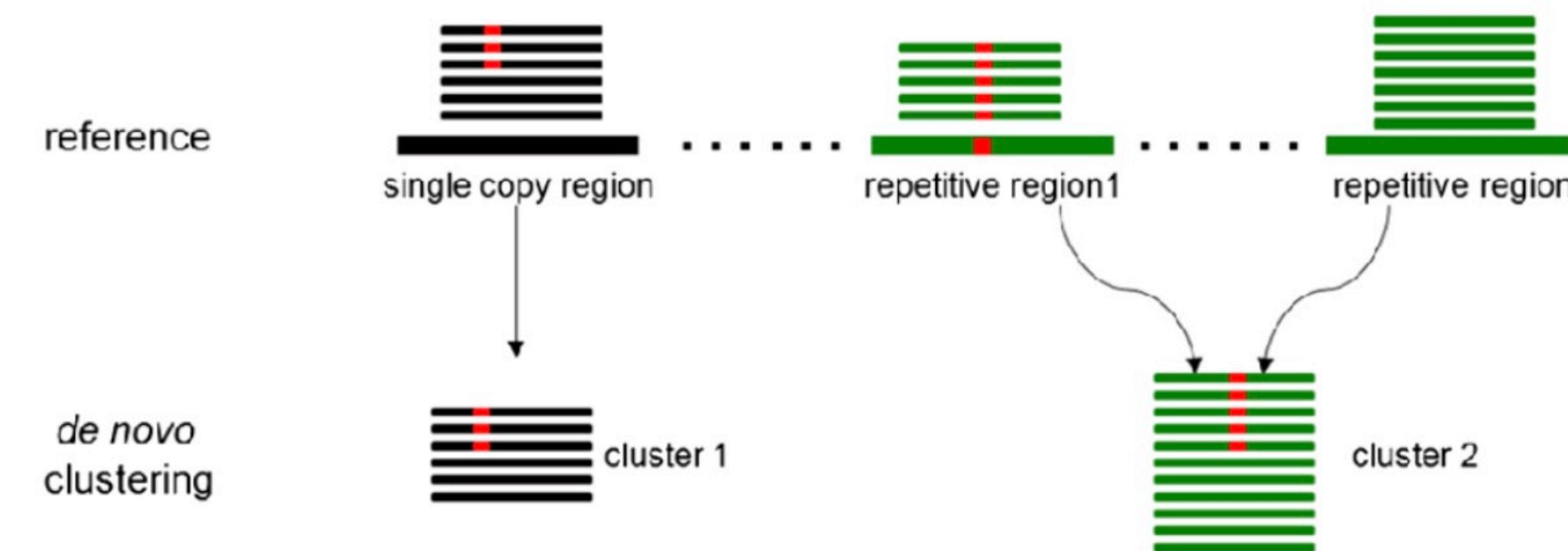
Variant tool	Version	Algorithm	Pipelines	Default filter	Reference
FreeBayes	v1.2.0–2	Haplotype-based	FreeBayes	^b 10, ^m 1	Garrison E, et al, 2012 [29]
		Bayesian			
GATK	4.0.11.0	Haplotype-based	MarkDuplicates	^b 10, ^m 20	DePristo M, et al, 2011 [27]
		significant test	BaseRecalibrator		
			HaplotypeCaller		
Platypus	0.8.1	Haplotype-based	Platypus callVariants	^b 20, ^m 20	Rimmer A, et al, 2014 [30]
		significant test			
Samtools /mpileup	1.9	Site align-based	Samtools/mpileup	^b 13, ^m 0	Li H, 2011 [28]
		gt likelihoods	bcftools call		
SNVer	0.5.3	Site align-based	SNVerIndividual	^b 17, ^m 20	Wei Z, et al, 2011 [31]
		MAF <i>p</i> -value		^f 0.25, ^r 1, ^P 0.05	
VarScan	v2.3.9	Site-based	Samtools/mpileup	^b 15, ^m 0	Koboldt D, et al, 2012 [33]
		allele frequency	mpileup2snp	^f 0.2, ^r 2, ^P 0.01	
VarDict	2018	Site-based	VarDict	^b 22.5, ^m 0	Lai Z, et al, 2016 [32]
		alleles Fisher's	var2vcf_valid	^f 0.01, ^r 2	

^aOnly default settings were listed. ^bBQ Base quality; ^mMQ Mapping quality; ^rVR Variant containing reads or total reads containing variants (TR); ^fVF Variant frequency; ^P P-value; ^dDP Depth coverage

Downstream analyses using SNPs

Population genetic studies

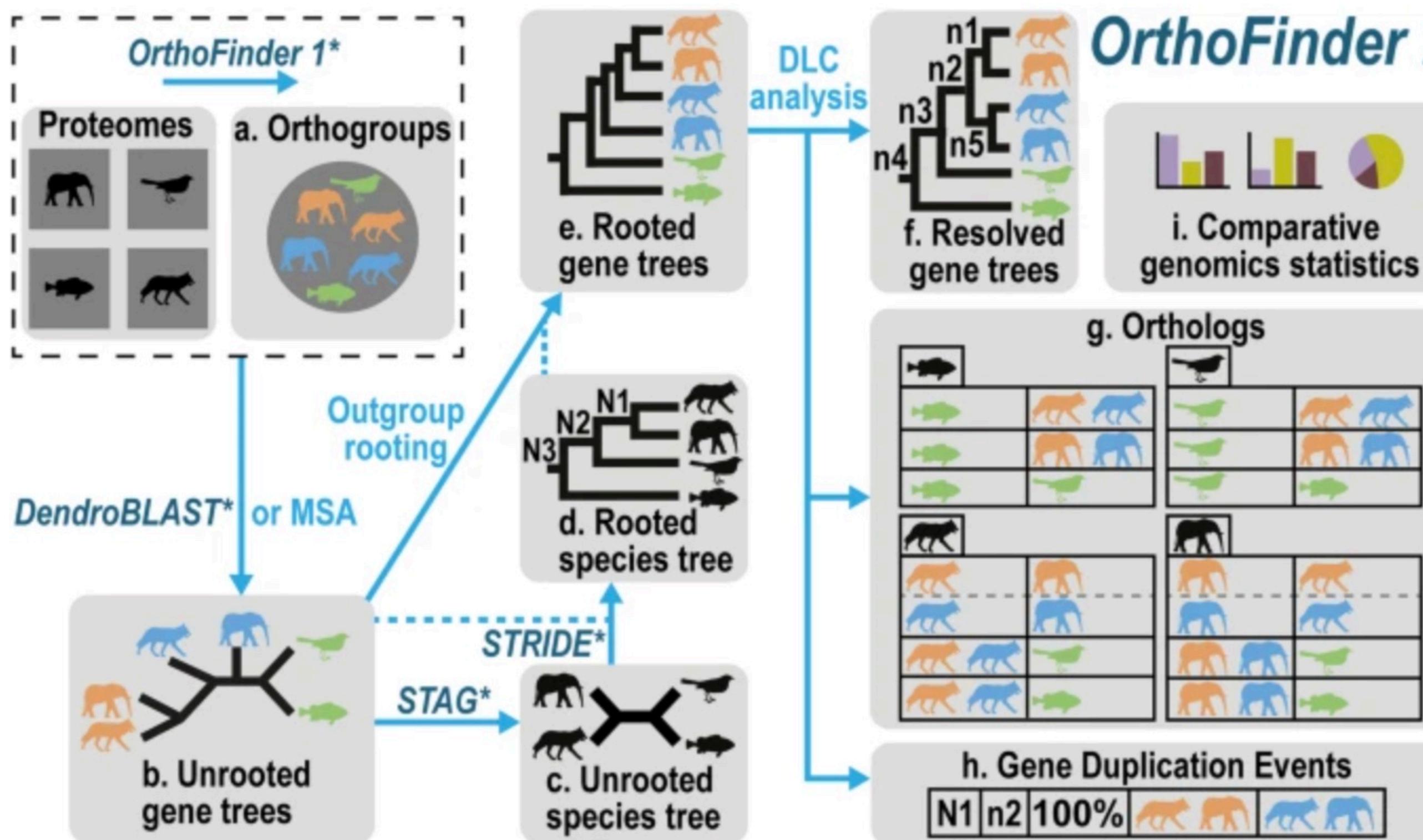
- Many options and different requirements for SNP calling
- RAD-Seq
 - Stacks, both reference guided or *de novo*
- Hyb-Seq, RNA-Seq or Whole Genome Sequencing
 - GATK, Freebayes, mpileup
- Must have a reference genome or at the very least something to map reads to



Phylogenomics

- How do sequences/genomes relate to each other?
- Align sequences
 - ClustalW
 - Muscle
 - MAFFT
- Build phylogenetic trees
 - Parsimony (PAUP)
 - Neighbor join (QuickTree)
 - Maximum likelihood (IQ-TREE, PHYLIP)
 - Bayesian (MrBayes)

Gene Families



The OrthoFinder workflow. The method used for each step is shown by the arrow. Published algorithms are shown in italics and are followed by an asterisk. A dotted blue line connecting with a solid arrow indicates additional data that are used in order to carry out the transformation indicated by the solid arrow. MSA, multiple sequence alignment-based tree inference; DLC, duplication-loss-coalescence. (a) Orthogroup inference using the original OrthoFinder algorithm (an orthogroup is the set of genes descended from a single gene in the last common ancestor of all the species under consideration). (b) Gene tree inference. (c) Species tree inference. (d) Species tree rooting (e) Gene tree rooting (f) Hybrid overlap + DLC analysis of rooted gene trees to infer orthologs and gene duplication events. (g) Illustration of the ortholog results table for the genes in each input species (four main boxes). The horizontal divisions within these show the orthologs for each individual species pair. (h) Illustration of the gene duplication event table showing the location of the gene duplication events mapped to the species tree, the location in the gene tree, the percent retention of the duplicate genes in the sampled species, and the genes descended from the gene duplication event. (i) Comparative genomics statistics

Single copy orthologs

CAFE: a computational tool for the study of gene family evolution FREE

Tijl De Bie, Nello Cristianini, Jeffery P. Demuth, Matthew W. Hahn ✉ Author Notes

Bioinformatics, Volume 22, Issue 10, 15 May 2006, Pages 1269–1271,

GWAS

- Genome-Wide Association Study
- Correlate SNPs with phenotypic effects
- Requires genotyped and phenotyped populations
- Allows identification of genes responsible for a phenotype

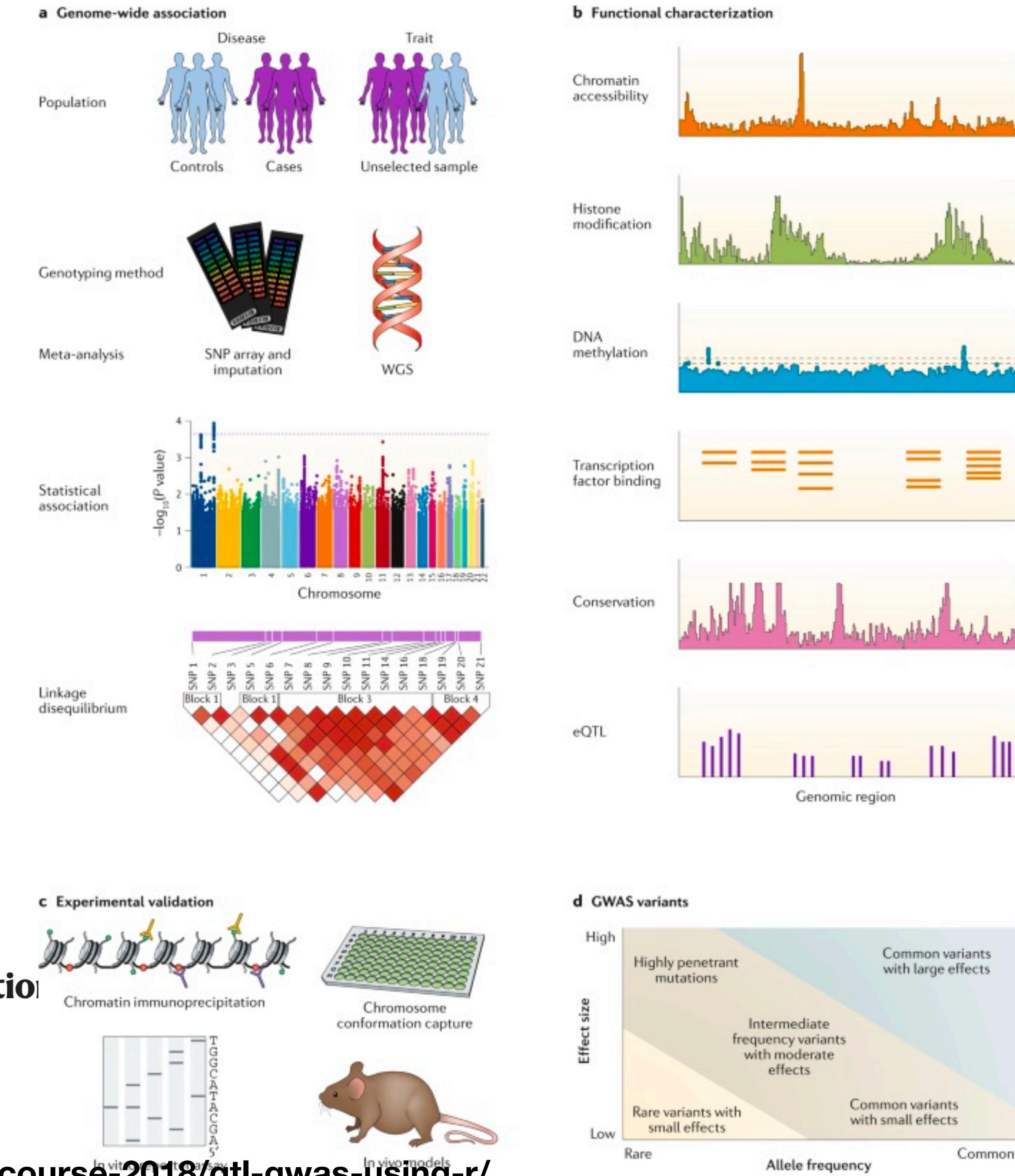
High LD = haplotype blocks

Review Article | Published: 08 May 2019

Benefits and limitations of genome-wide association studies

Vivian Tam, Nikunj Patel, Michelle Turcotte, Yohan Bossé, Guillaume Paré & David Meyre 

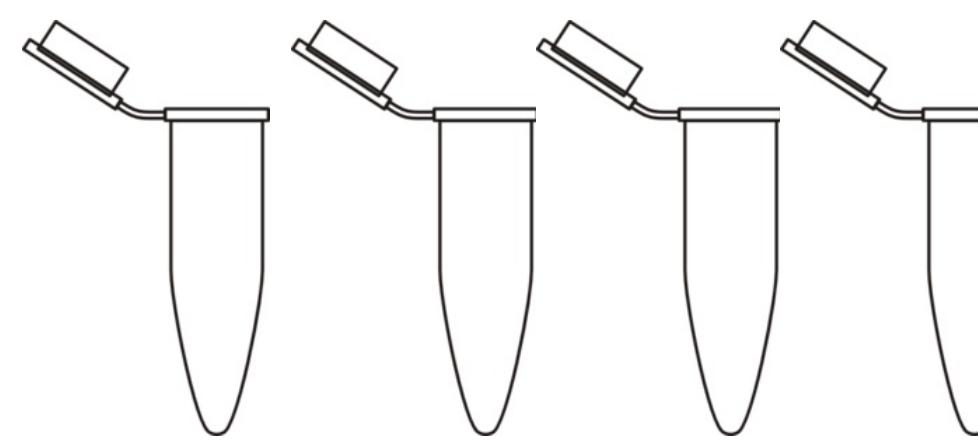
Nature Reviews Genetics 20, 467–484 (2019) | Cite this article



Transcriptomics

The study of transcriptomes

RNASeq



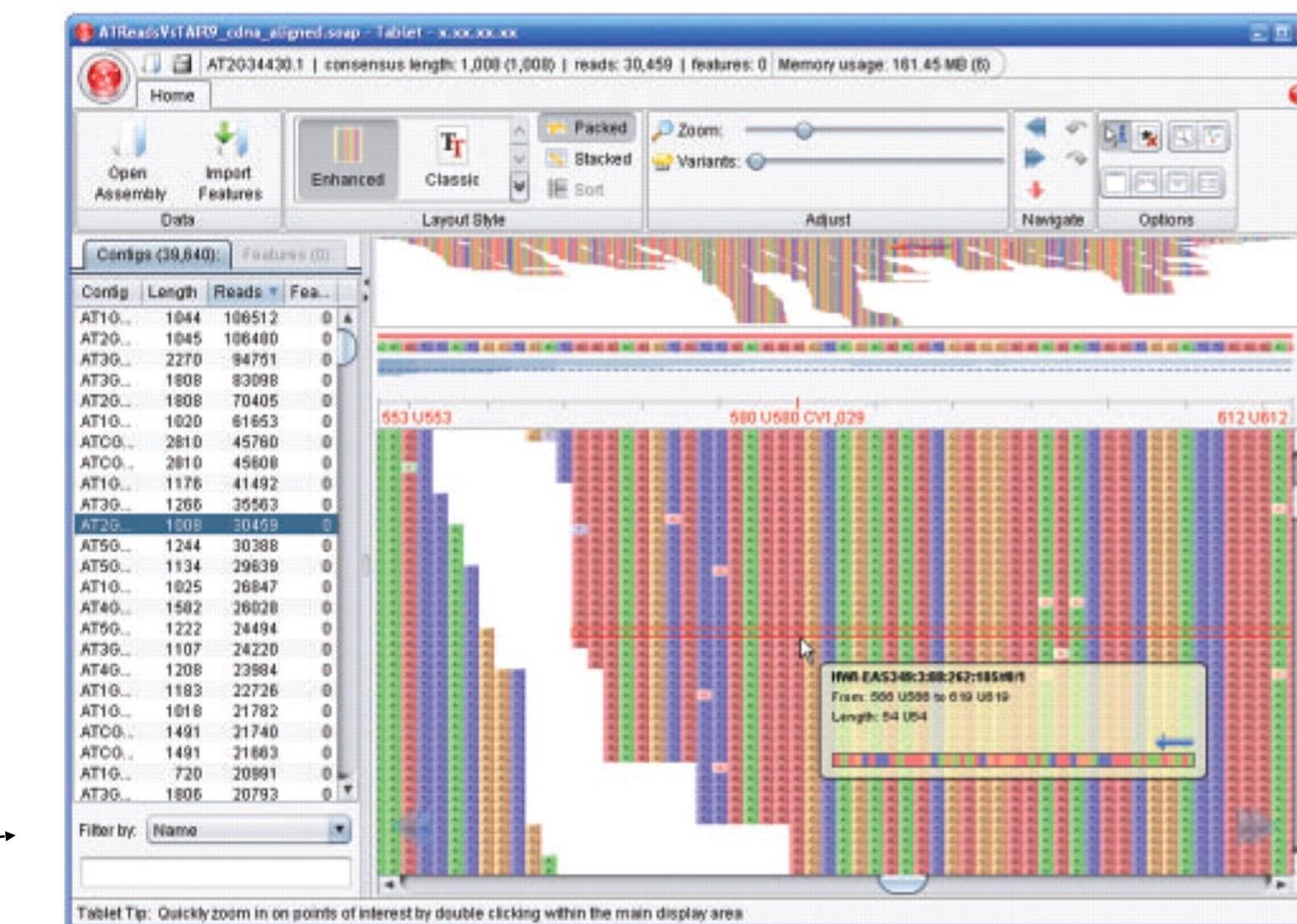
RNA samples



cDNA synthesis,
multiplex and sample
prep (bridge
amplification)



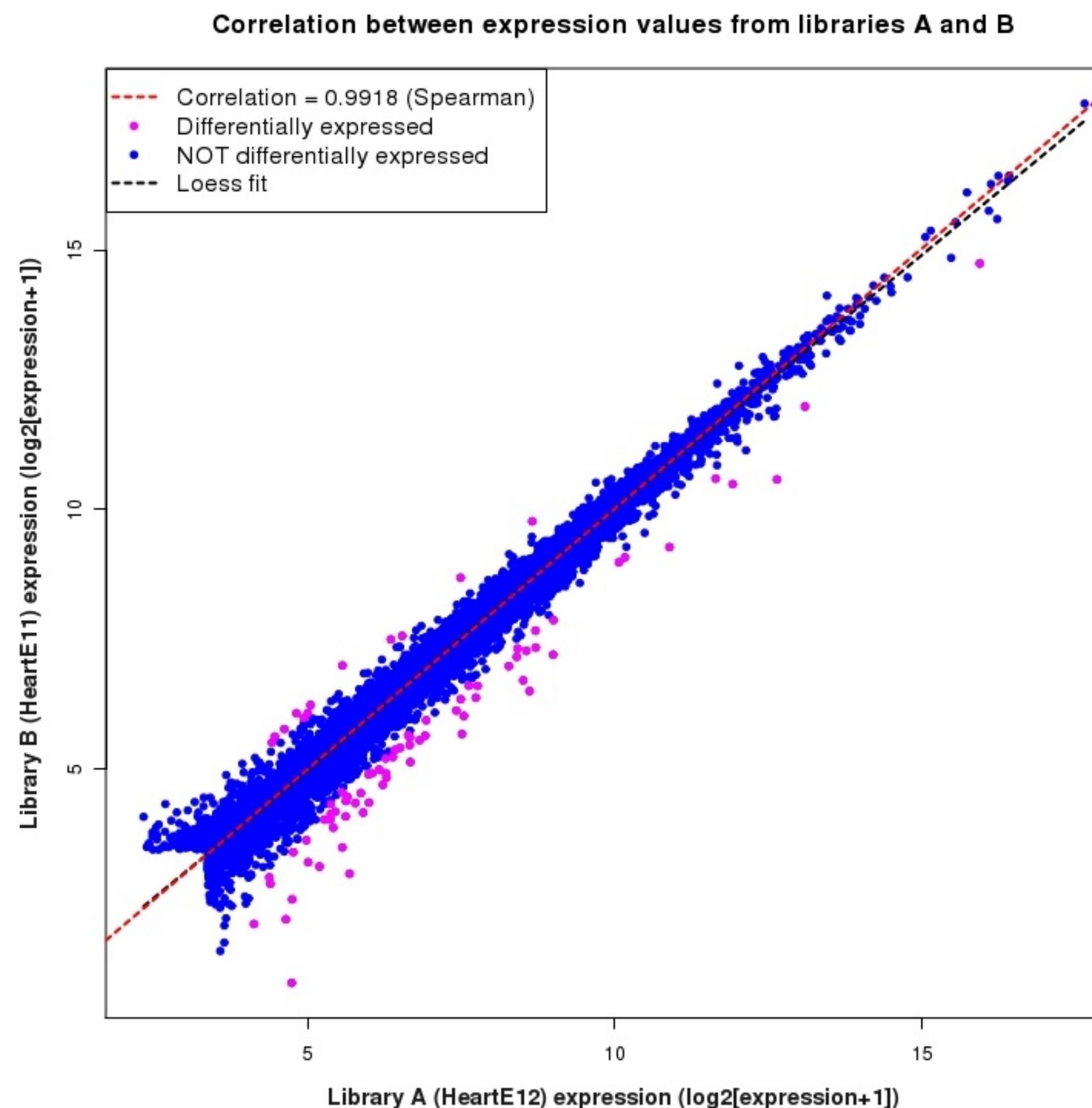
Align to reference
genome/transcriptome
Using bowtie2

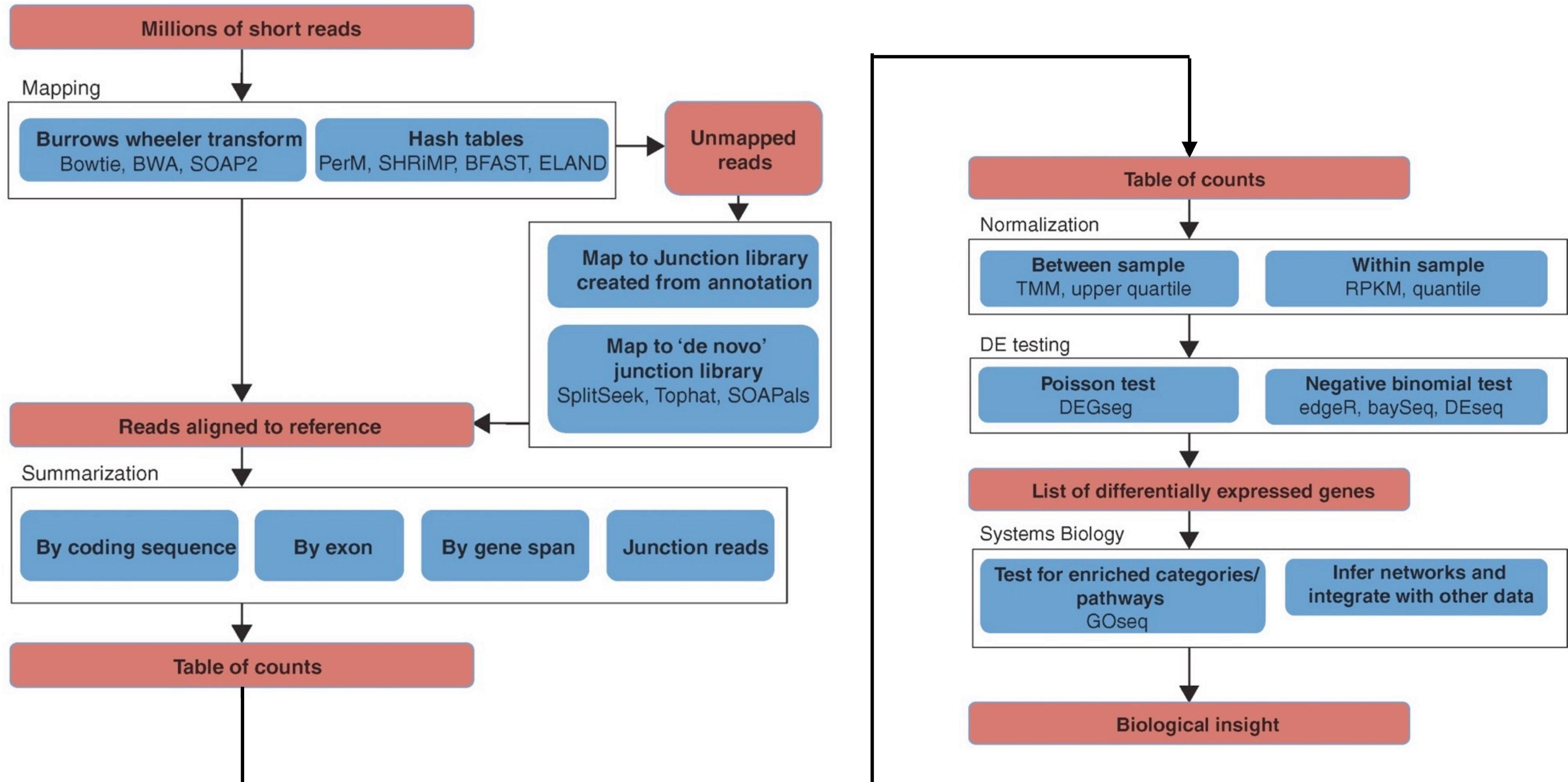


Visualize in Tablet
software, determine
rpkm values
(number of reads
per kilobase of exon
per million reads)

Differential Expression

- Identify differentially expressed genes





Differential Expression Analysis

edgeR: a Bioconductor package for differential expression analysis of digital gene expression data

Robinson MD, McCarthy DJ and Smyth GK (2010). “edgeR: a Bioconductor package for differential expression analysis of digital gene expression data.” *Bioinformatics*, 26, pp. -1.

<http://www.bioconductor.org/packages/release/bioc/vignettes/edgeR/inst/doc/edgeRUsersGuide.pdf>

DESeq

Differential expression analysis for sequence count data

Simon Anders  & Wolfgang Huber

Genome Biology 11, Article number: R106 (2010) | [Cite this article](#)

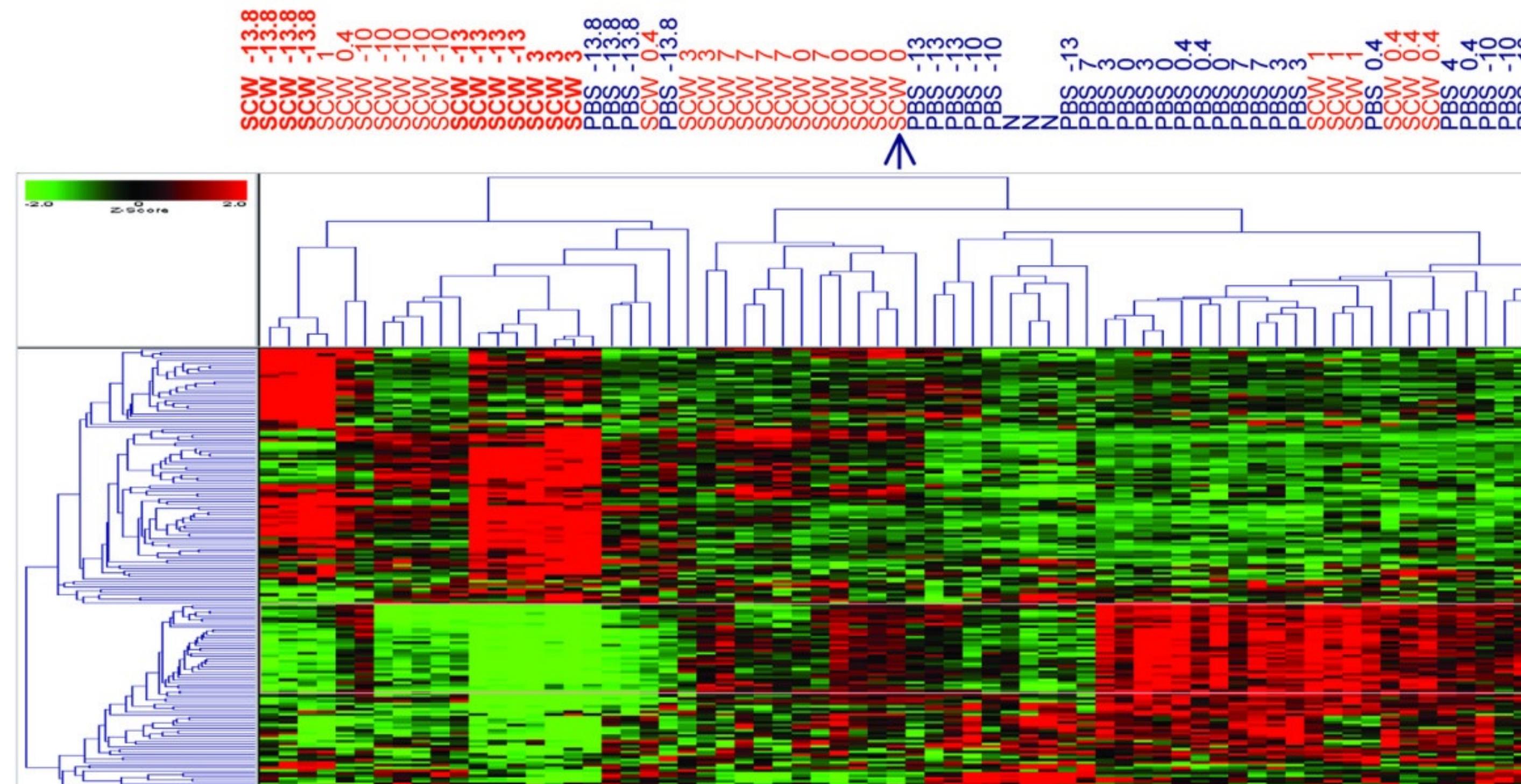
348k Accesses | 9676 Citations | 92 Altmetric | [Metrics](#)

A Beginner's Guide to Analysis of RNA Sequencing Data

Clarissa M. Koch ^{1*}, Stephen F. Chiu ^{1,2*}, Mahzad Akbarpour ², Ankit Bharat ^{1,2}, Karen M. Ridge ^{1,3}, Elizabeth T. Bartom ^{4‡}, and Deborah R. Winter ^{5‡}

Hierarchical Clustering

- Identify genes with similar expression profiles

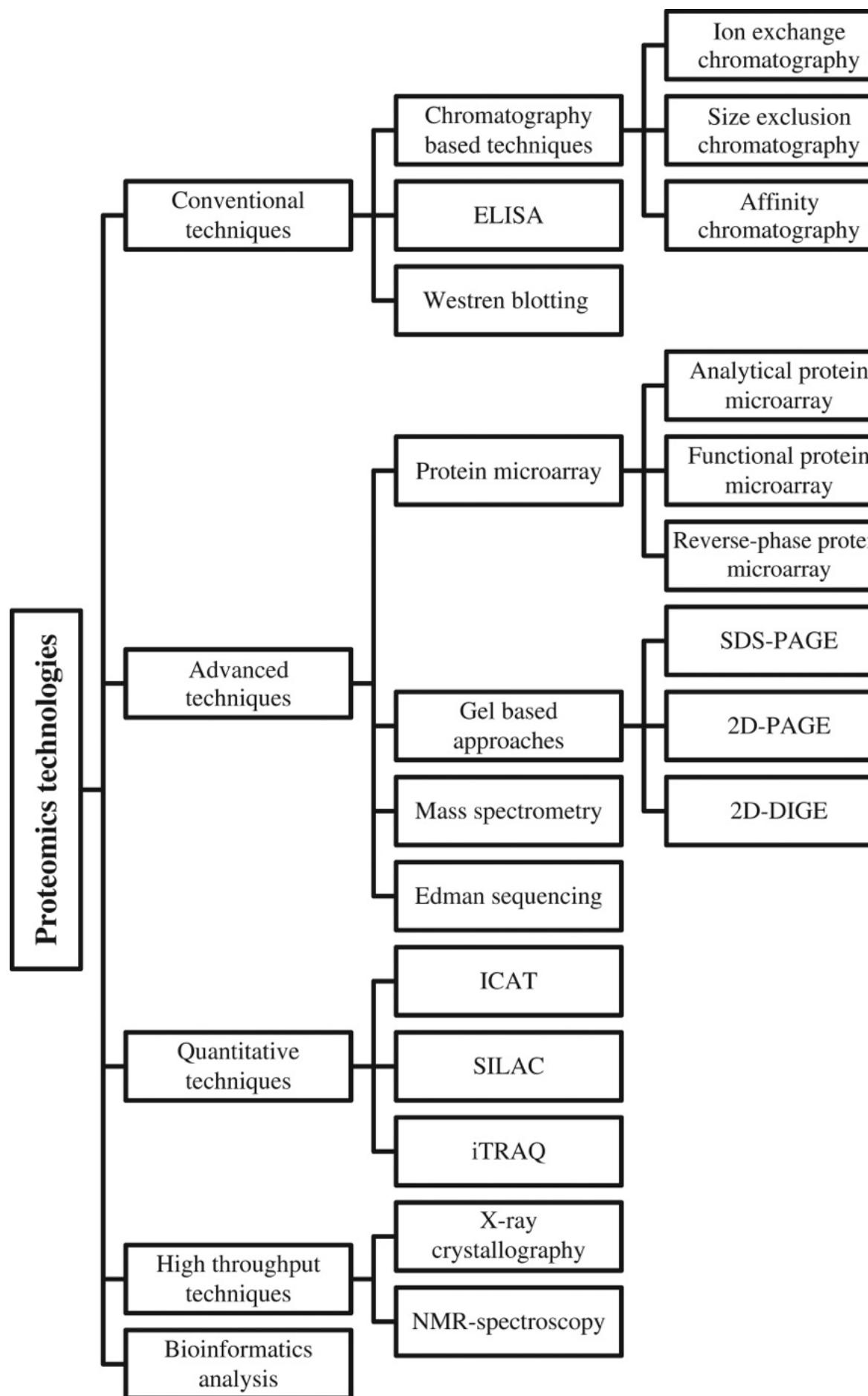


pheatmap2: Pretty and parallel heatmap
In [jokergoo/pheatmap2: Pretty Plus Parallel Heatmap](#)

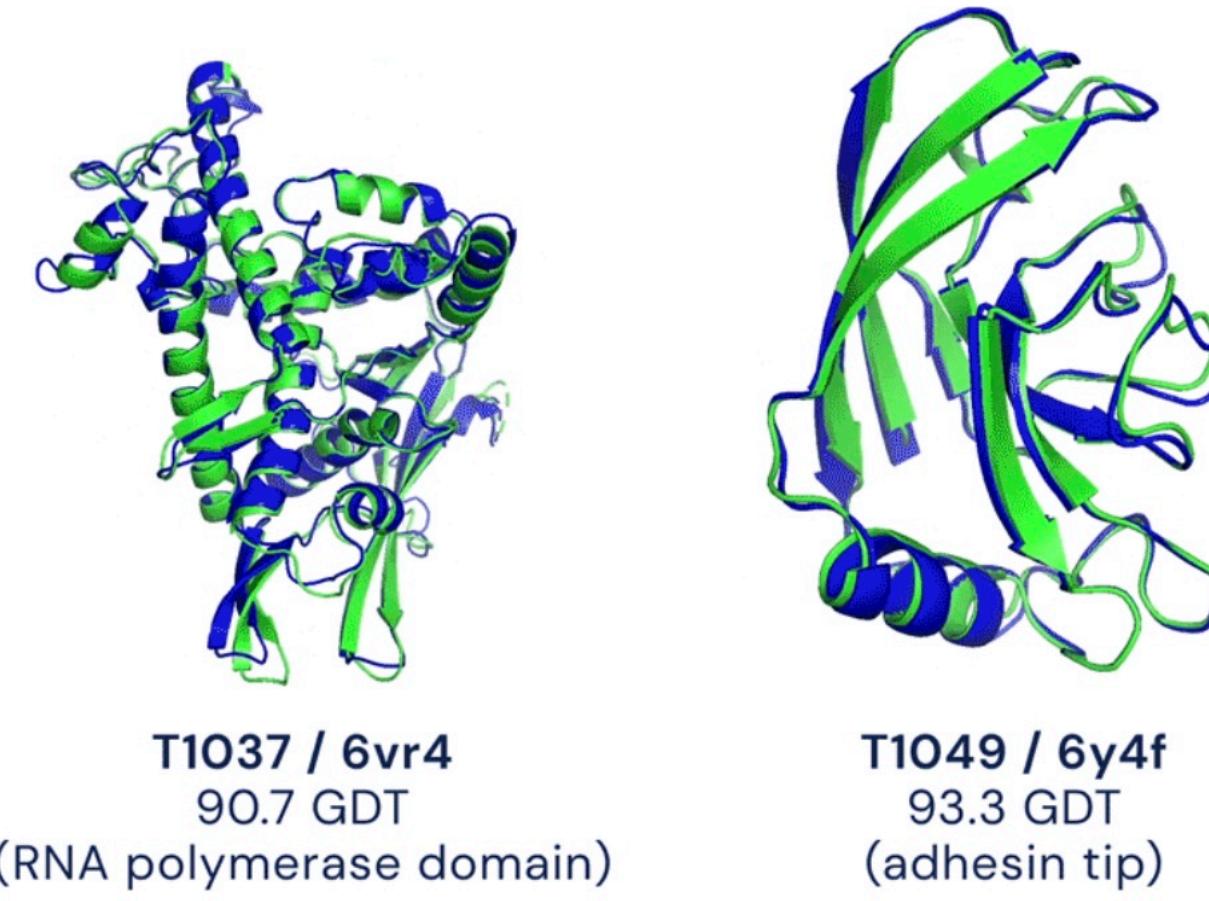
Proteomics

The study of proteomes

Figure 1. An overview of proteomics techniques.
Applications of proteomics techniques.
Figure 3.Schematic ...



AlphaFold2



● Experimental result
● Computational prediction

[nature](#) > [articles](#) > [article](#)

Article | [Open Access](#) | Published: 15 July 2021

Highly accurate protein structure prediction with AlphaFold

[John Jumper](#)✉, [Richard Evans](#), ... [Demis Hassabis](#) ✉ [+ Show authors](#)

ColabFold: AlphaFold2 using MMseqs2

Easy to use protein structure and complex prediction using [AlphaFold2](#) and [AlphaFold2-multimer](#). Sequence alignments/templates are generated through [MMseqs2](#) and [HHsearch](#). For more details, see [bottom](#) of the notebook, checkout the [ColabFold GitHub](#) and read our manuscript. Old versions: [v1.0](#), [v1.1](#), [v1.2](#), [v1.3](#)

[Mirdita M, Schütze K, Moriwaki Y, Heo L, Ovchinnikov S, Steinegger M. ColabFold: Making protein folding accessible to all. *Nature Methods*, 2022](#)



Metabolomics

- study of chemical processes involving **metabolites**, the small molecule substrates, intermediates, and products of cell metabolism.
- Detected by mass spectrophotometry

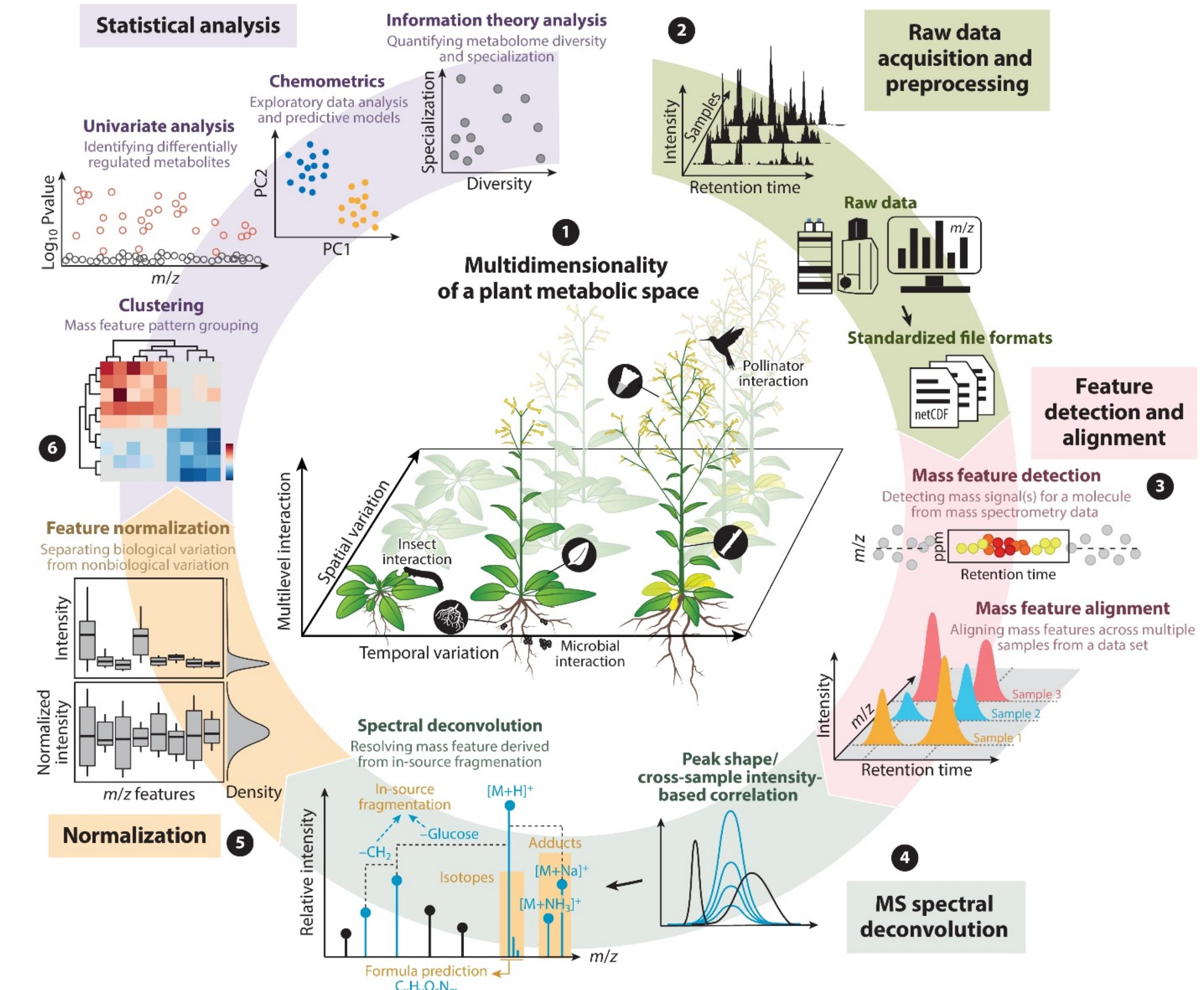
Next-Generation Mass Spectrometry Metabolomics Revives the Functional Analysis of Plant Metabolic Diversity

Annual Review of Plant Biology

Vol. 72:867-891 (Volume publication date June 2021)

First published as a Review in Advance on March 29, 2021

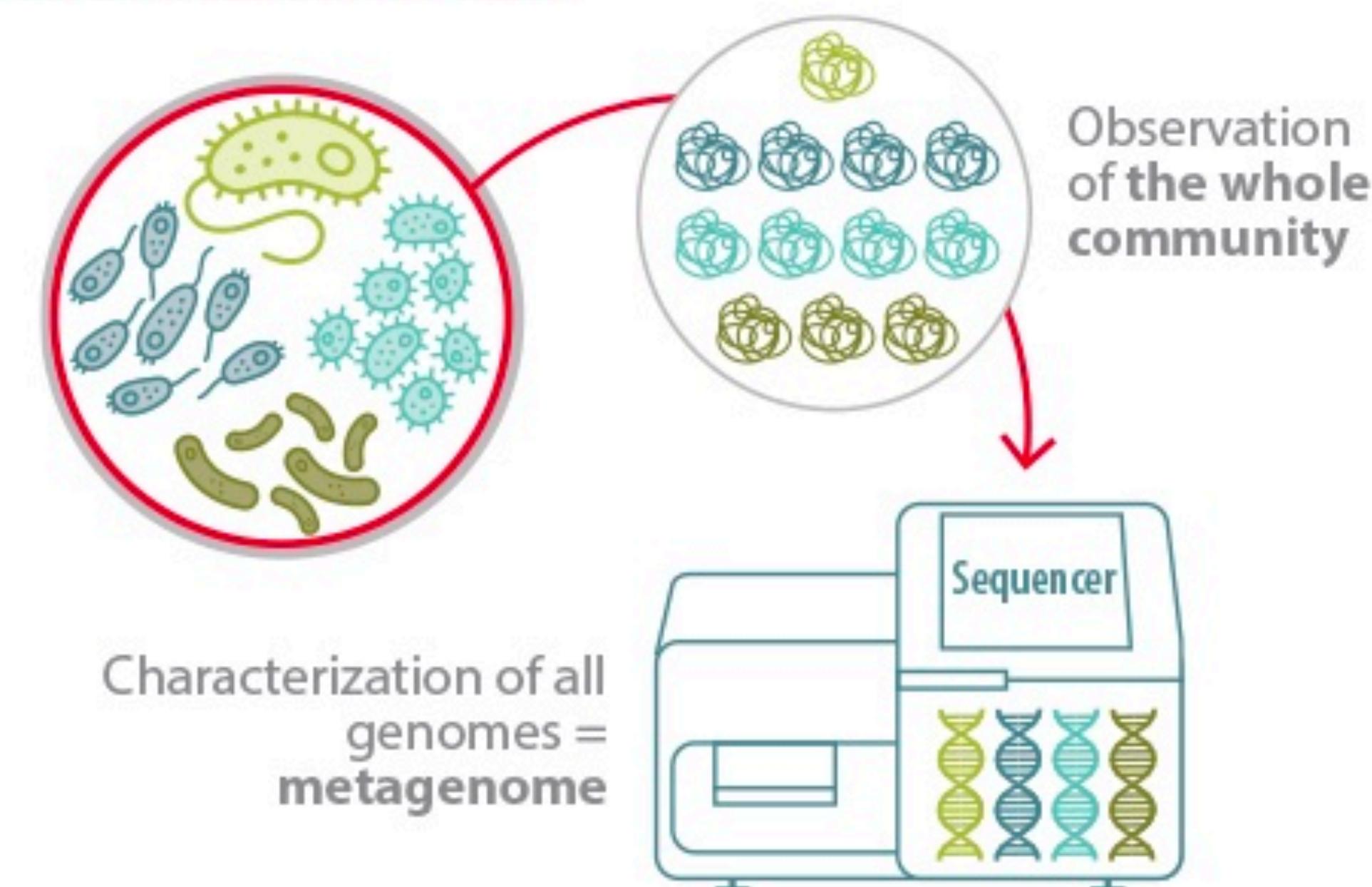
<https://doi.org/10.1146/annurev-arplant-071720-114836>



Metagenomics

the study of [genetic](#) material recovered directly from [environmental](#) samples.

What is metagenomics?



LALLEMAND

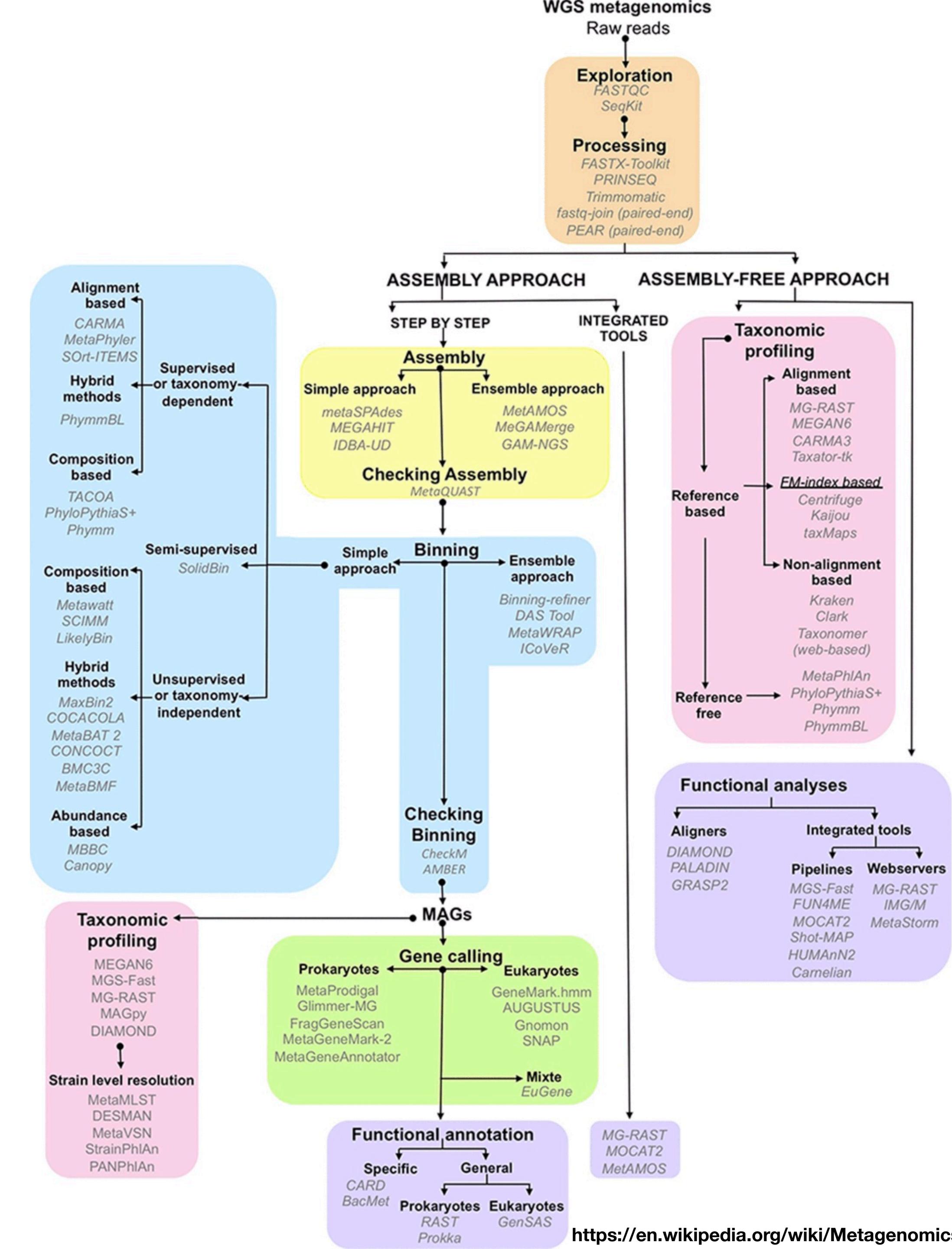


Table 1. A List of Benchmarked Classifiers and Their Various Characteristics

Type	Classifier	Custom Databases	Generates Abundance Profile	Memory Required	Time Required	Reference	Prophyle	yes	no	40 Gb	40 min	Břinda et al., 2017
DNA	Bracken	yes	yes	<1 Gb	<1 min	Lu et al., 2017	taxMaps	yes	yes	65 Gb	25 min	Corvelo et al., 2018
	Centrifuge	yes	yes	20 Gb	7 min	Kim et al., 2016	Kaiju	yes	yes	25 Gb	1 min	Menzel et al., 2016
	CLARK	yes	yes	80 Gb	2 min	Ounit et al., 2015	MMseqs2	yes	no	85 Gb	9 h	Steinegger and Söding, 2017
	CLARK-S	yes	yes	170 Gb	40 min	Ounit and Lonardi, 2016	Markers	MetaPhlAn2	no	2 Gb	1 min	Truong et al., 2015
	Kraken	yes	yes	190 Gb	1 min	Wood and Salzberg, 2014	mOTUs2	no	yes	2 Gb	1 min	Milanese et al., 2019
	Kraken2	yes	yes	36 Gb	1 min	Wood and Salzberg, 2014						
	KrakenUniq	yes	yes	200 Gb	1 min	Breitwieser et al., 2018						
	k-SLAM	yes	yes	130 Gb	2 h	Ainsworth et al., 2017						
	MegaBLAST	yes	no	61 Gb	4 h	Morgulis et al., 2008						
	metaOthello	no	no	30 Gb	1 min	Liu et al., 2018						
	PathSeq	yes ^a	no	140 Gb	5 min	Walker et al.,						

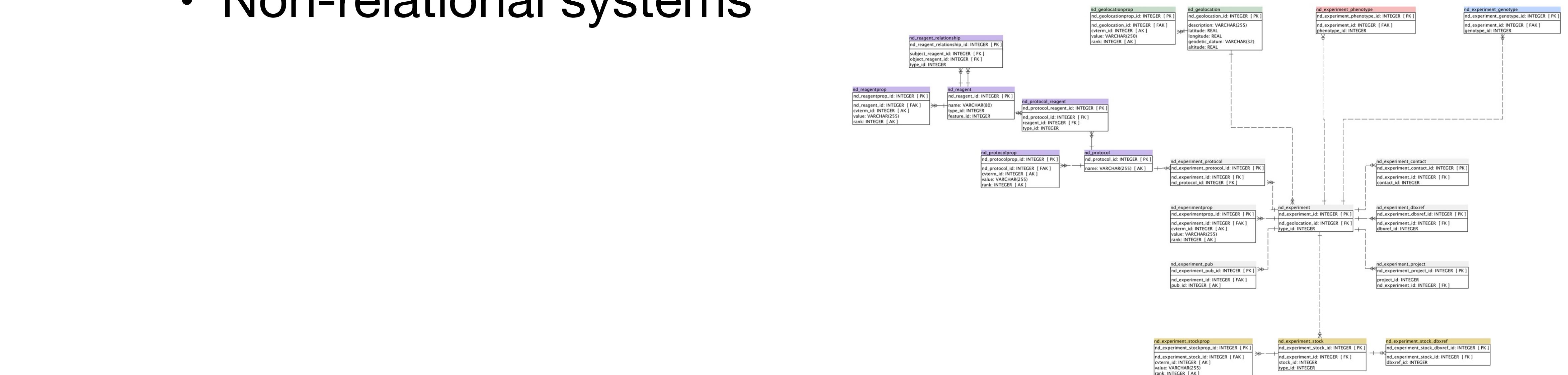
Benchmarking Metagenomics Tools for Taxonomic Classification

Simon H. Ye^{1, 2} , Katherine J. Siddle^{2, 3}, Daniel J. Park², Pardis C. Sabeti^{2, 3, 4, 5}

3. Databases

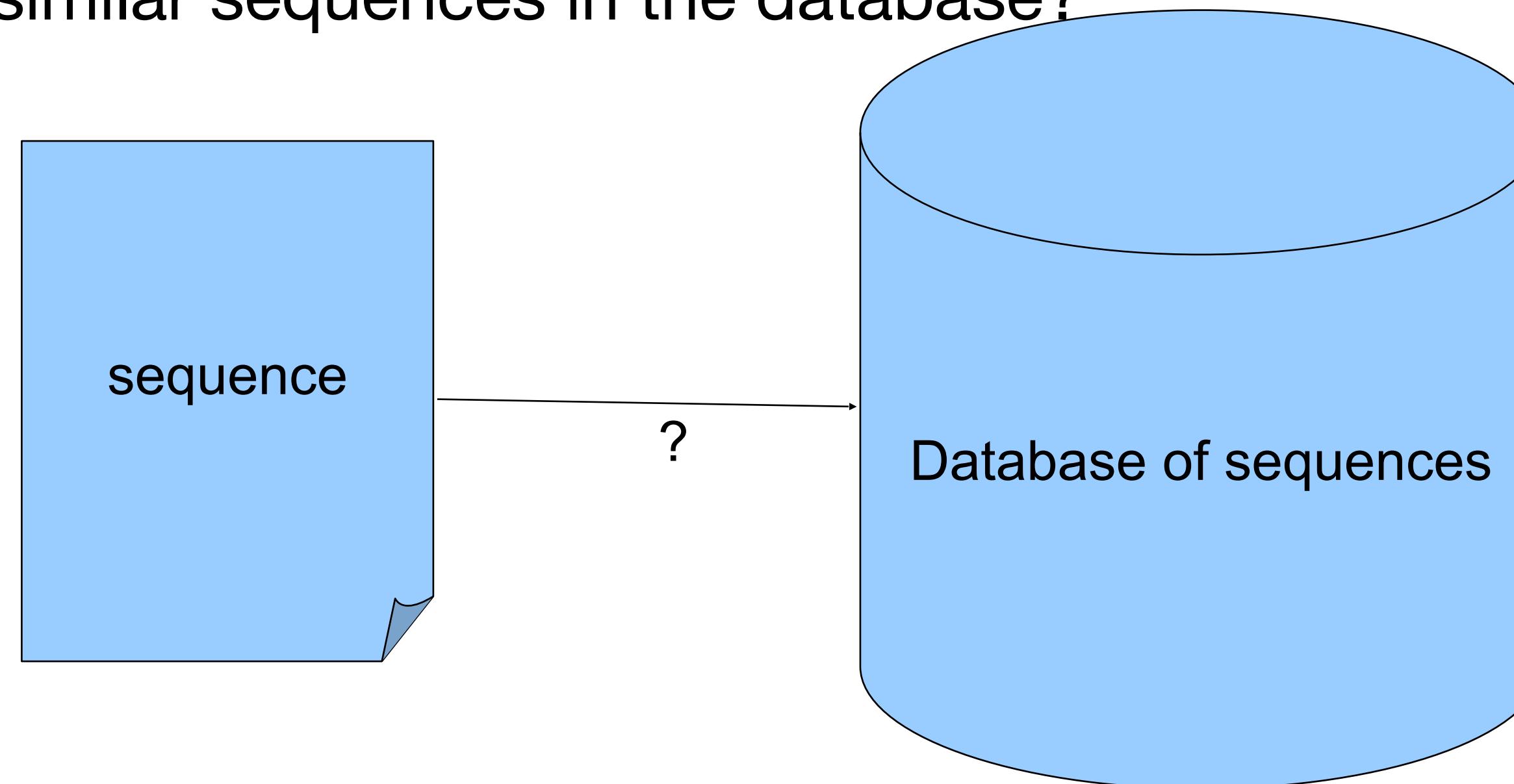
Databases

- Need to store and query data
- Biological data is highly structured
- Relational database systems (postgres, mysql)
- Database schemas - normalization
- Non-relational systems



Identify similar sequences

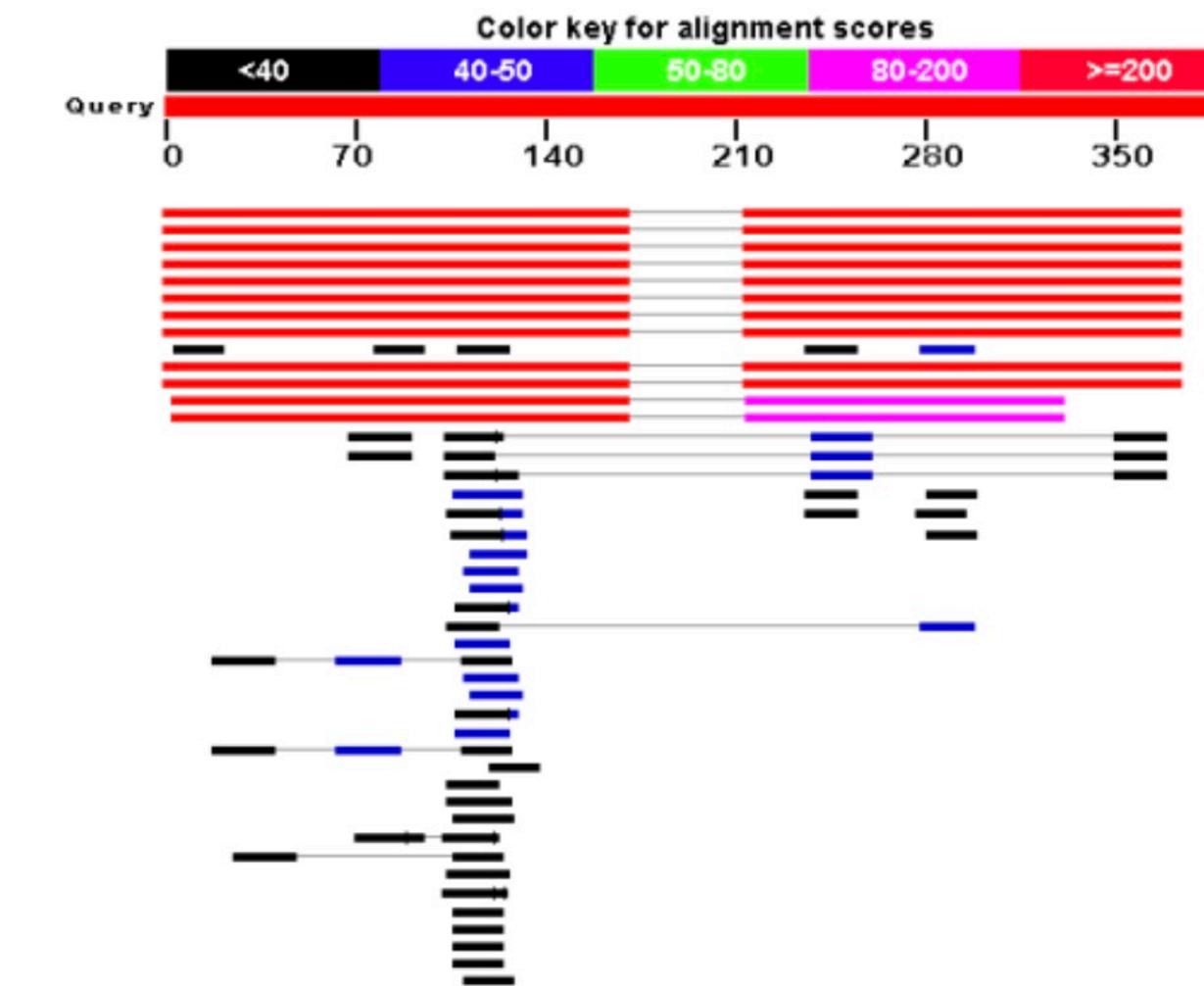
- Is the sequence (query) contained within a list of known sequences (database)?
- Are the similar sequences in the database?



- Compare the query sequence to each sequence in the database and calculate a similarity score, report sequences above a certain cutoff.

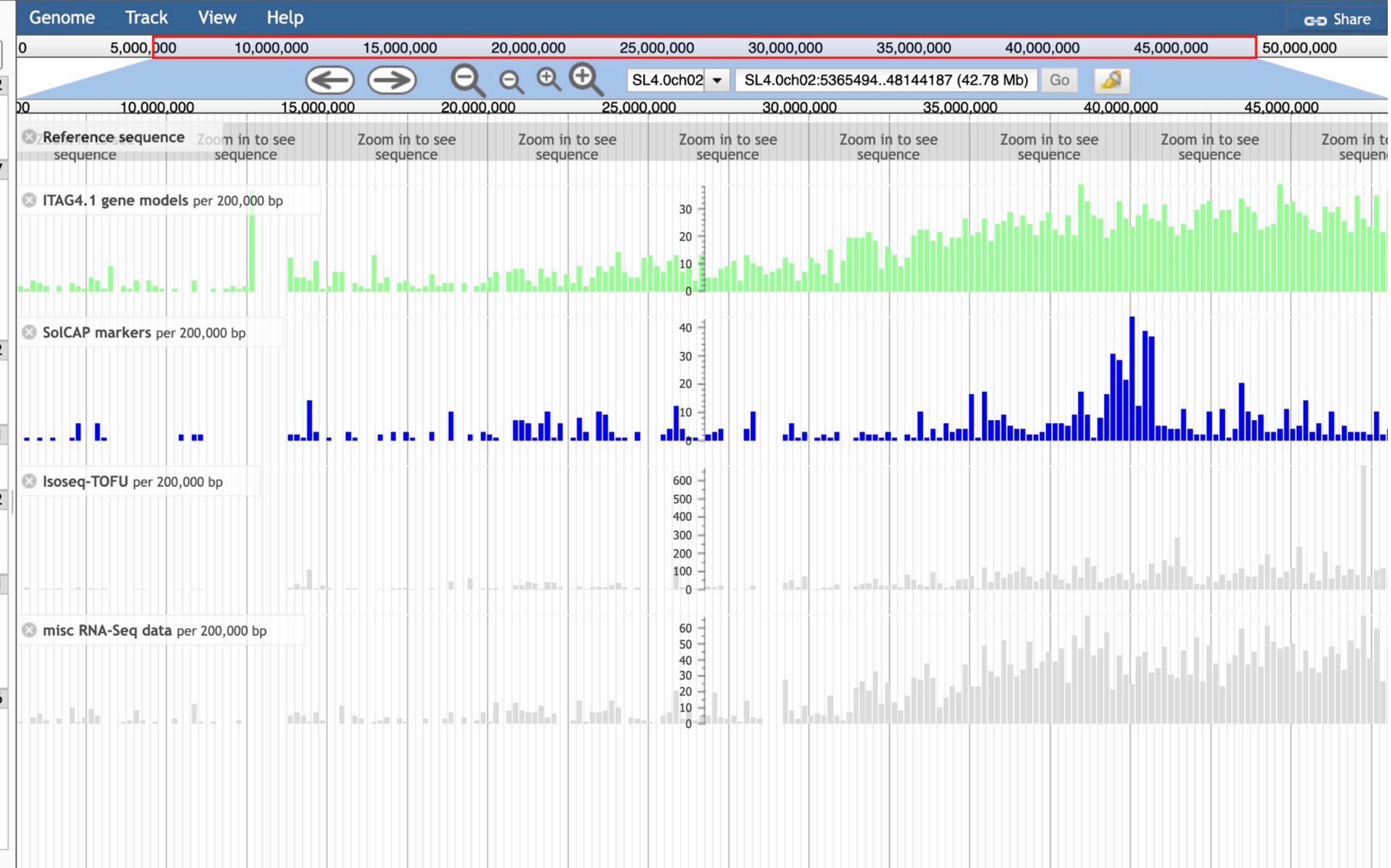
BLAST

- Basic Local Alignment and Search Tool
- Created by Altschul et al. at NCBI
- Homology search
 - Find similar sequences
 - Ranked by homology
 - Reports an e-value (expected number of hits by chance)
- A sequence *query* is run against many other sequences which form the *database*
- Offered by many websites



Name	Type	Web links
GenBank	Database	http://www.ncbi.nih.gov/entrez/query.fcgi?db=protein
RefSeq	Database	https://www.ncbi.nlm.nih.gov/refseq/
nr	Database	http://www.ncbi.nlm.nih.gov/BLAST/
UniProt	Database	http://www.pir.uniprot.org/
UniRef	Database	http://www.pir.uniprot.org/database/nref.shtml
UniParc	Database	http://www.pir.uniprot.org/database/archive.shtml
TrEMBL	Database	http://kr.expasy.org/sprot/
SwissProt	Database	http://kr.expasy.org/sprot/
PIR	Database	http://pir.georgetown.edu/
OWL	Database	http://www.bioinf.man.ac.uk/dbbrowser/OWL/

Genome browsers



4. Operating systems used in bioinformatics

Why use Unix?

- Free to use.
- Stable and secure
- Easily maintained
- Can edit and enhance source code (Open Source)
- World-wide community of developers
- Can be installed on just about anything
- Many bioinformatics programs will not run on Windows
- Many analysis servers run Linux

Linux Distributions

- Distributions have been created around the Linux kernel
- Many distros available (distrowatch.com)
- Examples:
 - RedHat
 - Debian (we will use this)
 - Ubuntu

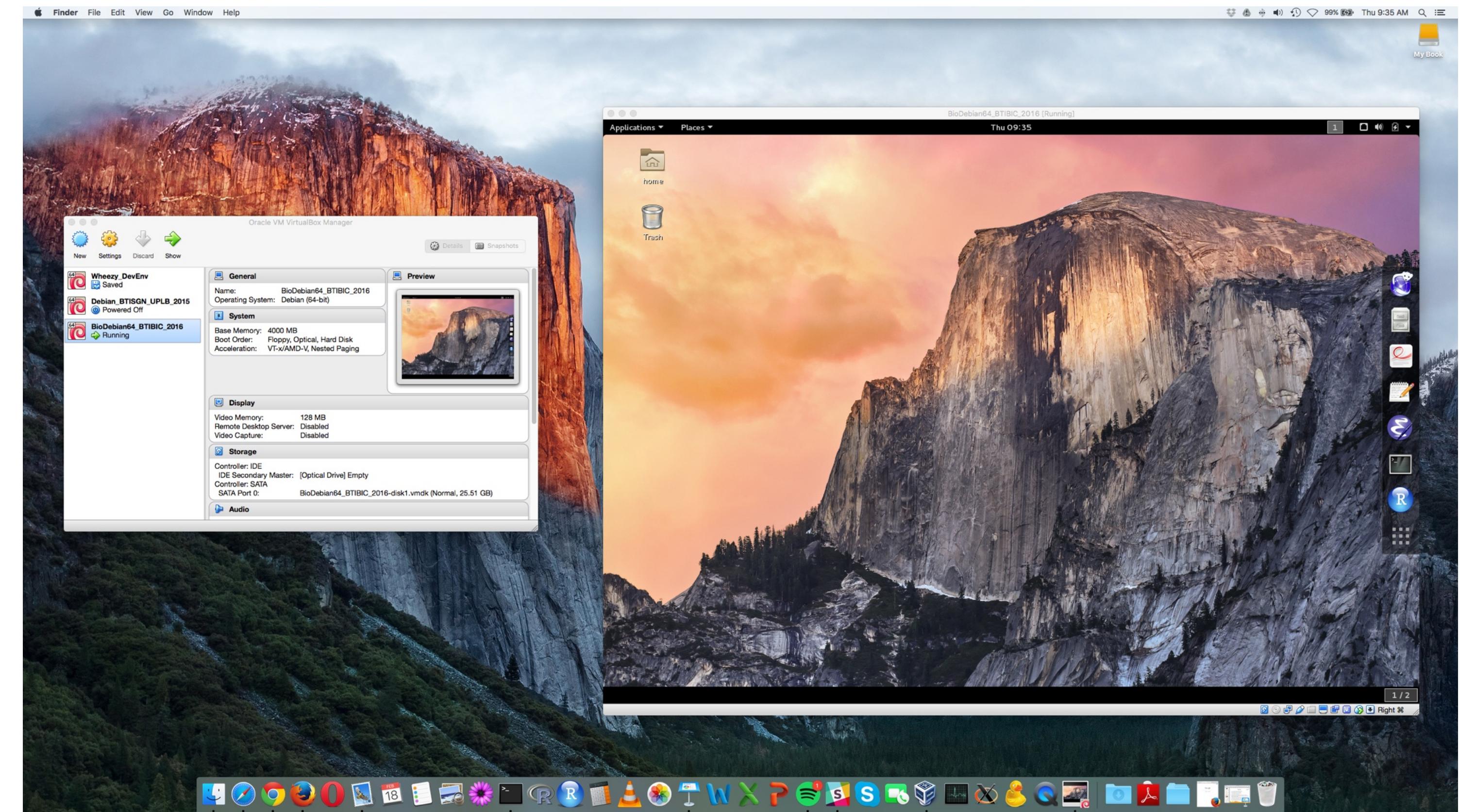
Fall 2020

Page Hit Ranking			Page Hit Ranking		
Data span:			Data span:		
Last 6 months			Last 6 months		
Rank	Distribution	HPD*	Rank	Distribution	HPD*
1	MX Linux	3823▼	1	MX Linux	3161▼
2	Manjaro	2640▼	2	EndeavourOS	2627▲
3	Mint	2368▼	3	Manjaro	2229▼
4	Ubuntu	1670▼	4	Mint	1928▲
5	Pop!_OS	1398▲	5	Pop!_OS	1542-
6	Debian	1377▼	6	Ubuntu	1327▼
7	elementary	1341▼	7	Debian	1259-
8	Solus	1050▼	8	Garuda	1152▲
9	Fedora	1000▼	9	elementary	1129-
10	Zorin	921▼	10	Fedora	958▲
11	KDE neon	904▼	11	Zorin	859▲
12	deepin	893▼	12	openSUSE	819▼
13	openSUSE	782-	13	KDE neon	692▲
14	EndeavourOS	764▲	14	Solus	607▼
15	Ubuntu Kylin	723-	15	antiX	543▲
16	Arch	701▼	16	Arch	503▼
17	antiX	662▼	17	Slackware	463▲
18	CentOS	599▼	18	Lite	459▼
19	Linuxfx	563▲	19	Artix	438-
20	ArcoLinux	552▼	20	PCLinuxOS	429-
21	PCLinuxOS	531-	21	Kali	412-
22	Puppy	519▲	22	Puppy	404-
23	Kali	509-	23	deepin	398▼

<https://distrowatch.com/>

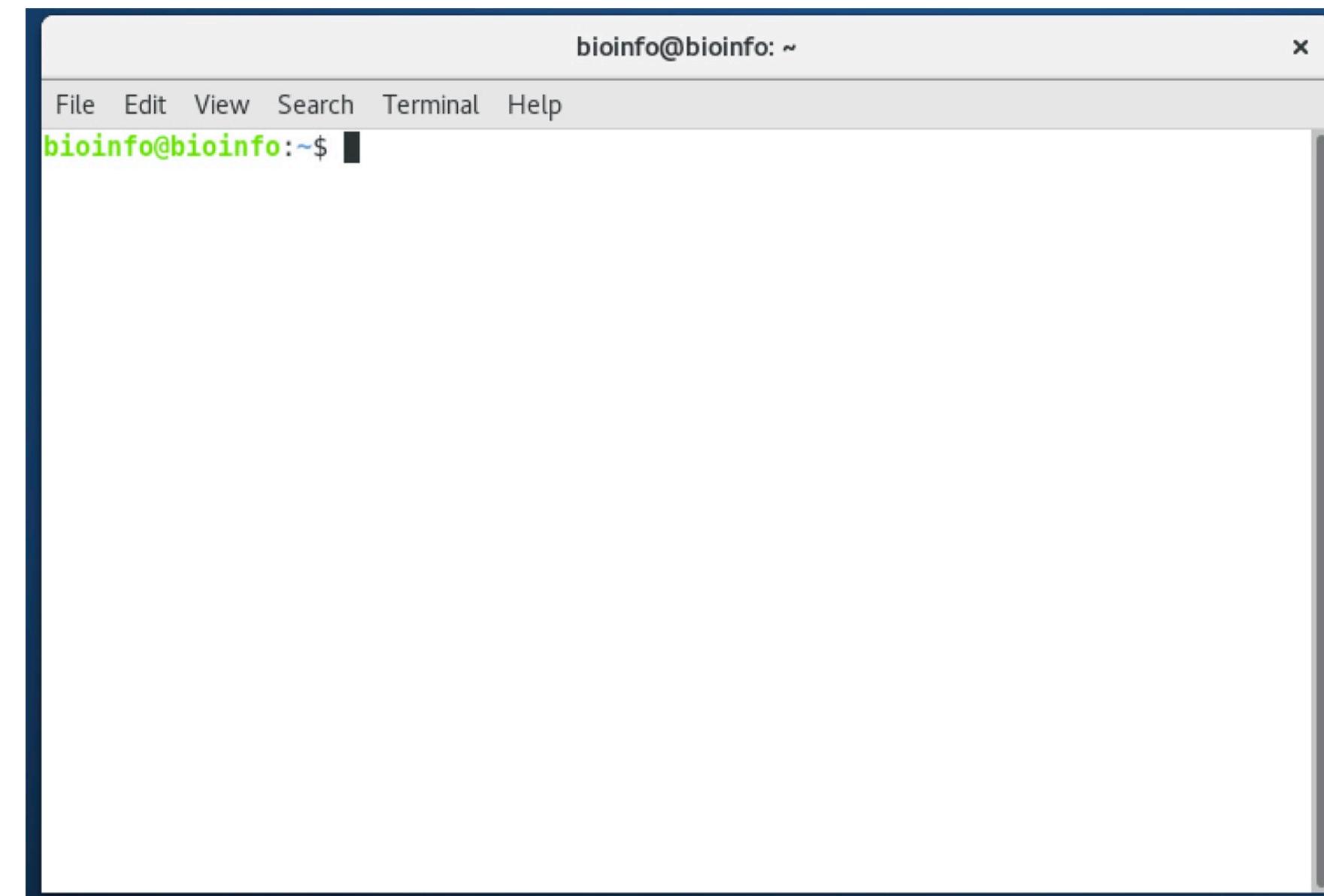
What is a Virtual Machine?

- A VM is an operating system which has been installed inside a simulated environment.
- Unlike emulators, virtual machines interface with real hardware.



What is a Terminal?

- A terminal is a textual interface for interacting with a computer.
- Using the terminal, one can issue powerful and concise command-line instructions for the computer to follow.

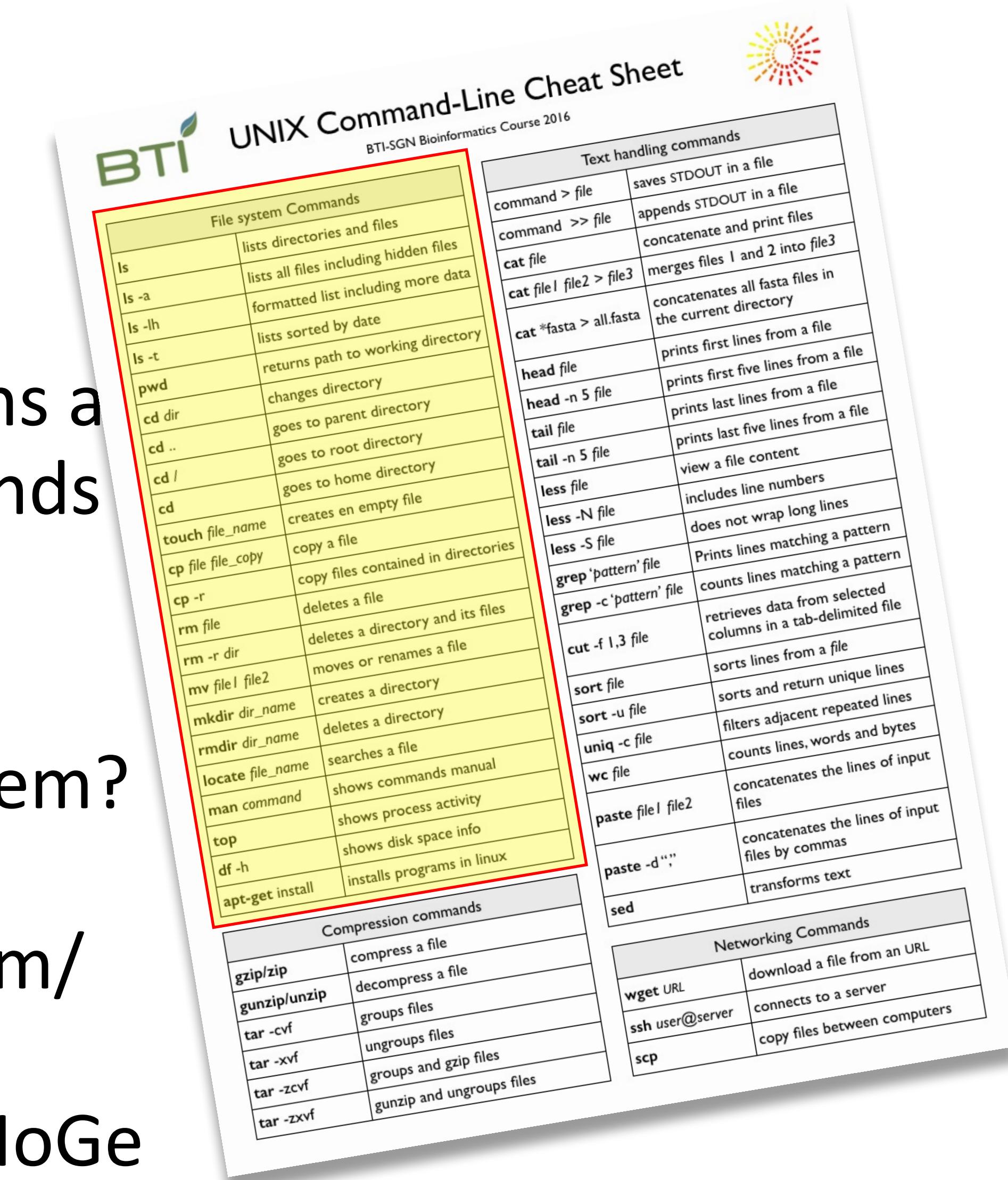


Why use the command-line (terminal)?

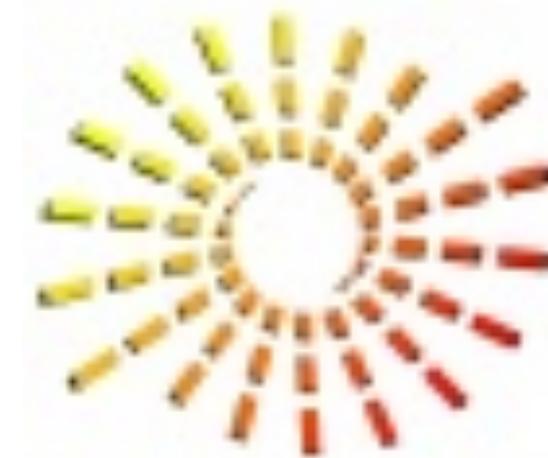
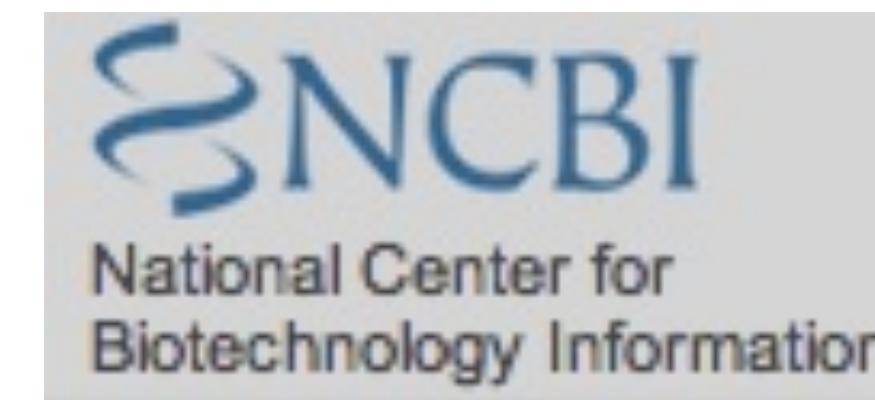
- Most software for biological data analysis is used through UNIX command-line operations.
- Most of the servers for biological data analysis use Linux/Unix as their operating system.
- Data analysis on calculation servers are much faster since we can use more CPUs and RAM than in a PC or laptop (*e.g.* BTI's "Boyce" server has 64 cores and 1TB RAM)
- Large NGS data files can not be opened or loaded in most graphical software and web sites.

Command-line File System Navigation.

- The cheat sheet contains a list of common commands for navigating the file system.
- But what *is* the file system?
- <https://drive.google.com/file/d/1Z8ryxpO-xyYhulfrWUfVsCS0YuZHoGeR/view>



Web-based bioinformatics

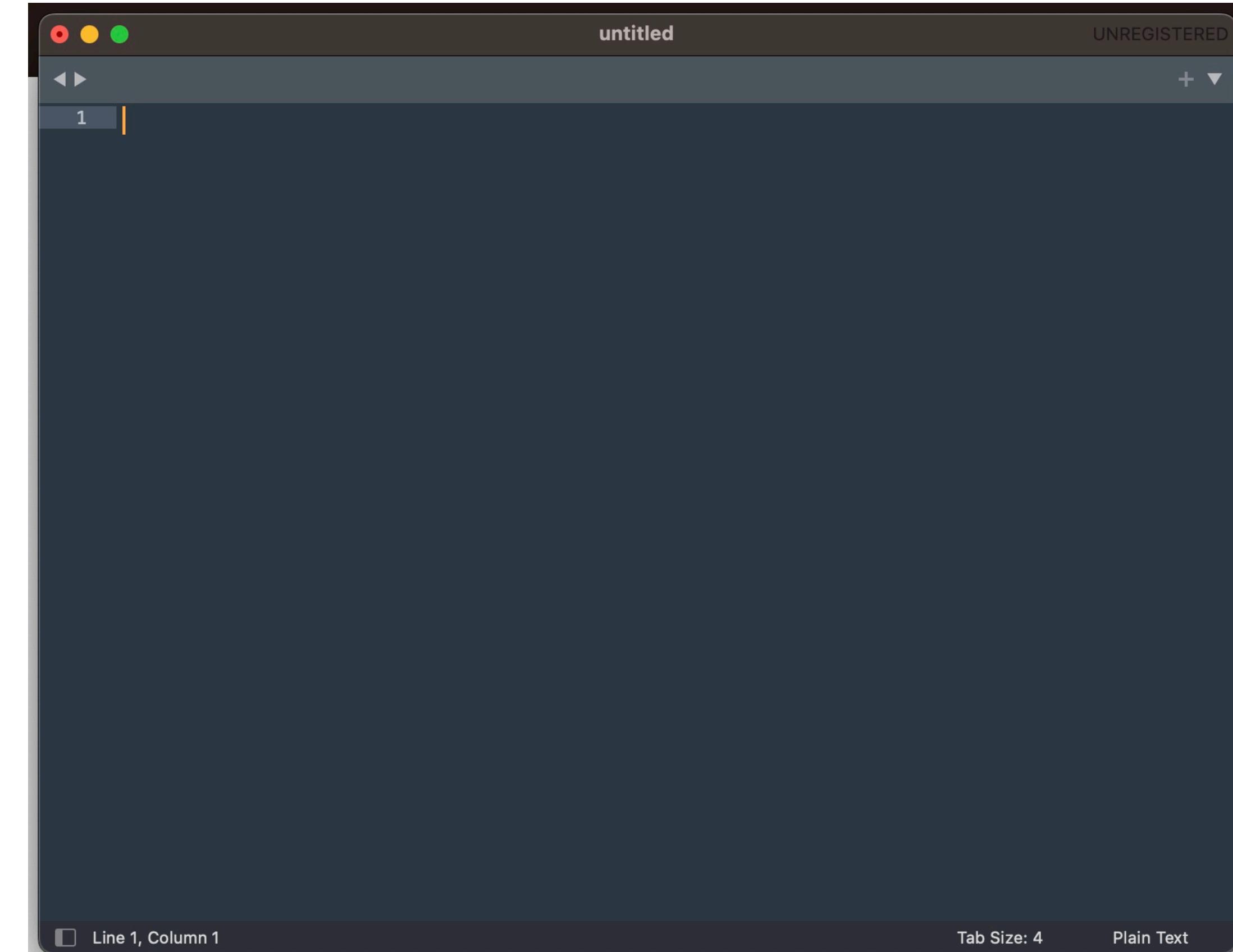


*Web-based not as powerful

5. Programming languages used in bioinformatics

Text Editors for Programming

- Don't write code in Word!
- Vi, Vim, Emacs
- Sublime
- Visual Studio



What are scripts?

- A text document that contains instructions to be executed by a program
- Usually combine existing components
- Usually use an interpreter instead of compiler to translate to machine code
 - **Bash, Python, Perl, JavaScript, PHP, Ruby**

```
1 #!/bin/sh
2 #####
3 # pipeline for running hisat2 on paired end files #
4 #
5 # usage:
6 #
7 #      hisat_pe_annot.sh $base_dir $CPU
8 #
9 #####
10
11 cd $1
12
13 #make a dir for your output and copy or symlink some files there
14 mkdir /scratch/annotation_output
15 cd /scratch/annotation_output
16 cp /scratch/Botany2020NMGWorkshop/annotation/2transfer/contig_15.fasta .
17 ln -s /scratch/Botany2020NMGWorkshop/annotation/2transfer/*.fastq .
18
19 #index reference fasta file; We will use only contig_15 for demo purposes
20 /opt/hisat-genotype-top/hisat2-build contig_15.fasta contig_15
21
22 #map RNA-seq reads to reference genome fasta file
23 for file in `dir -d *_1.fastq` ; do
```

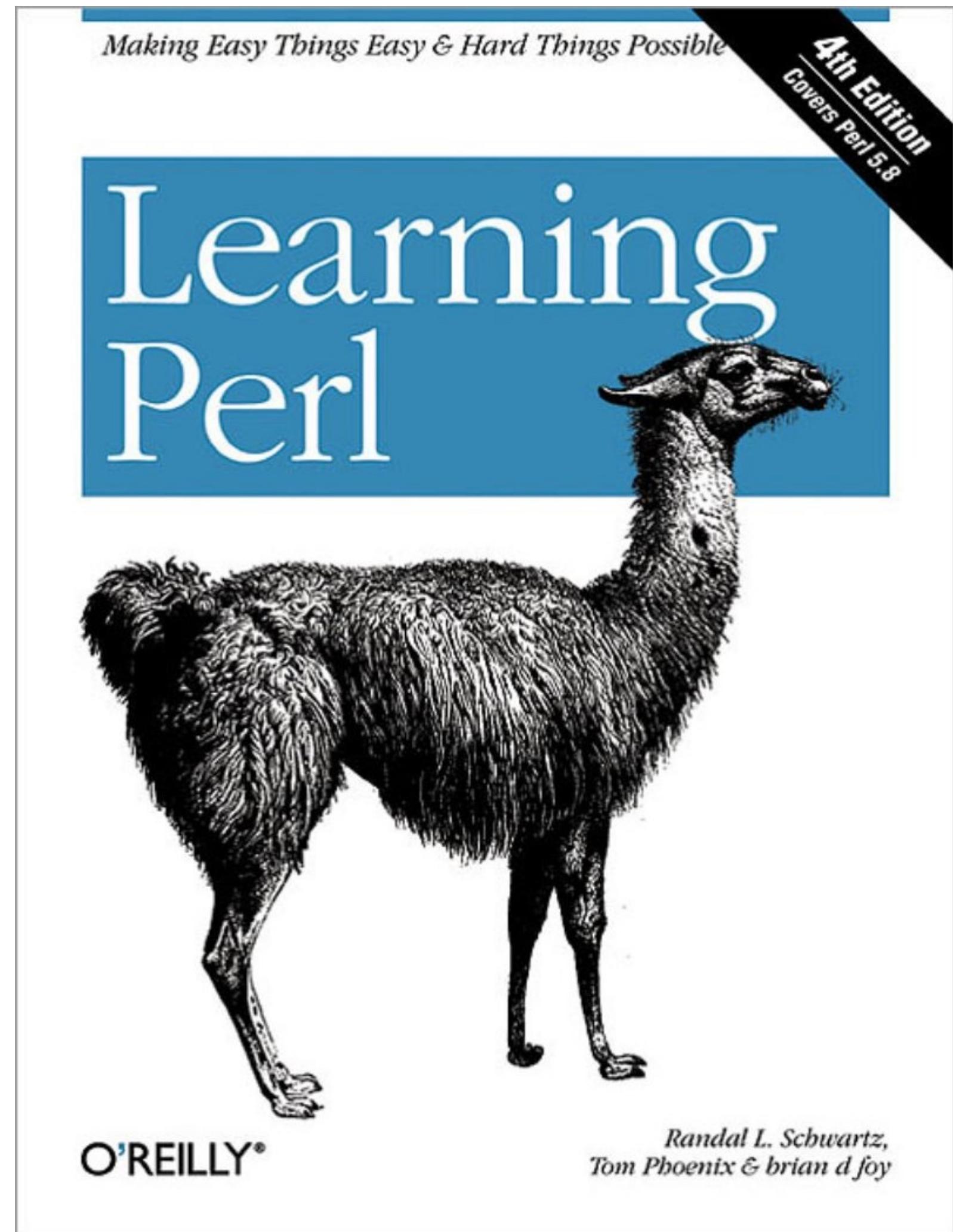
Why use bash?

- It is useful for chaining other programs together.
- This allows one to automate tasks and build pipelines.

```
1 #!/bin/sh
2
3 #move to annotation directory
4 cd /opt/annotation
5
6 #get some files
7 cp /scratch/Botany2020NMGWorkshop/annotation/2transfer/list.txt .
8 cp /scratch/Botany2020NMGWorkshop/annotation/2transfer/uniprot_sprot_plants.fasta .
9
10 #run Portcullis: https://bioconda.github.io/recipes/portcullis/README.html
11 #docker pull maplesond/portcullis:stable #already done with Adrian
12 docker run --rm -v /scratch/annotation_output:/data maplesond/portcullis:stable portcullis full -t 7 -v /data/contig_15.fasta /data/SRR5046448_contig15.sort.bam
13
14 #run Mikado pipeline: https://bioconda.github.io/recipes/mikado/README.html
15 #docker pull cyverseuk/mikado #already done with Adrian
16 docker run --rm -v /scratch/annotation_output:/data cyverseuk/mikado mikado configure --list list.txt --reference data/contig_15.fasta --mode permissive --scoring-method jaccard
17 docker run --rm -v /scratch/annotation_output:/data cyverseuk/mikado mikado prepare --json-conf /data/configuration.yaml
18
19 makeblastdb -in uniprot_sprot_plants.fasta -dbtype prot
20 blastx -max_target_seqs 5 -num_threads 7 -query mikado_prepared.fasta -outfmt 5 -db uniprot_sprot_plants.fasta -evalue 0.000001 2> blast.log | sed '/^$/d' | gzip
21
22 screen -L /opt/TransDecoder-TransDecoder-v5.5.0/TransDecoder.LongOrfs -t mikado_prepared.fasta
23
24 screen -L hmmscan --cpu 7 --domtblout pfam.domtblout Pfam-A.hmm mikado_prepared.fasta.transdecoder_dir/longest_orfs.pep
25
26 blastp -query mikado_prepared.fasta.transdecoder_dir/longest_orfs.pep -db uniprot_sprot_plants.fasta -max_target_seqs 1 -outfmt 6 -evalue 1e-5 -num_threads 7 >
27
28 /opt/TransDecoder-TransDecoder-v5.5.0/TransDecoder.Predict -t mikado_prepared.fasta --retain_blastp_hits blastp.outfmt6 --cpu 7 --retain_pfam_hits pfam.domtblout
29
30 docker run --rm -v /scratch/annotation_output:/data cyverseuk/mikado mikado serialise --json-conf /data/configuration.yaml --xml /data/mikado.blast.xml.gz --orf
```

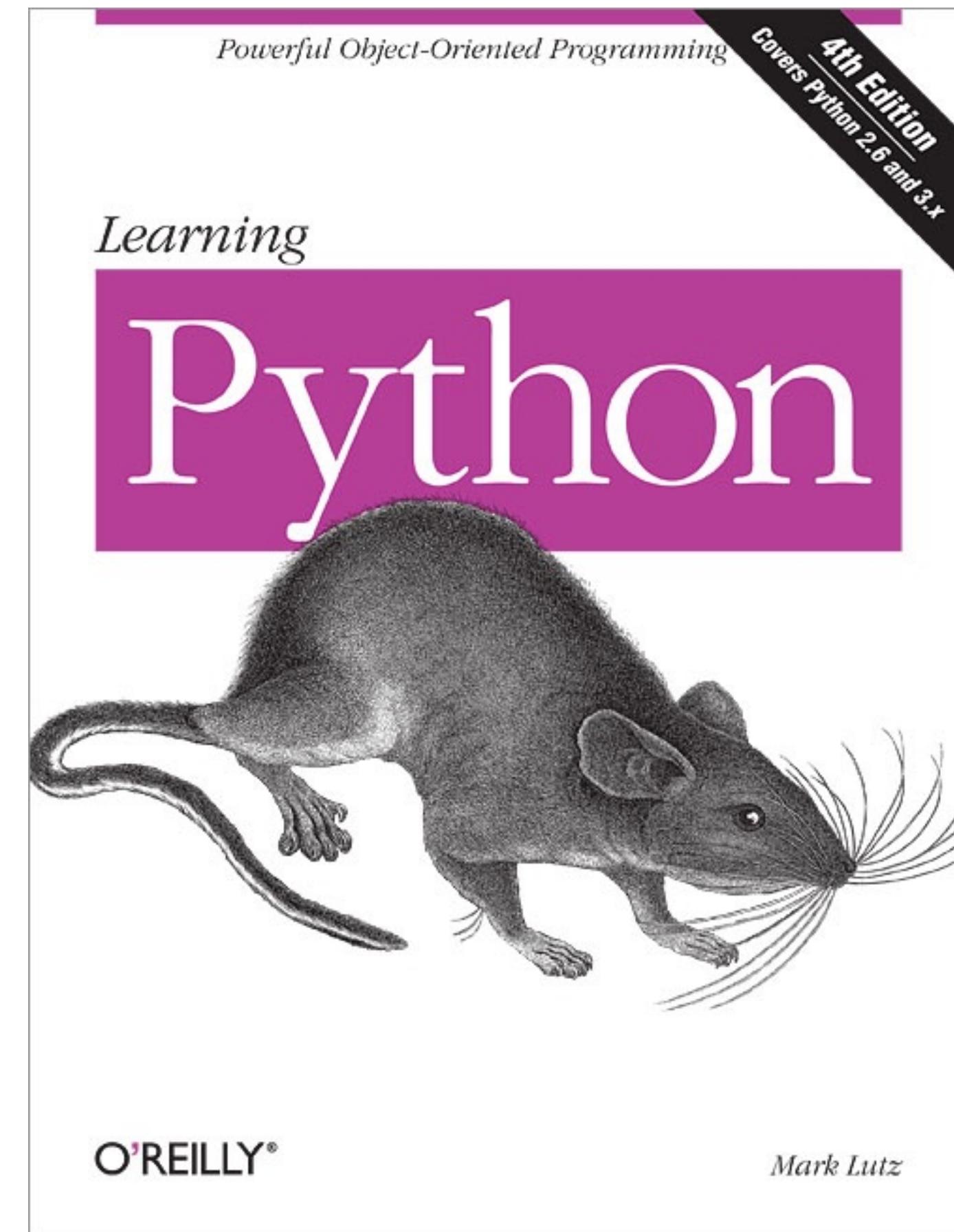
Perl

- Developed since 1980s by Larry Wall
- Useful for bioinformatics and web development
- Support for *objects*
- Excellent integration of *regular expressions* (text handling language)
- Vast open source code library (<http://cpan.org/>)
 - BioPerl (<http://bioperl.org/>)
- Easy to learn
- <http://www.perl.org/>



Python

- Created by Guido van Rossum in 1989
- Very elegant language
- BioPython libraries
- The “new” popular language
- Many frameworks (Django for web etc.)



Virtual Environments

- Virtual environments are isolated from system directories
- Each has its own Python binary and has its own independent set of packages
- Example:

```
(base) bioinfo@bioinfo:~/Scripts/python$ conda create --name testpy3 python=3.8
(base) bioinfo@bioinfo:~/Scripts/python$ conda activate testpy3
(testpy3) bioinfo@bioinfo:~/Scripts/python$ less ~/.conda/environments.txt
```

Python Package managers

- pip

```
pip install "SomeProject==1.4"
```

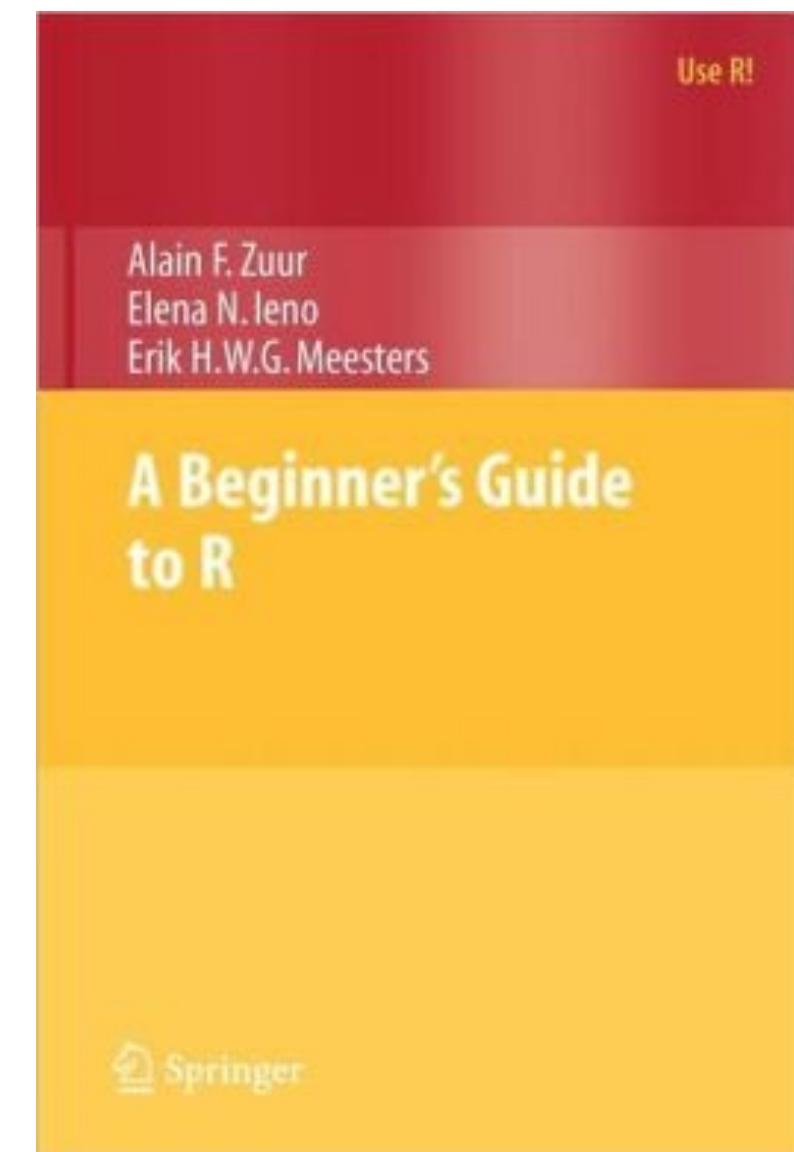
- Conda

```
conda install package-name=2.3.4
```

- Virtual environments



- Language designed for statistics
- Support for matrix calculations, graphics
- Expression analysis, Next-Gen sequence analysis, Graphics, genome annotation statistics, phylogeny
- Interactive
- “Bioconductor” package



R Studio

The screenshot shows the RStudio interface with four main panes:

- Source:** The top-left pane shows an untitled R script file with a single line of code: "1".
- Console:** The bottom-left pane displays the R console output, including the welcome message and various help texts.
- Environment:** The top-right pane shows the Global Environment tab with a list of objects.
- Packages:** The bottom-right pane shows the Packages tab listing available packages in the User Library.

Annotations:

- 1. Source: where you write your code.** Points to the Source pane.
- 2. Console: where your code is evaluated.** Points to the Console pane.
- 3. Environment/History** Points to the Environment pane.
- 4. Files/Plots/Packages/Help** Points to the Packages pane.

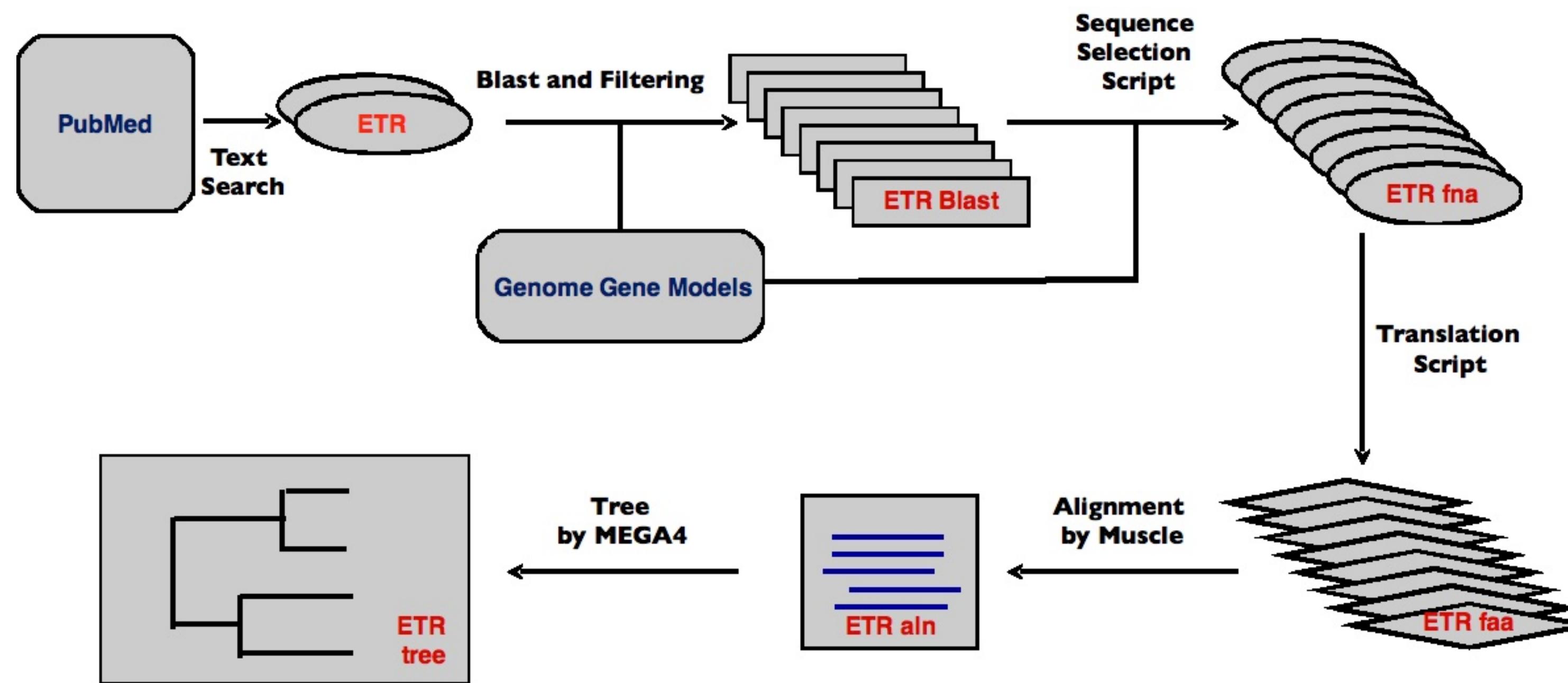
Name	Description	Version
assertthat	Easy-to-use assertion functions	0.2.0
BH	Boost C++ Header Files	1.66.0-1
Biobase	Biobase: Base functions for Bioconductor	2.34.0
BiocGenerics	S4 generic functions for Bioconductor	0.20.0
BiocInstaller	Install/Update Bioconductor, CRAN, and github Packages	1.24.0
BiocParallel	Bioconductor facilities for parallel evaluation	1.8.2
Biostrings	String objects representing biological sequences, and matching algorithms	2.42.1
bit	A class for vectors of 1-bit booleans	1.1-12
bit64	A S3 Class for Vectors of 64bit Integers	0.9-7
bitops	Bitwise Operations	1.0-6
blob	A Simple S3 Class for Representing Vectors of Binary Data ('BLOBS')	1.1.0
cli	Helpers for Developing Command Line Interfaces	1.0.0
crayon	Colored Terminal Output	1.3.4
DBI	R Database Interface	0.8

6. Data Best Practices

Data Best Practices

- Data should be stored in the appropriate repository
 - EX: SRA/Genbank, TreeBase
- Organizing schemas for your own data
- Backup storage, redundancy
- Compress large files to save disk space
- rsync, md5sum, diff, zless

Running “pipelines”



Code Libraries

- Modularizing code
 - Breaking apart problem into individual units
 - Easier to maintain
 - Reusable
- Modules, Packages - grouping of modules, Library - a collection of packages
- Libraries we will use:
 - BioPython, NumPy, Pandas in later classes
- Package Repos
 - Python Package Index (PyPi) and Anaconda

Jupyter Notebooks

- Starting the notebook server

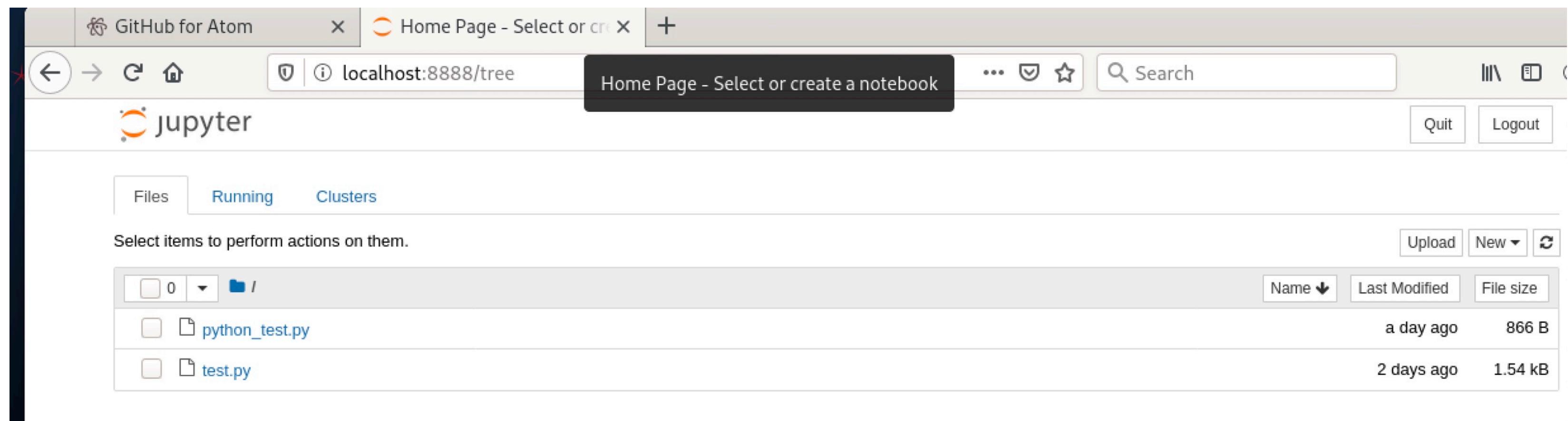
mkdir ~/Scripts/python

cd Scripts/python/

jupyter notebook

```
bioinfo@bioinfo: ~/Scripts/python
File Edit View Search Terminal Help
(base) bioinfo@bioinfo:~$ cd Scripts/python/
(base) bioinfo@bioinfo:~/Scripts/python$ jupyter
usage: jupyter [-h] [--version] [--config-dir] [--data-dir] [--runtime-dir]
                [--paths] [--json]
                [subcommand]
jupyter: error: one of the arguments --version subcommand --config-dir --data-di
r --runtime-dir --paths is required
(base) bioinfo@bioinfo:~/Scripts/python$ jupyter notebook
[I 20:05:44.744 NotebookApp] Writing notebook server cookie secret to /home/bioi
nfo/.local/share/jupyter/runtime/notebook_cookie_secret
[T 20:05:45.292 NotebookApp] JupyterLab extension loaded from /home/bioinfo/min
```

- Notebook Dashboard





GitHub

- GitHub is a repository hosting service.
- Has a command line tools and GUI options
- Wikis can be used for software documentation

- Collaborative coding
- Version control

A screenshot of the GitHub homepage. At the top, there's a navigation bar with links for Pull requests, Issues, Marketplace, and Explore. Below the navigation, a user profile for 'srs218' is shown, along with a search bar and a 'New' button. A list of repositories owned by 'srs218' is displayed, including 'bcbc-group/Botany2020NMGWorkshop', 'bcbc-group/PLSci7202', 'srs218/PLSci7202_test', 'srs218/jupyter-scipy-cc2020', 'srs218/manuscripts', and 'srs218/phylogeny'. On the right side of the page, there's a large, semi-transparent overlay for a 'Hello World' guide. The overlay features the text 'Learn Git and GitHub without any code!', 'Using the Hello World guide, you'll create a repository, start a branch, write comments, and open a pull request.', a green 'Read the guide' button, and a white 'Start a project' button.

Versioning

The screenshot shows a GitHub repository page for `bcbc-group / Botany2021NMGWorkshop`. The repository is public and has 1 branch and 0 tags. The main branch is `main`, which has 46 commits. A red arrow points to the `Add file` button in the top right corner of the commit list. Another red arrow points to the bottom right corner of the page, indicating where to click to view more commits.

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main 1 branch 0 tags Go to file Add file Code

Commits

srs218 Merge branch 'main' of https://github.com/bcbc-group/Botany2021NMGWor... 6f4a90d on Jul 18 46 commits

1.Introduction Create Introduction.pdf 3 months ago

2.BaselineSkills Add files via upload 3 months ago

3.ExperimentalDesign Create Botany_Workshop_2021.pdf 3 months ago

4.GenomeAssembly Folder name change

5.Annotation added papers

6.DownstreamAnalyses Add files via upload

.DS_Store files added

.gitattributes Initial commit

BTI_BCBC_Bioinfo_unix_command_... Add files via upload

README.md Update README.md

Schedule.pdf Added schedule as pdf

~\$hedula.docx Added schedule as pdf

Merge branch 'main' of https://github.com/bcbc-group/Botany2021NMGWor... srs218 committed on Jul 18

Create Botany_Workshop_2021.pdf srs218 committed on Jul 18

Add files via upload afpowell committed on Jul 18

Add files via upload afpowell committed on Jul 18

Add files via upload afpowell committed on Jul 18

Add files via upload afpowell committed on Jul 18

Add files via upload afpowell committed on Jul 18

added papers srs218 committed on Jul 18

Merge branch 'main' of https://github.com/bcbc-group/Botany2021NMGWor... srs218 committed on Jul 18

6f4a90d a948f12 f4a1f5c 422cca3 873b976 77fd76c 4088253 531d089

Reproducible Research

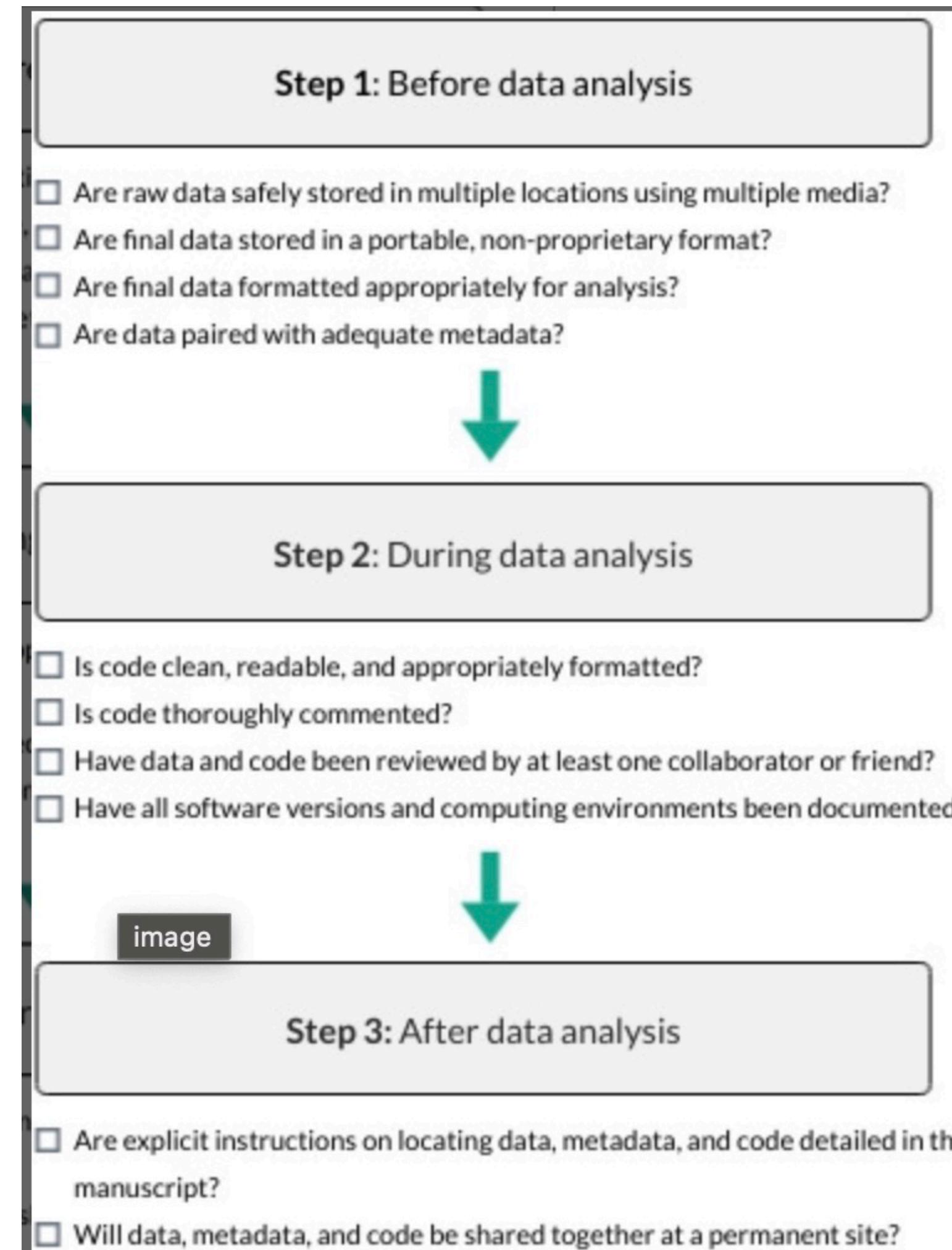
- No “custom scripts”
- Versioning of software and datasets
- System architecture and software
- Making code publicly available
- Raw data is available
- Everything is provided to produce a consistent result
- Reproducible research increases dependability of results and make it easier and more efficient to continue the research.

Reproducible Research

	Free	Open source	Website	
Data and code management				
Version control				
GitHub	Y [†]	N	https://github.com	
BitBucket	Y [†]	N	https://bitbucket.com	
GitLab	Y [†]	Y	https://www.gitlab.com	
Make				
GNU Make	Y	Y	https://www.gnu.org/software/make/	
Software containers and virtual machines				
Docker	Y	Y	https://docker.com	
Singularity	Y [†]	Y	https://syslabs.io	
Oracle VM VirtualBox	Y	Y	https://virtualbox.org	
Sharing research				
Preprint servers				
ArXiv	Y		https://arxiv.org/	
bioRxiv	Y		https://www.biorxiv.org/	
Preprint servers				
ArXiv			Y	https://arxiv.org/
bioRxiv			Y	https://www.biorxiv.org/
EcoEvoRxiv			Y	https://ecoevorxiv.org/
Manuscript creation				
Overleaf		Y [†]	Y	https://overleaf.com
TeXstudio		Y	Y	https://www.texstudio.org/
Rstudio		Y	Y	https://rstudio.org
Data repositories				
Dryad			N	https://datadryad.org/
Figshare			Y [†]	https://figshare.com/
Zenodo			Y	https://zenodo.org/
Open Science Framework			Y	https://osf.io/

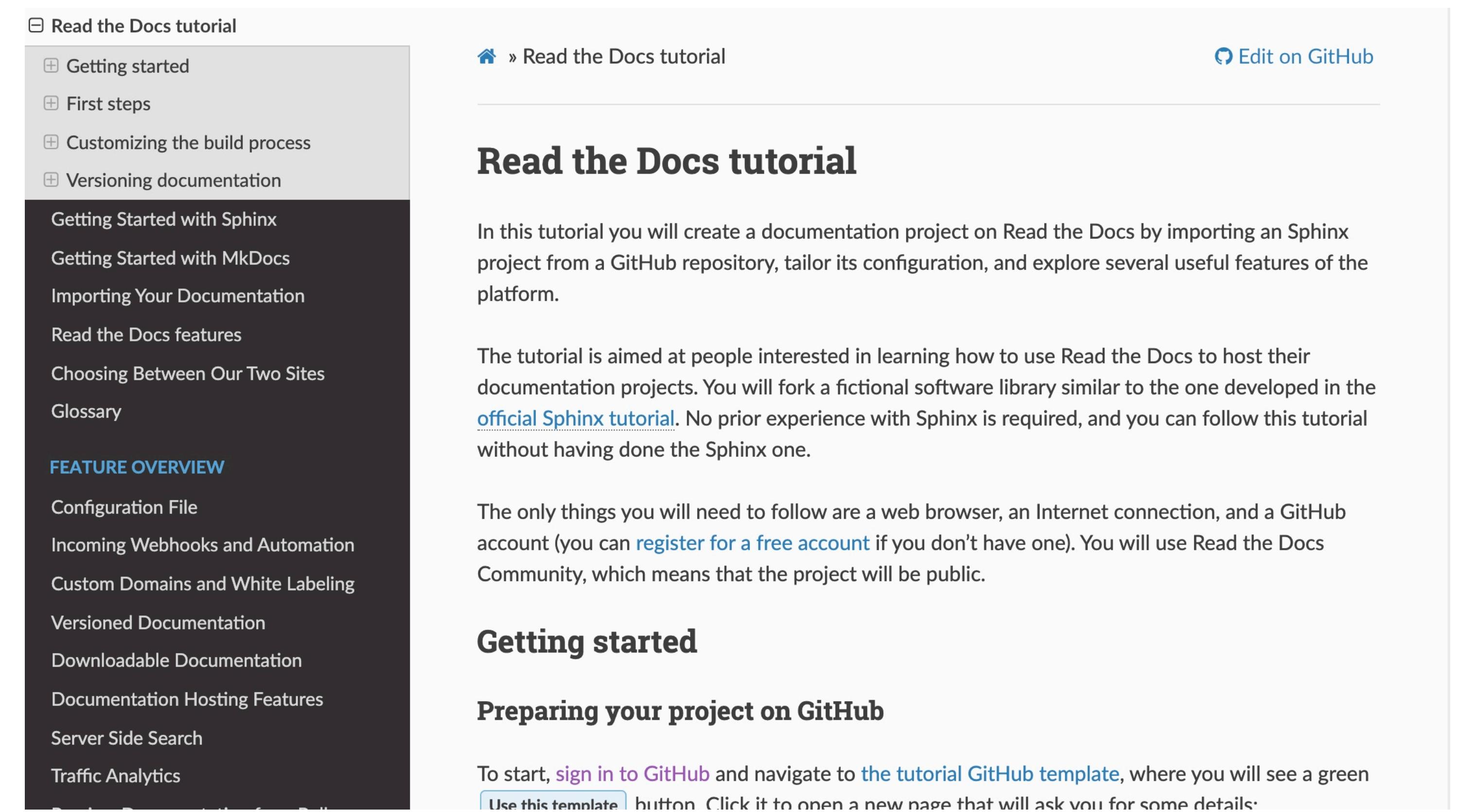
Alston and Rick, 2021

Reproducible Research



Documentation

- ReadTheDocs
- Jupyter / Markdown
- Commenting code
- README
- Documenting data - where it came from



The screenshot shows a documentation page for a "Read the Docs tutorial". The left sidebar contains a navigation menu with sections like "Read the Docs tutorial", "Getting started", "First steps", "Customizing the build process", "Versioning documentation", "Getting Started with Sphinx", "Getting Started with MkDocs", "Importing Your Documentation", "Read the Docs features", "Choosing Between Our Two Sites", and "Glossary". Below this is a "FEATURE OVERVIEW" section with links to "Configuration File", "Incoming Webhooks and Automation", "Custom Domains and White Labeling", "Versioned Documentation", "Downloadable Documentation", "Documentation Hosting Features", "Server Side Search", and "Traffic Analytics". The main content area has a breadcrumb trail ("Read the Docs tutorial") and a "Edit on GitHub" button. The title "Read the Docs tutorial" is displayed prominently. The text in the main content area describes the purpose of the tutorial: "In this tutorial you will create a documentation project on Read the Docs by importing an Sphinx project from a GitHub repository, tailor its configuration, and explore several useful features of the platform." It also states that the tutorial is aimed at people interested in learning how to use Read the Docs to host their documentation projects, mentioning the official Sphinx tutorial as a resource. The text concludes by noting that no prior experience with Sphinx is required and that the tutorial can be followed without having done the Sphinx one. The sidebar also includes a "Read the Docs tutorial" link.

Read the Docs tutorial

In this tutorial you will create a documentation project on Read the Docs by importing an Sphinx project from a GitHub repository, tailor its configuration, and explore several useful features of the platform.

The tutorial is aimed at people interested in learning how to use Read the Docs to host their documentation projects. You will fork a fictional software library similar to the one developed in the [official Sphinx tutorial](#). No prior experience with Sphinx is required, and you can follow this tutorial without having done the Sphinx one.

The only things you will need to follow are a web browser, an Internet connection, and a GitHub account (you can [register for a free account](#) if you don't have one). You will use Read the Docs Community, which means that the project will be public.

Getting started

Preparing your project on GitHub

To start, [sign in to GitHub](#) and navigate to [the tutorial GitHub template](#), where you will see a green [Use this template](#) button. Click it to open a new page that will ask you for some details:

Containers

- What is a container?

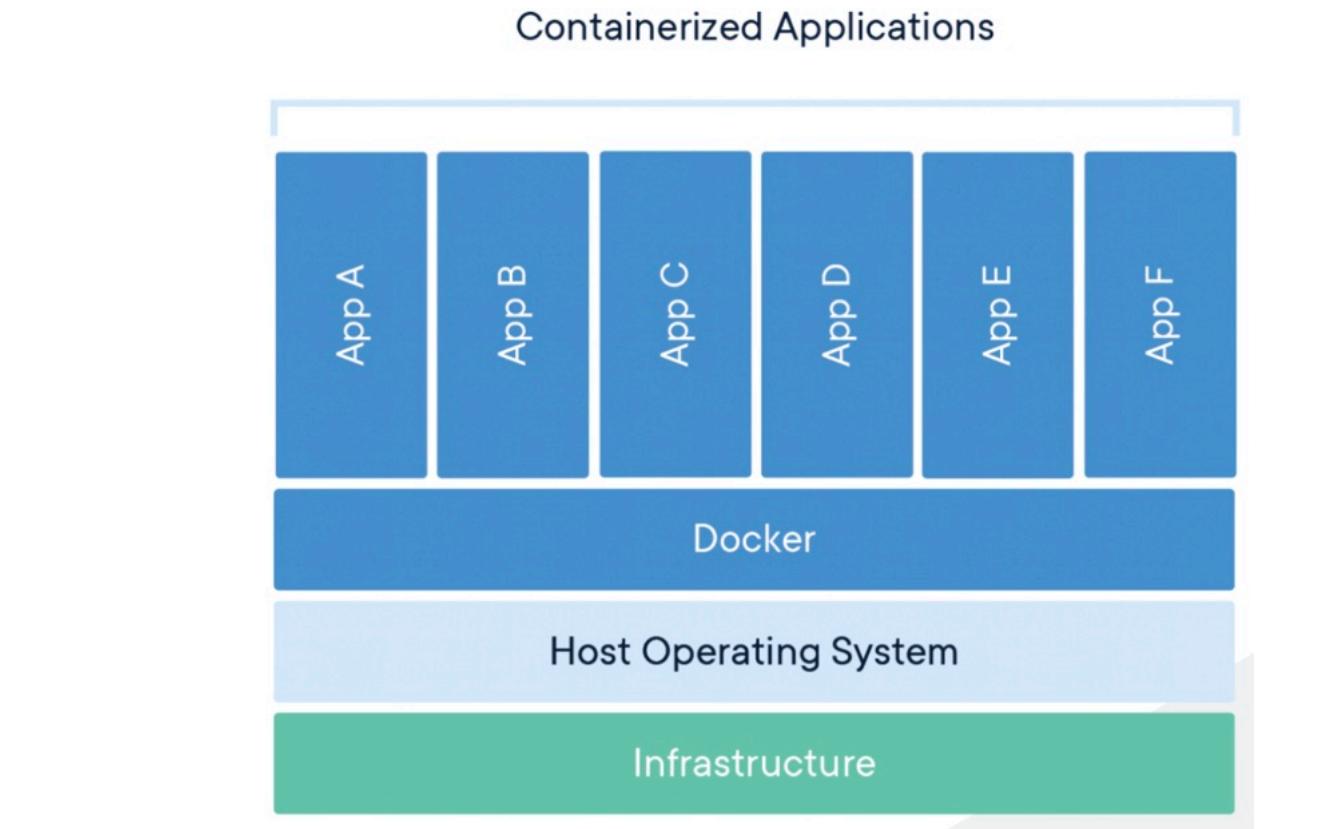
- A standardized unit of software that packages code and all of its dependencies so the software can run reliably across platforms (<https://www.docker.com/resources/what-container>)



- docker

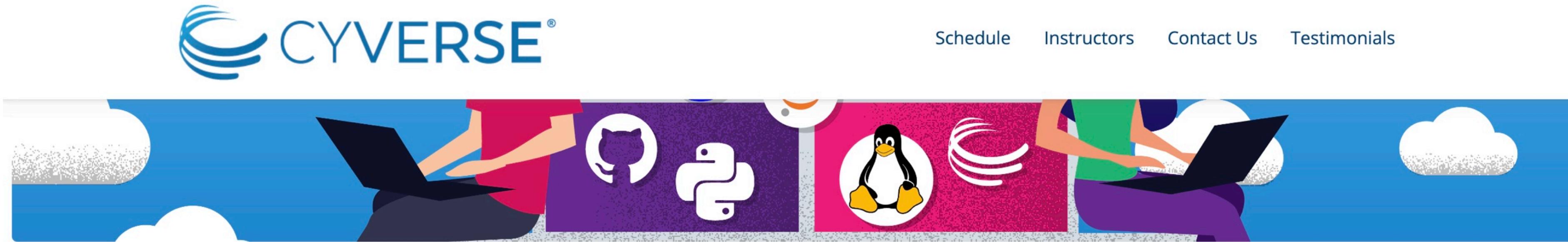
- Cheatsheet: <https://www.docker.com/sites/default/files/d8/2019-09/docker-cheat-sheet.pdf>

- dockerhub



A screenshot of the Docker Hub website. The top navigation bar includes "Explore", "Repositories", and "Organizations". The search bar contains the text "srs57". Below the search bar, there are two search input fields: "Search for great content (e.g., mysql)" and "Search by repository name...". A "Create Repository" button is located in the top right corner. The main content area displays a list of repositories owned by "srs57":
1. srs57 / snpbiner (Updated 7 months ago) - 0 stars, 6 downloads, Public
2. srs57 / pi-estimator (Updated 7 months ago) - 0 stars, 7 downloads, Public
3. srs57 / jupyter-scipy (Updated 7 months ago) - 0 stars, 15 downloads, Public
4. srs57 / lolcow (Updated 7 months ago) - 0 stars, 4 downloads, Public

Additional Resources



Foundational Open Science Skills

CyVerse's 8-week virtual workshop teaches you the principles, practices, and how-tos for doing collaborative open science using cutting-edge, open source cyberinfrastructure, in a collaborative, hands-on setting. To see how our FOSS workshop can support your work, check out the [curriculum](#).

September 7 - November 2, 2023

8 weekly virtual sessions on Thursdays; 11 am - 1 pm Pacific Daylight Time (or MST)
September 7 is an optional, pre-session brush-up on Git and Unix with an intro to ChatGPT

[REGISTER NOW](#)

7. Computational infrastructure

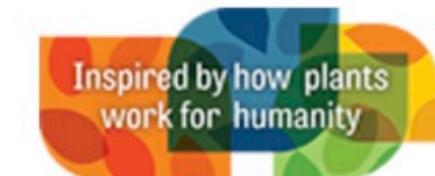


- Powerful servers with large amounts of memory, compute cores, and disk

8. Careers

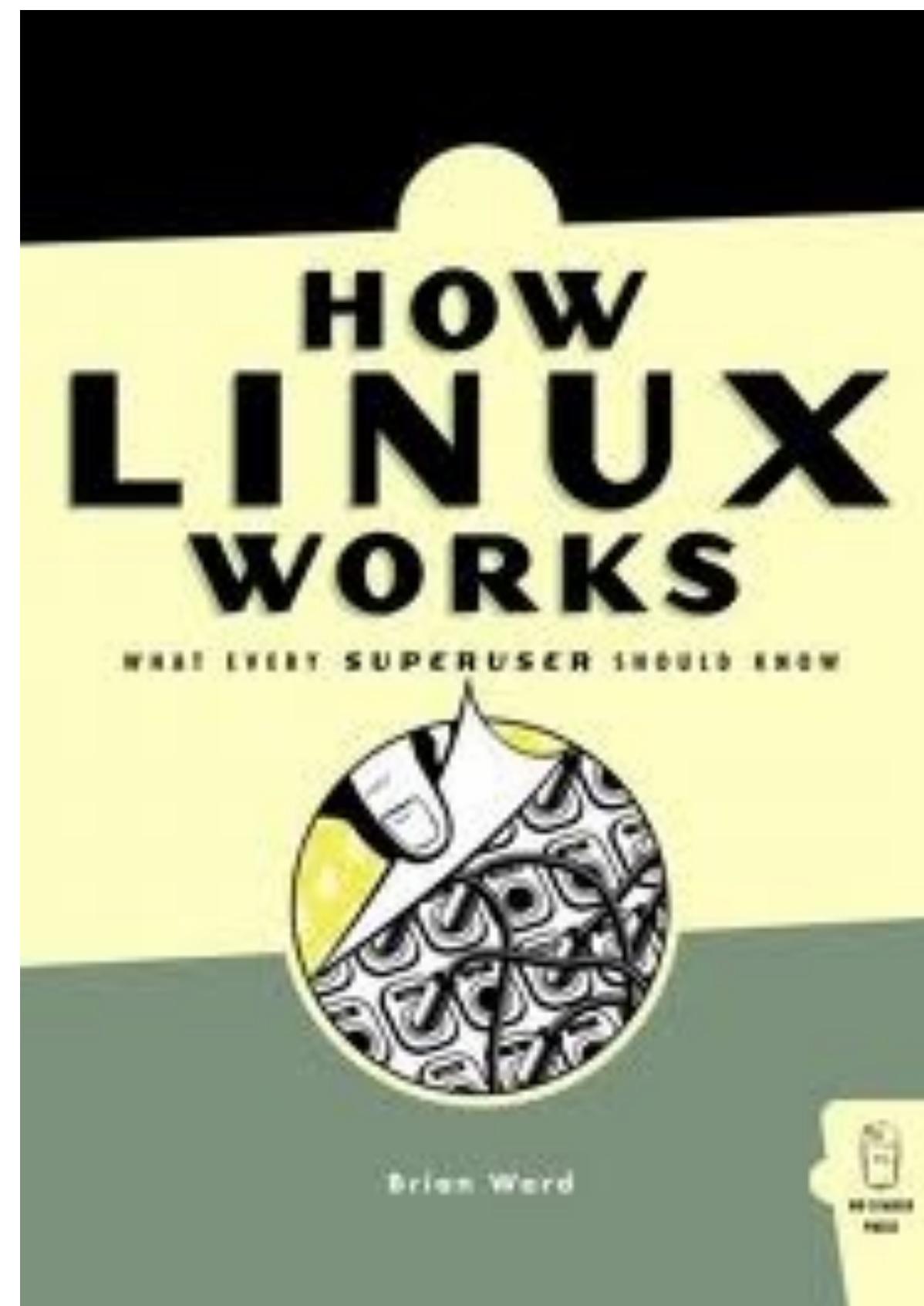
- Academia/Industry/Government
- BS/MS/PhD
- Bioinformatics core
- Professor
- Scientist
- Biomedical/Pharma
- Plant Breeders
- Botanic Gardens ;)
- Tends to have higher pay than other biology jobs
- Many positions offer remote work option
- <https://www.biology.pitt.edu/sites/default/files/publication-images/BINF%20Career.pdf>

More on Bioinformatics

[Explore](#)[Our Research](#)[Education & Outreach](#)[Get Involved](#)[Staff Hub](#)[Home](#)[Our Research](#)[BTI Computational Biology Center](#)[What is
Bioinformatics?](#)[Research](#)[Consulting &
Training](#)[Infrastructure](#)[Seminars](#)

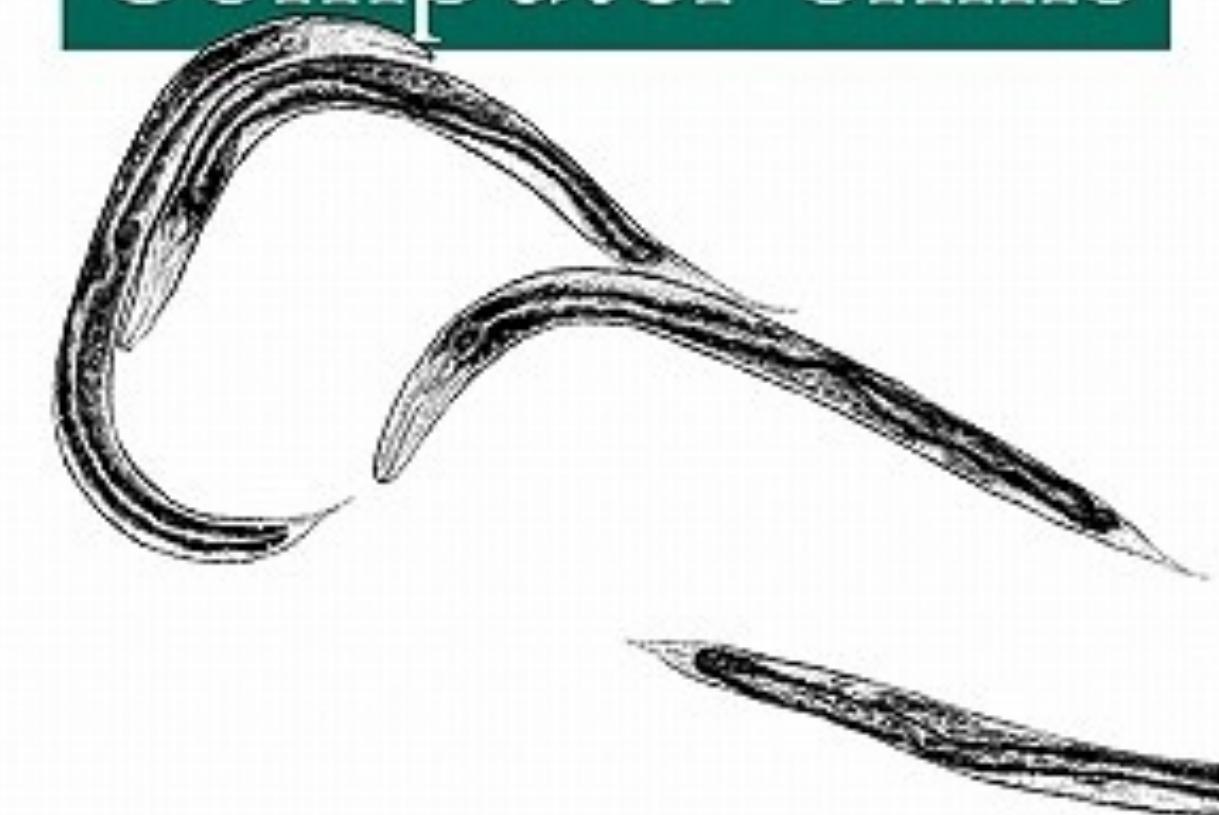
- <https://btiscience.org/our-research/computational-biology/bcbc-symposium/bcbc-2018-symposium/>

Further Reading



An Introduction to Software Tools for Biological Applications

*Developing
Bioinformatics
Computer Skills*



O'REILLY®

Cynthia Gibas & Per Jambekk

Empowering Beginners in Bioinformatics with ChatGPT

[Evelyn Shue](#),¹ [Li Liu](#),^{2,3} [Bingxin Li](#),⁴ [Zifeng Feng](#),⁵ [Xin Li](#),⁶ and [Gangqing Hu](#)^{1,*}

► [Author information](#) ► [Copyright and License information](#) [Disclaimer](#)

The complete version history of this preprint is available at [bioRxiv](#).

Next time

- Who would like to attend?
- You will need a Linux Virtual Machine
 - Instructions for installing it here: <https://btiscience.org/our-research/computational-biology/educationtraining/intern-2023-bcbc-bioinformatics-course/>
 - Please try installing it today and let me know by end of day if it doesn't work for you