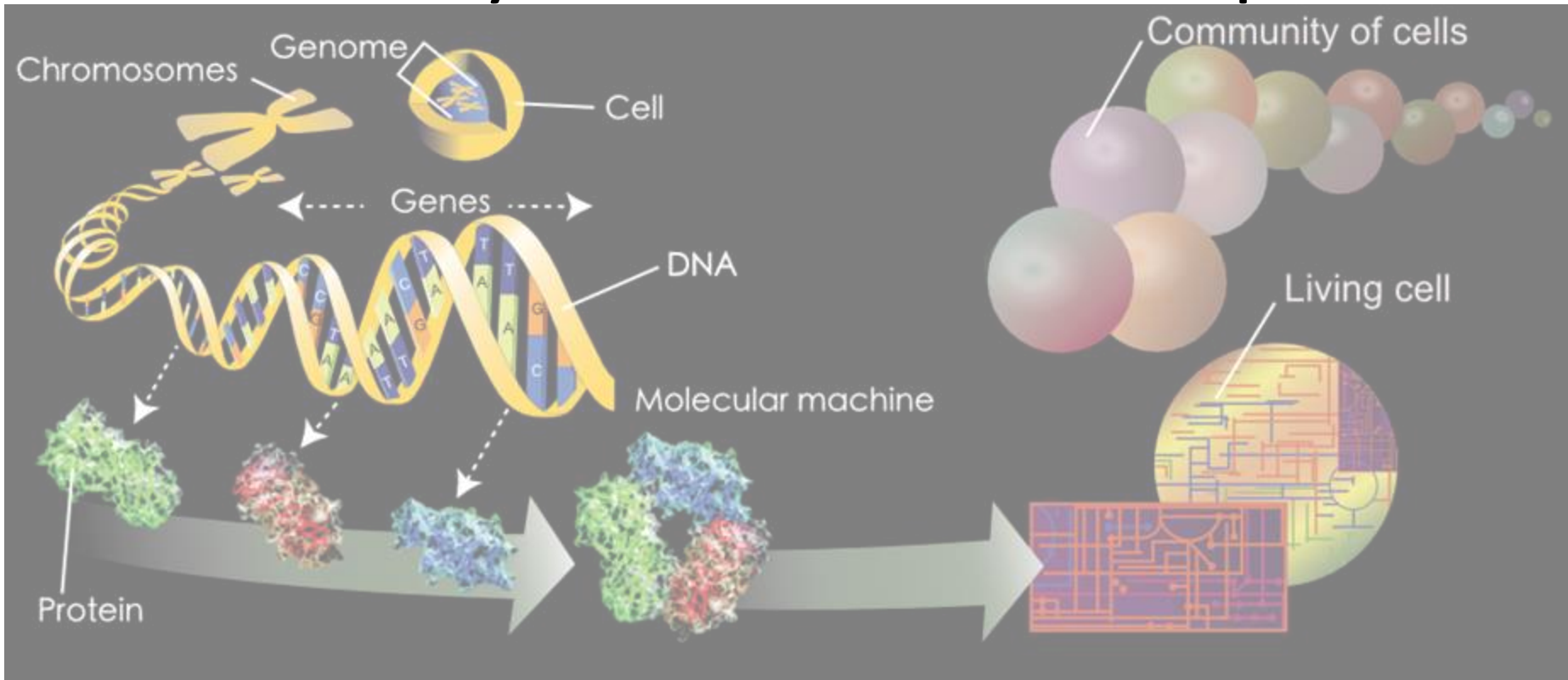


# From Molecules to Ecosystems: An Introductory Tour of "OMICS" Pipelines



Making sense of 'big' biological data

Emma Rocke | SEEC Seminar | University of Cape Town

# What do we mean by “Big Data” in Omics?



- **Volume** – billions of sequences
- **Complexity** – many genes, organisms, environments
- **Computation** – requires pipelines to extract biological meaning



**Bioinformatics Pipeline**

QC → Assembly    Mapping    Annotation

**Illumina**  
Paired-End Sequencing  
(~200-300 bp)

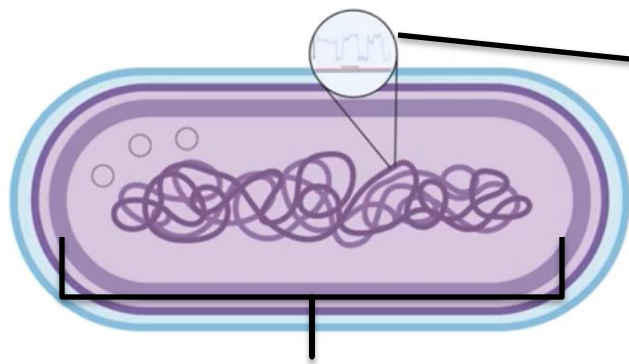
**Oxford Nanopore**  
Long Reads (from kb to Mb)

**PacBio**  
Long Reads (20+ kb)

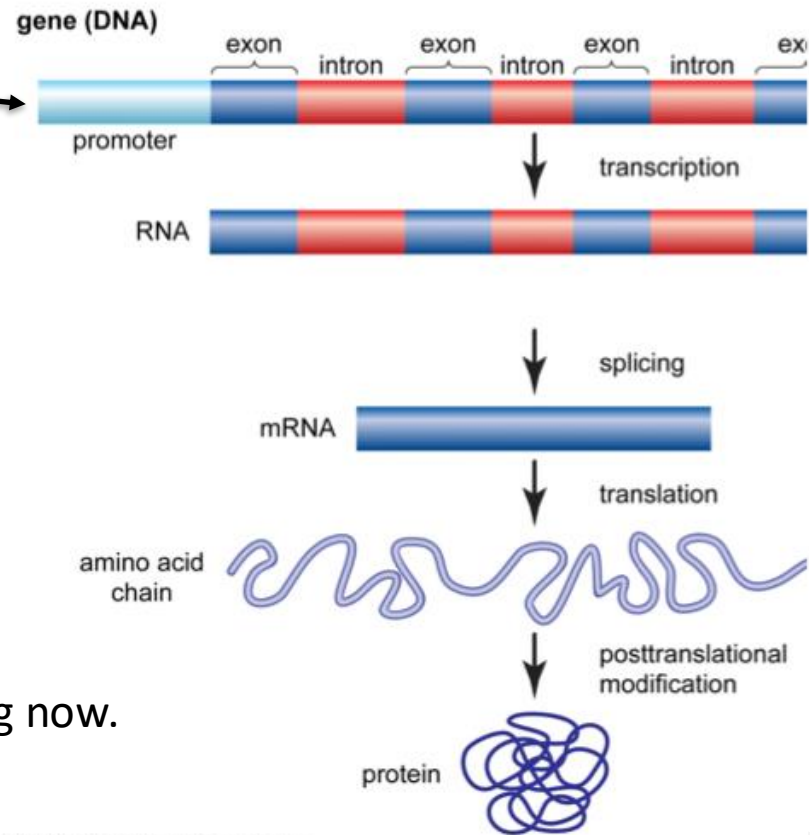
**PacBio**  
Long Reads (20+ kb)

- **High-throughput** sequencing of DNA or RNA samples
- **Different technologies** yield varying read lengths and outputs
- **Massive** amounts of **data** generated from a single run

# What is a **gene**, a **genome**, a **transcript**, or a **protein**?



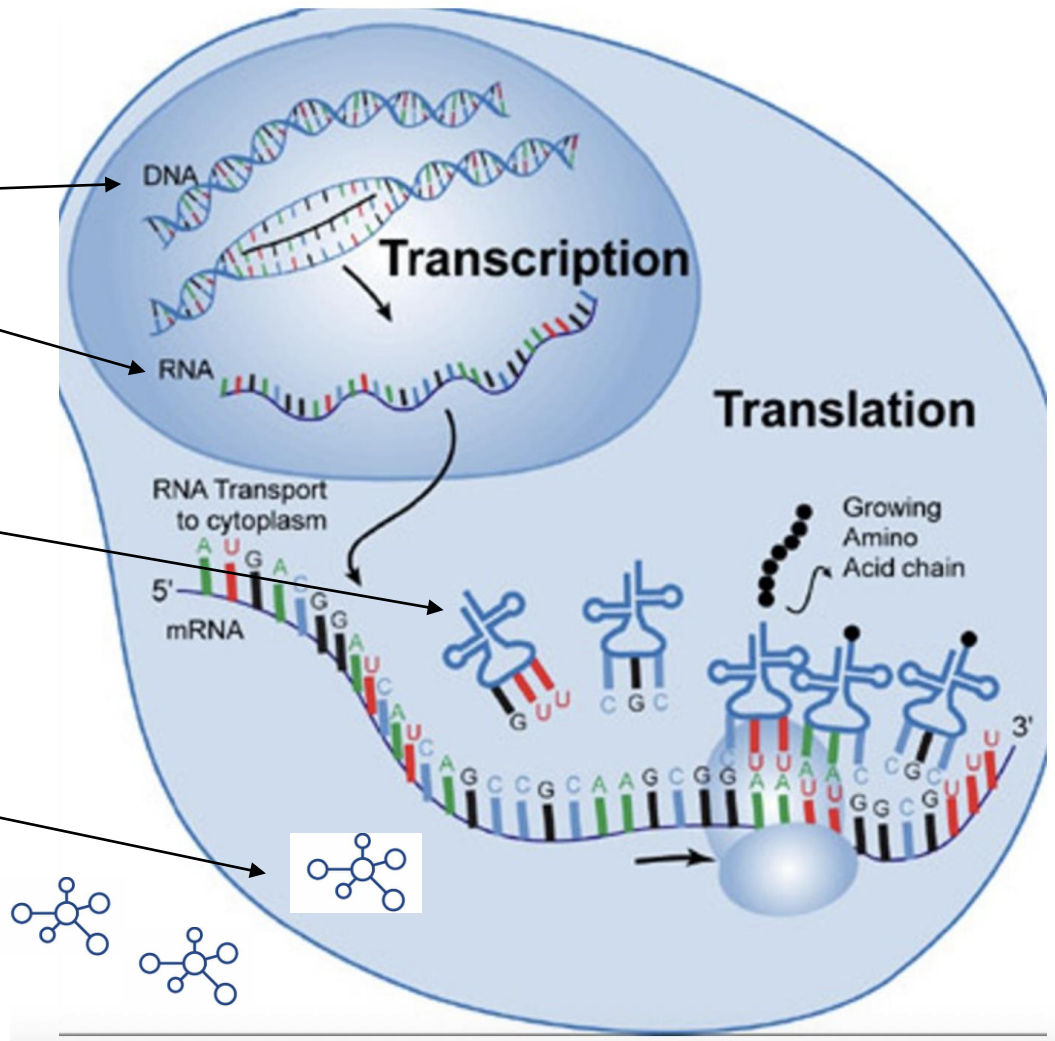
**Genome**



DNA=relatively stable blueprint  
RNA = what's being used  
Proteins/metabolites: What's happening now.

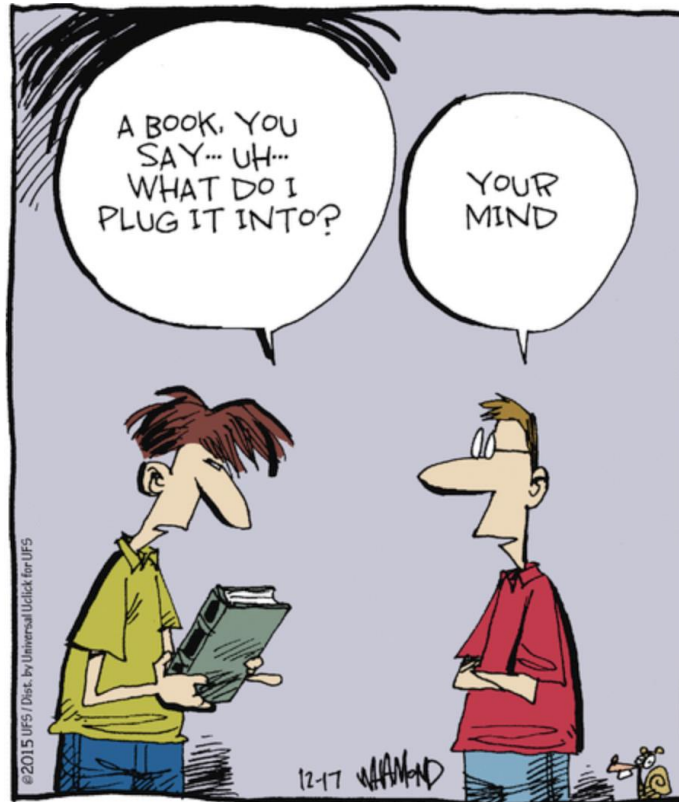
# What Do We Mean by 'OMICS'?

- **Genomics** → DNA: Who's there?\*
- **Transcriptomics** → RNA: What's active?
- **Proteomics** → Proteins: What's functioning?
- **Metabolomics** → Metabolites: What's produced or consumed?



\*amplicon/eDNA: Who is present in an environment

# Reality check!!



Omics data can feel like a black box of codes and acronyms...

# Resources – there are lots!!

**DIPL**  **MICS**

START HERE!!

<https://www.diplomics.org.za/contact>

Blogs:

<https://merenlab.org/posts/>

<https://bigomics.ch/blog/>

EMBL-EBI | MGnify

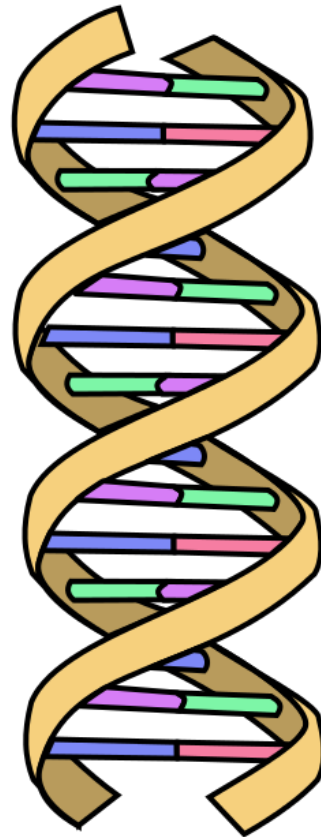
# MGnify


Submit, analyse, discover and compare microbiome data





Example searches: Tara oceans, MGYS00000410, Human Gut


# DNA




 = Adenine

 = Thymine

 = Cytosine

 = Guanine

 = Phosphate  
backbone

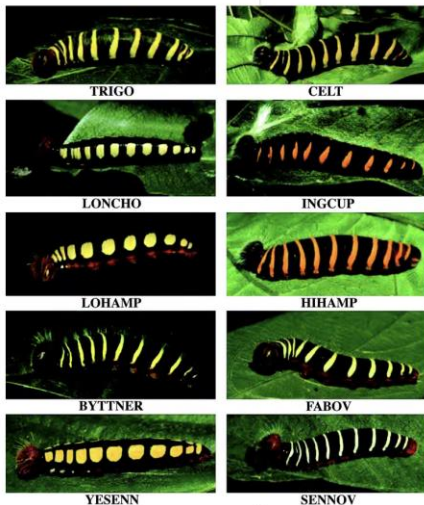
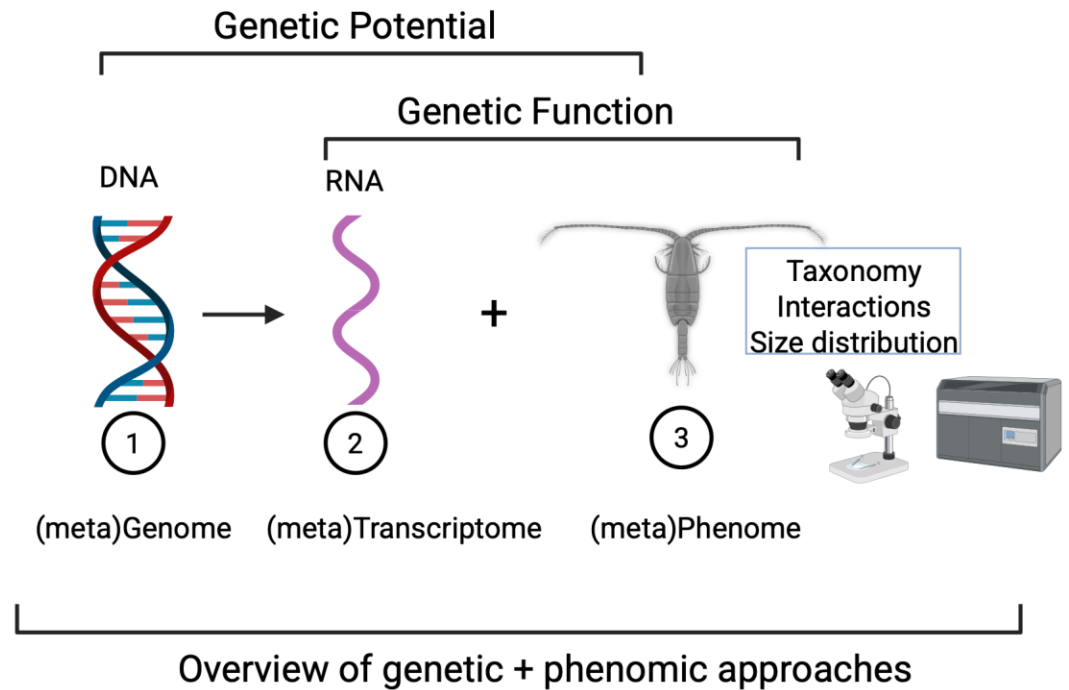
DNA

# Why OMICS?

Reveals hidden biodiversity  
(uncultured/cryptic species)

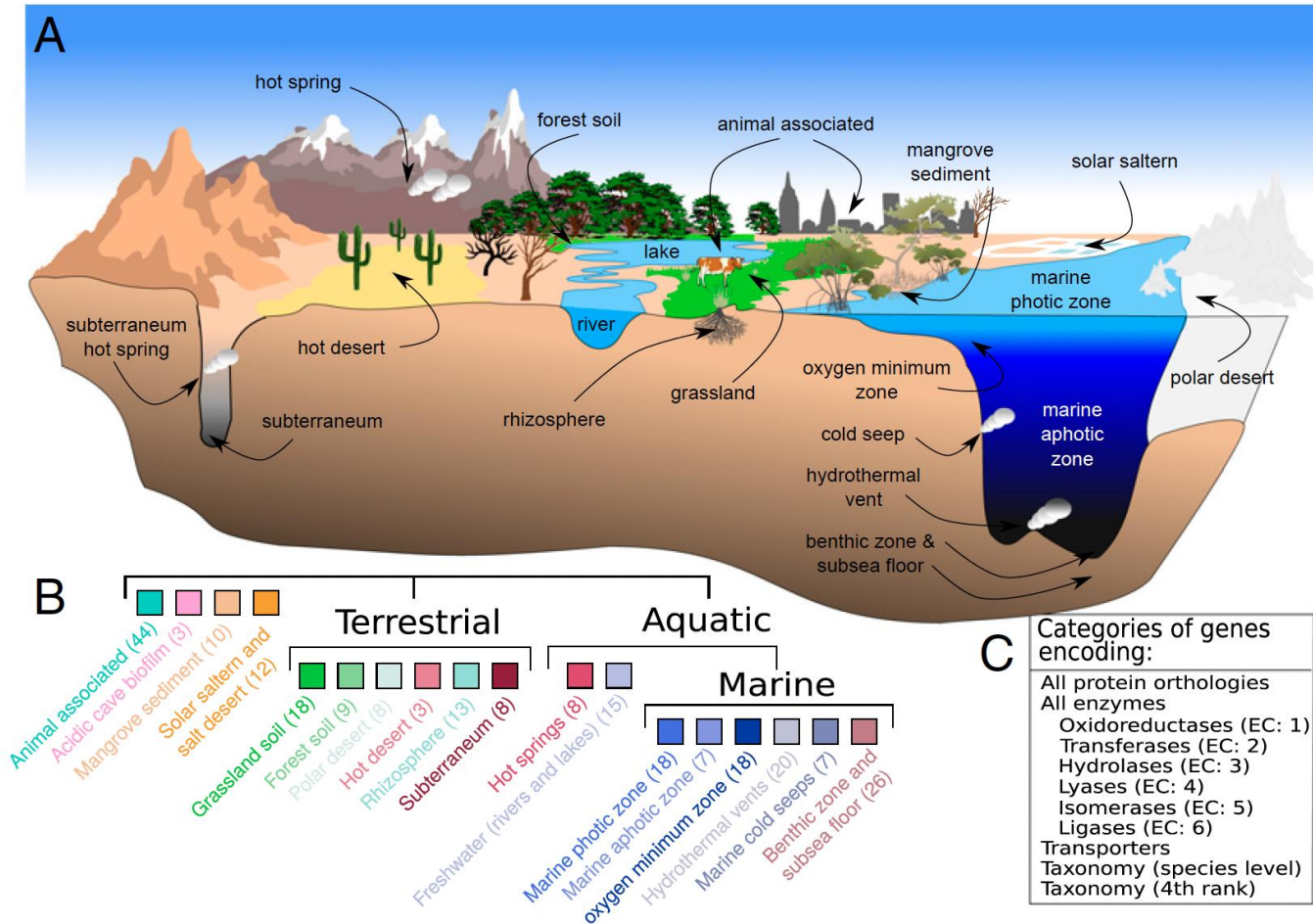
Links molecular function to  
ecosystem/organism processes

Enables early detection and  
monitoring (ie: AMR, stress  
biomarkers as a case study)



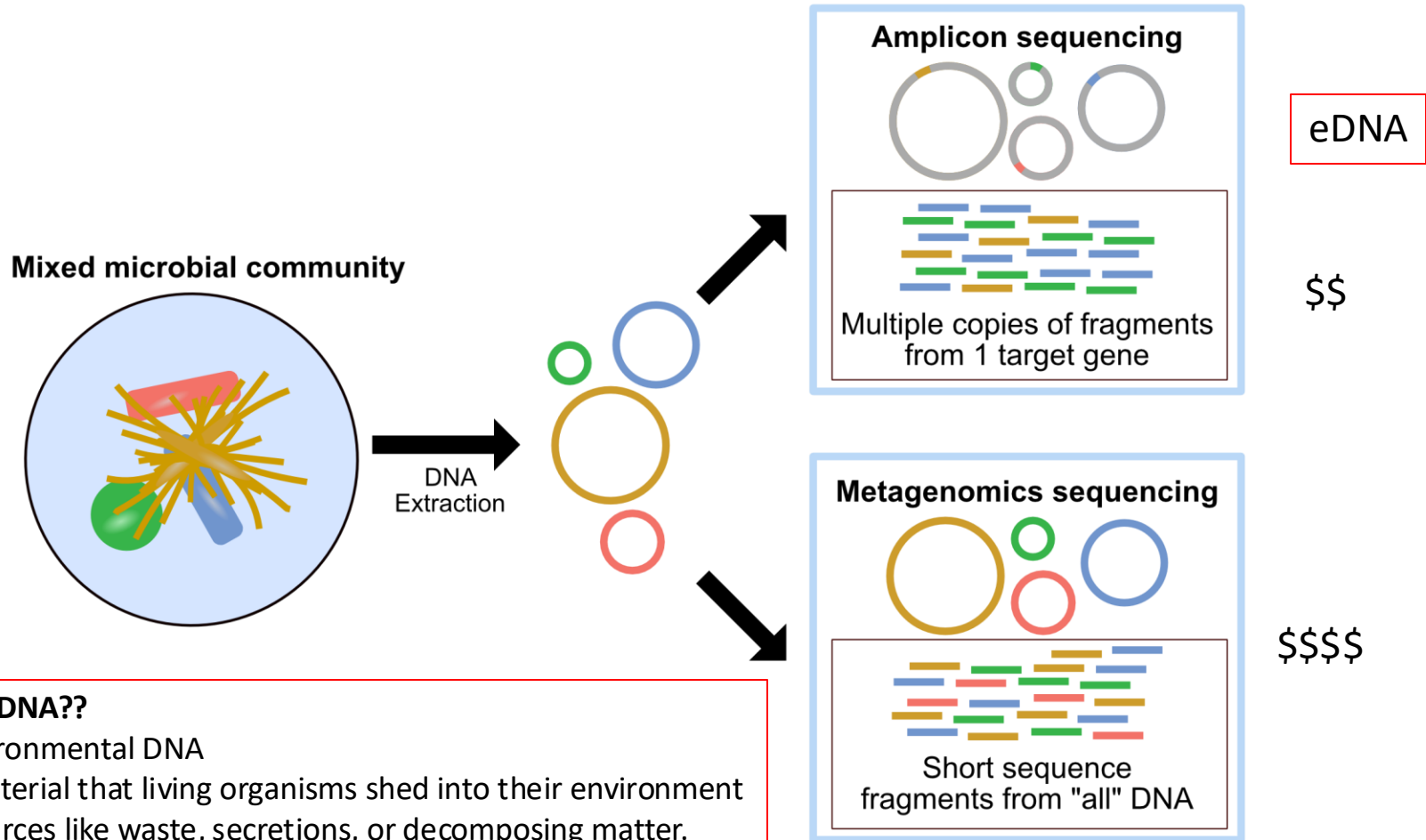
**Omics turns your organisms and ecosystem into quantitative datasets you can plug into the models you already use.**

# Why OMICS?



Ramirez-Flandes, PNAS, 2019

# What data product should I use?

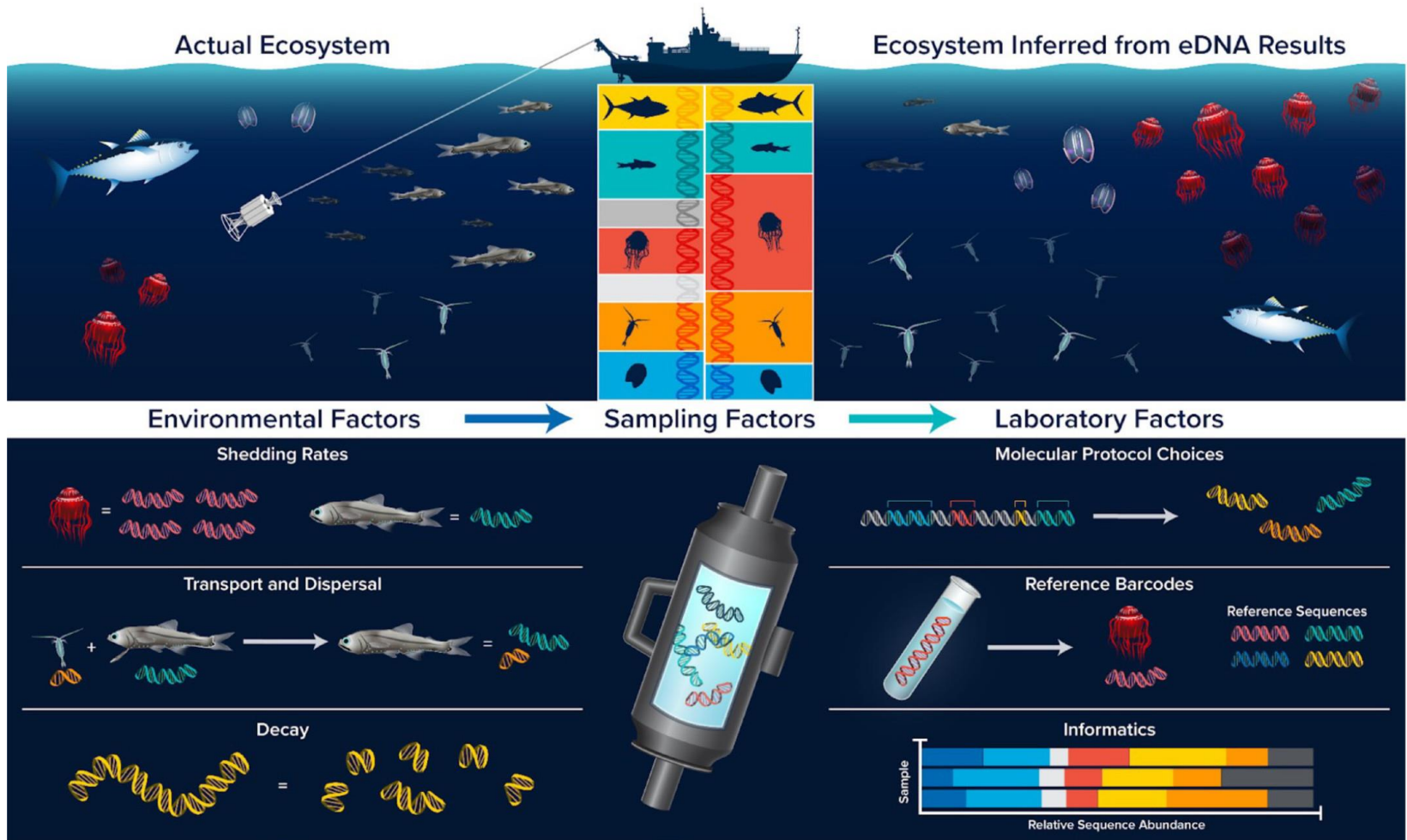


## What about eDNA??

eDNA, or environmental DNA

- Genetic material that living organisms shed into their environment through sources like waste, secretions, or decomposing matter.
- Allows researchers to identify the types of species present in an area without needing to sample the organisms directly.
- Cost-effective and non-invasive tool for biodiversity monitoring, invasive species detection, and wildlife surveillance

# eDNA realities



**Detection limits, degradation and transport mean 'present' = 'local and abundant' + your species ID is only as good as your chosen database and primers!**

# The OMICS Workflow — The Big Picture

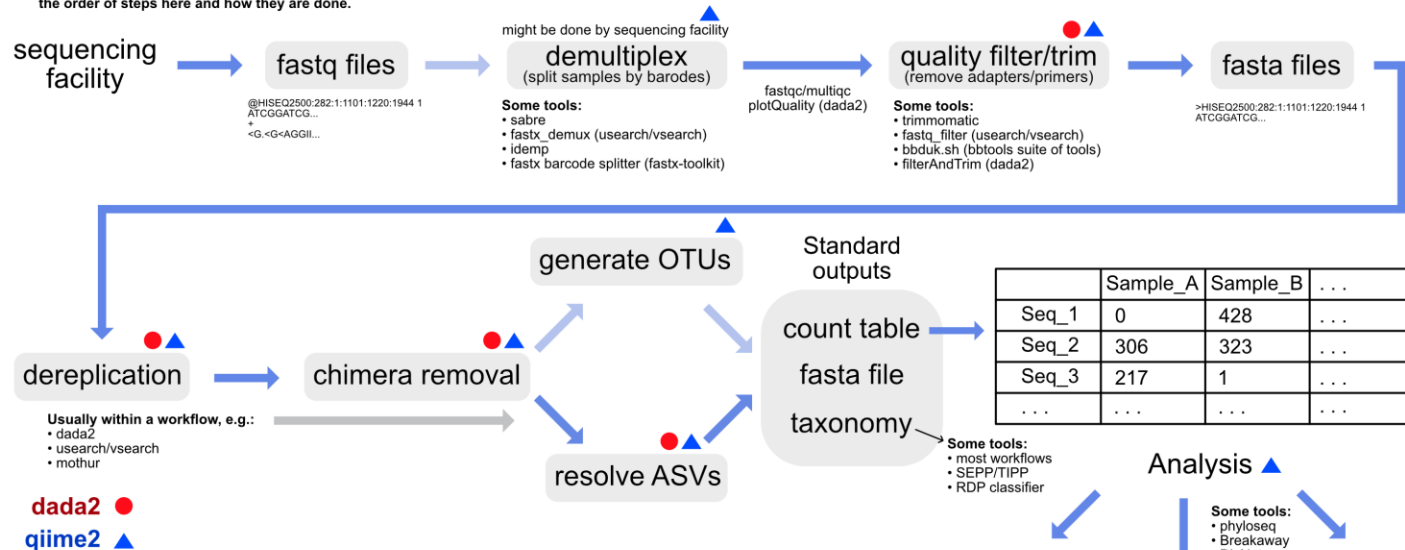
Sampling → Data generation → **Processing** →  
**Annotation** → **Statistical analysis** →  
**Interpretation**

Same logic across ALL OMICS types

# Amplicon/eDNA overview:

## Overview of generic\* amplicon workflow

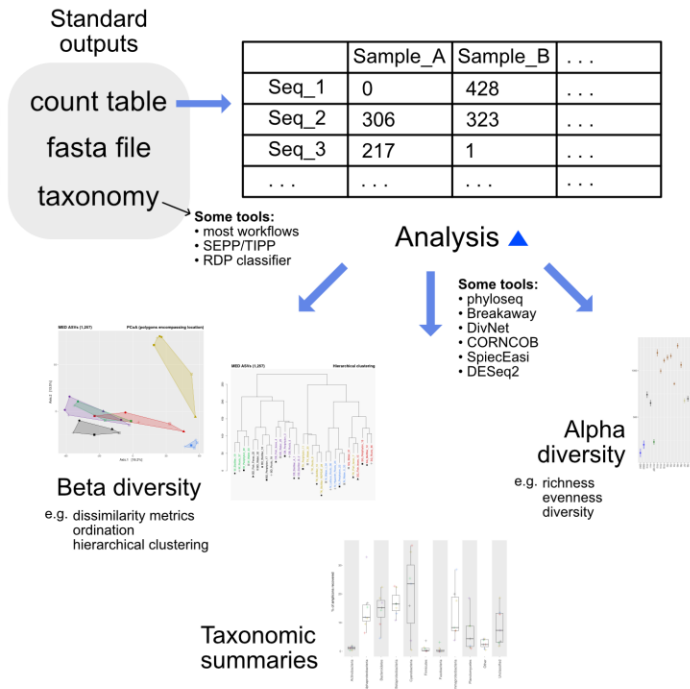
\*This is generic; specific workflows can vary on the order of steps here and how they are done.



When working with your own data you should never follow any pipeline blindly. There can be critical differences based on your data.

**Some tools that provide whole workflows:**  
**dada2** runs within R (ASVs)  
**usearch/vsearch** runs at the command line (ASVs and OTUs)  
**mothur** runs at the command line (OTUs only currently)

**qiime2** provides a multi-interface environment that employs processing tools like those above, infrastructure for easily documenting all processing performed, and interactive visualizations



# File types...

Fastq files = raw sequence reads from the machine

Fasta file: Your sequences!

FASTA file:

```
>NC_000006.12:151654148-152129619 Homo sapiens chromosome 6, GRCh38.p13 Primary Assembly  
TATTGATTTTGGTAAACATGGTTTGTATATCTATAACGAGAAGCTCAAGTCATACTGTAATCCTAT  
TTTGAAAACGACTTTTTCCTTTATCAGTATATCAAGATTATTTCCACATCATTGACATTTTTCT  
ACAGTGAATTTAATGGCTACATGTTTCTATCCTATGAATATATCAAACCTATTTCTAAAAACCTA  
CTCAGGGATTTAAAAAATAAAAAAGATGTTTAAATATTAAAGATTCAAGTGAAGGATATTTCTATACG  
TACACATTTCTAAGGTTTGAGTTCTTACAAGATGCTGAAGCTAGCTAAGACTACTGGTCTCATCTGTAC  
ATAGGGAAAAATTAGAAGGAAAAATCAAGATTGGAAAAATCTGTGAGAATTGTTTGCATTAGTGT  
GTAGGTGTGTGTGTGGGGTGGTGGCTGCAGCTTGGGGCAGAGGCTCAGGTGTGGCTGTGGAGTGATCA  
GATAGAGTTTTGGAGTTCGGCTTTTCCAGGACACTTGGTGCCTGCCCCAGAGCTGCAGCCAGAA  
GGCCGTTCTCAGAGGTGAAGTCCAGGCAGTGAGGAGCTGTGTCAGTAGGCAGTTGAAGAAAAAATG  
AGCTAGAGGAAAAAACAATAAAATCTCCCTTCTAATGCTGCAGGCTGCCGGAGCTGGAATGGA  
AGCACTGACAGGAGTGGGATTTTCATGGTGAAGGAATAATCAACTGGTTTTTTGGTACCAAGACTTT  
CCACCTTCCACACACACATGAGATGCTTTGAAATAAGATAGTCACTTGAAGTTGTTGAC  
ATAAAAATAGAGAAATACCAAAGAATACAAAAAGGAAACTTCGTTAATATTATCAGACTTAAATTC  
CAGATTGTATCAACATTAAGGGGGTGTGATGAAACATGGGAGAAAGCCAGGGACGTGAGATCGGGCTCA  
ATTCCTGACTGCTGGGGGAAAGGTATCAACACAGAACTTTAAGAATTAGAAGCATTAAAAAGAAATAG  
AAATCCTGAATCAAATGAAACAGTAAAAATAAATAGTCCAAGATGTGTAATATATCACTATCAAAAT
```

→ HEADER

→ SEQUENCE

Fasta file extensions:

Extension	Meaning	Notes
fasta, fa <sup>[9]</sup>	generic FASTA	Any generic fasta file. See below for other common FASTA file extensions
fna	FASTA nucleic acid	Used generically to specify nucleic acids.
ffn	FASTA nucleotide of gene regions	Contains coding regions for a genome.
faa	FASTA amino acid	Contains amino acid sequences. A multiple protein fasta file can have the more specific extension mpfa.
frn	FASTA non-coding RNA	Contains non-coding RNA regions for a genome, in DNA alphabet e.g. tRNA, rRNA

# Ampliseq: one stop shop:

<https://nf-co.re/ampliseq/2.6.1/>

```
nextflow run nf-core/ampliseq \  
-profile singularity \  
--input "samplesheet.tsv" \  
--FW_primer GTGYCAGCMGCCGCGGTAA \  
--RV_primer GGACTACNVGGGTWTCTAAT \  
--metadata "data/Metadata.tsv" \  
--outdir "./results"
```

*How many genes regions in your sample?*

*Spreadsheet indicating file path to your sequences*

*Primer set that was used*

Standardised, documented pipeline = less time reinventing wheels (github is a **great** resource here)

Produces summary reports, diversity metrics, taxonomic tables ready for R

You will need: Raw sequence files, primer info and a sample metadata sheet

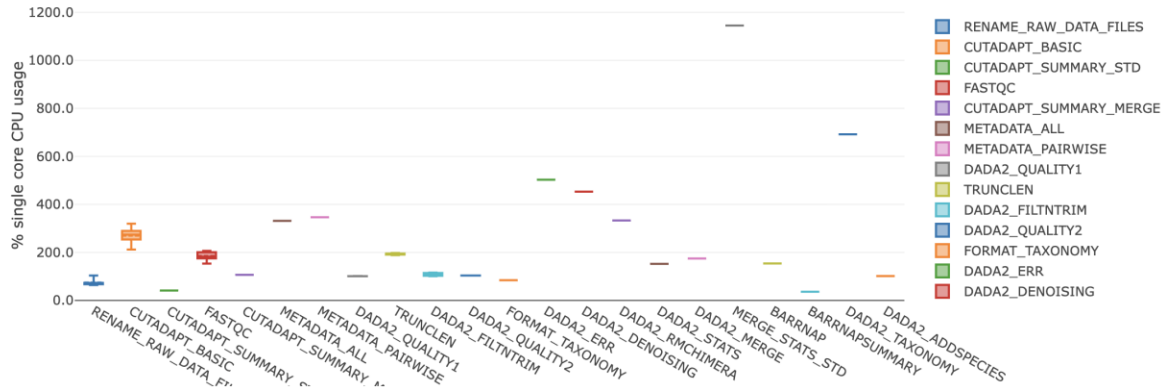
**R packages:**

Phyloseq

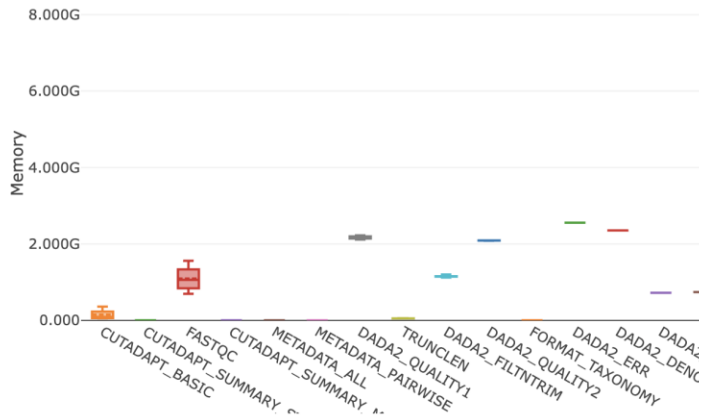
Microeco

# QC Summary reports: important to check

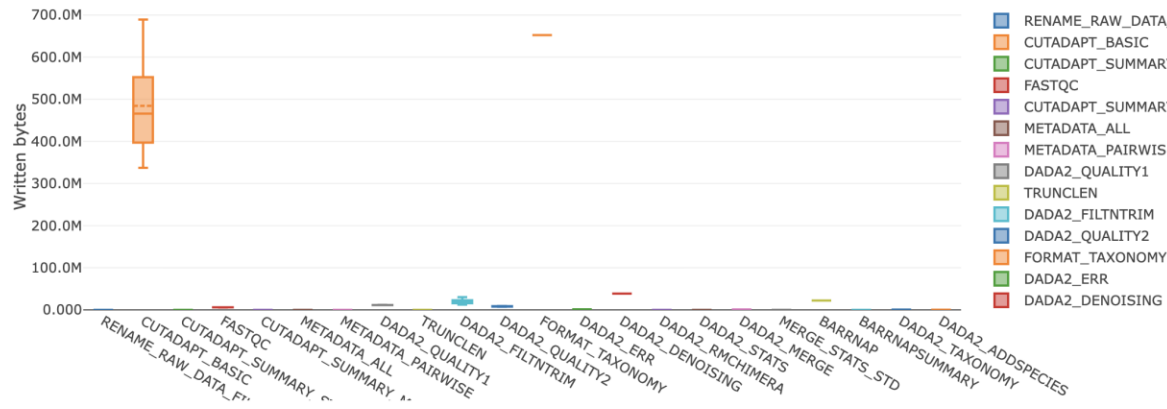
CPU Usage



Physical Memory U



Number of bytes written



# 1. Quality filtering: QC

FASTQ file---FASTA file:

Every new entry starts with an "@" sign at the start of a line followed by an ID

IDs are not always unique in the file. If they are not, the order of sequences in the file is important.

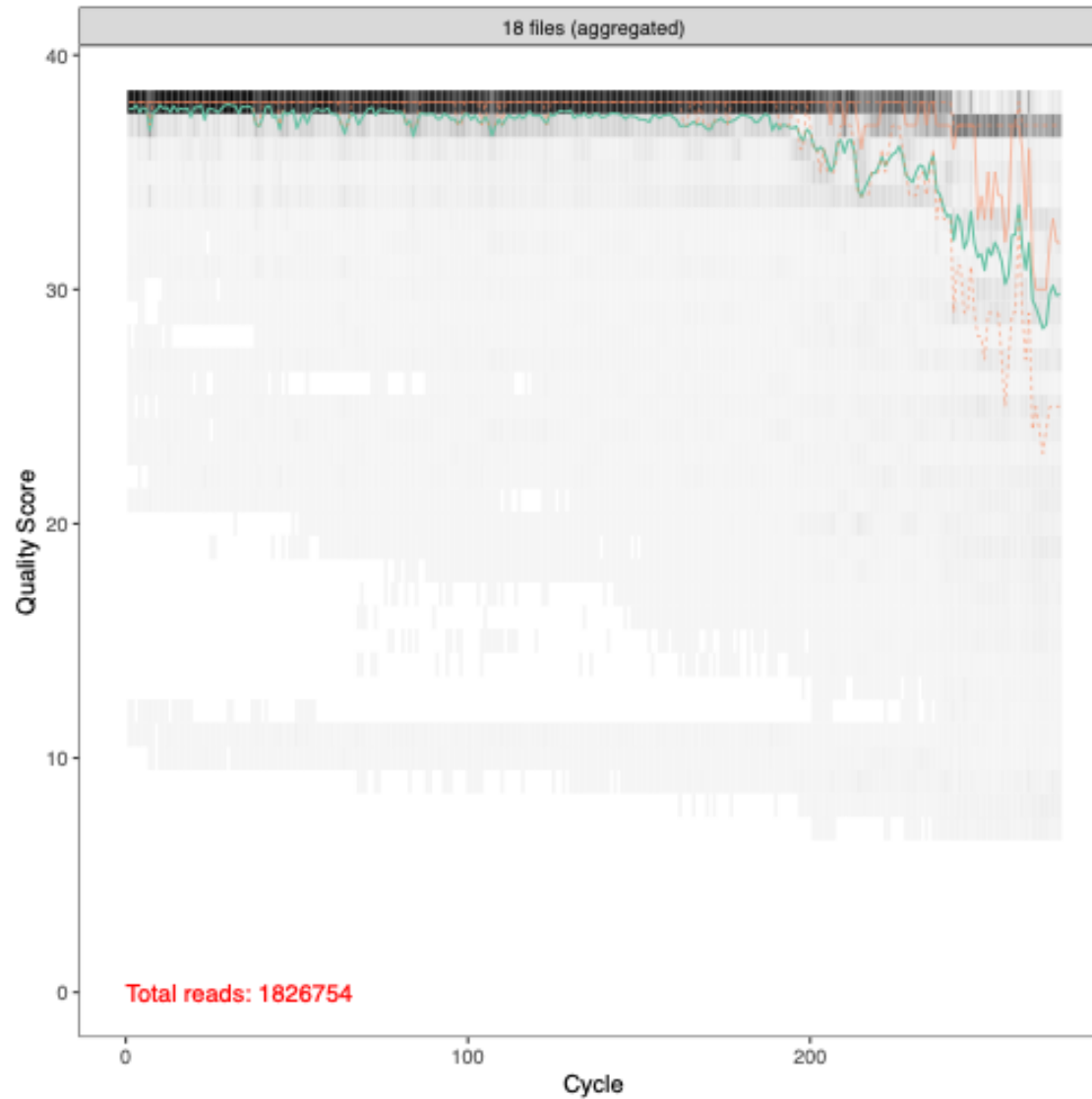
Every third line starts with a "+" and may or may not repeat the ID

```
@M01965:5:000000000-A9228:1:1101:10116:1028 1:N:0:14
NCCCTGCATGATTGTCTCCATCTTAAGCTCTGAGGAGTGATGCTCTATCCACTGACTTA
+
#8BCCFGFGGGGGGEFFGGGDGGGGFFFGF8FGCCFFCFC9FCFGGGF9F6CFGDGGFF9
@M01965:5:000000000-A9228:1:1101:13369:1030 1:N:0:14
NTTTATAGTTGTATTCATTTTTTATAATCAACAAATTTTGTGATAAAGGCTTCTTAGTG
+
#8ACCGGGGGGGGGGGGGGGGGGGGGGGFGG@FF9AE@,EFFGGGGGGGFFGGFFCFGGFEAFG
```

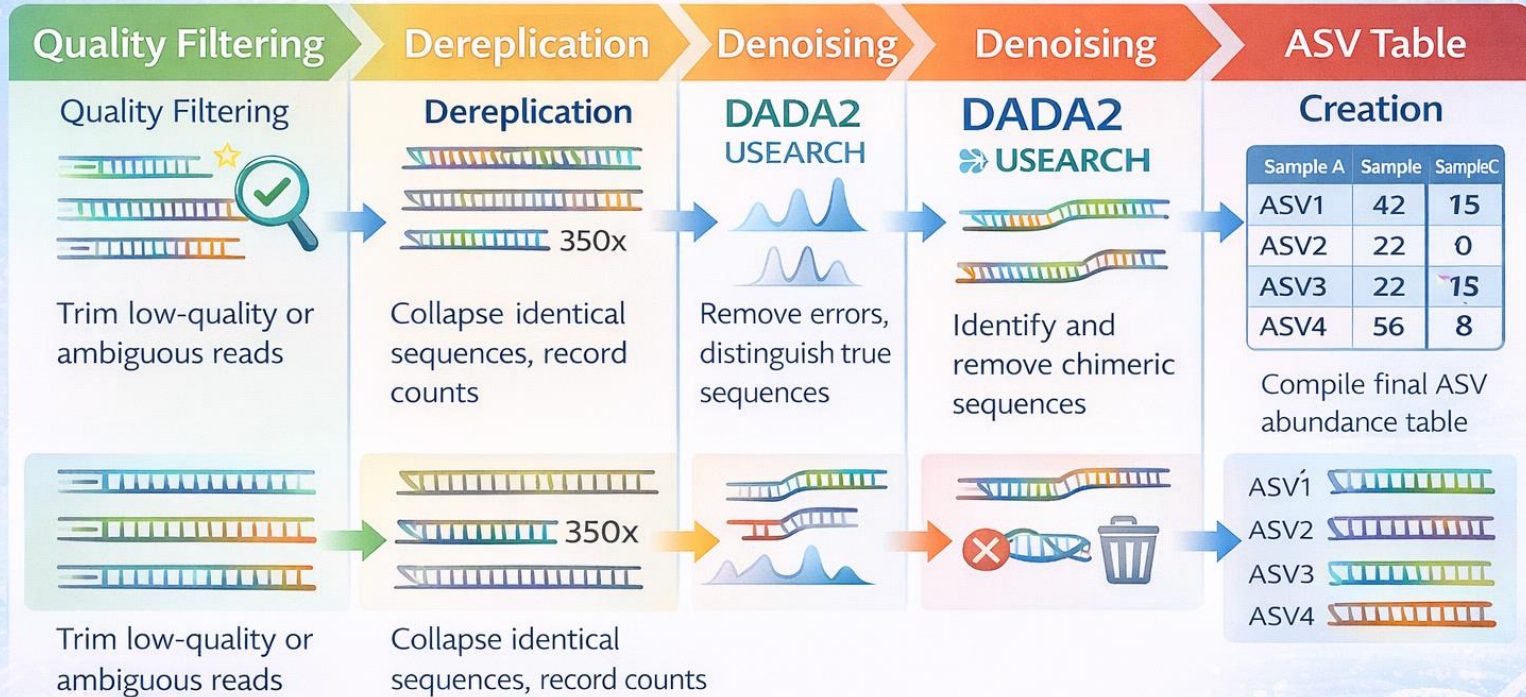
Per-nucleotide quality scores are coded in ASCII, often from ! to J (Phred score 0-41)

Every entry consists of four lines: identifier, the nucleotide sequence, a line starting with +, and per-nucleotide quality scores

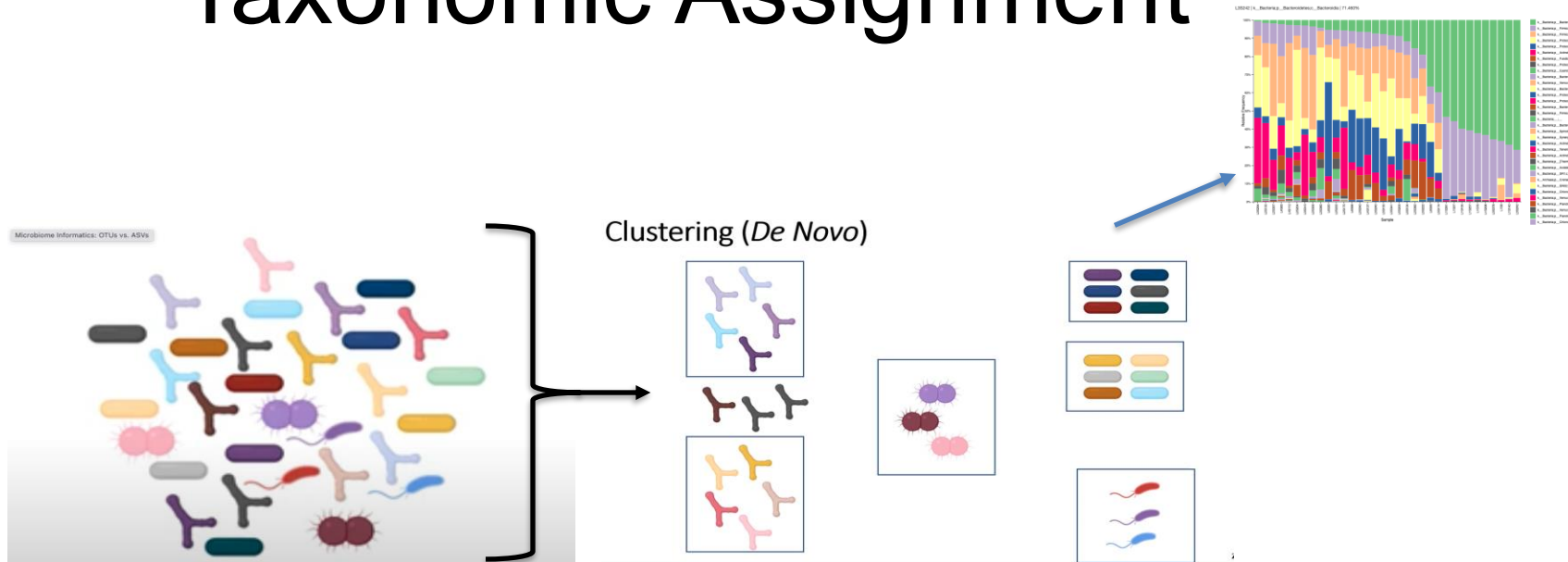
# Preprocessing – Quality scores



# Steps in ASV Generation



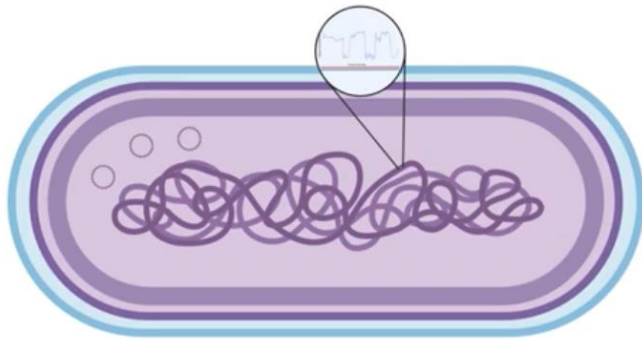
# Under the hood: ASV's and Taxonomic Assignment



ASV= Amplicon sequence variant (statistical inference of true biological sequences)

# A few warnings...

Amplicon sequence variant




A Quick Note

## 16S Sequencing Challenges

**Targeted sequence with a few bases differentiating species**


- Sequencing is imperfect
  - Illumina usually makes some base call errors
  - Nanopore makes more
  - Errors are not necessarily evenly-distributed
- We do not want errors to be confused with real diversity/new species

**FOR OUR PURPOSES:**

 = 1 sequence and one species

**REAL BACTERIA**

 One species can have multiple, different copies of a gene

 Two similar species may share an identical sequence

# ASV inference

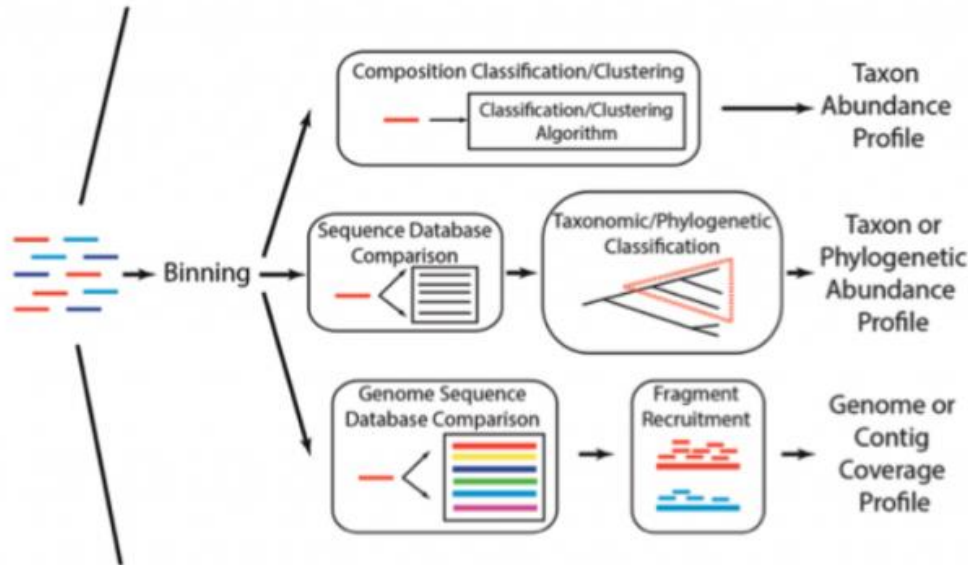
ASV\_table

ASV_ID	E10a	E10b	E11a	E11b	E13a	E13b	E14b	E14c	E4a	E4astar	E4b	E4bstar	E5a	E5b
cae609ea893f535d61720405d255b8ad	4285	9764	6399	5404	1506	1307	26962	35226	10672	5626	5455	3756	18177	17938
9f2865c7bbbce298473d4ae532a335e6	18204	15465	107	27	7611	8863	17	23	122	23631	22379	27871	131	59
0002c03993330198b02457afef4d00e5	170	150	5446	5811	843	567	310	283	6422	559	227	108	10923	3508
d11d8a29479c2d36943b2b68362221f4	4088	1621	1919	2420	6356	5455	122	147	580	1675	5628	1376	660	260
ffb5ad37299ba5fbe1946a5639f722cd	0	0	29	22	195	13	24000	8291	181	0	0	0	246	3357
eca00533a595089a33766b2f14a1961c	112	46	4382	3259	199	195	468	417	3492	272	152	161	5069	3681
5b7ab8aab9e49afadb578146ef26e5a7	33	0	3975	2728	134	0	3500	2506	1663	0	39	0	2518	2335
9fe4b245dc54657f85440234d48a02ed	1919	1321	125	0	2540	3304	0	0	34	3140	2352	1622	61	0
0df0db8002857b2a9bfe61fd662710da	7	0	931	445	0	0	17	31	4132	12	4	0	6162	4099
e3ea1a931ec40e4c21752dac82031b8d	4153	1523	0	0	598	289	0	6	0	680	5578	592	0	0



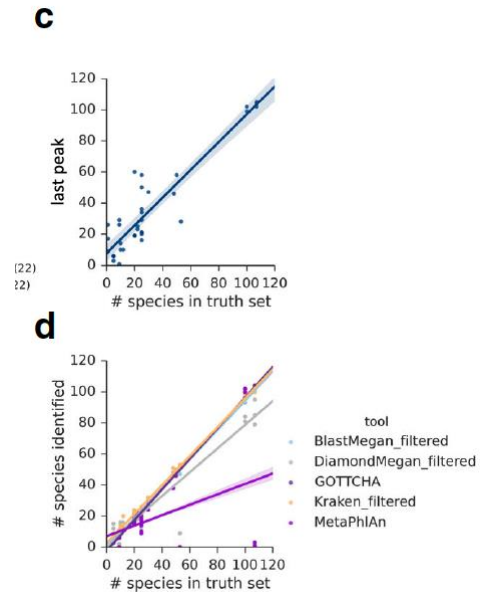
Note: Normalisation. Relative abundance/ Hellinger... have a look at your dataset !

# Taxonomy/gene calling



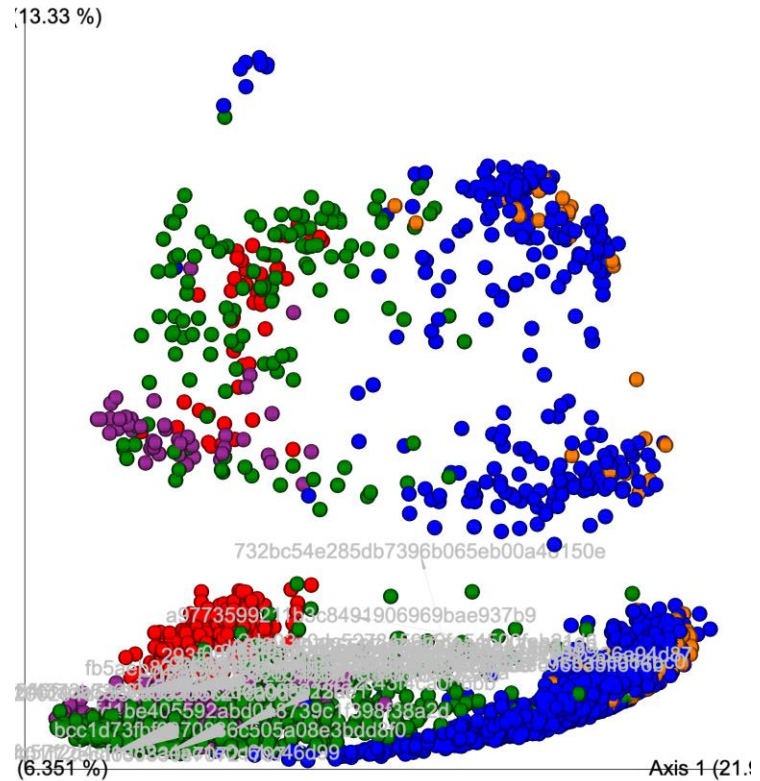
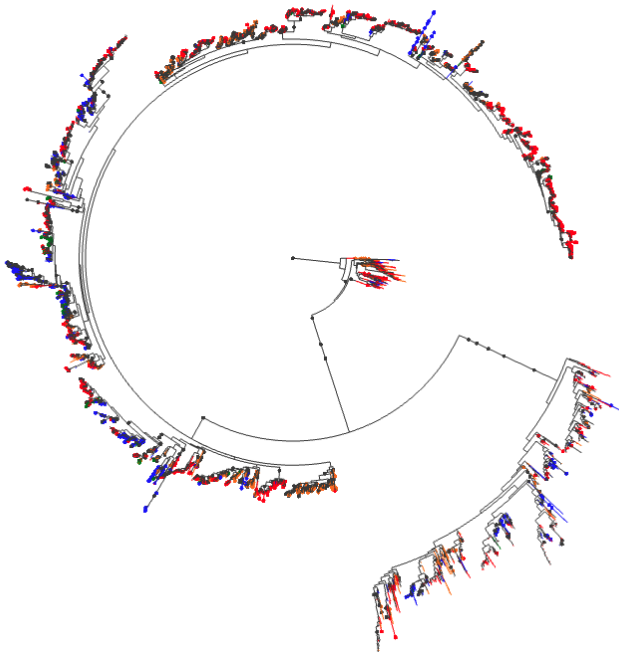
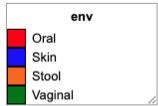
## Comprehensive benchmarking and ensemble approaches for metagenomic classifiers

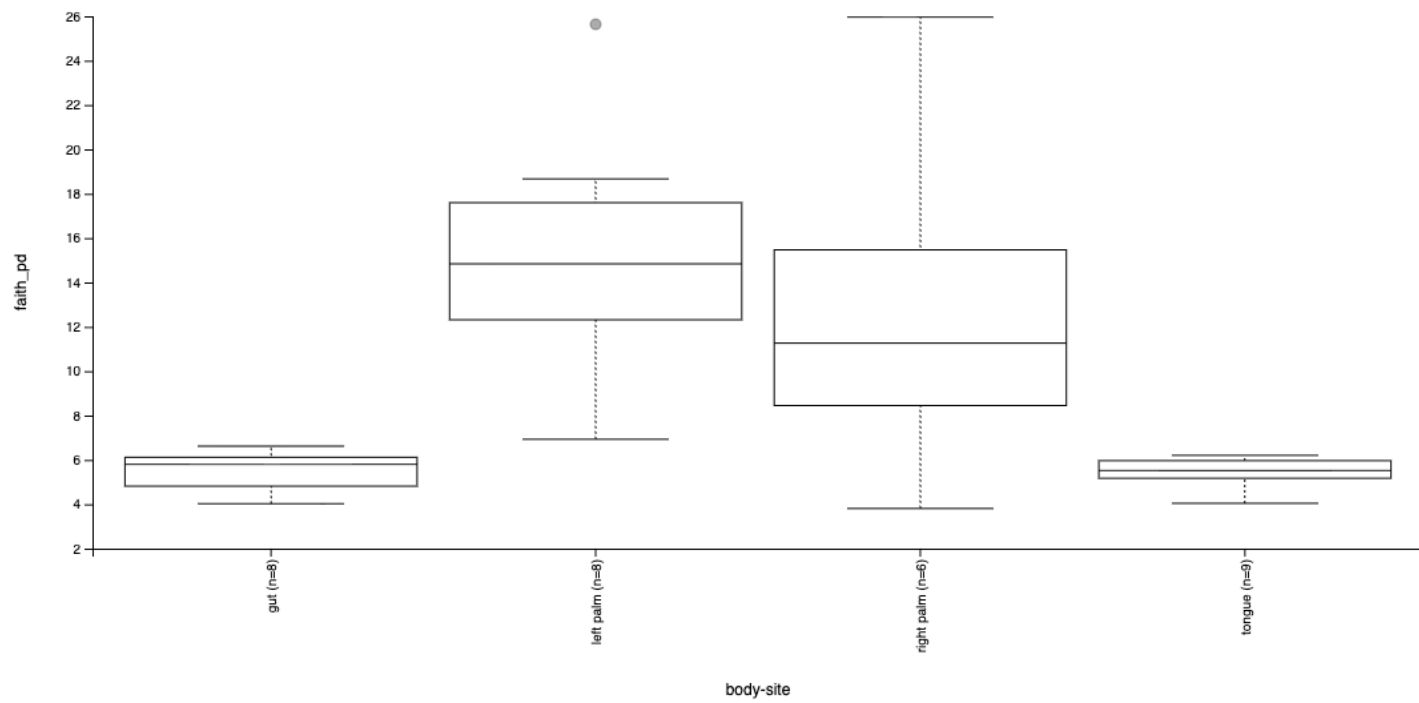
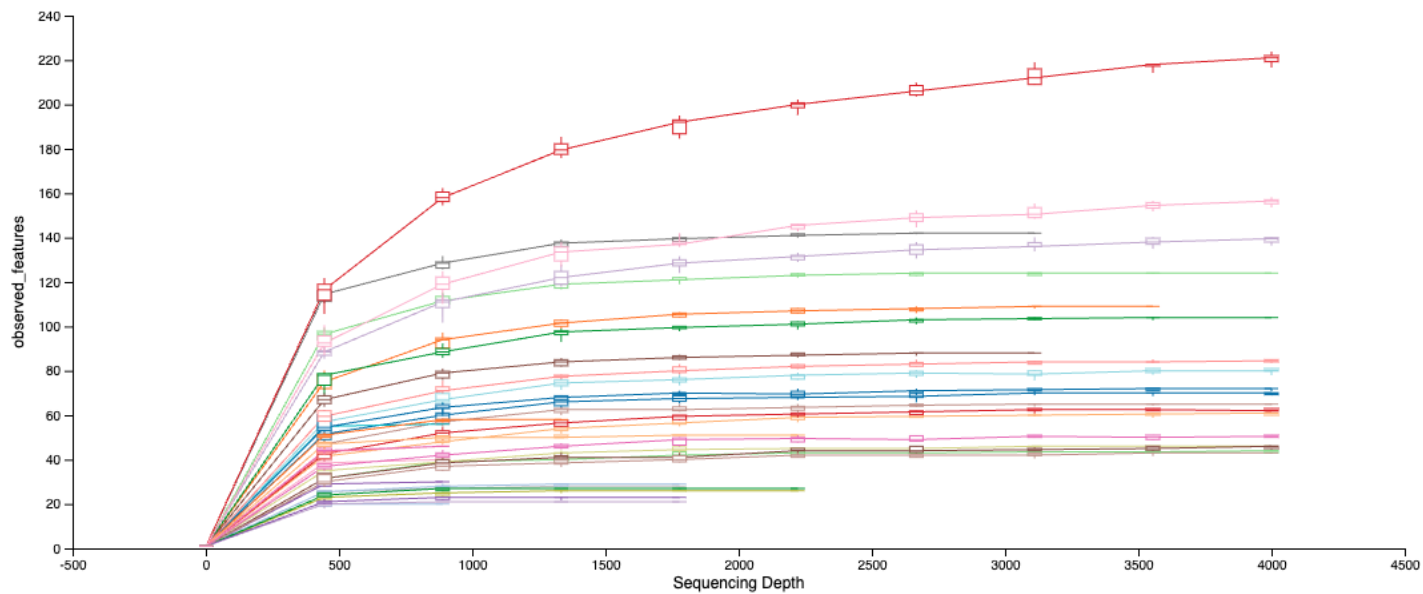
Alexa B. R. McIntyre<sup>1,2,3</sup>, Rachid Ounit<sup>4</sup>, Ebrahim Afshinnekoo<sup>2,3,5</sup>, Robert J. Prill<sup>6</sup>, Elizabeth Hénaff<sup>2,3</sup>, Noah Alexander<sup>2,3</sup>, Samuel S. Minot<sup>7</sup>, David Danko<sup>1,2,3</sup>, Jonathan Foox<sup>2,3</sup>, Sofia Ahsanuddin<sup>2,3</sup>, Scott Tighe<sup>8</sup>, Nur A. Hasan<sup>9,10</sup>, Poorani Subramanian<sup>9</sup>, Kelly Moffat<sup>9</sup>, Shawn Levy<sup>11</sup>, Stefano Lonard<sup>8</sup>, Nick Greenfield<sup>7</sup>, Rita R. Colwell<sup>12</sup>, Gail L. Rosen<sup>13</sup> and Christopher E. Mason<sup>2,3,14\*</sup>



A mock community is key here

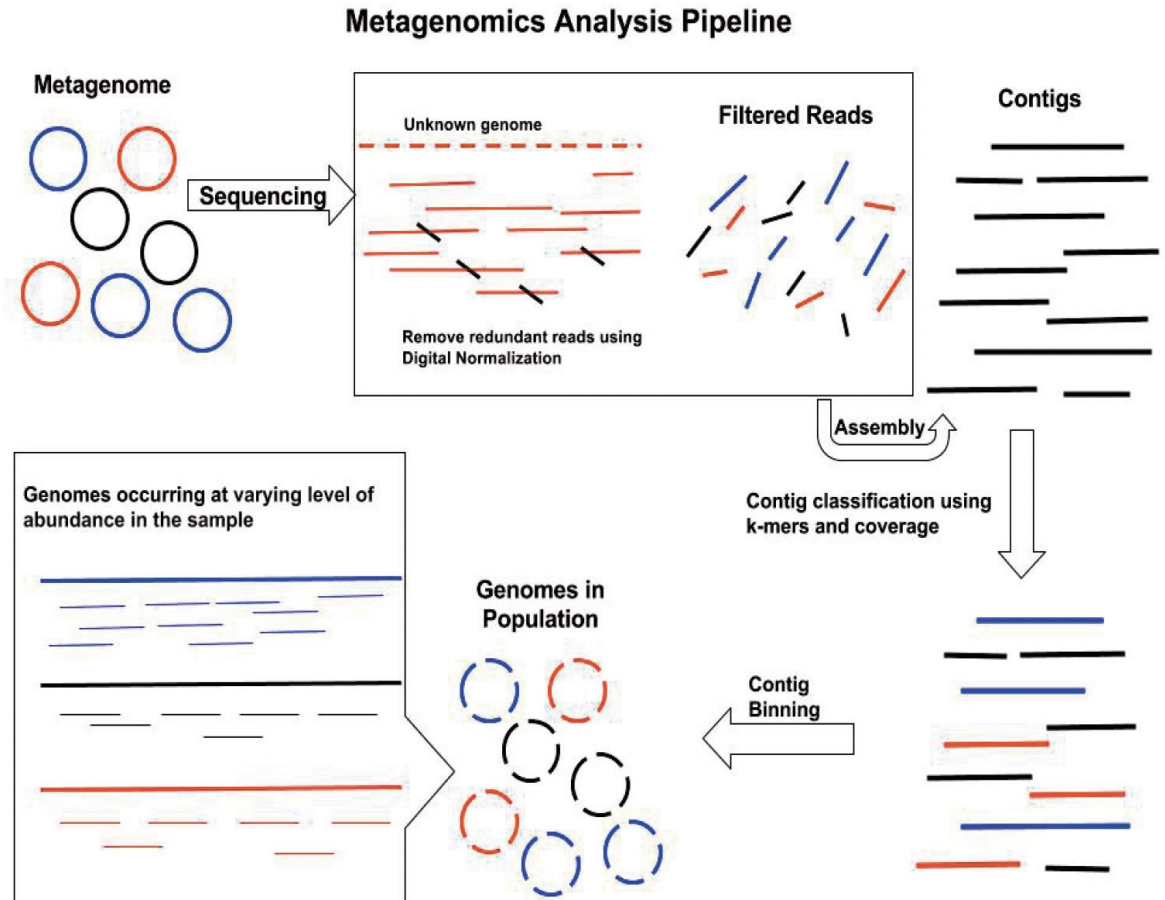
# What more can we do?





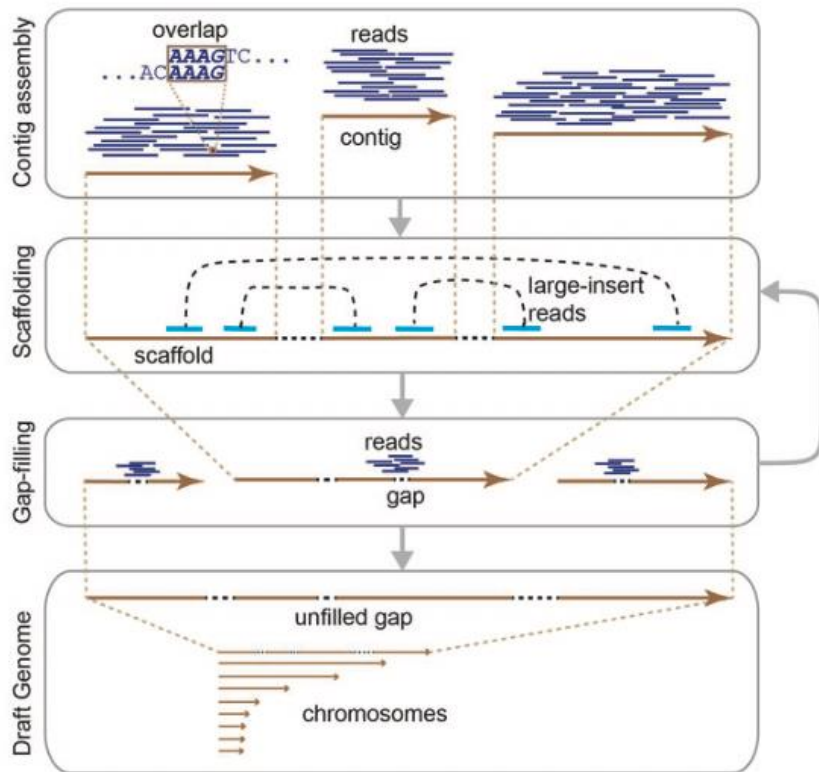
# Genomics & Metagenomics

- Sequencing DNA from isolates or communities
- Produces FASTQ files → **assembly** → annotation → gene counts
- Tools: FastQC, MEGAHIT, MetaSpades, MetaBAT, eggNOG-mapper



# ASSEMBLY....

‘Denovo’ versus ‘Mapping’

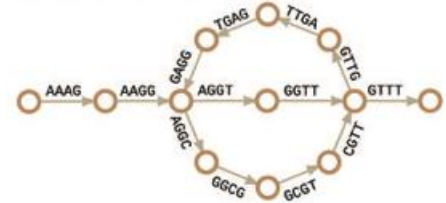


A. Short read to  $k$ -mers ( $k=4$ )

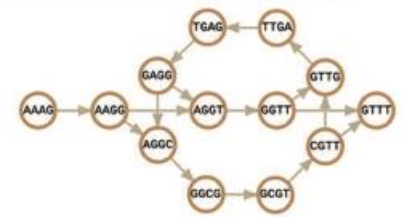
AAAGGCGTTGAGGTT

AAAG  
AAGG  
AGGC  
GGCG  
GCGT  
CGTT  
GTTG  
TTGA  
TGAG  
GAGG  
AGGT  
GGTT

B. Eulerian de Bruijn graph



C. Hamiltonian de Bruijn graph

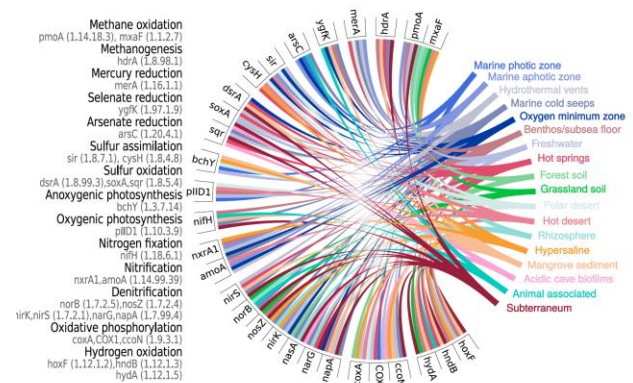


# Transcriptomics (Who's Active?)

- RNA → cDNA → sequencing
- Key step: differential expression across conditions (careful!!how do we normalize??:  
<https://bigomics.ch/blog/why-how-normalize-rna-seq-data/> )
- Analogy: same instrument (genome), different songs (transcripts) will tell a different story

# Statistical Thinking in OMICS

- High-dimensional abundance tables
- PCA, NMDS, clustering, regression, random forests
- Ecological parallels: taxa  $\rightleftharpoons$  genes; community  $\rightleftharpoons$  pathways



## End product:

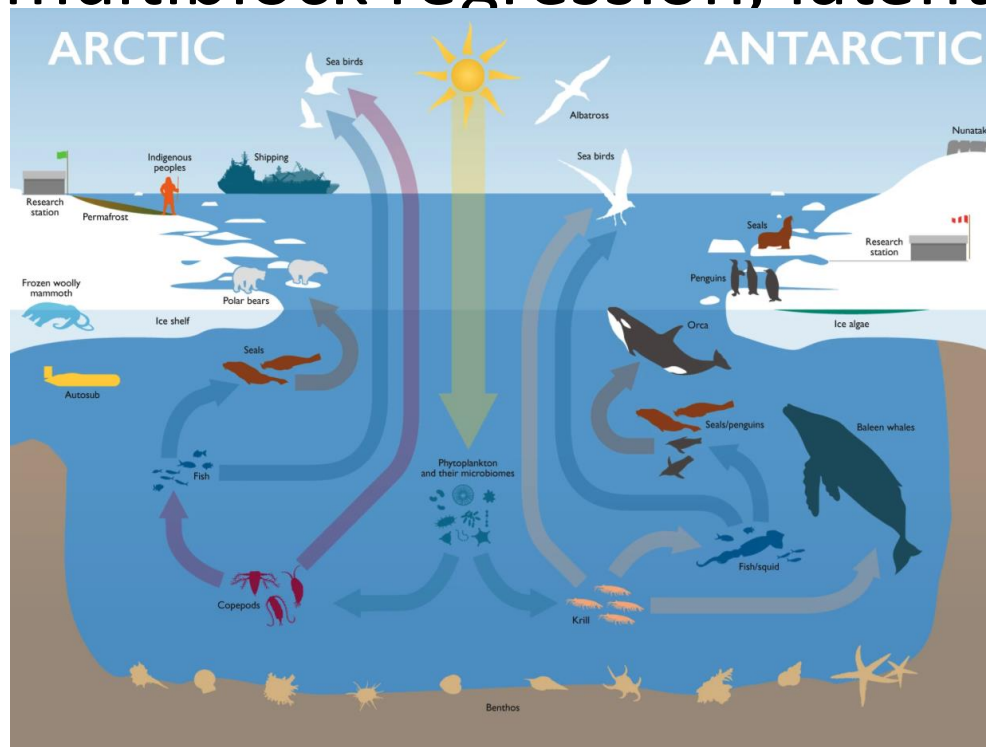
Omics gives you big multivariate matrices; your PCA/NMDS/regression/... skills still apply here.

# Typical Analytical Questions

- Diversity: how rich and even is the community?
- Function: which pathways dominate where?
- Differential features: which respond to stress?
- Linkages: environment → function → ecosystem process

# Integrating Multiple 'OMES'

- Multi-omics connects layers for causality
- Methods: MixOmics, MOFA, DIABLO
- View: multiblock regression, latent variables



# Common Pitfalls

- **Compositionality** & false correlations
- **Batch effects** / unequal coverage
- **Unknown genes**
- **Over-interpretation** of small effects

These are design and analysis issues – not reasons to avoid omics!

# Computational & Ethical Landscape

- HPC, cloud pipelines, reproducibility (Nextflow, Snakemake, Docker)
- FAIR + CARE data principles:
- FAIR + CARE principles help ensure data are re-usable and that communities benefit, not only high-income labs.

# Future Horizons

- Portable sequencing (Nanopore/Pacbio field kits)
- AI/ML for annotation & pattern discovery
- Predictive eco-models linking OMICS → fluxes → climate

# Summary

- OMICS = quantitative molecular ecology
- Pipelines bridge raw data to ecological insight
- Integration & reproducibility are the next frontier

# Acknowledgements & Further Resources

MARiS | SEEC | DIPLOMICS

## **Key references:**

Knight et al. 2018 – Nature Biotech (Microbiome analysis)

Misra et al. 2022 – Nature Methods (Multi-omics)

Love et al. 2014 – DESeq2; Rohart et al. 2017 – mixOmics

## Pipelines:

<https://www.geneious.com>

<https://mothur.org>

<https://qiime2.org>

<https://nf-co.re/ampliseq/2.6.1/>

## DATABASES



**National Library of Medicine**  
*National Center for Biotechnology Information*

<https://www.ncbi.nlm.nih.gov>

**ENA**  
European Nucleotide Archive



<https://www.ddbj.nig.ac.jp/index-e.html>

**BOLD** SYSTEMS

[Boldsystems.org](https://boldsystems.org)