

# Exploring microbiomes with cultivation-independent genome-resolved metagenomics

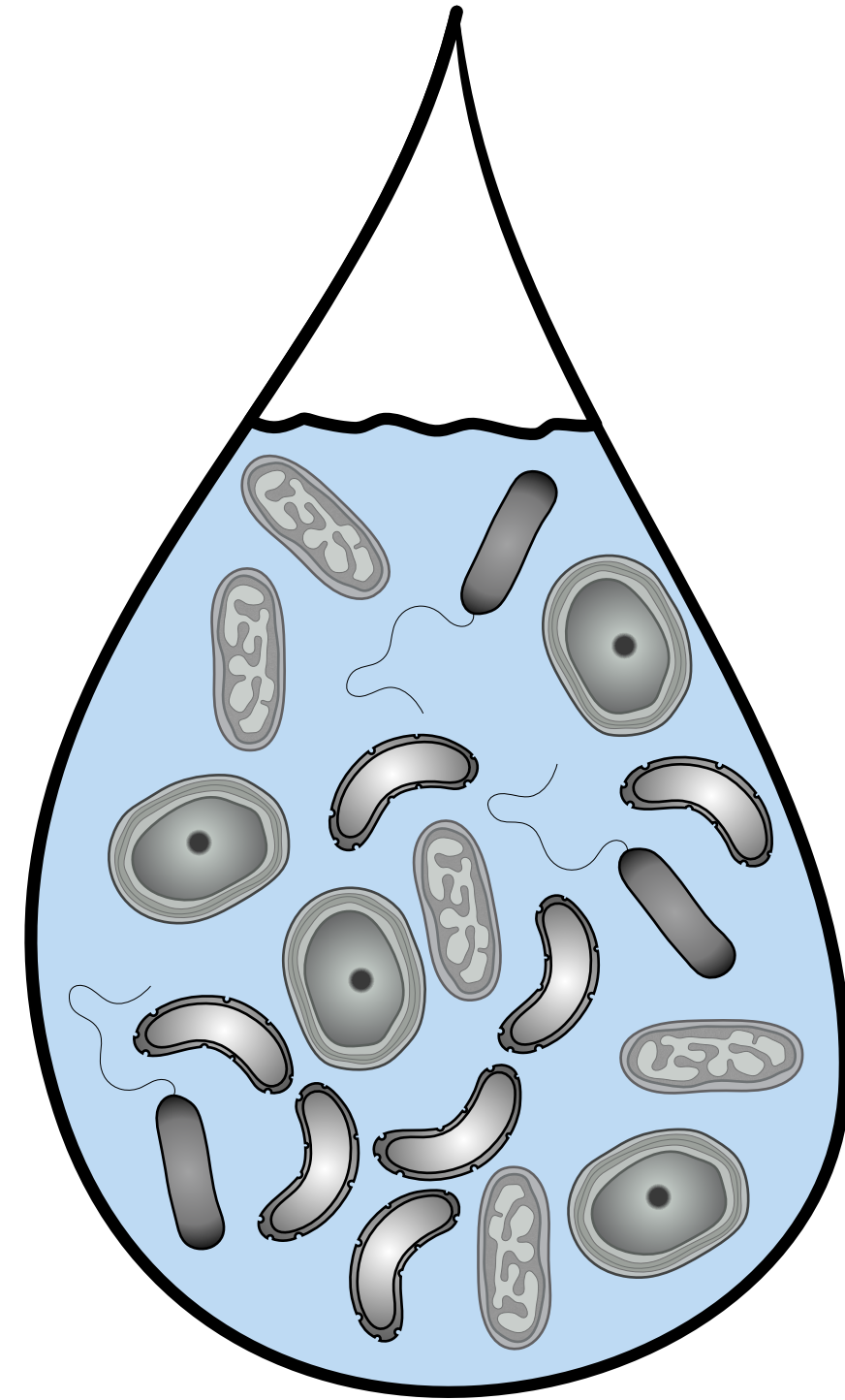
# Why explore microbiomes?

- Ubiquitous across earth's ecosystems
- Support global food webs
- Underpin biogeochemical cycles
- Determine Host's health and disease
- ...
- Untapped metabolic diversity

**Oceans cover >70% of the planet**



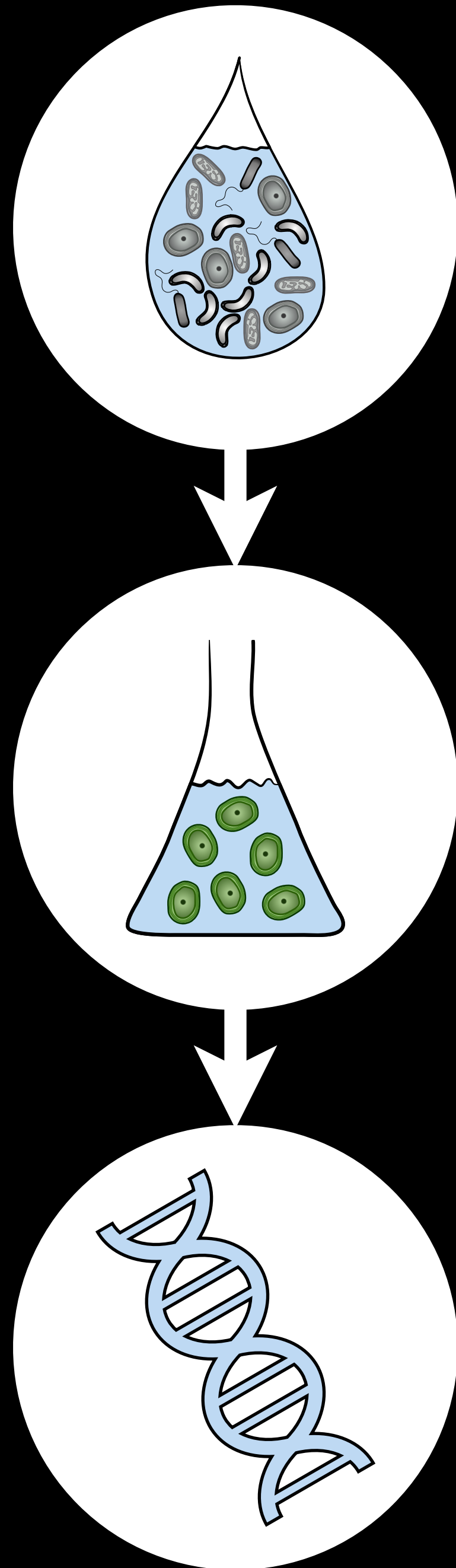
- > **500,000 microbial cell per mL**
- > **50% of the oxygen production**



# Traditional microbiology

# Traditional microbiology

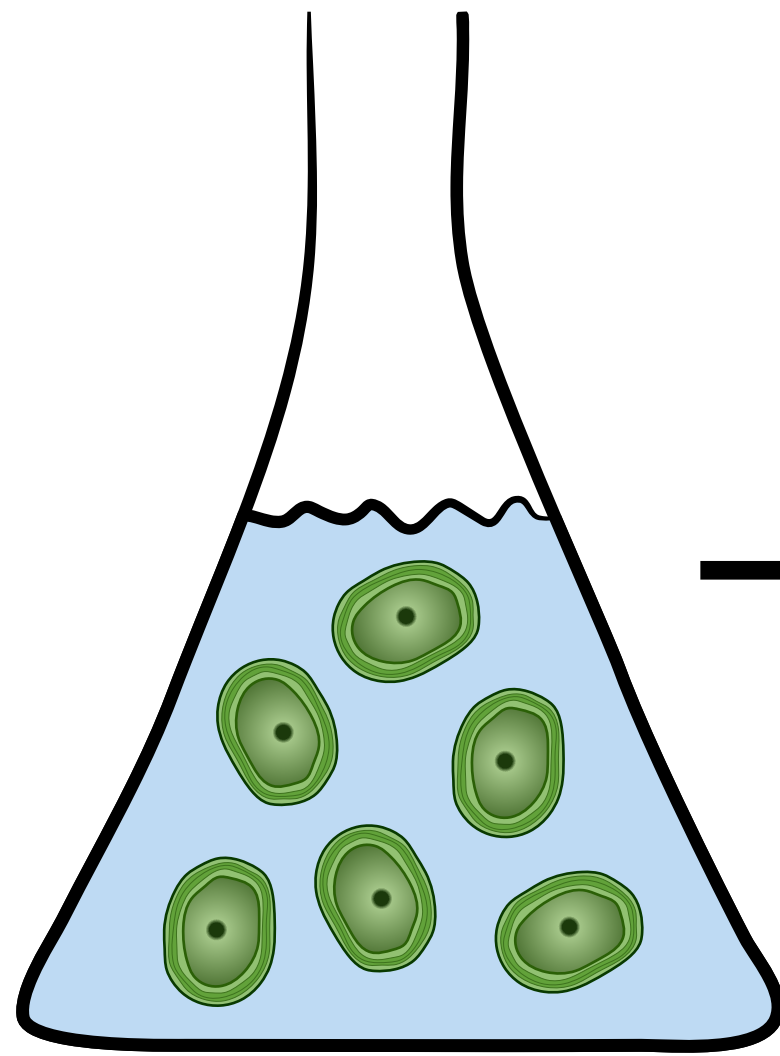
## Cultivation-based analysis of microbiomes



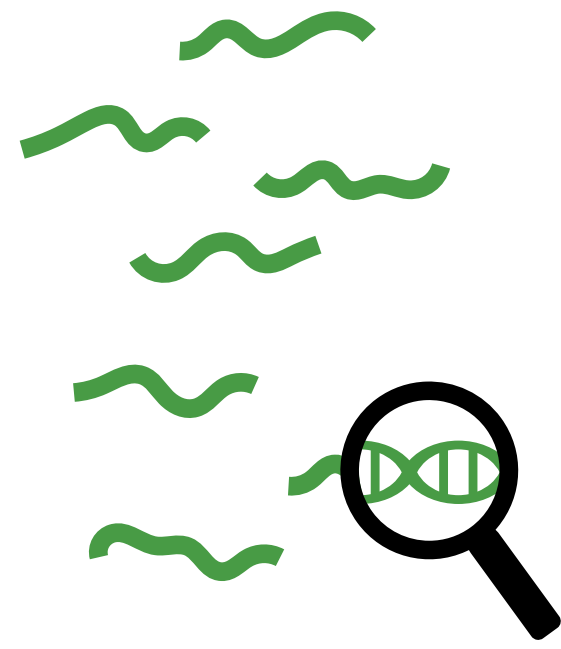
- Isolate
- Cultivate
- Sequence



Cultivate



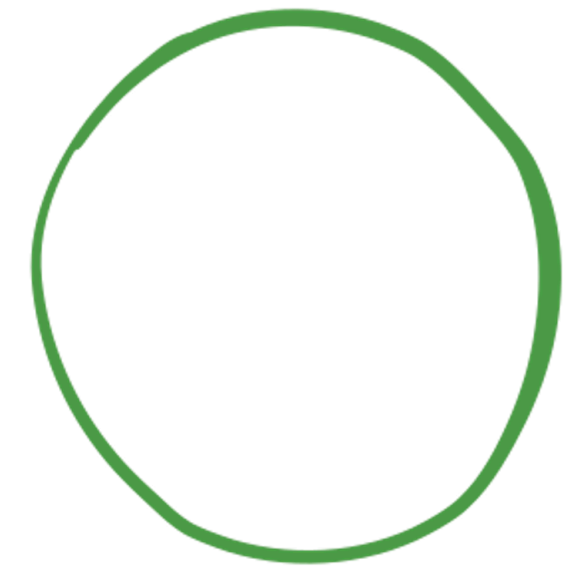
Extract DNA



Sequence fragments



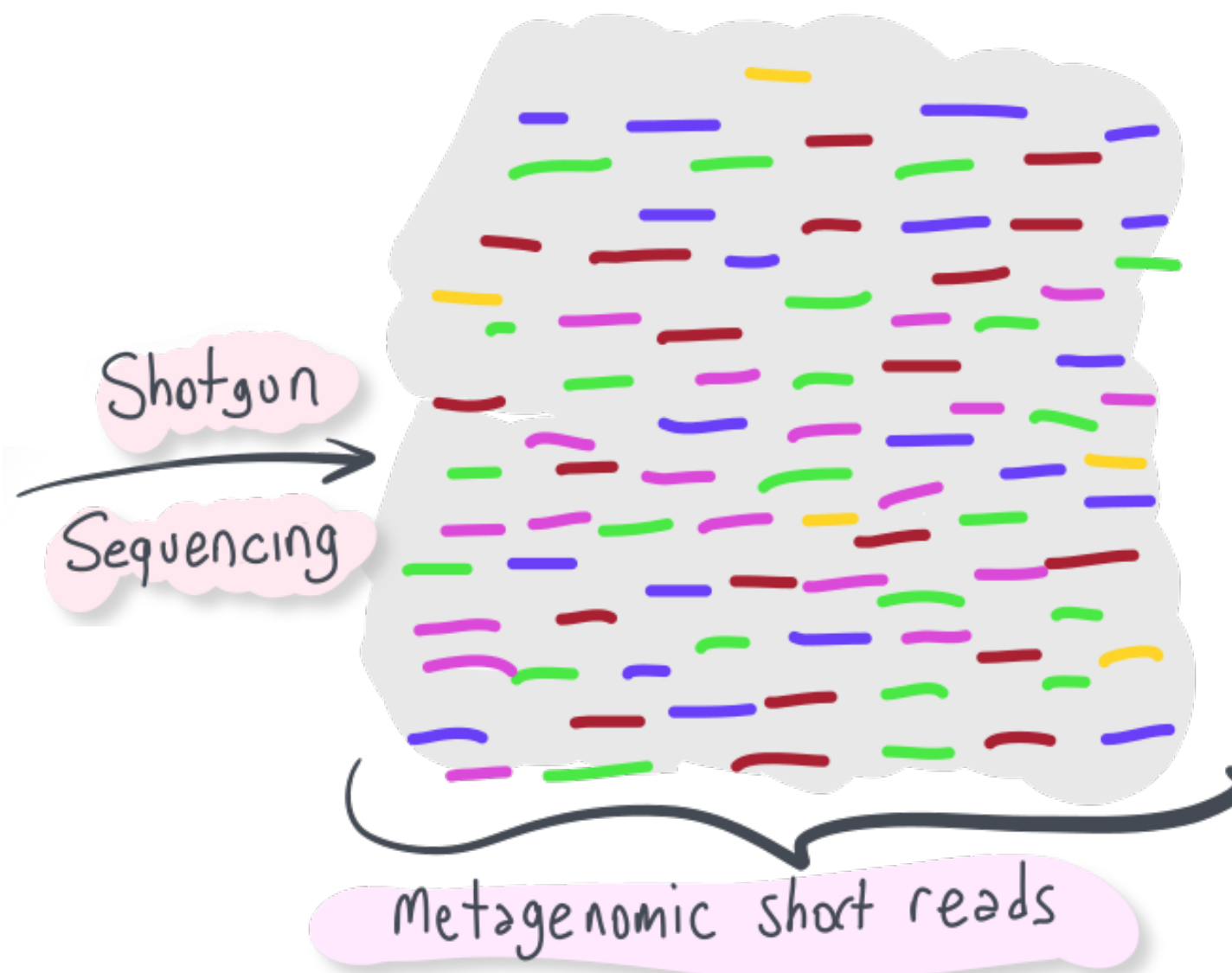
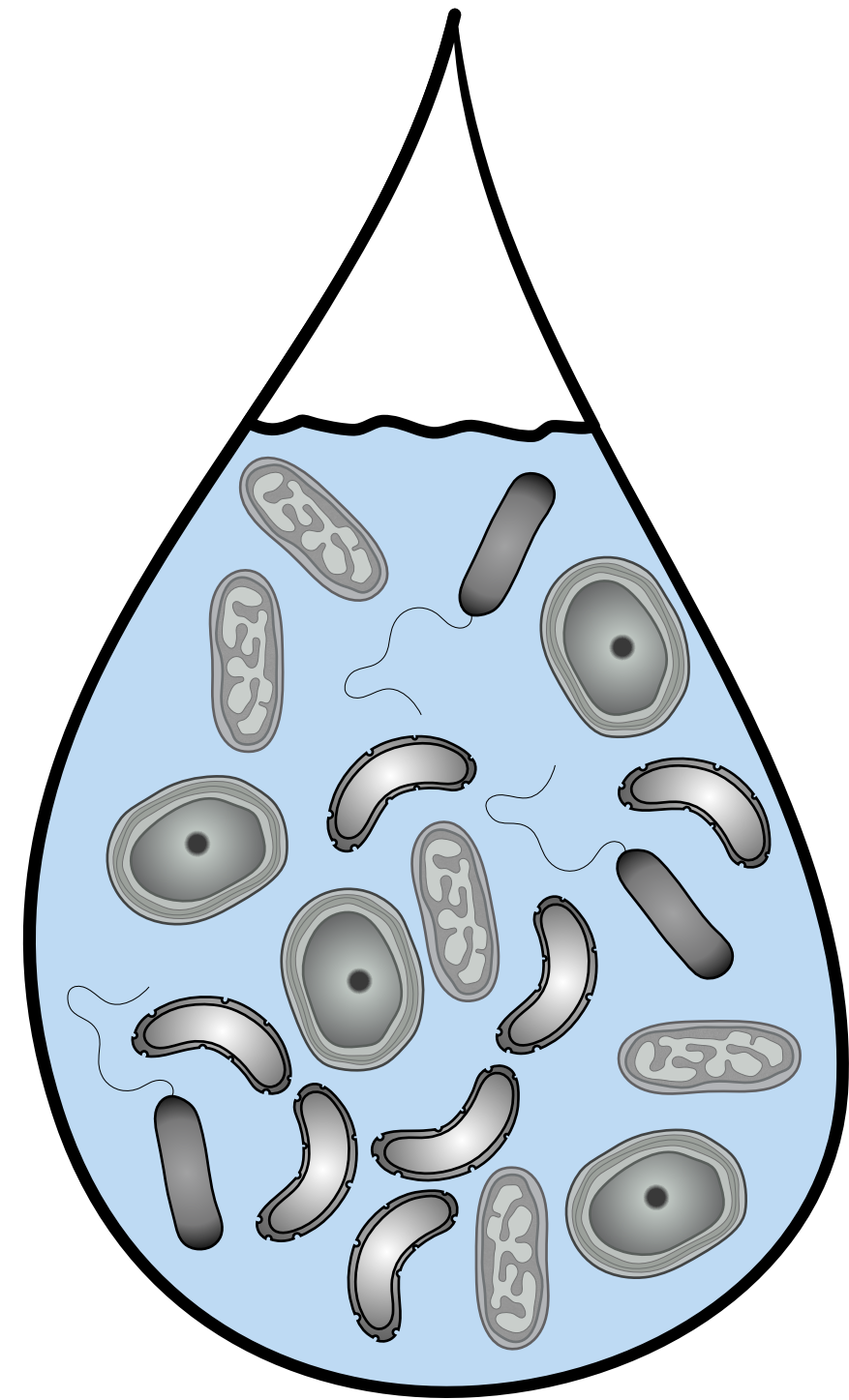
Reconstruct genome



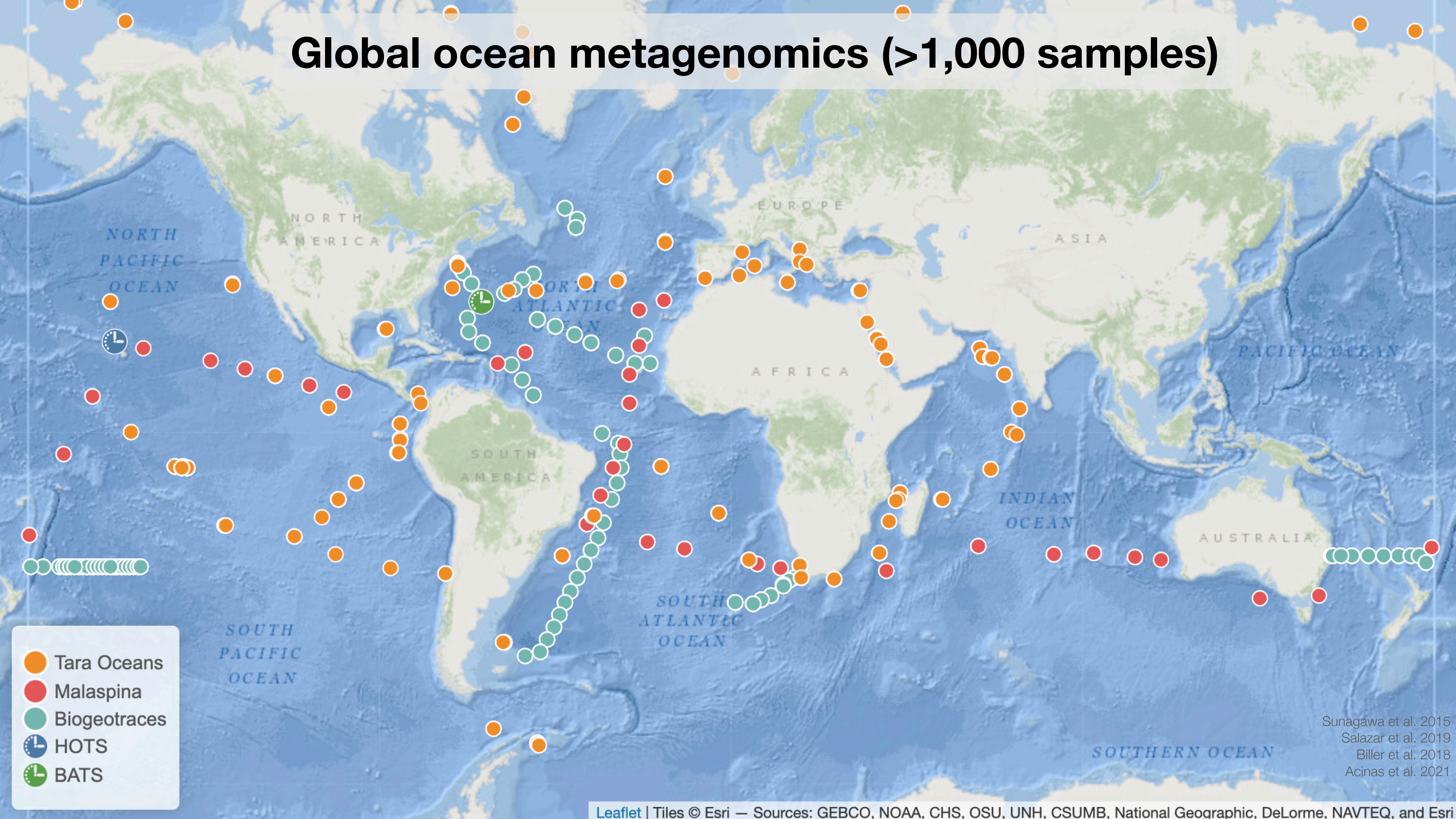


# Metagenomics

Accessing the genomic content of microbiomes

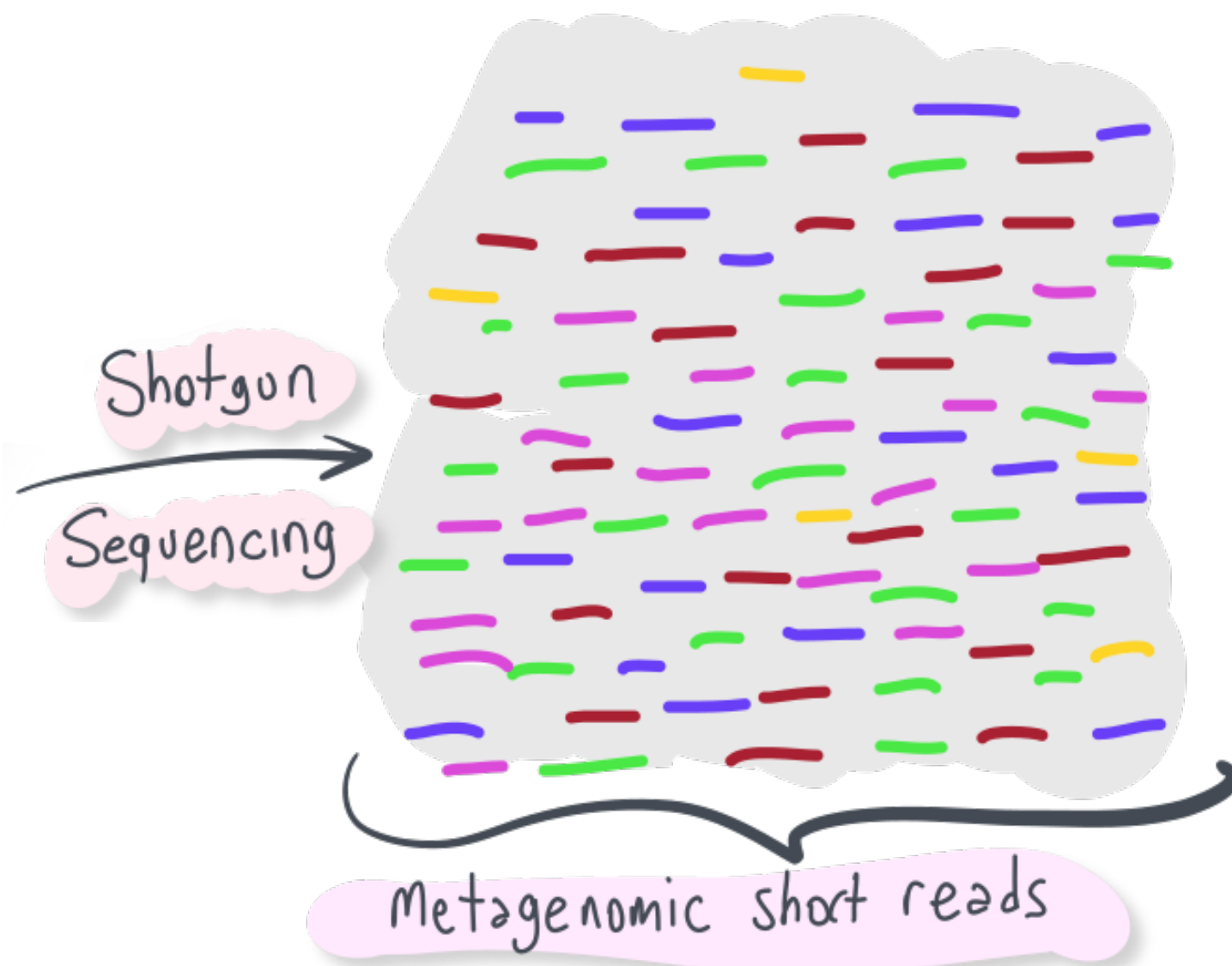
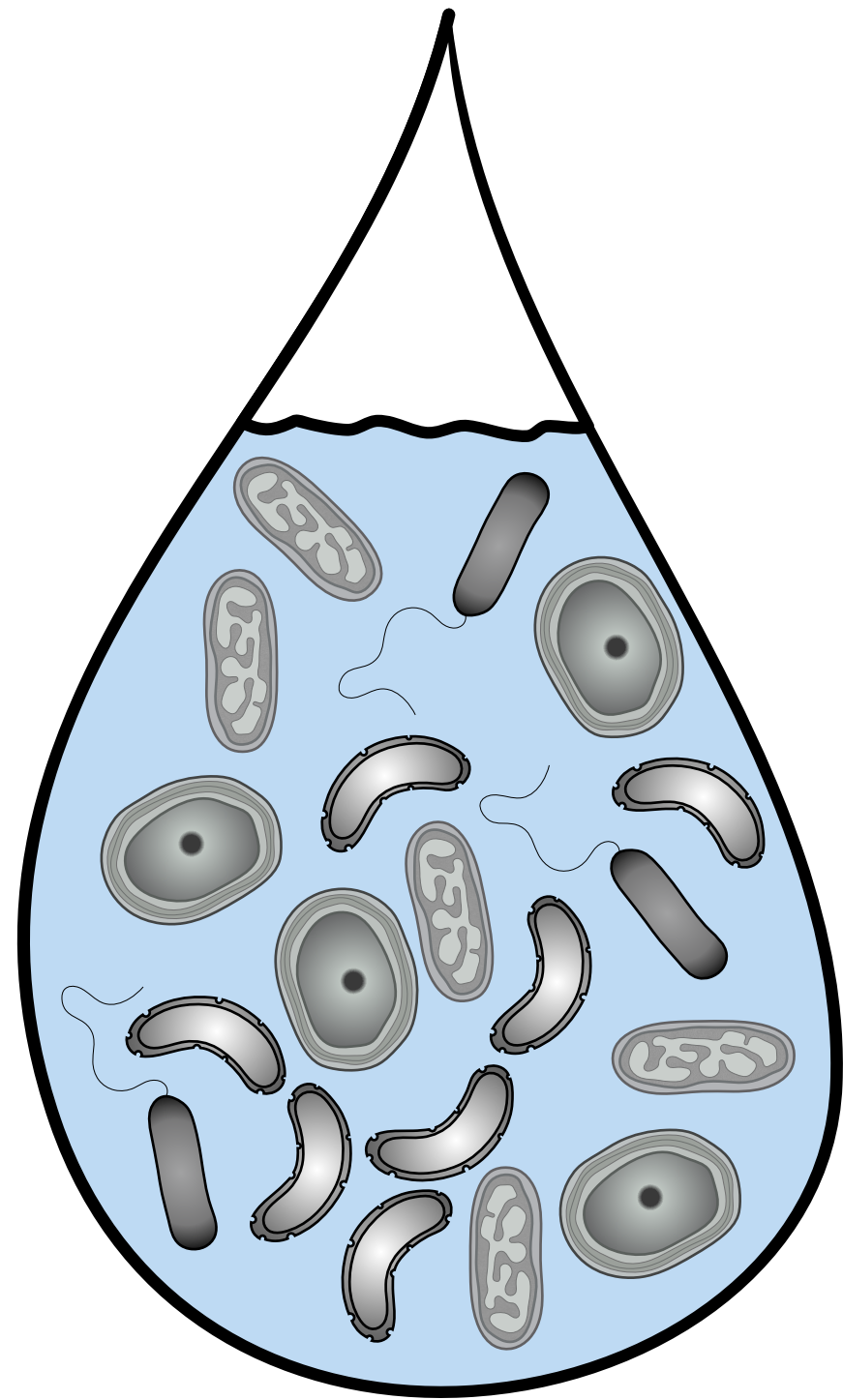


# Global ocean metagenomics (>1,000 samples)



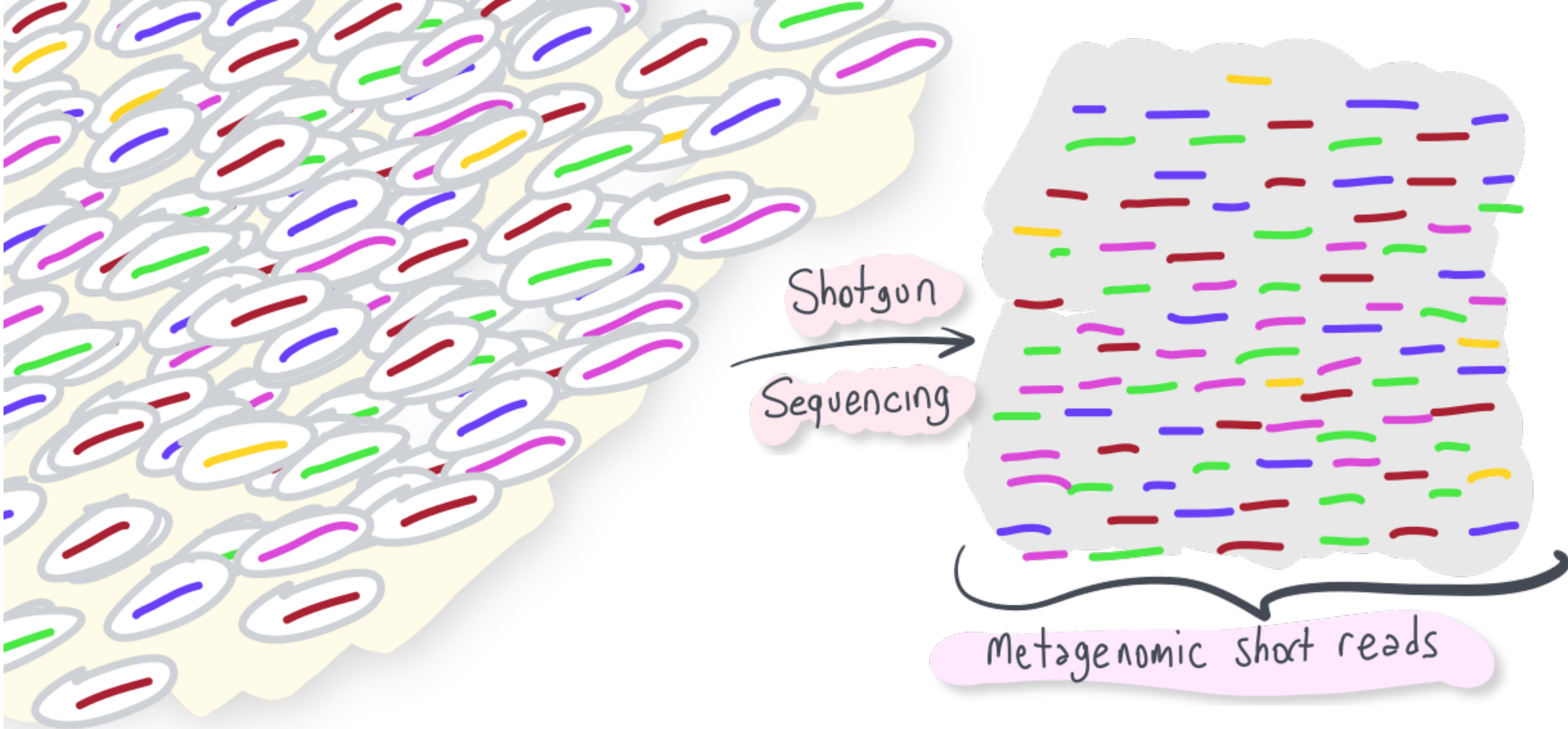
- Tara Oceans
- Malaspina
- Biogeotraces
- HOTS
- BATS

Sunagawa et al. 2015  
Salazar et al. 2019  
Biller et al. 2018  
Acinas et al. 2021

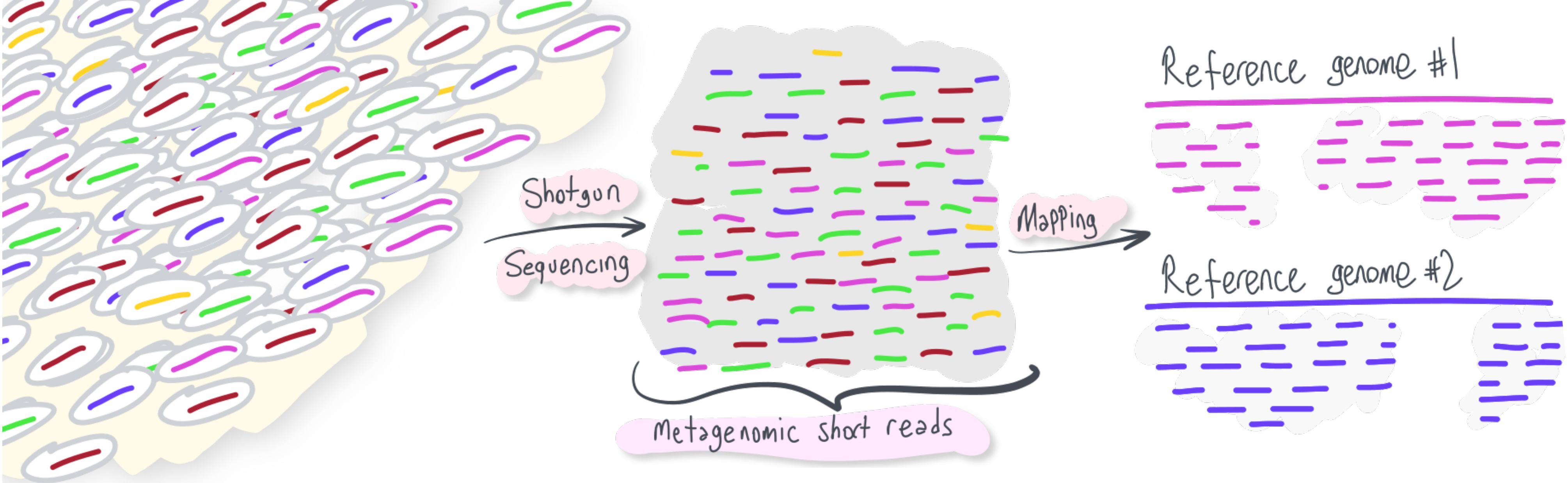


# Genome-resolved Metagenomics

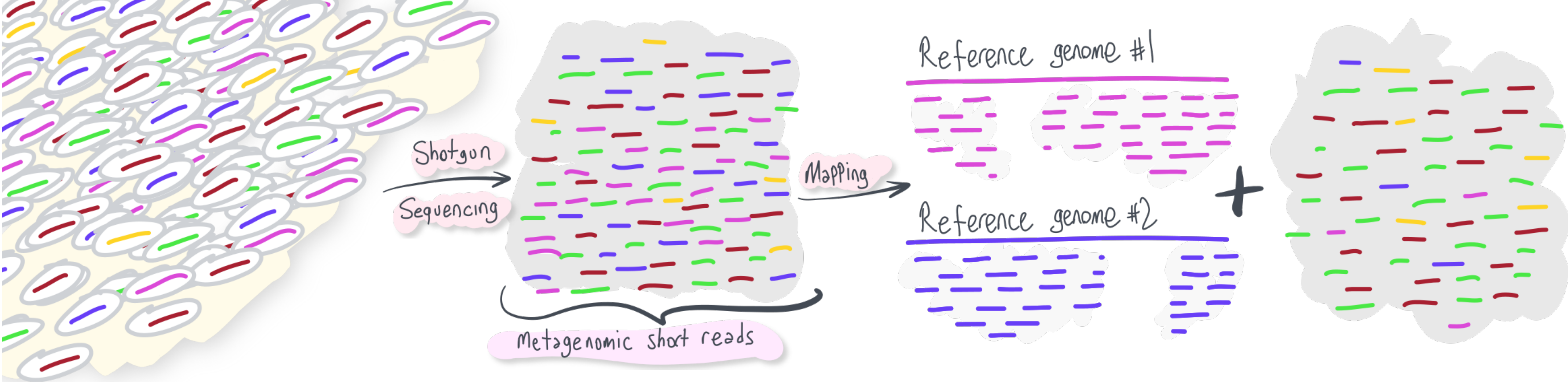
Contextualizing the genomic content of microbiomes



# GENOME RESOLVED METAGENOMICS



# GENOME RESOLVED METAGENOMICS

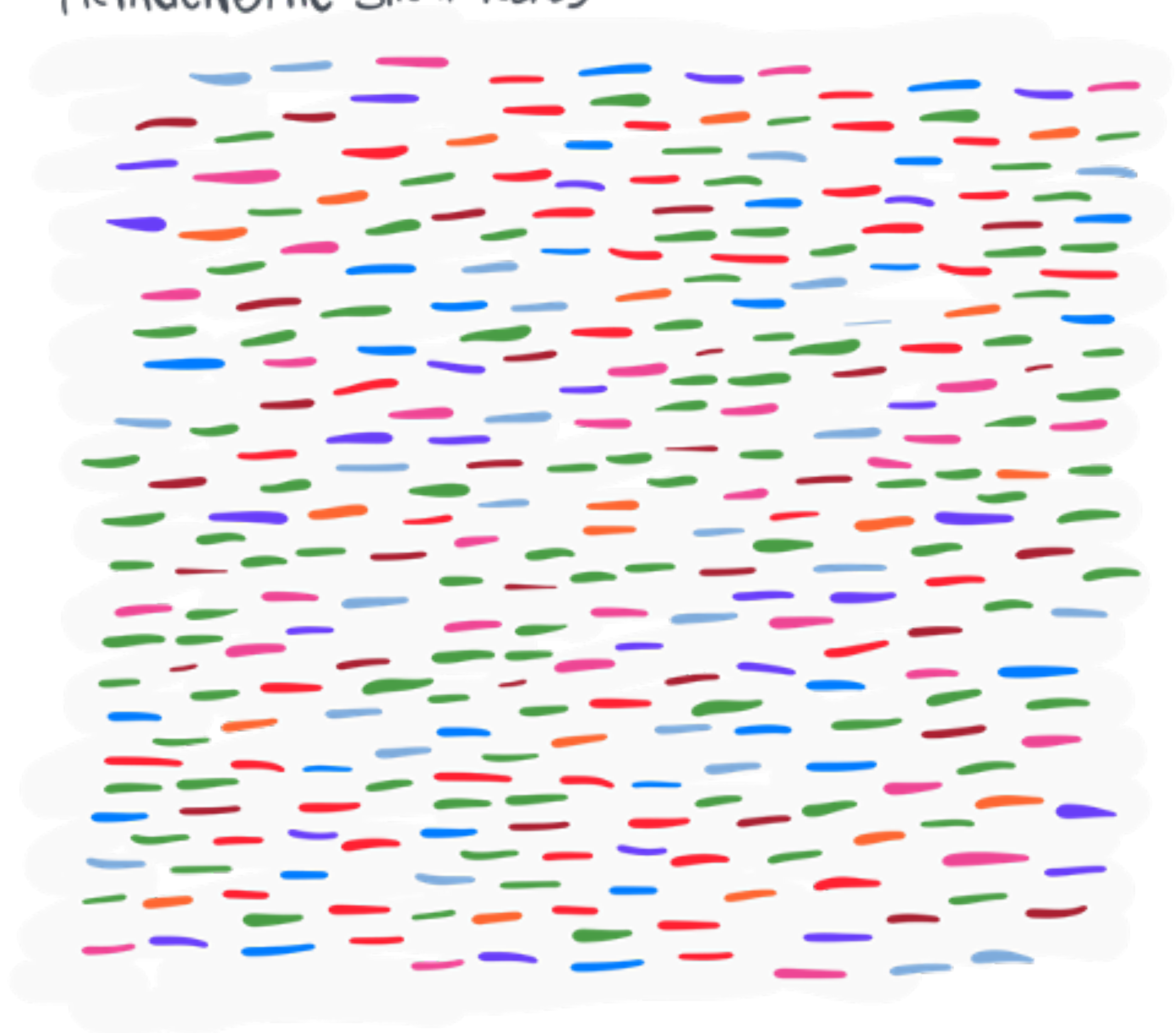


# GENOME RESOLVED METAGENOMICS

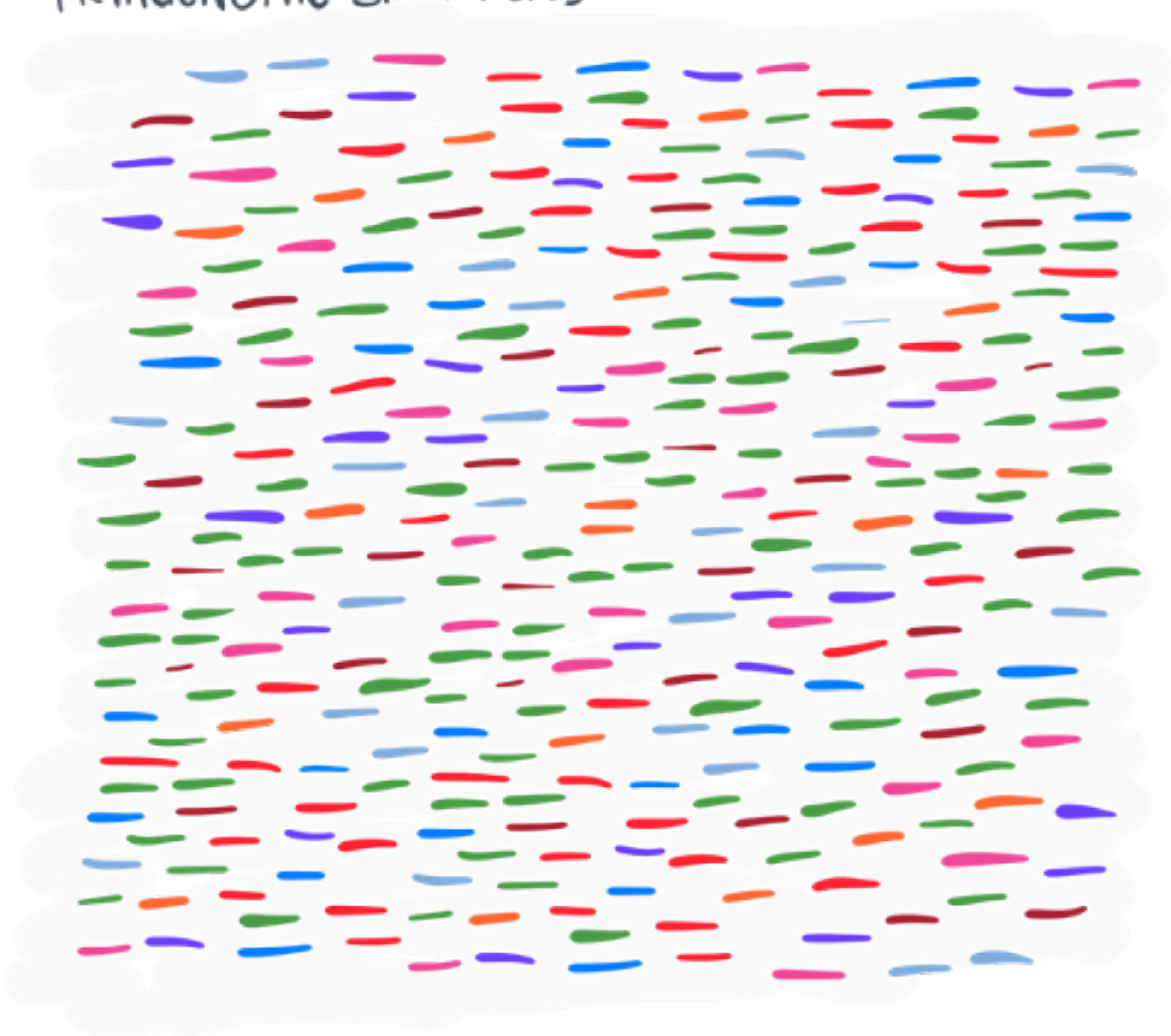


Mapping a metagenome  
Alignment of reads to a reference

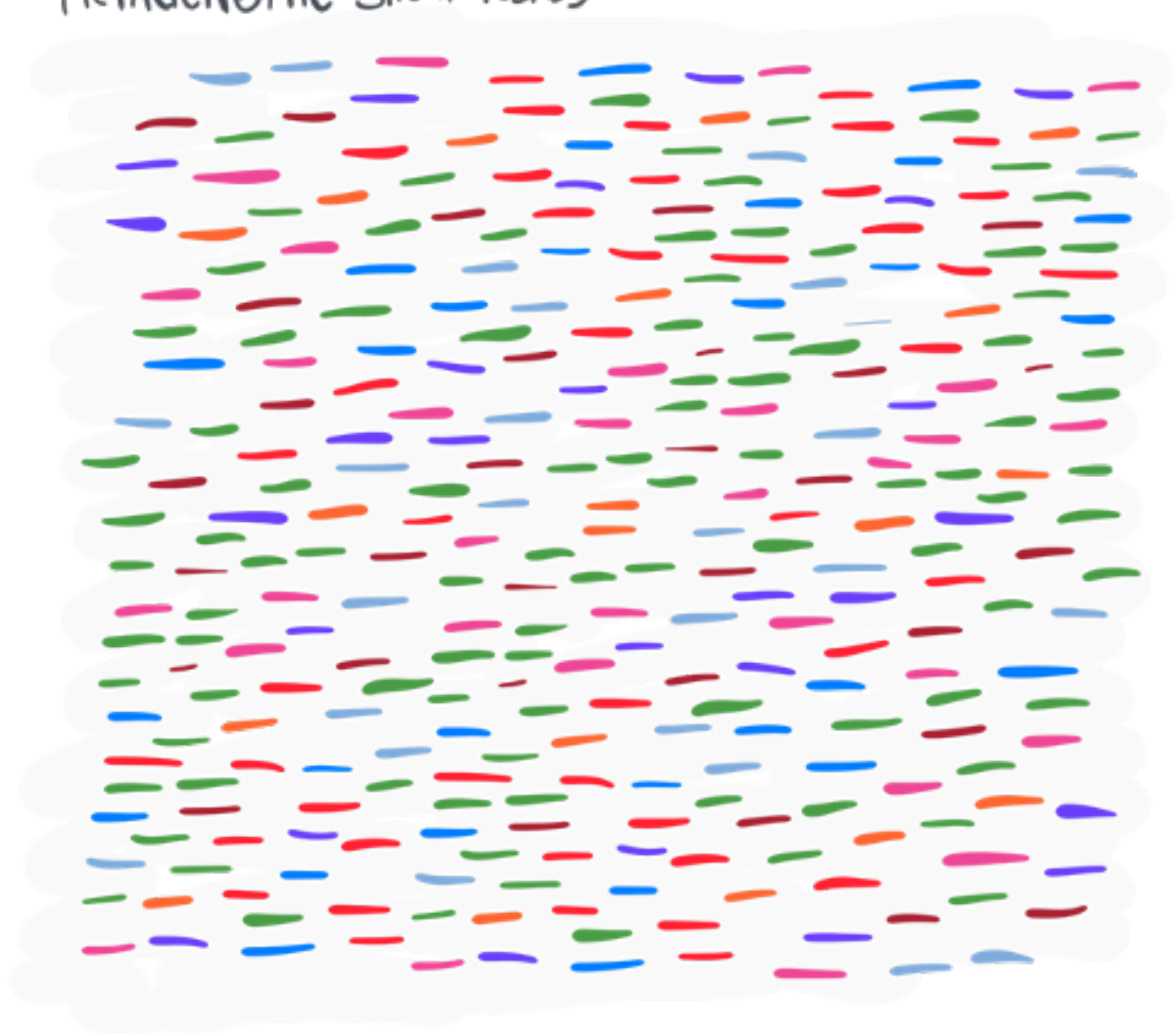
# METAGENOMIC SHORT READS



# METAGENOMIC SHORT READS



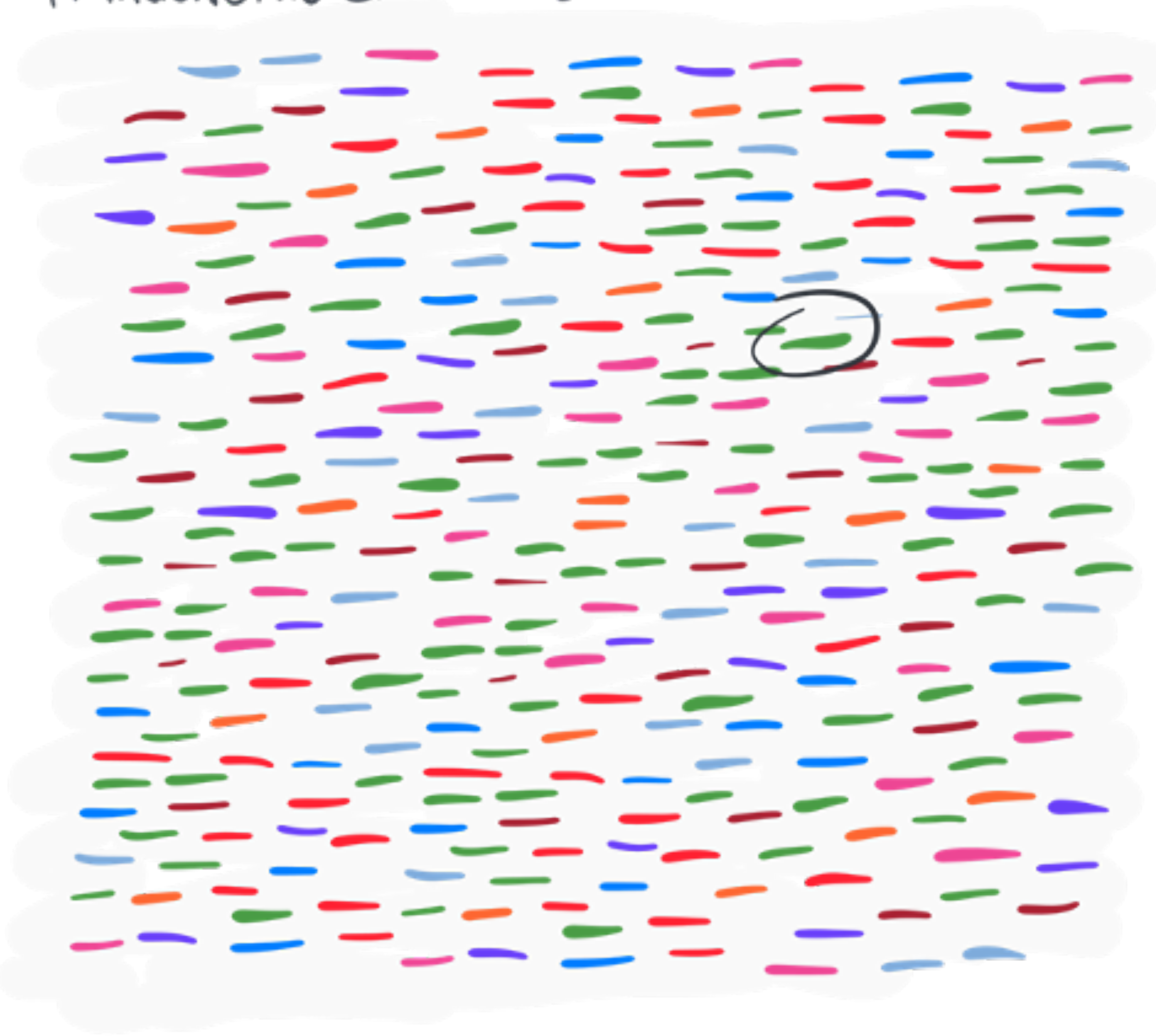
# METAGENOMIC SHORT READS



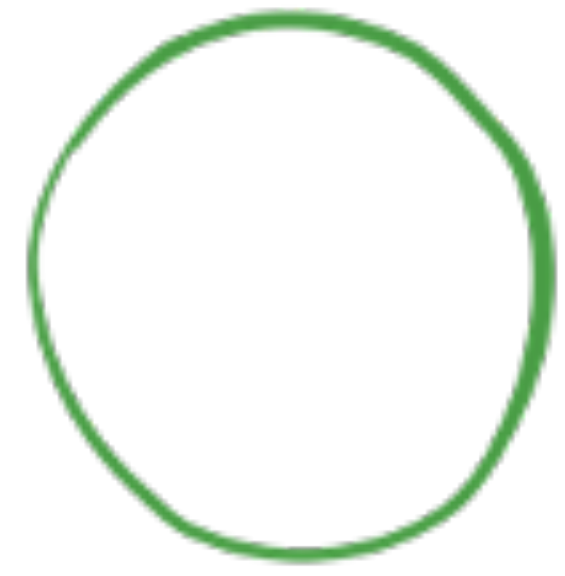
READ  
RECRUITMENT →



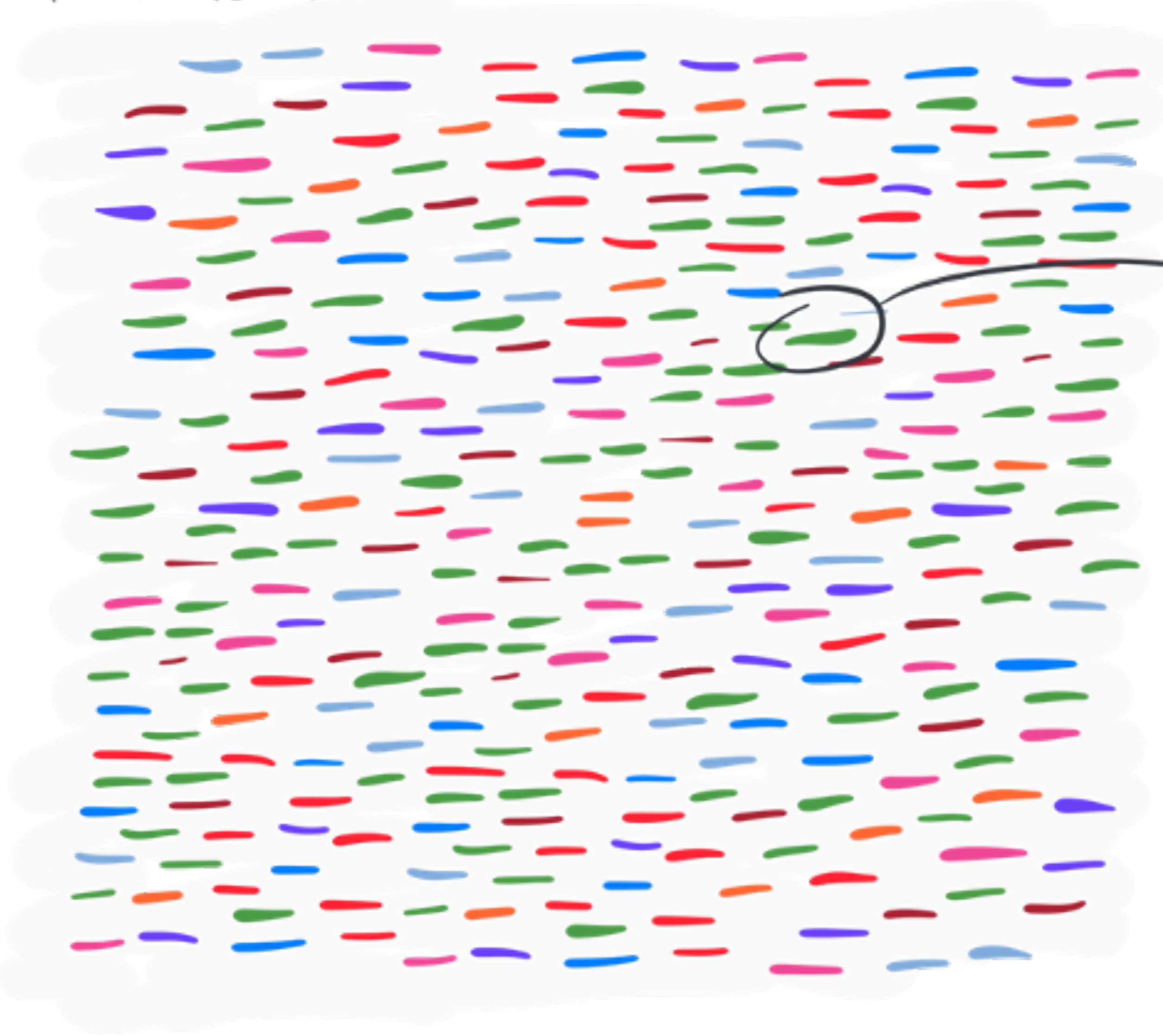
# METAGENOMIC SHORT READS



READ  
RECRUITMENT →



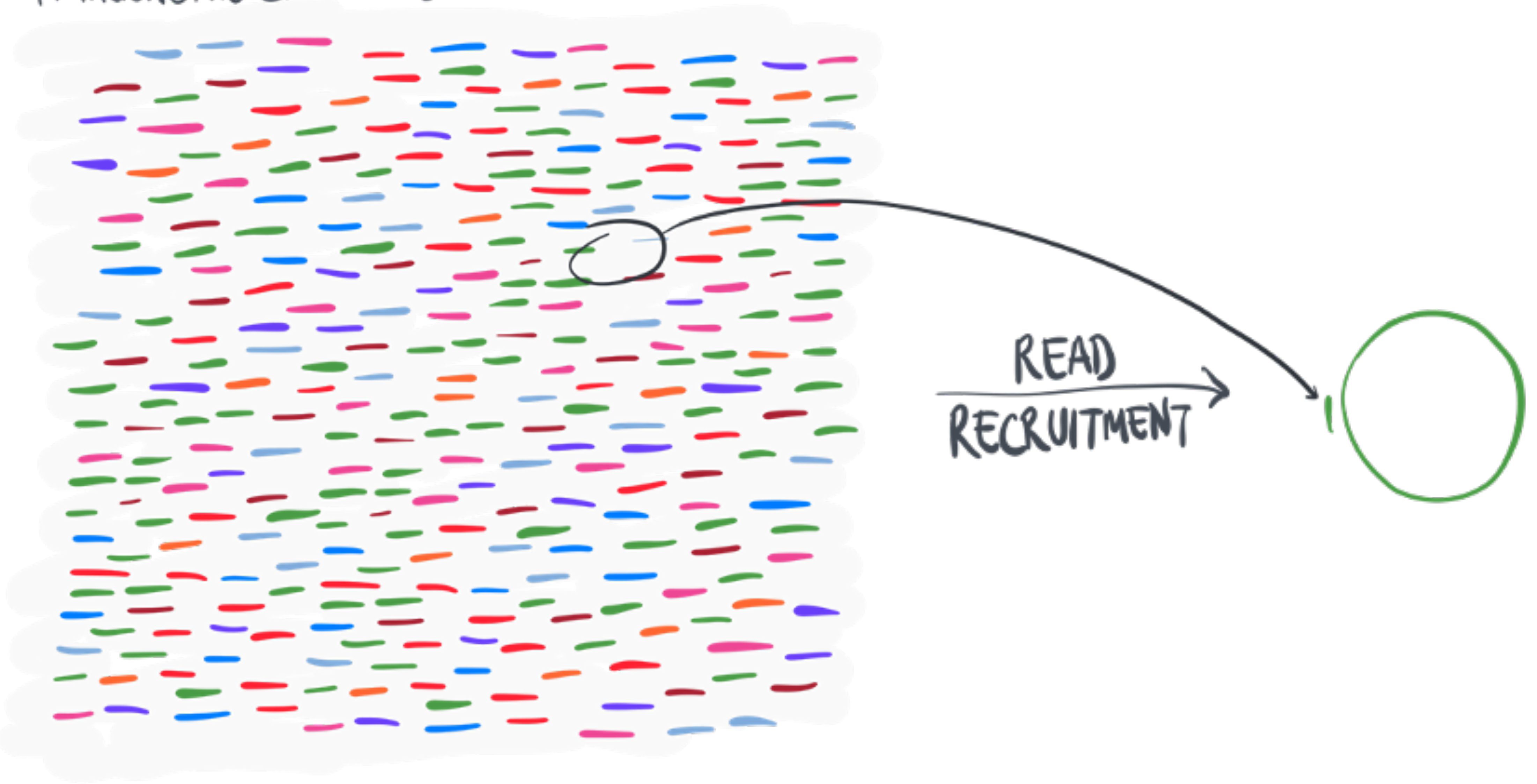
# METAGENOMIC SHORT READS



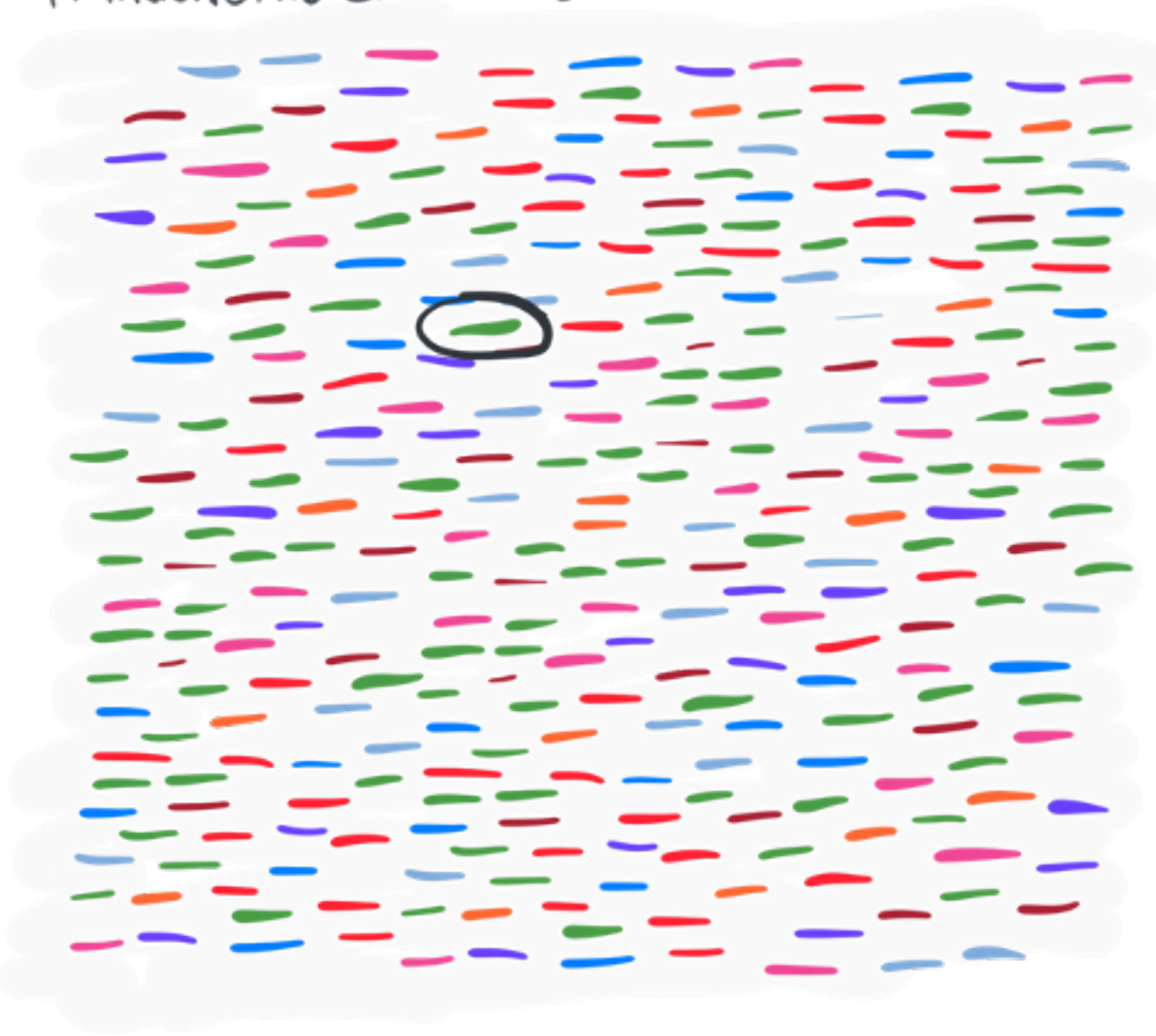
READ  
RECRUITMENT →



# METAGENOMIC SHORT READS



# METAGENOMIC SHORT READS

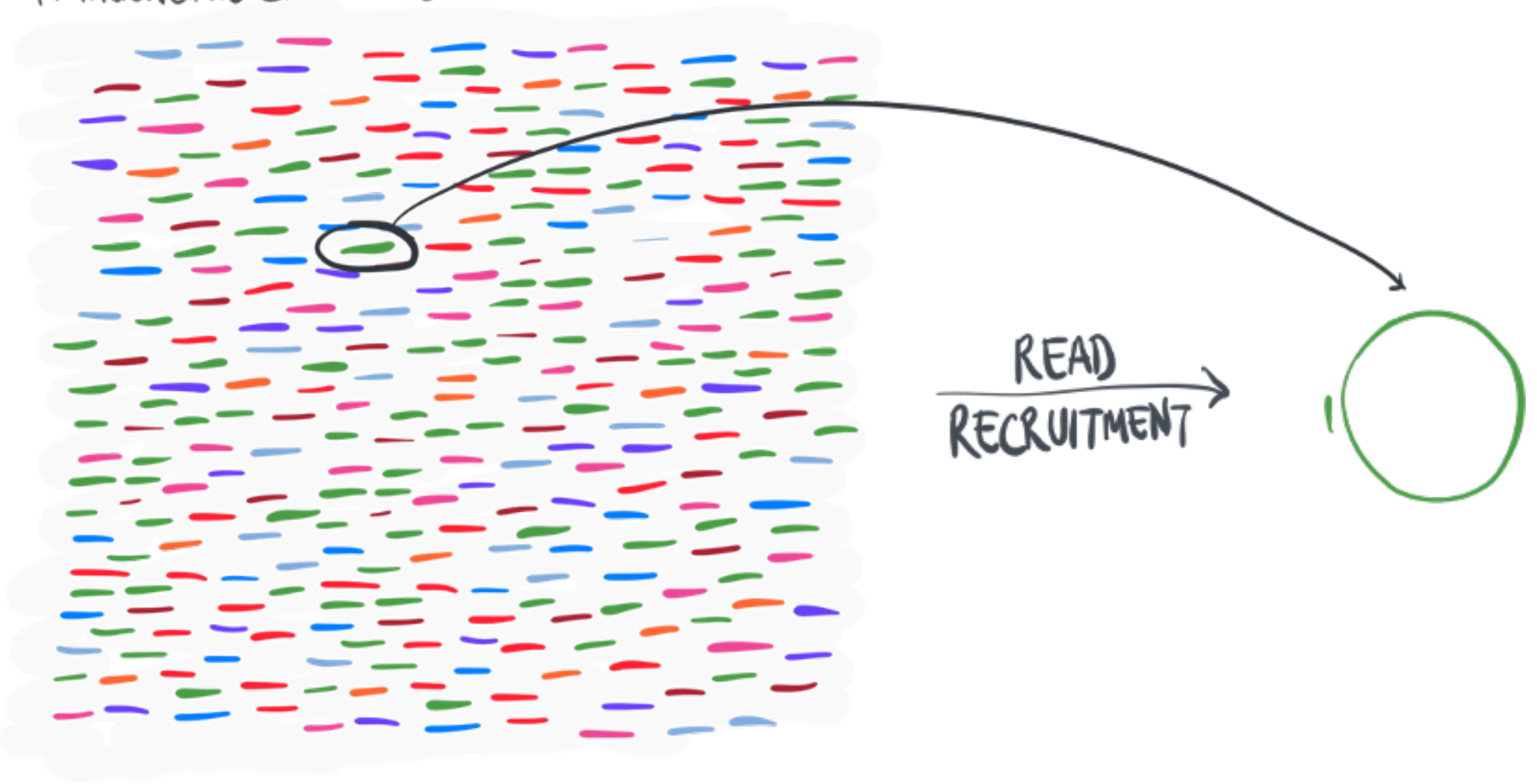


READ  
RECRUITMENT →

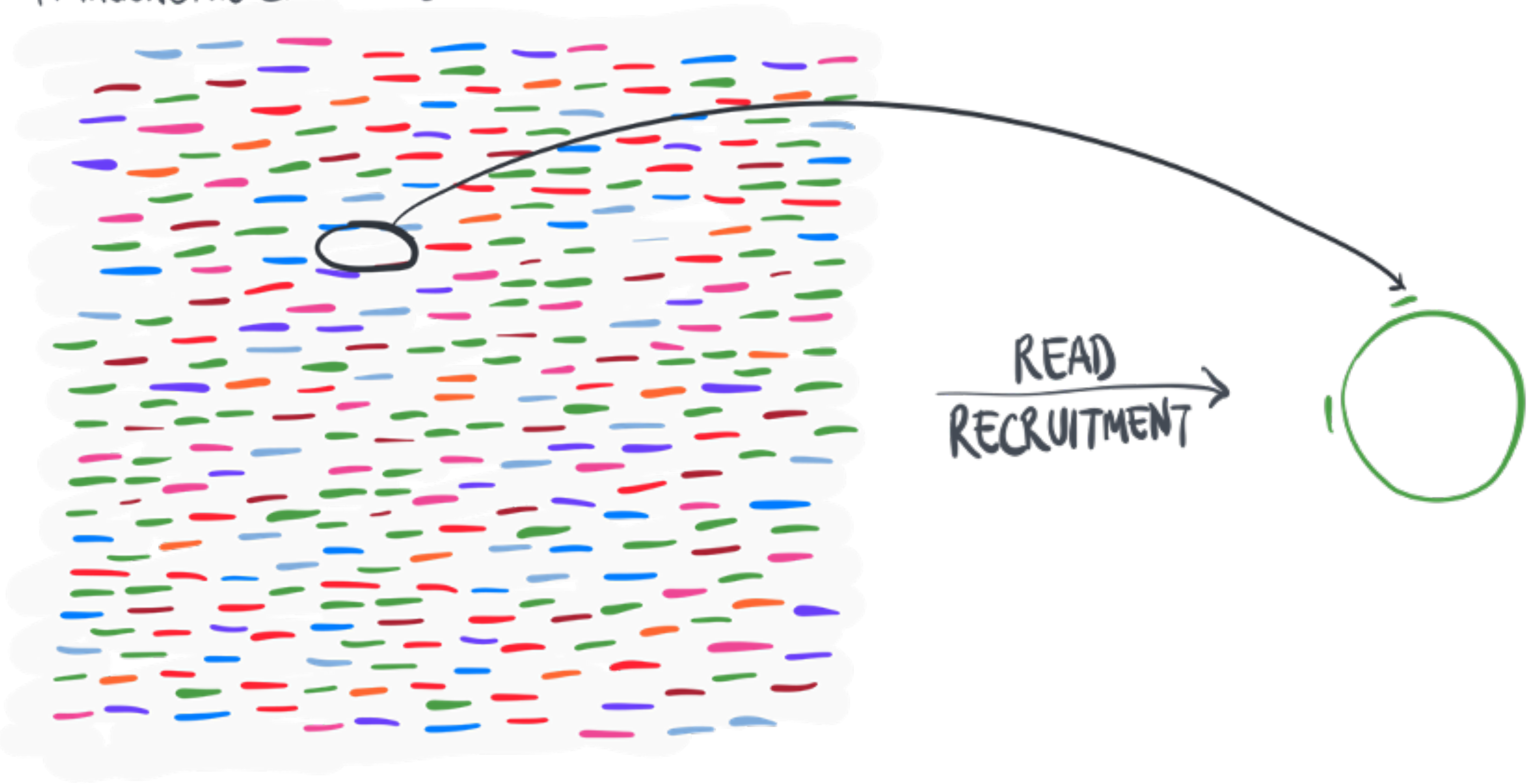




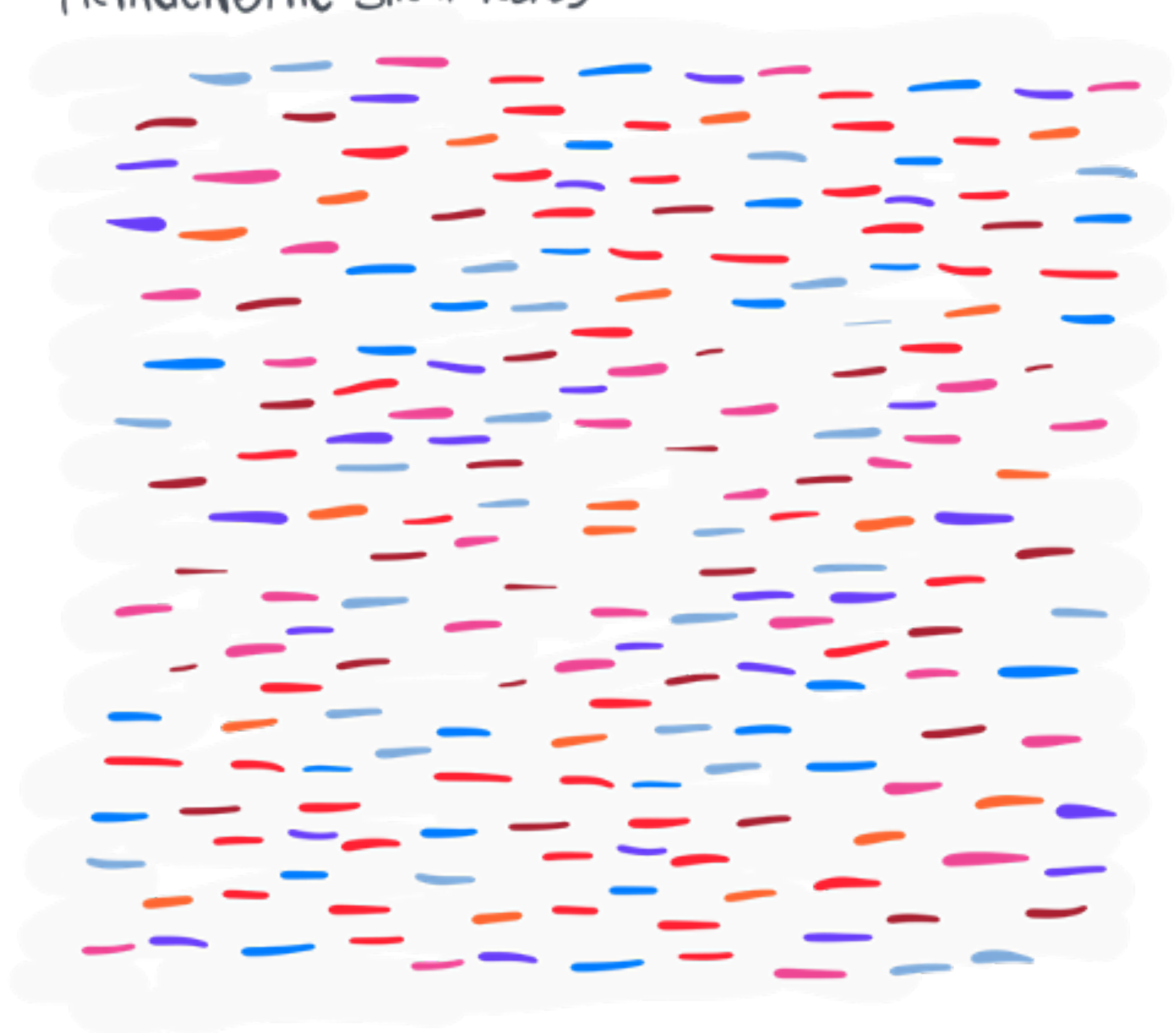
# METAGENOMIC SHORT READS



# METAGENOMIC SHORT READS



# METAGENOMIC SHORT READS

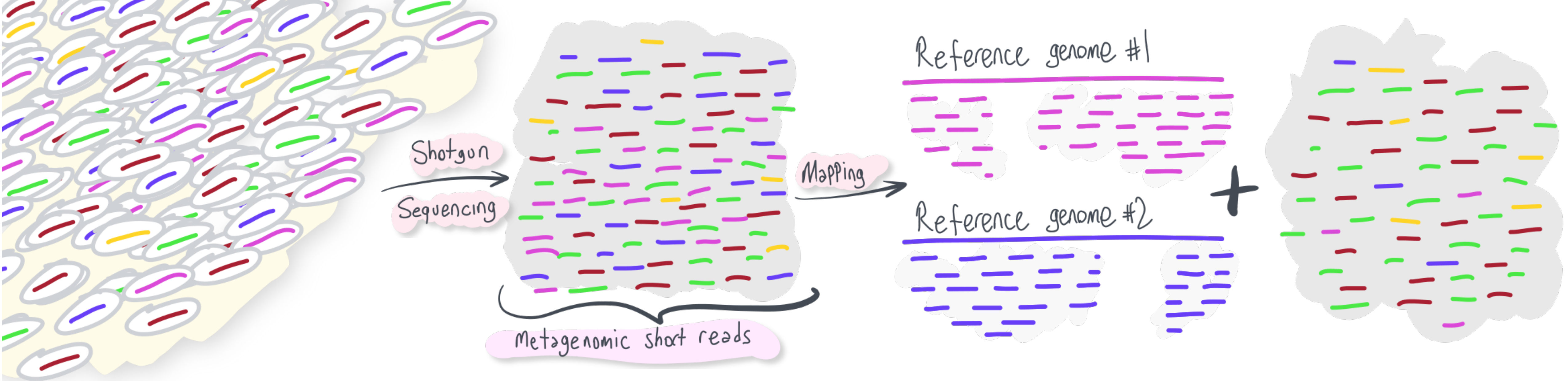


READ  
RECRUITMENT →



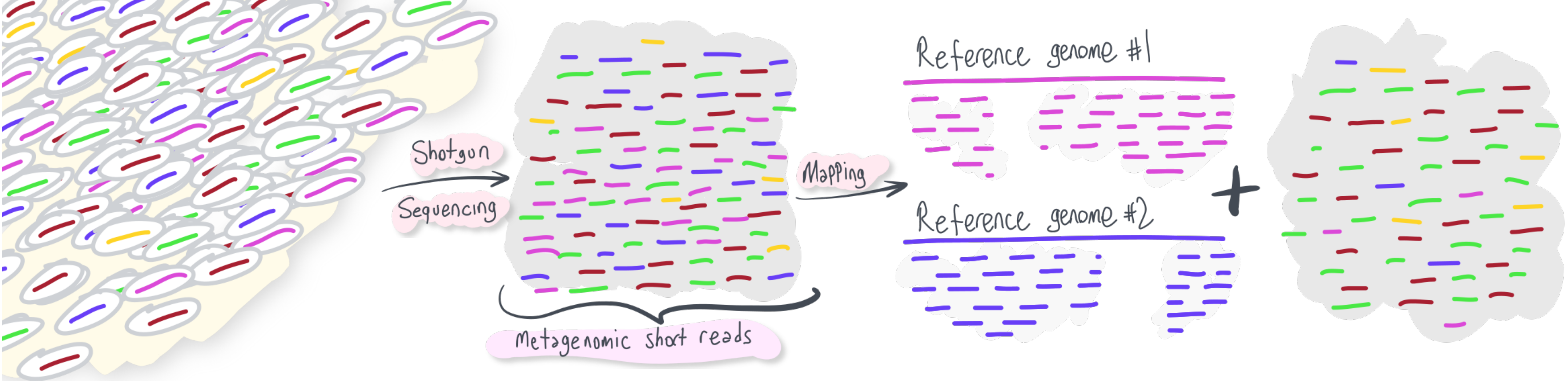
# Mapping rates

How much can we put in the reference context?

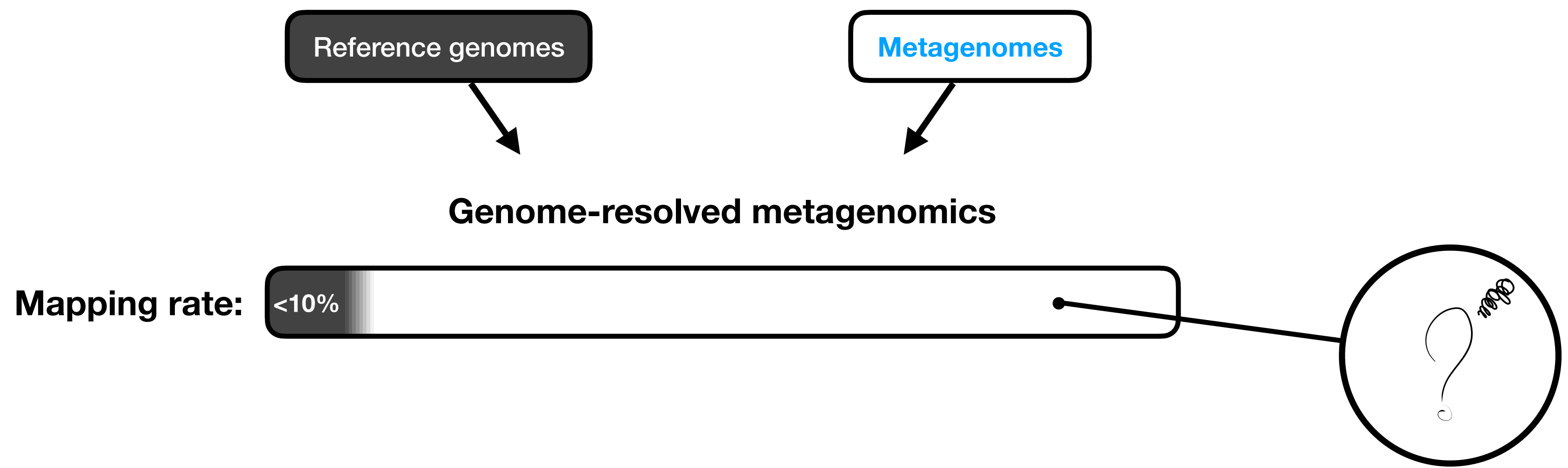


# GENOME RESOLVED METAGENOMICS



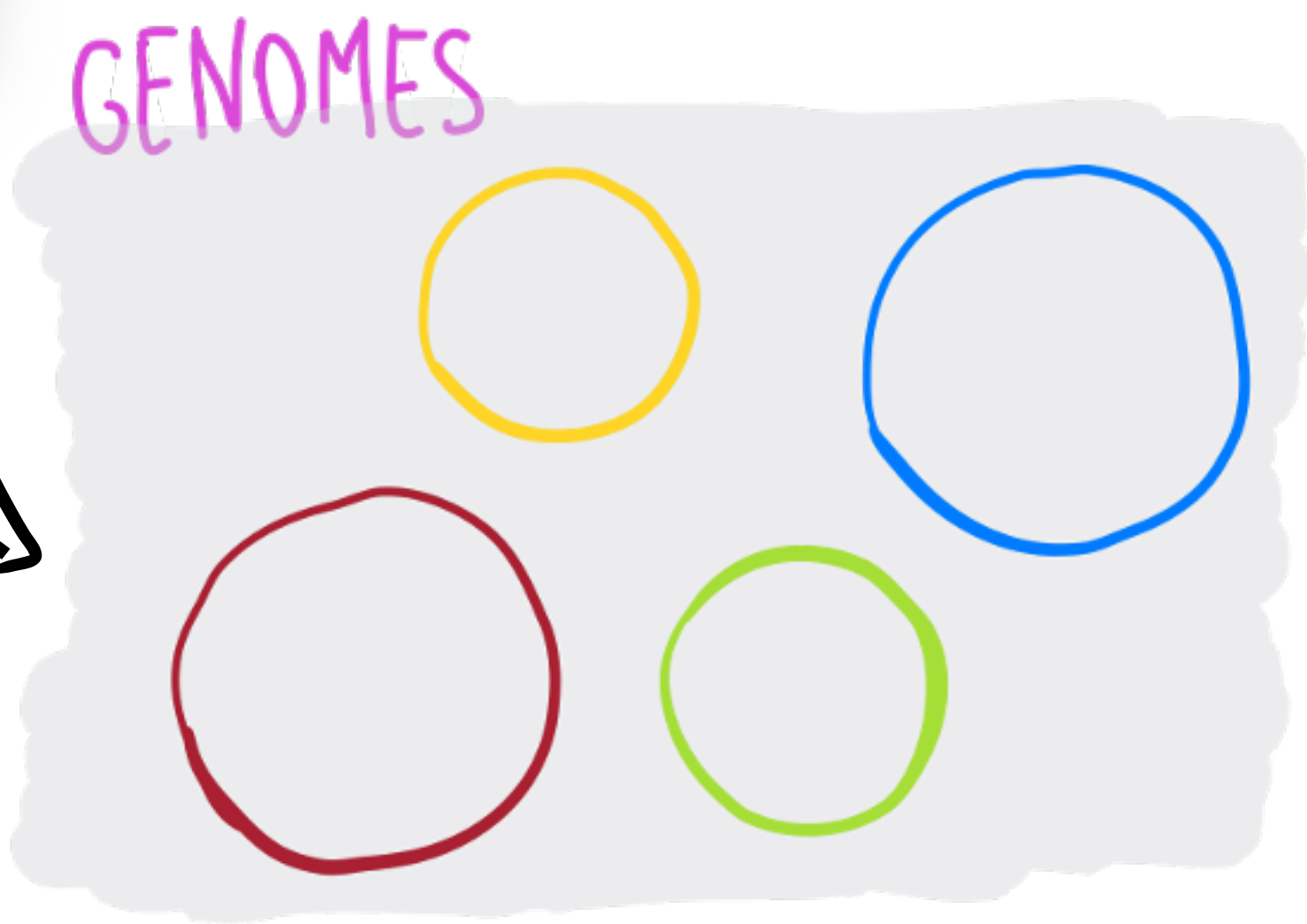
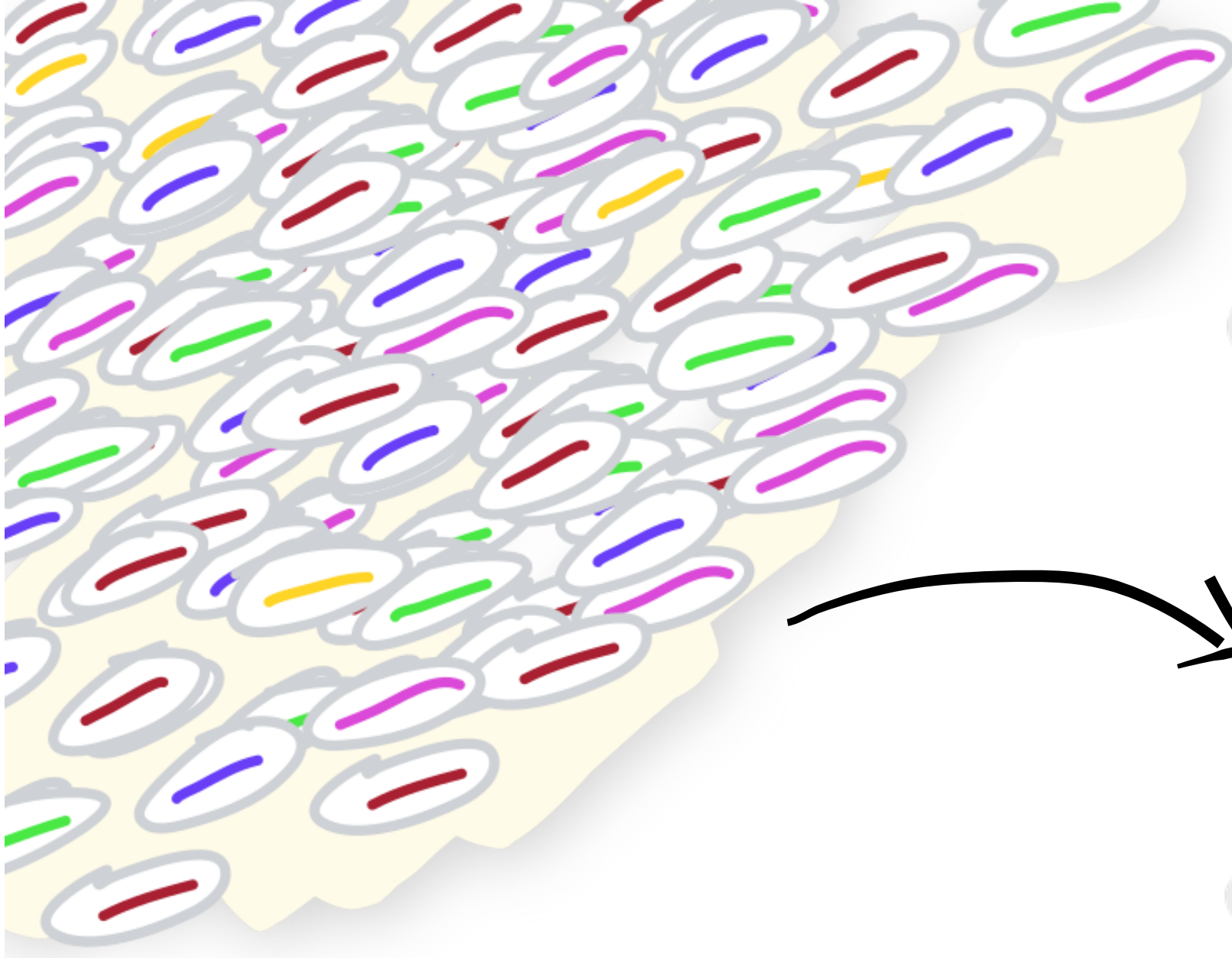


# GENOME RESOLVED METAGENOMICS

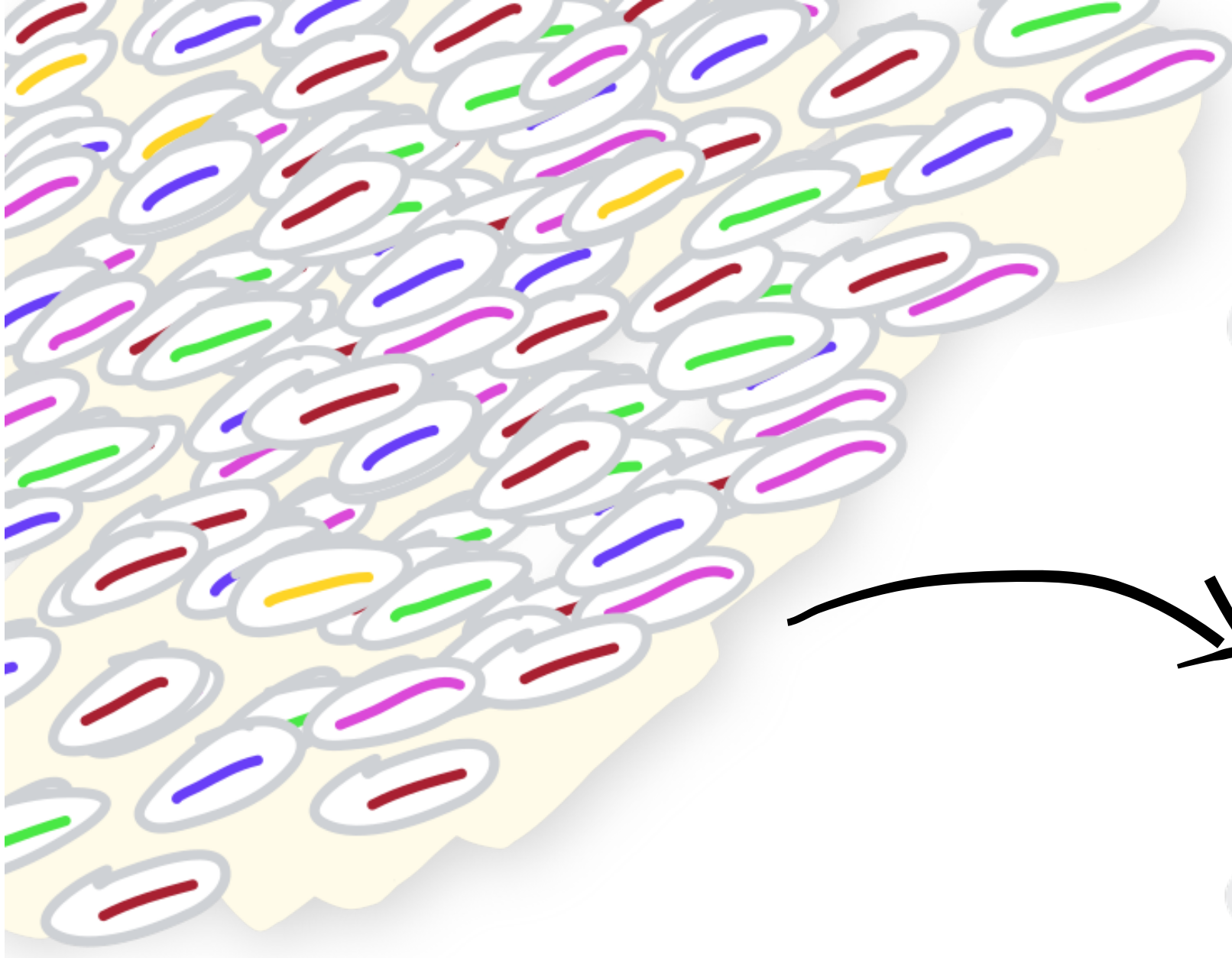


# Accessing the missing genomic content

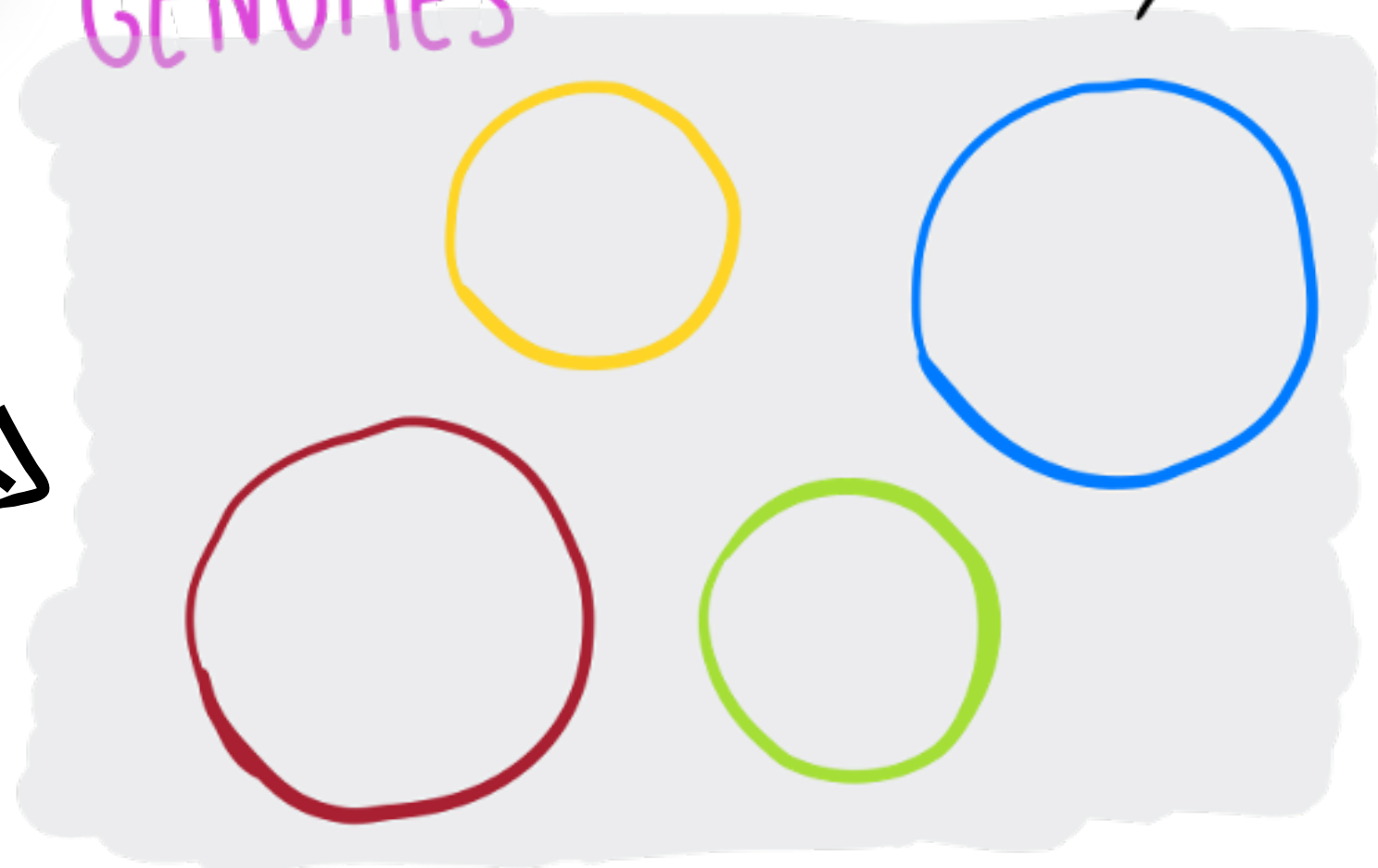
## Cultivation-independent methods







GENOMES

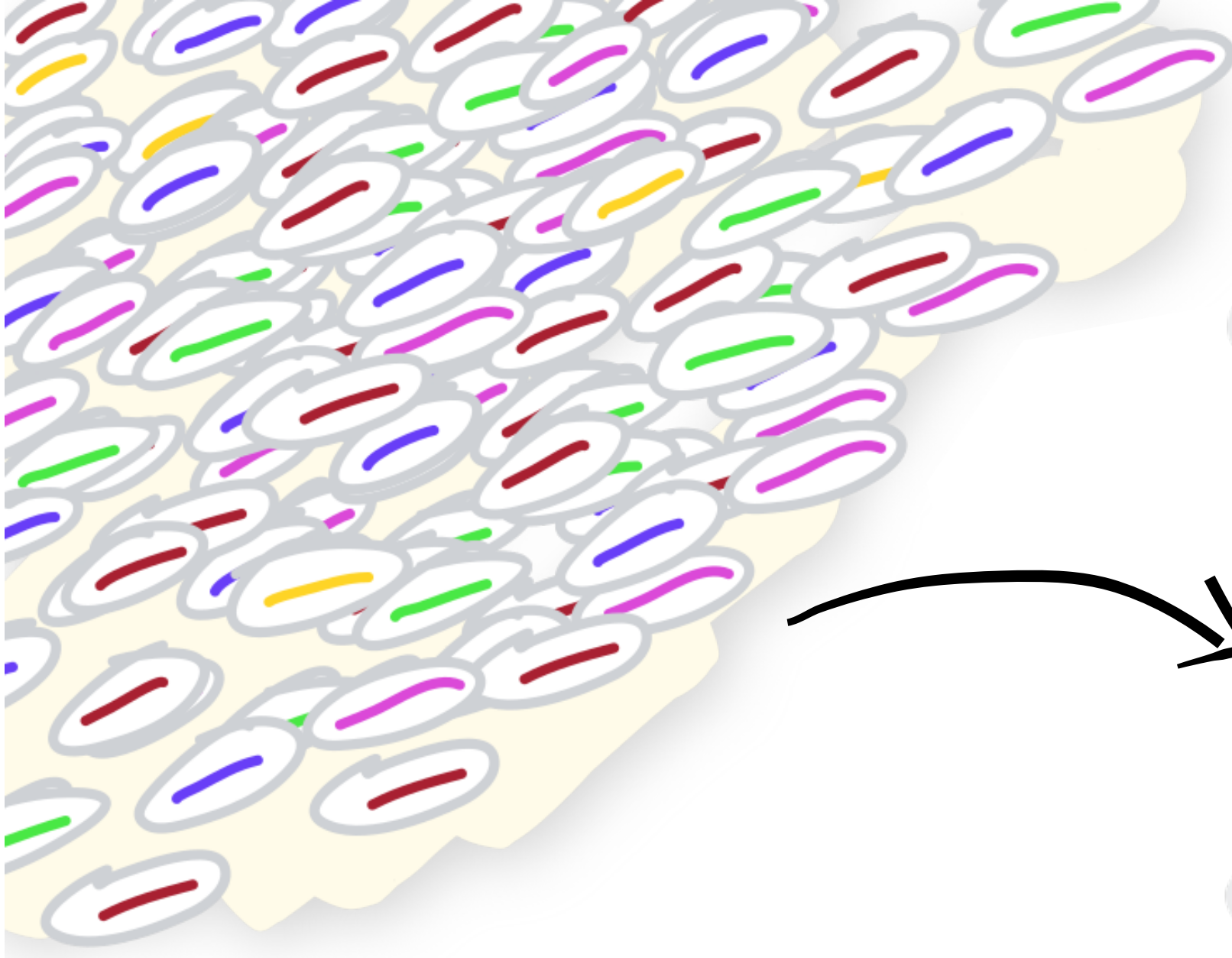


SHOTGUN SEQUENCING

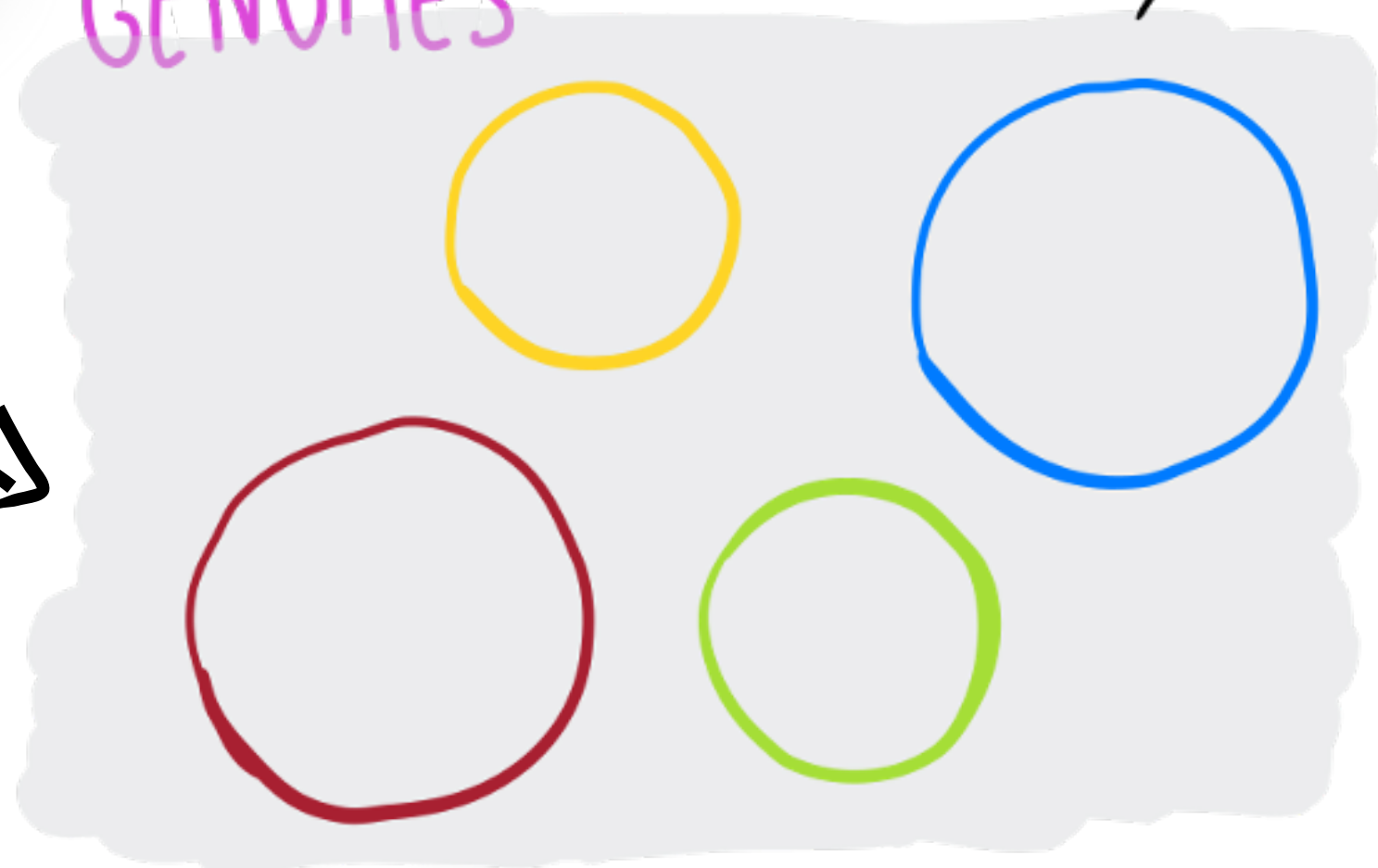


SHORT READS





GENOMES

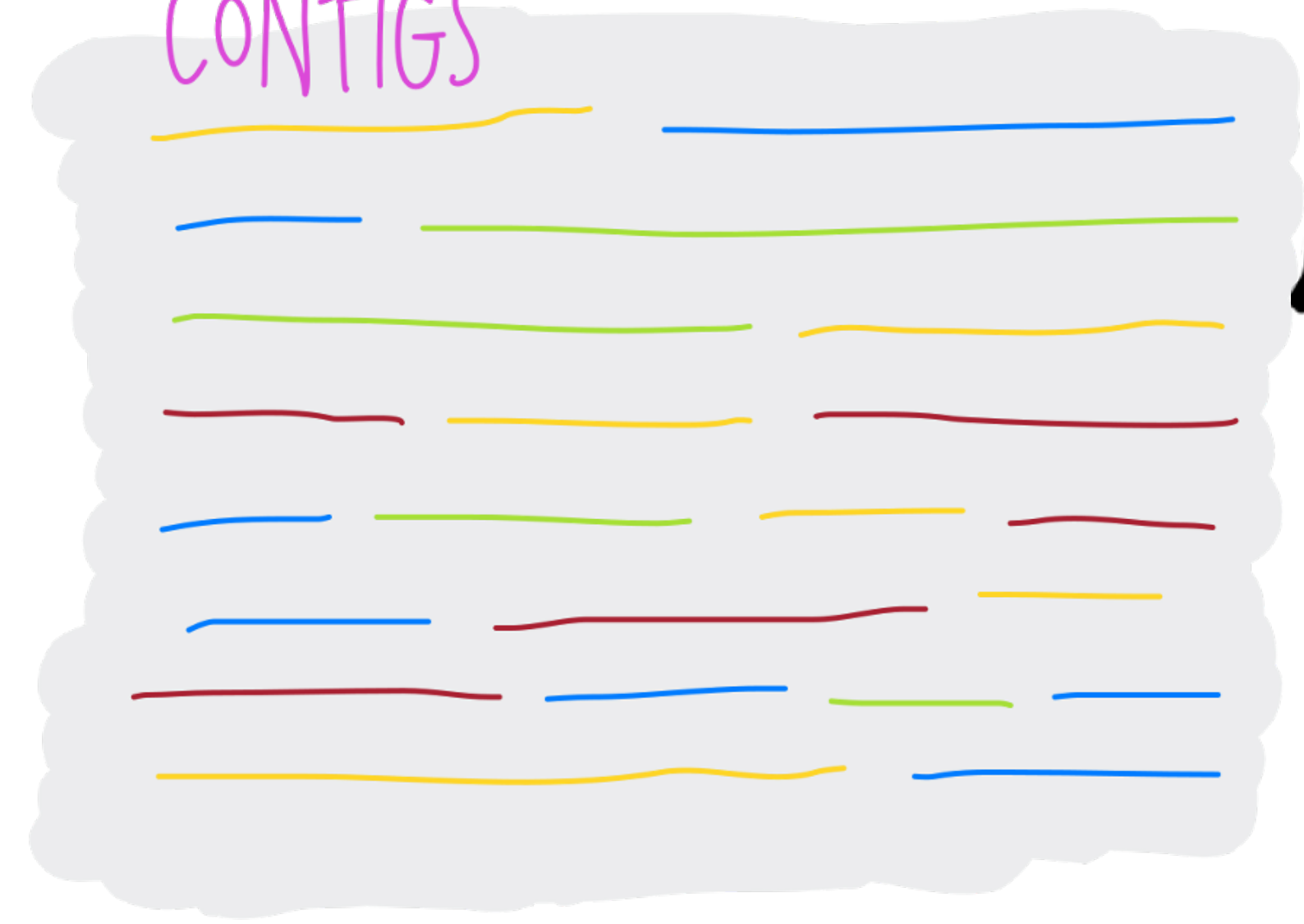


SHOTGUN SEQUENCING

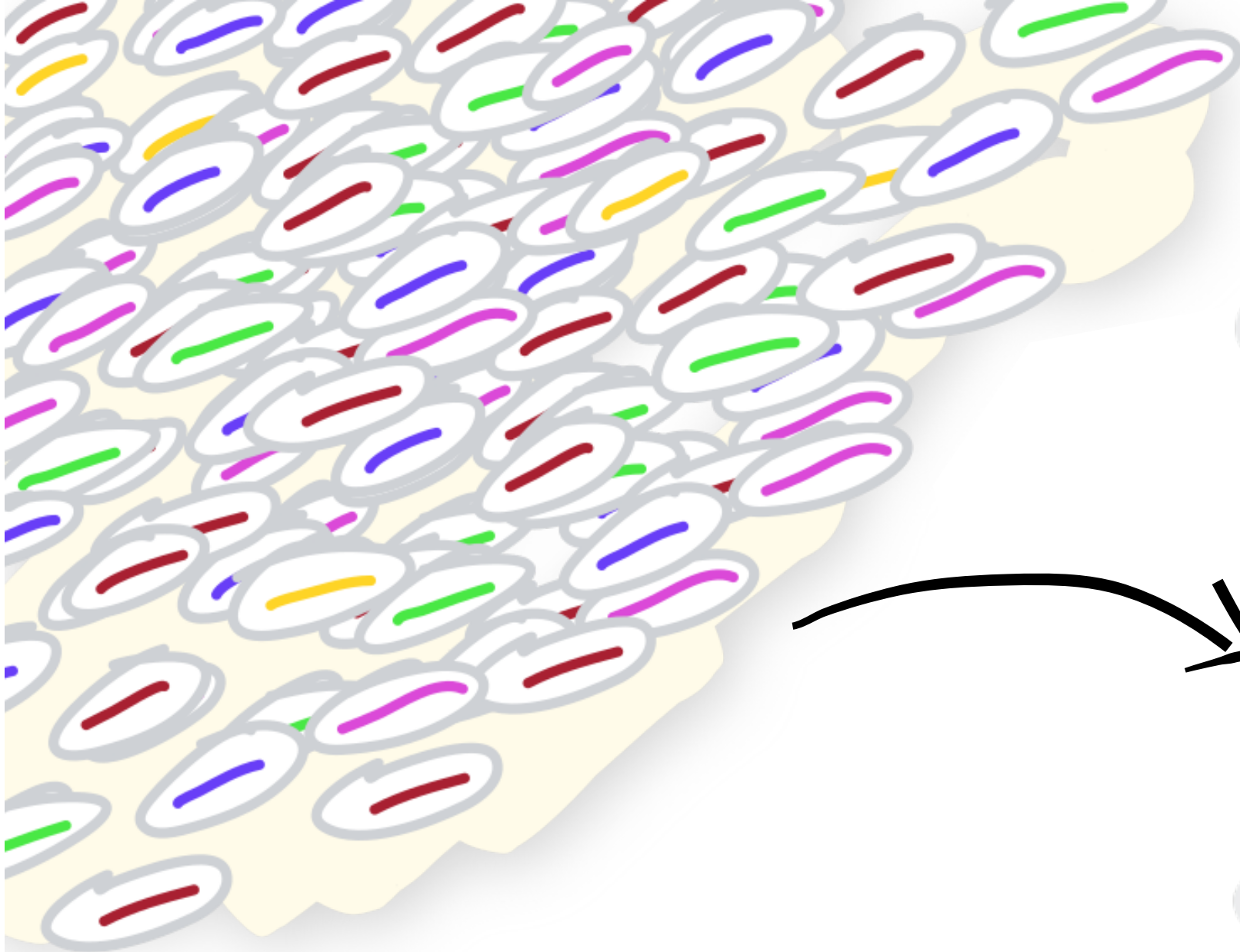
SHORT READS



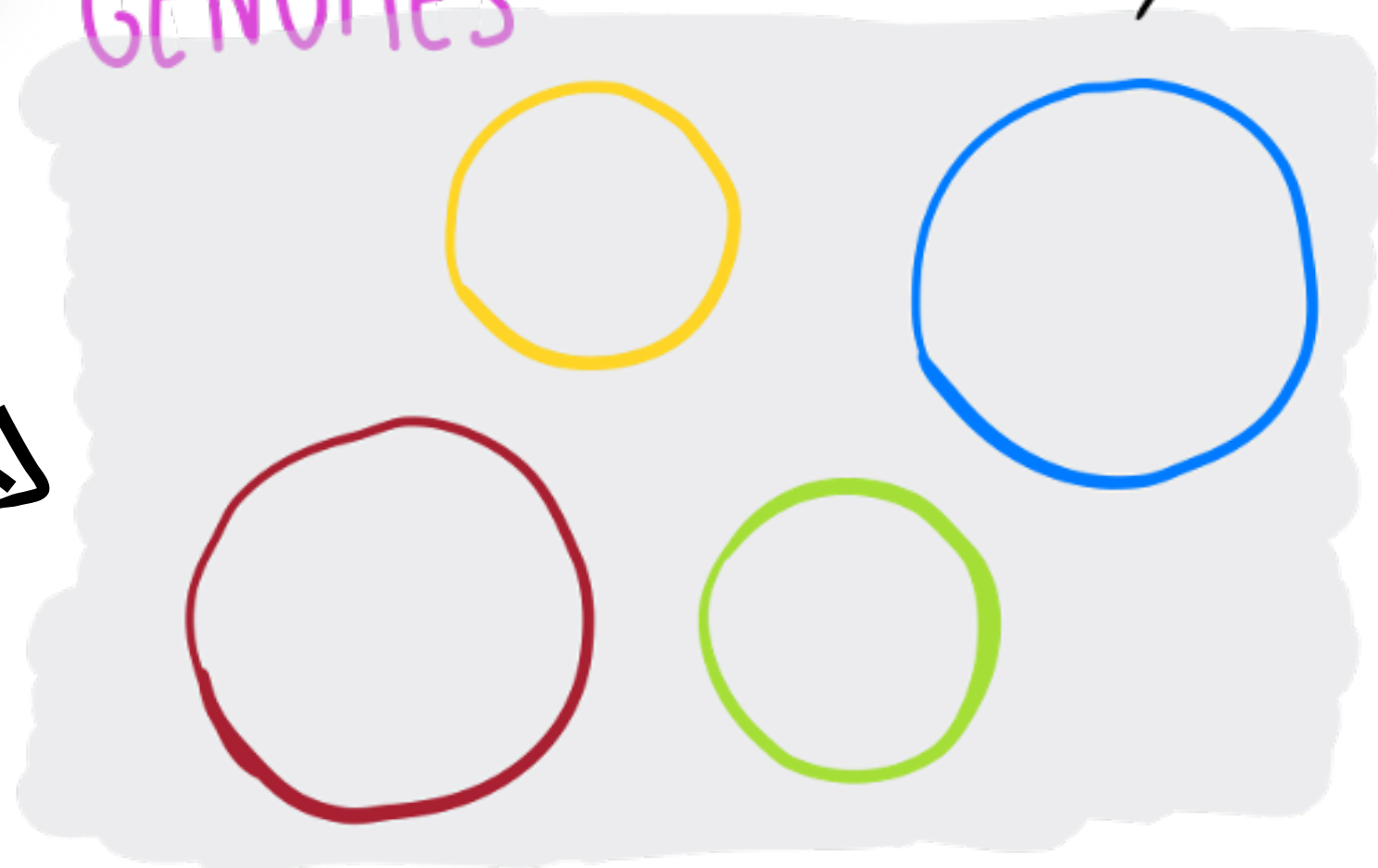
CONTIGS



DE-NOVO ASSEMBLY

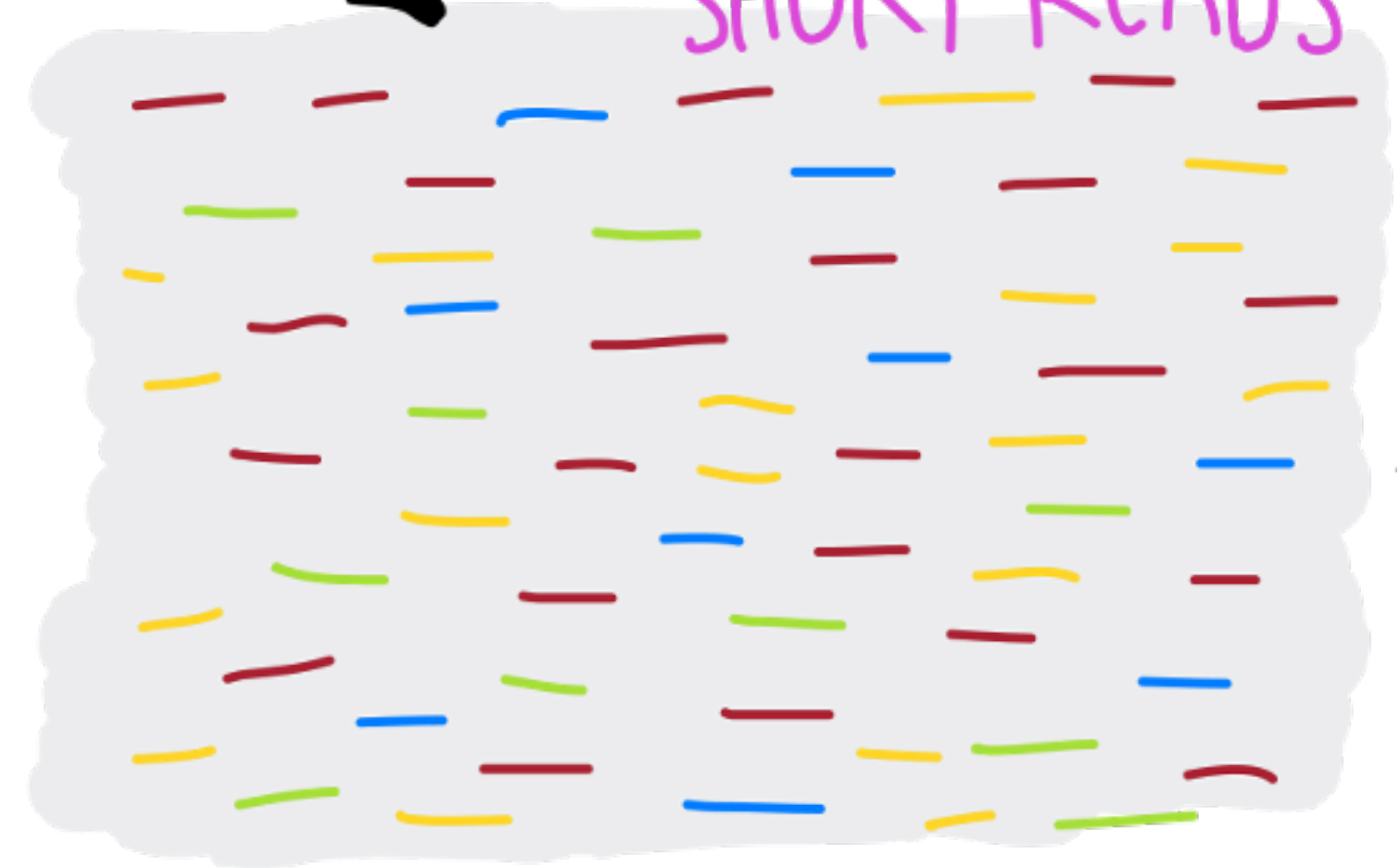


GENOMES



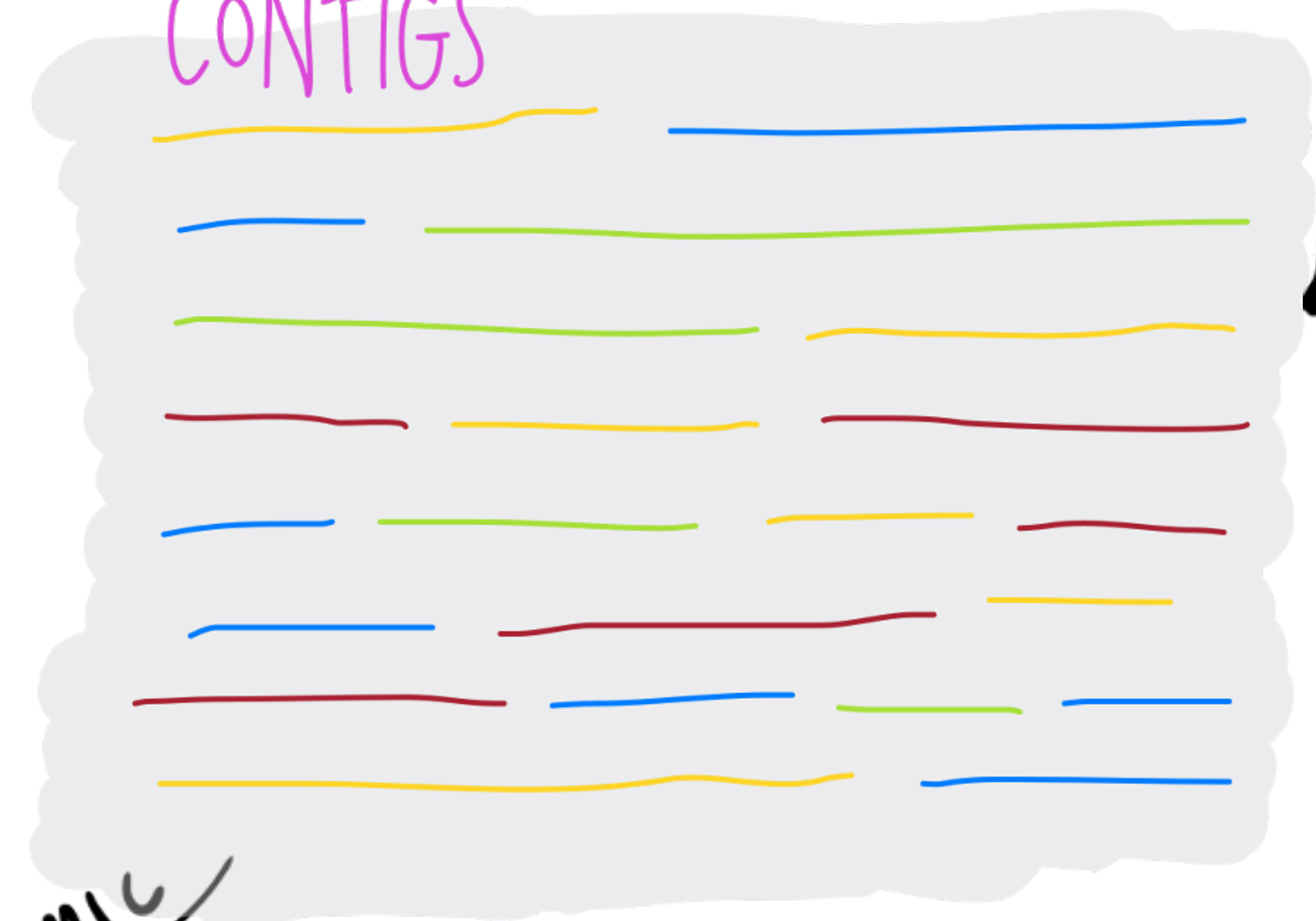
SHOTGUN SEQUENCING

SHORT READS

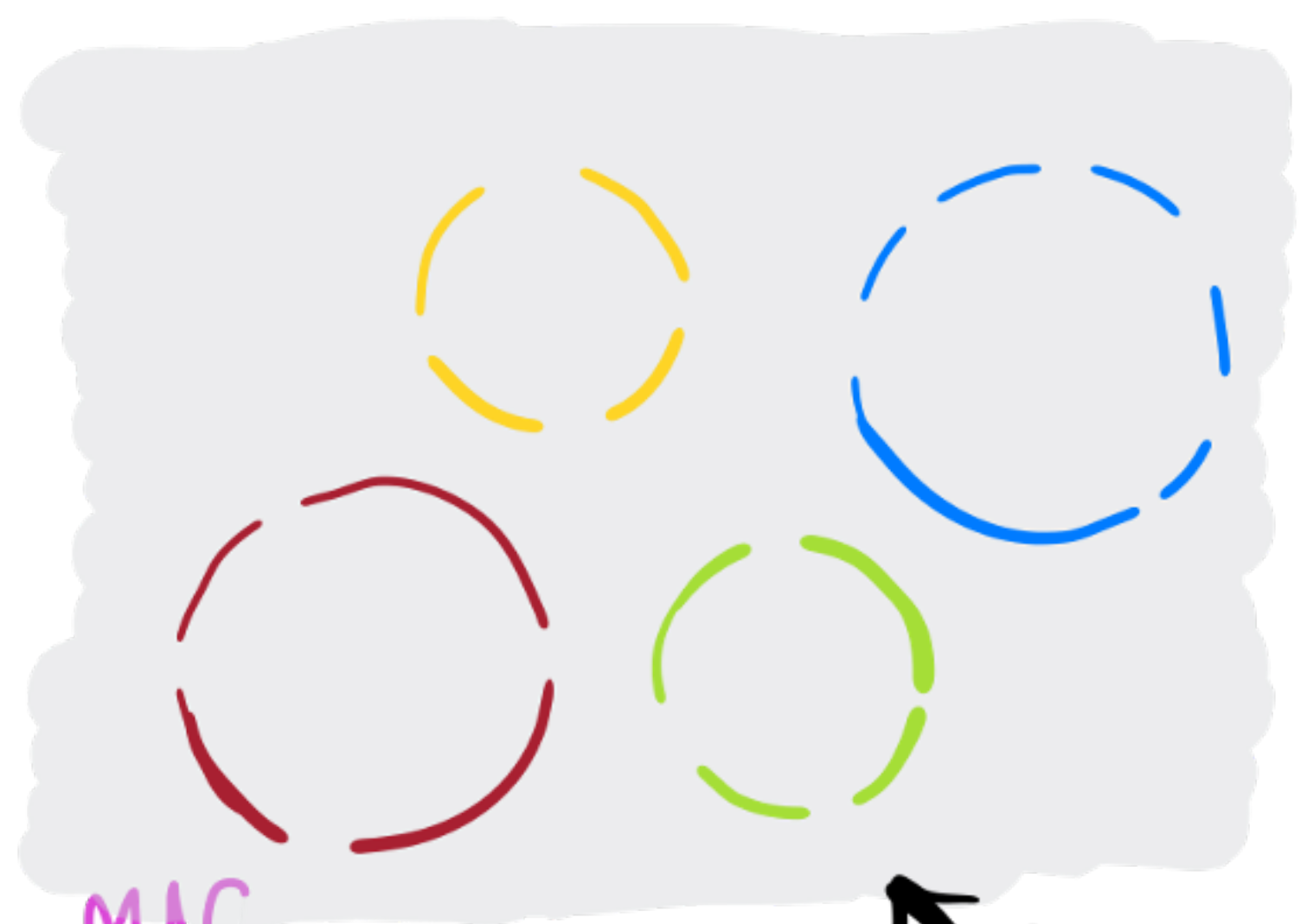


DE-NOVO ASSEMBLY

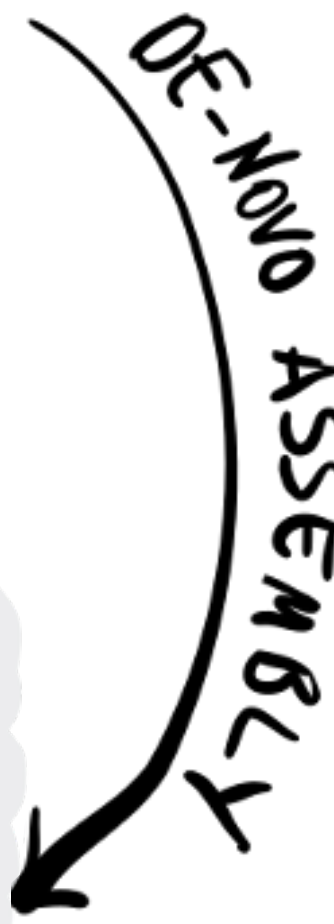
CONTIGS

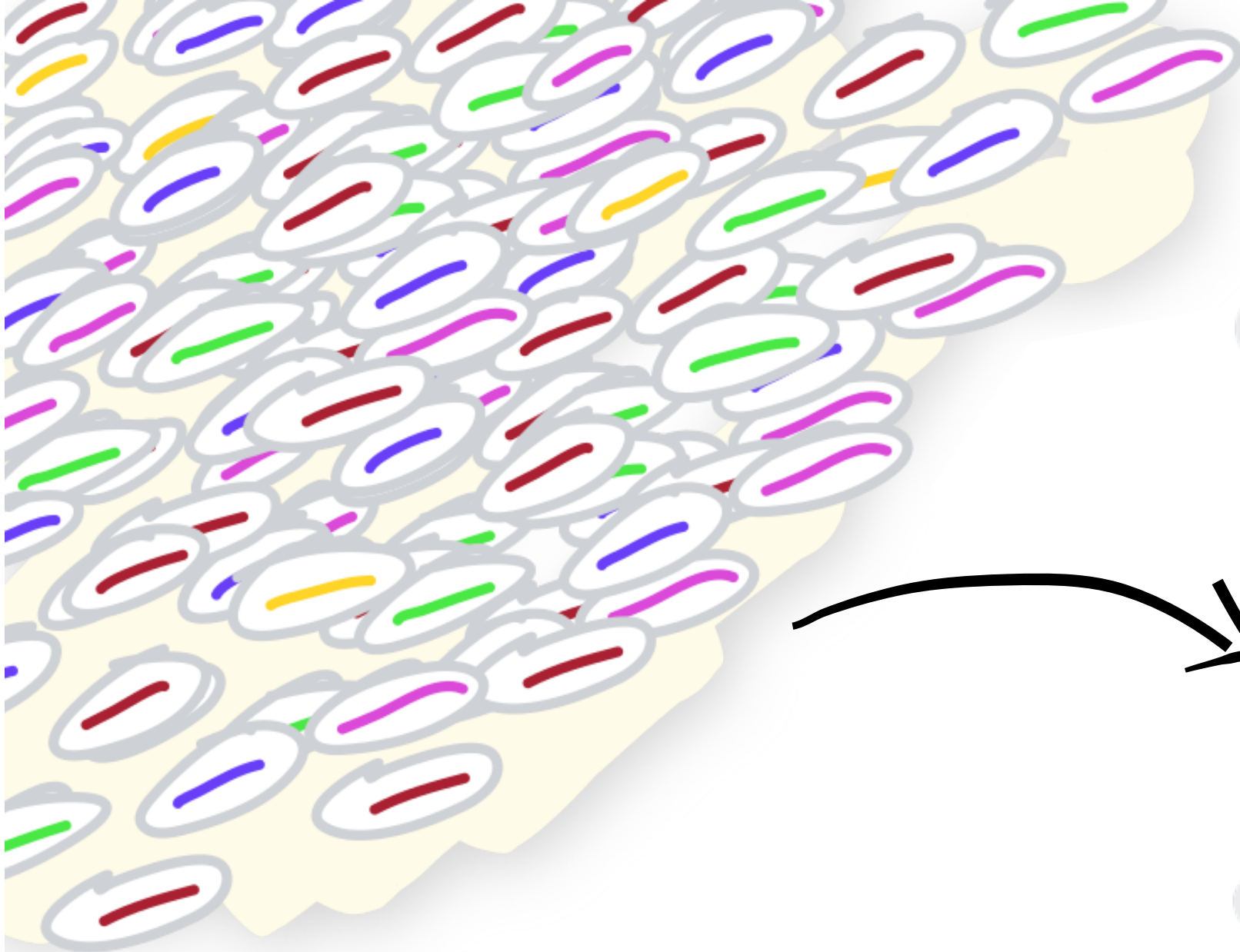


MAGs

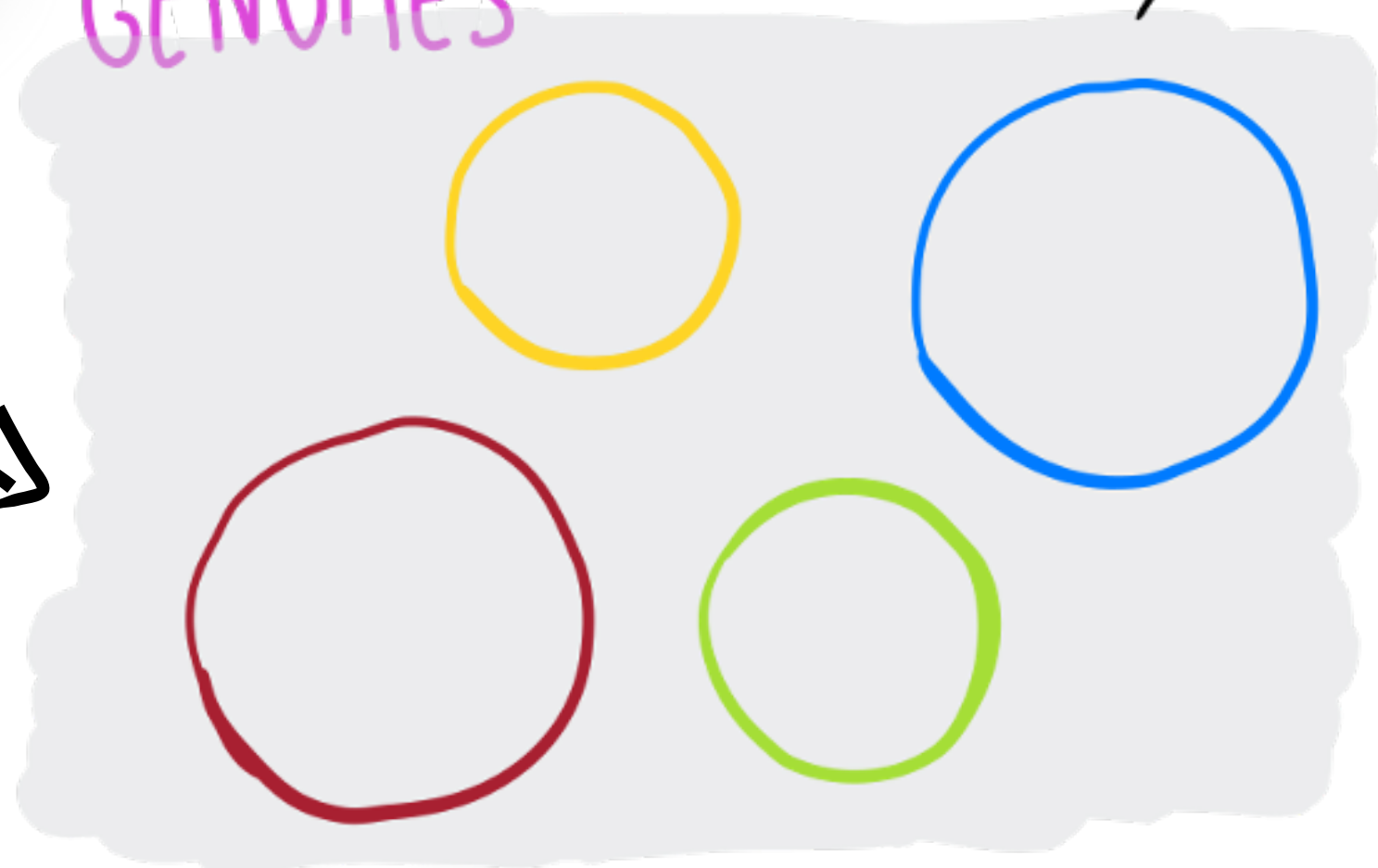


METAGENOMIC BINNING



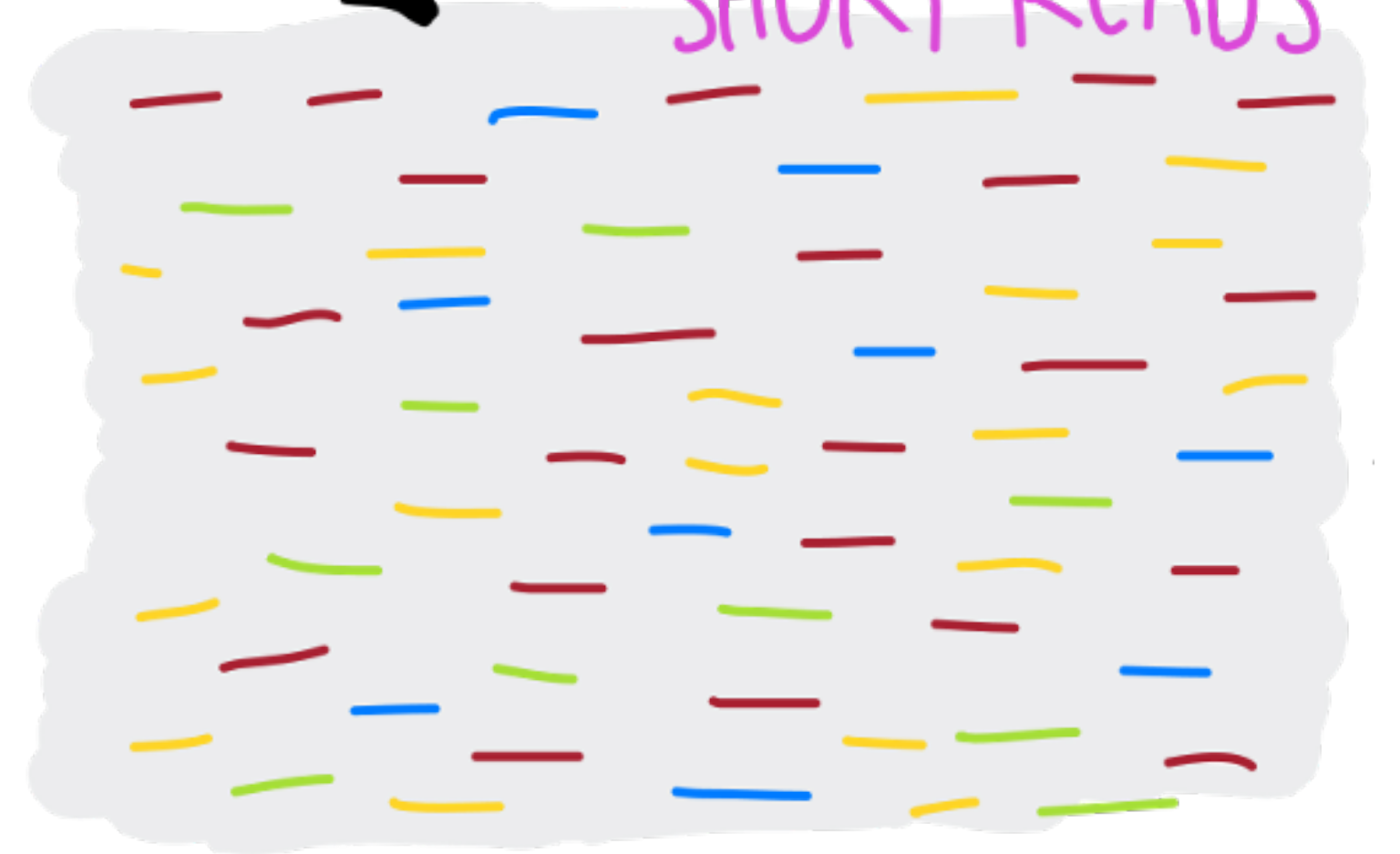


GENOMES



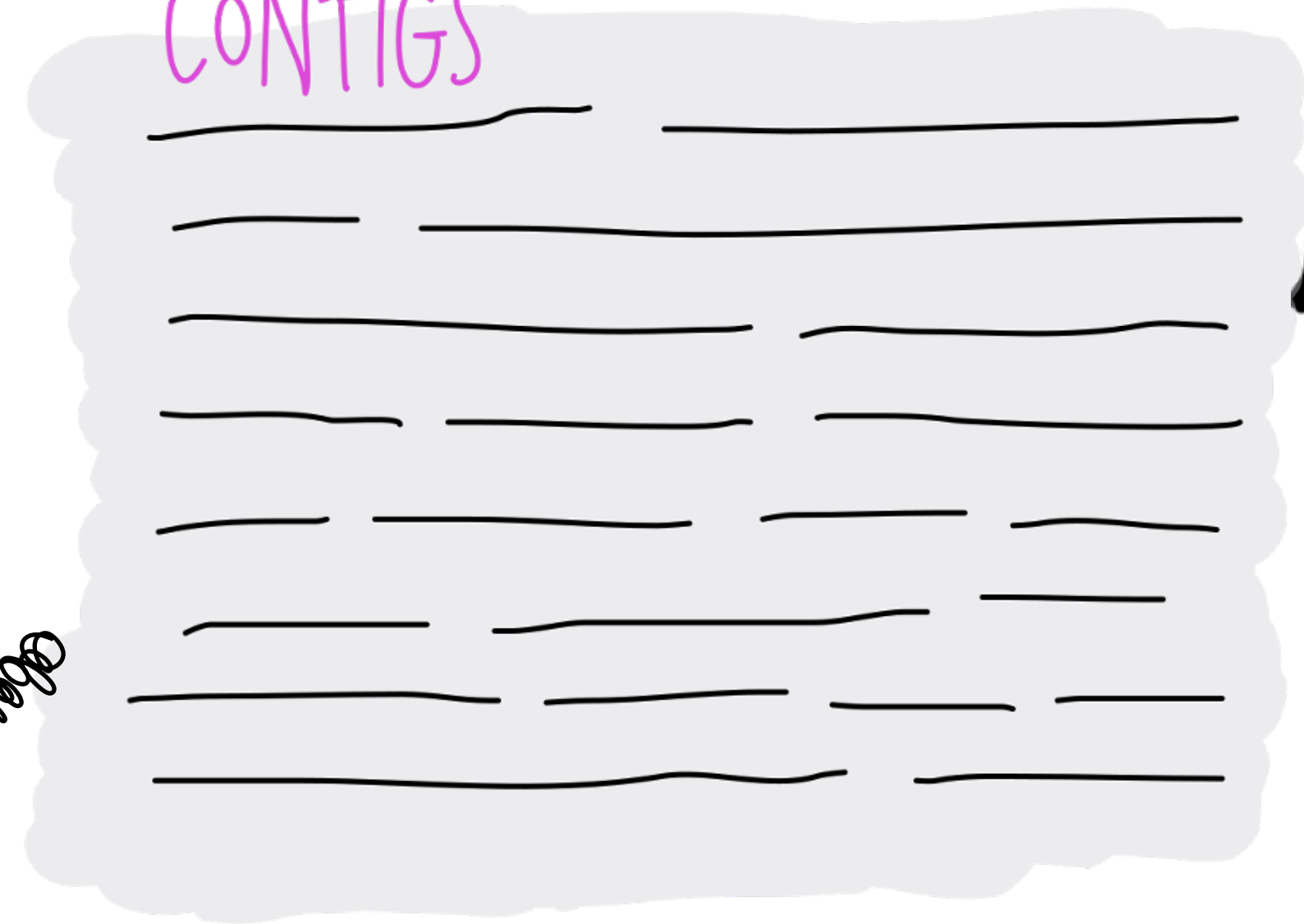
SHOTGUN SEQUENCING

SHORT READS

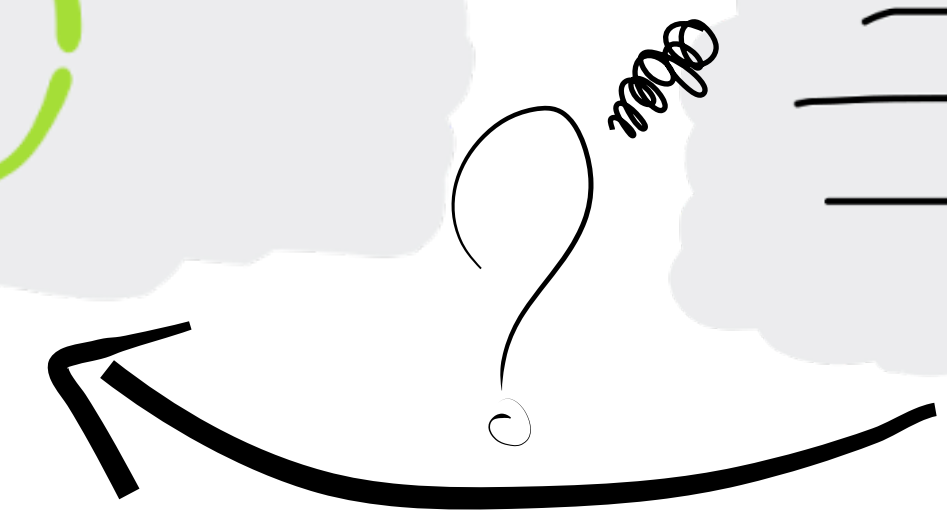
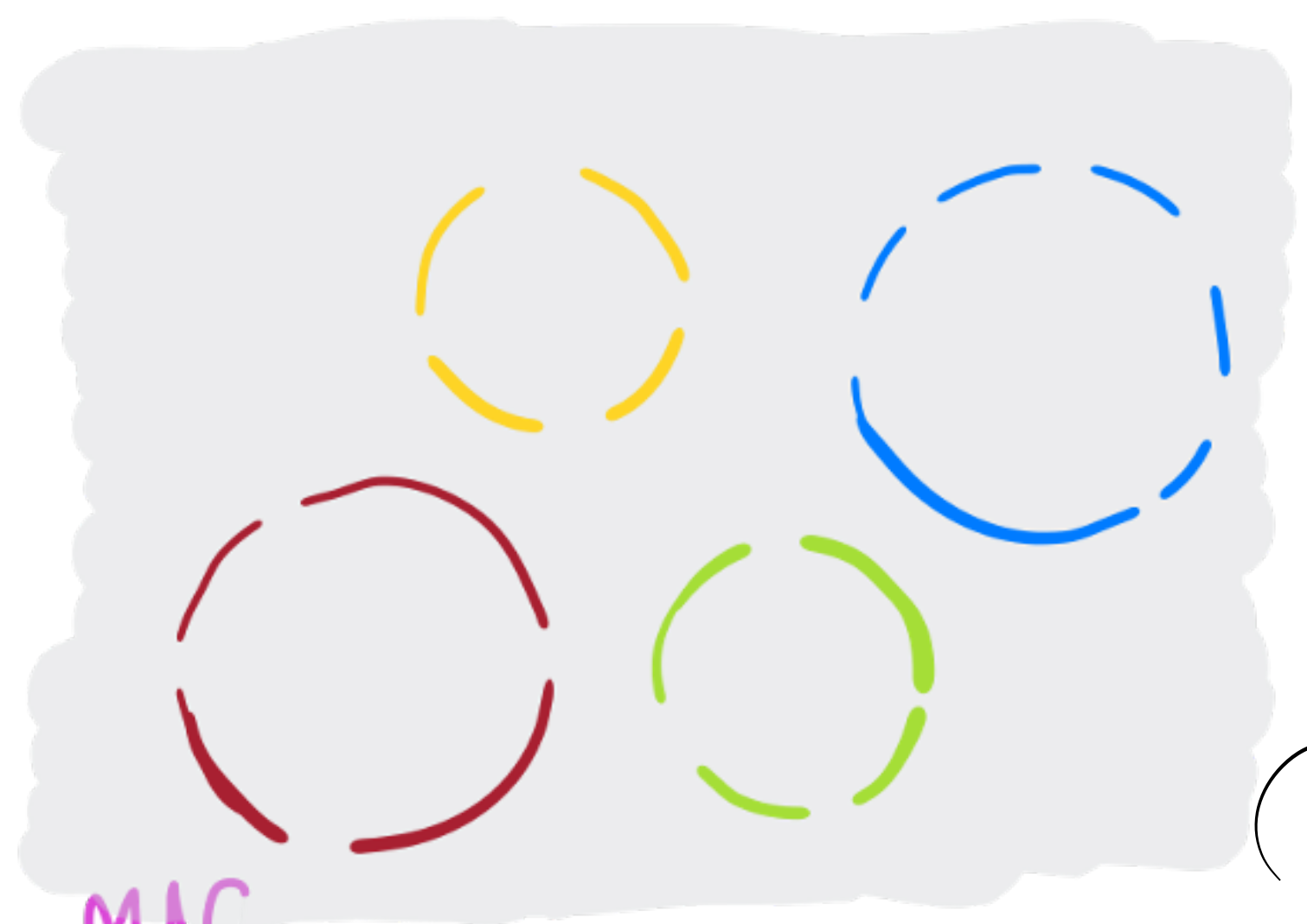


DE-NOVO ASSEMBLY

CONTIGS



MAGs



Sequence composition  
Computing k-mer frequencies

GTTTGGCATGATTAAGGAGTTTCTTTGTGCTTC

k=2

GTTTGGCATGATTAAGGAGTTTCTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT

k=2

GT TTTGGCATGATTAAGGAGTTTCTTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0

k=2



CTT TGGCATGATTAAGGAGTTTCTTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	1

k=2

GTTTTGGCATGATTAAGGAGTTTCTTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

k=2

GTTTTGGCATGATTAAGGAGTTTTCTTTTGTGCTTC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
1	0	2	2	1	0	0	2	2	2	2	3	1	2	4	10

GAAGCACAAAAGAAACTCCTTTAATCATGCCAAAAC

AA	AC	AG	AT	CA	CC	CG	CT	GA	GC	GG	GT	TA	TC	TG	TT
10	3	2	2	4	2	0	2	2	2	0	0	1	2	1	1

GTTTTGGCATGATTAAGGAGTTTTCTTTTGTGCTTC  
GAAGCACAAAAGAAACTCCTTTAATCATGCCAAAAC

AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
11	3	4	4	5	2	0	2	2	1

→ PALINDROMES :)

k=2

GTTTGGCATGATTAAGGAGTTTCTTTGTGCTTC

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y										
Z										
L										
K										
M										

k=2

ACTTCCGCAGTCGGGCATTACGCGTTGTGGAATGA

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z										
L										
K										
M										

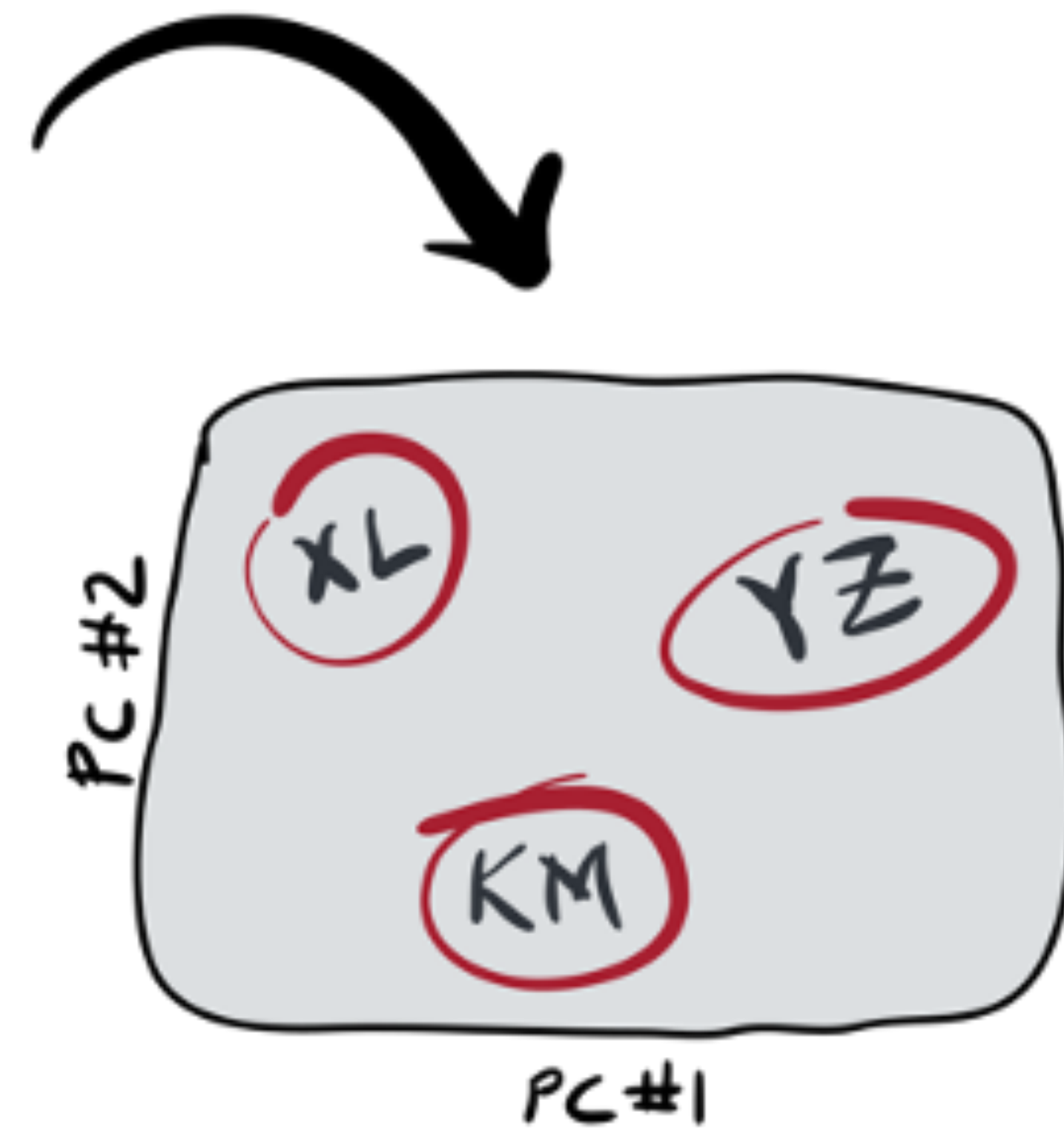
k=2

GGGCCTGCGCCGGTCCAGTCACCCGGCTGCGACCT

	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0

k=2

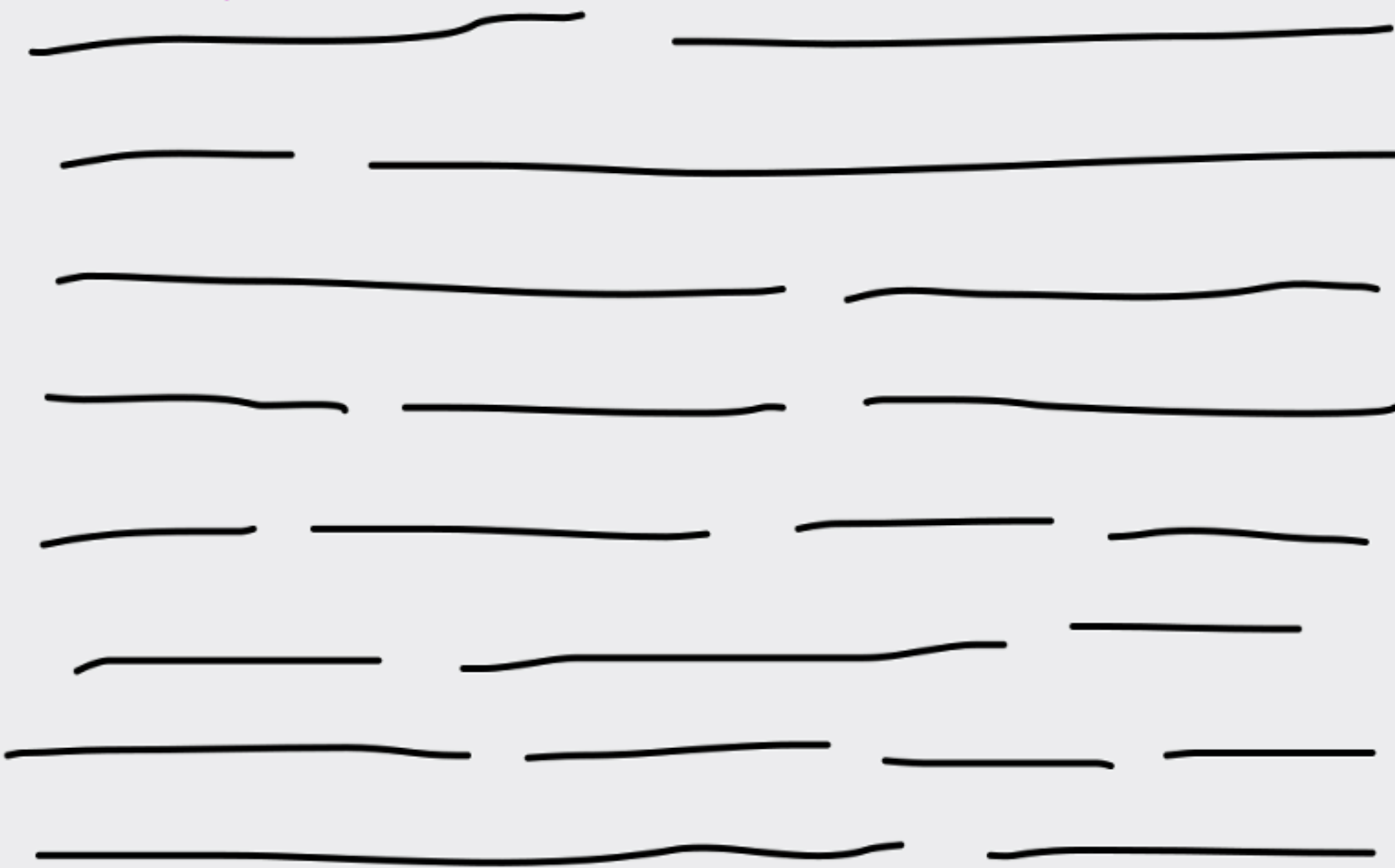
	AA	AC	AG	GA	CA	CC	CG	GC	AT	TA
X	11	3	4	4	5	2	0	2	2	1
Y	4	5	2	4	5	4	4	3	2	1
Z	4	5	3	2	4	1	5	5	2	3
L	11	6	3	2	2	3	2	1	1	4
K	1	1	2	2	1	8	9	10	0	0
M	0	4	4	3	4	10	4	5	0	0



k=2

SEQUENCE COMPOSITION

CONTIGS



MAGs



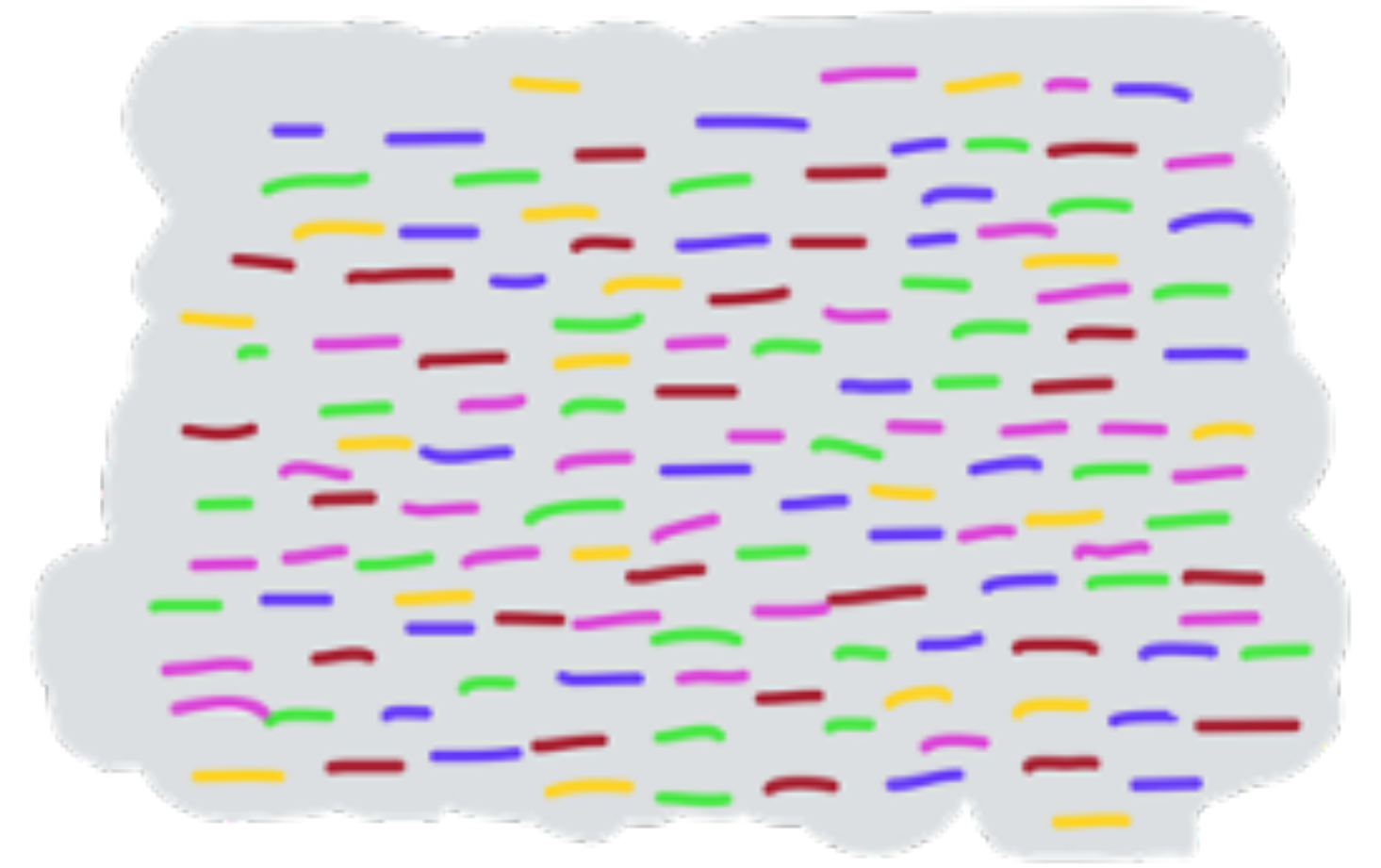


Abundance correlation  
Counting stuff across samples

CONTIG #1

CONTIG #2

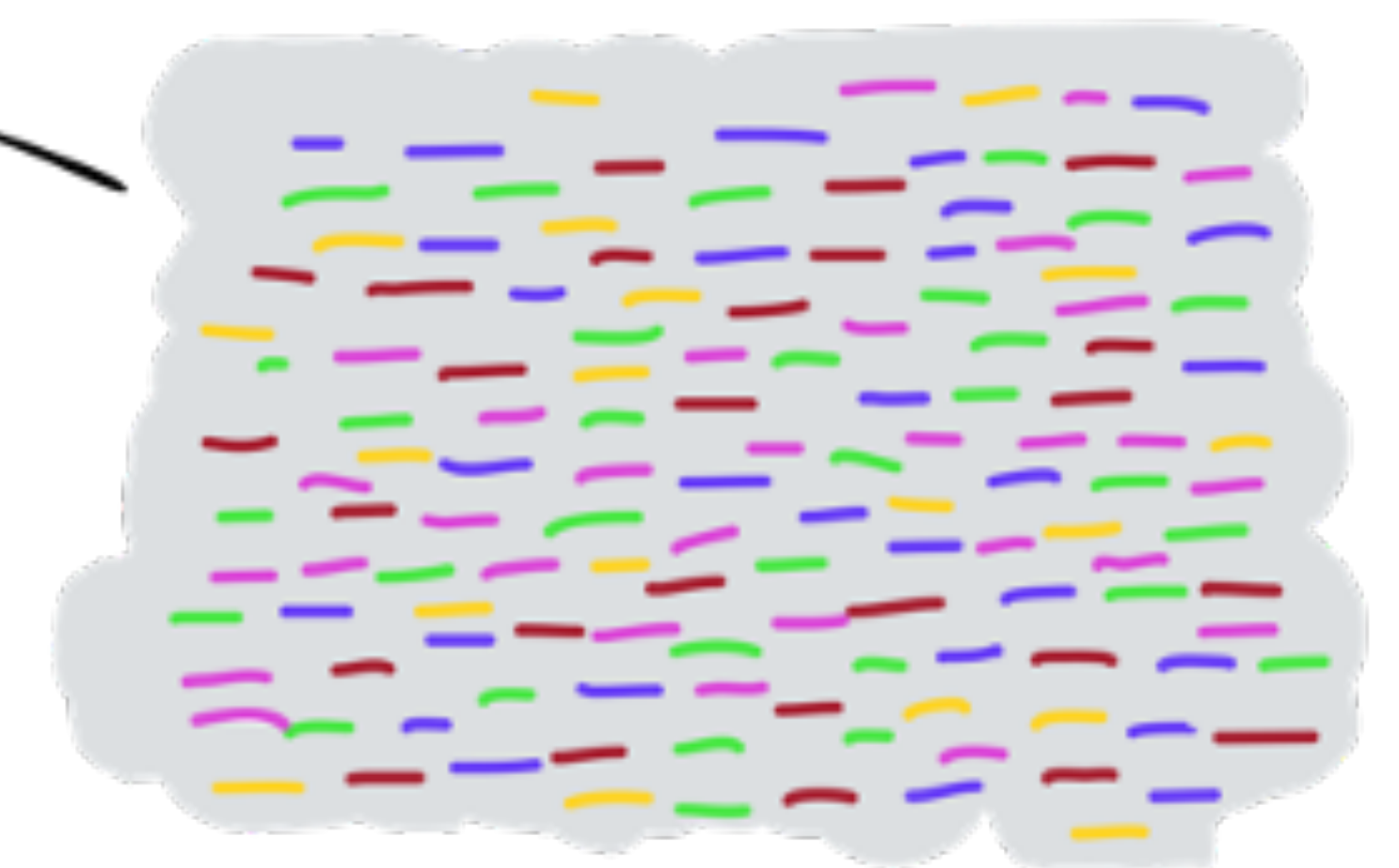
CONTIG #1



METAGENOMIC READS

CONTIG #2

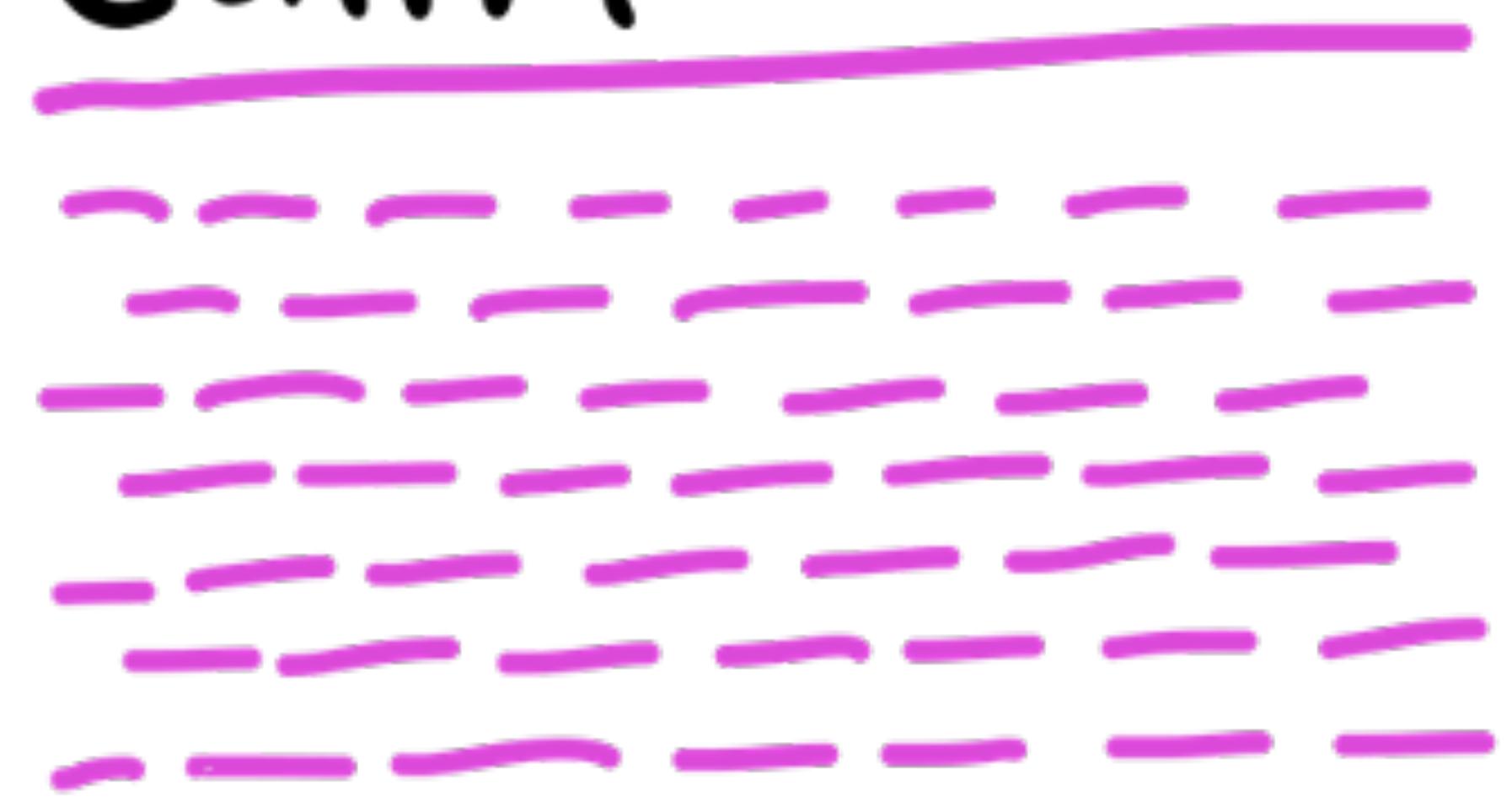
CONTIG #1



METAGENOMIC READS

CONTIG #2

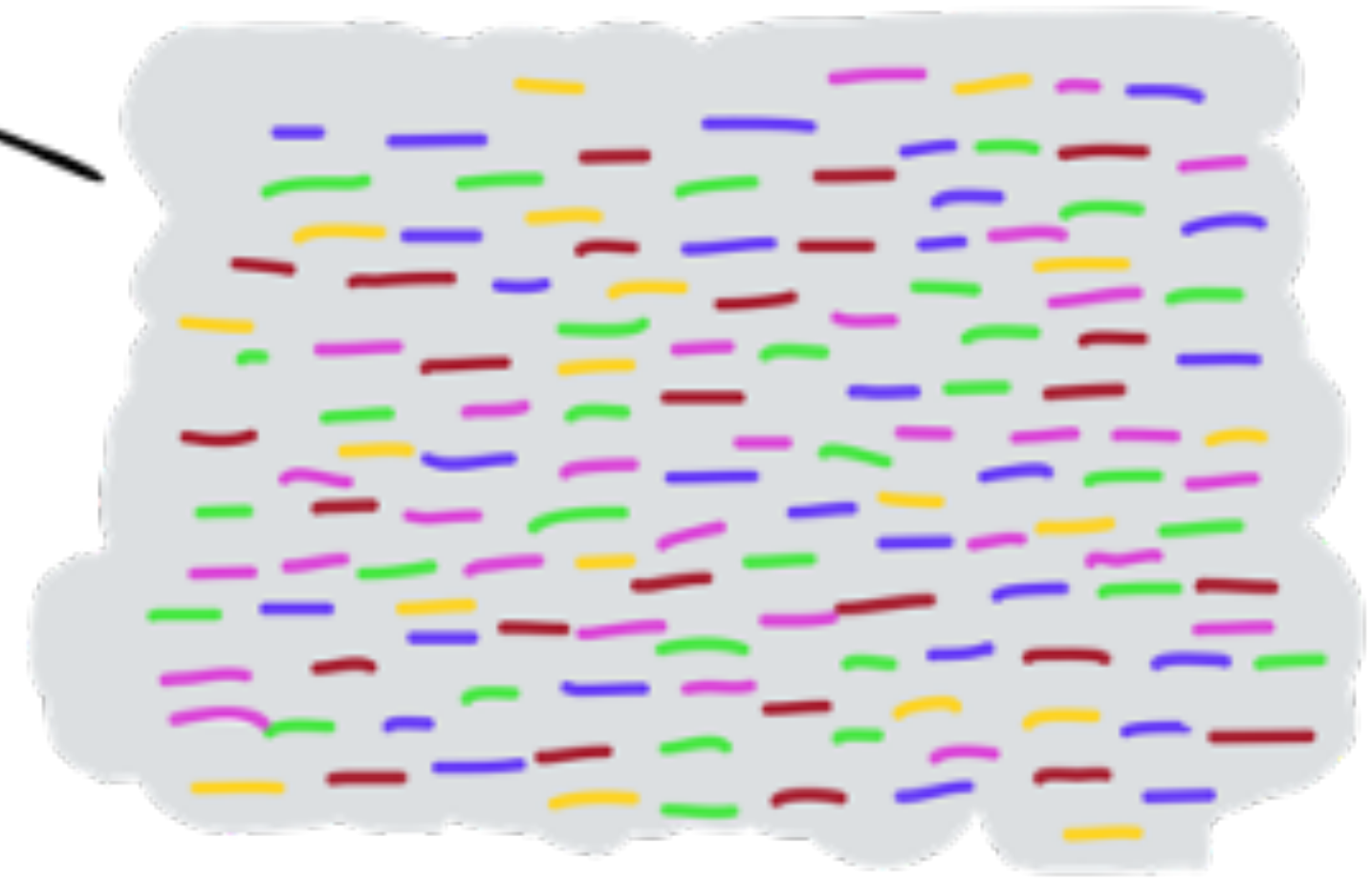
CONTIG #1



CONTIG #2

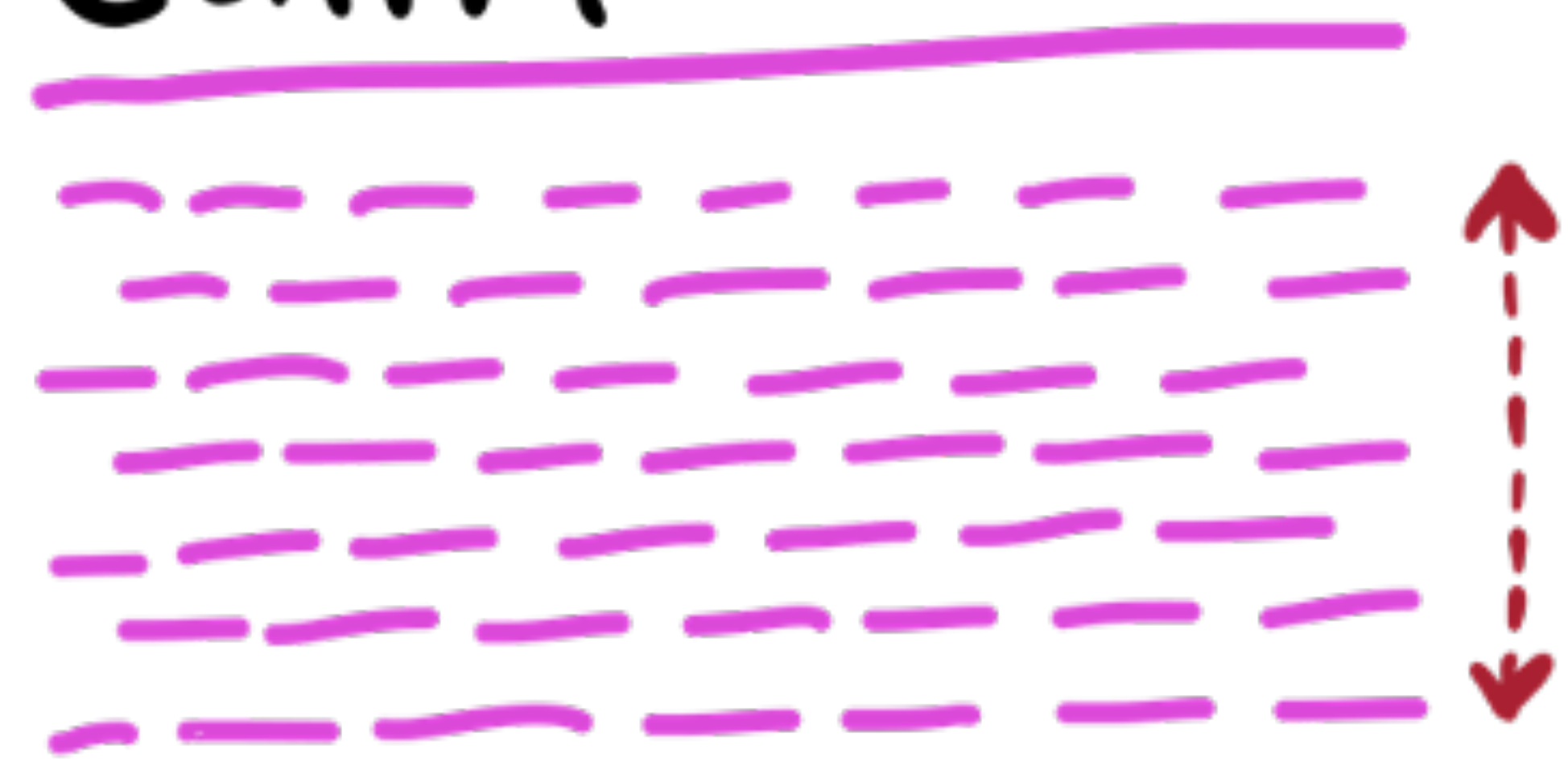


MAPPING

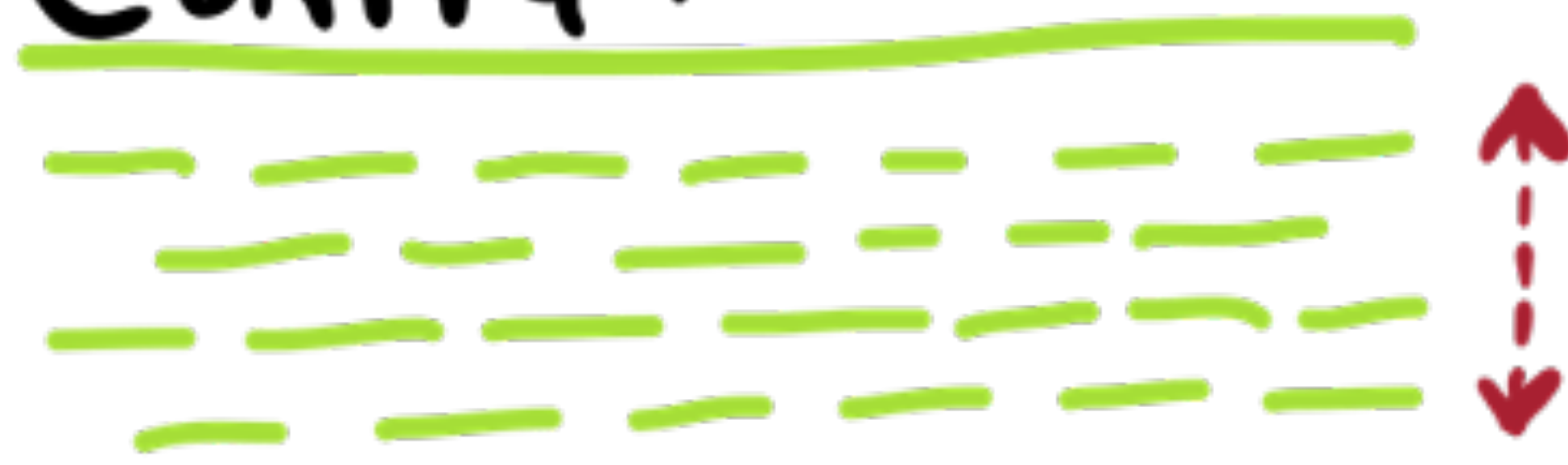


METAGENOMIC READS

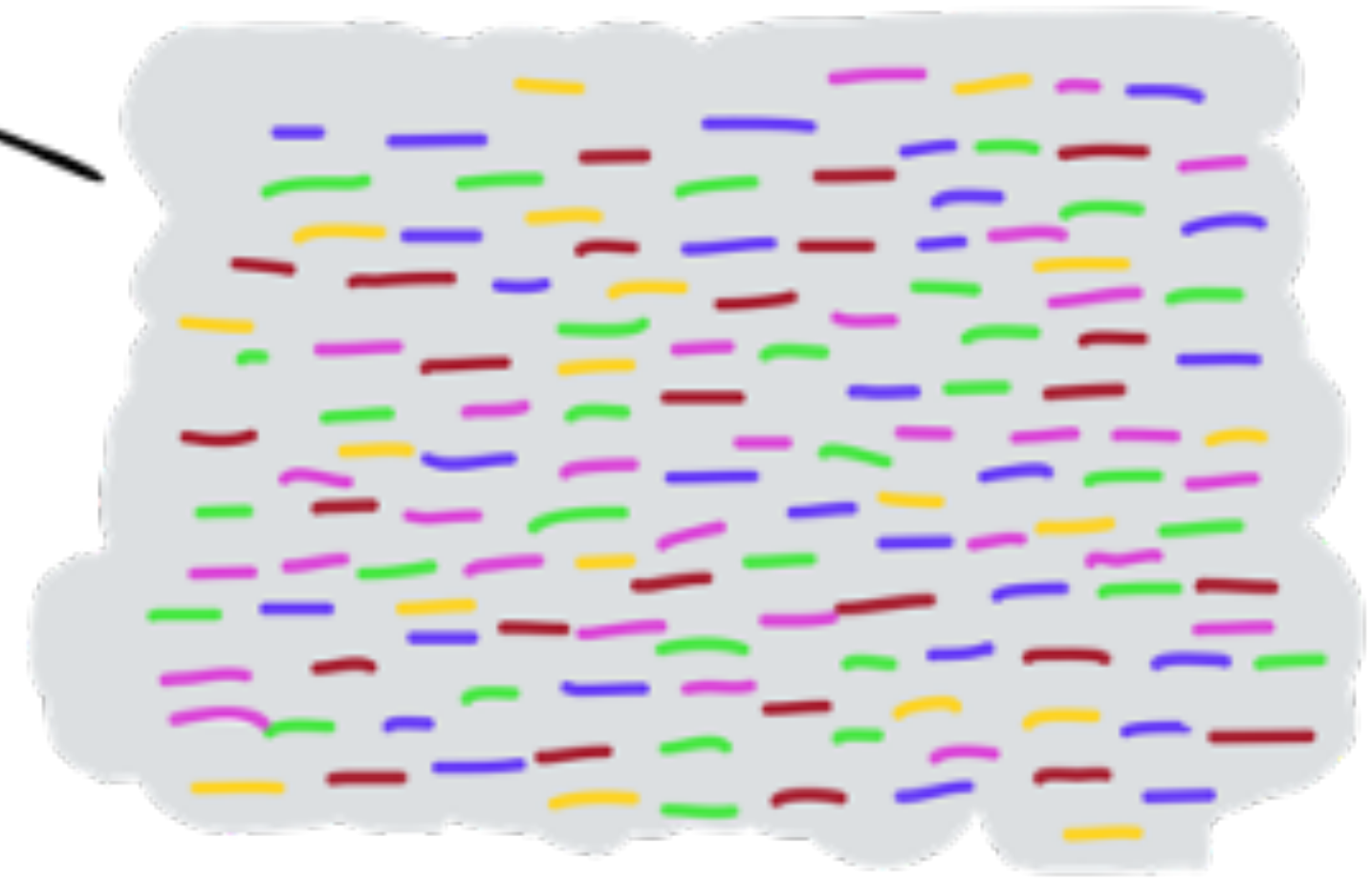
CONTIG #1



CONTIG #2

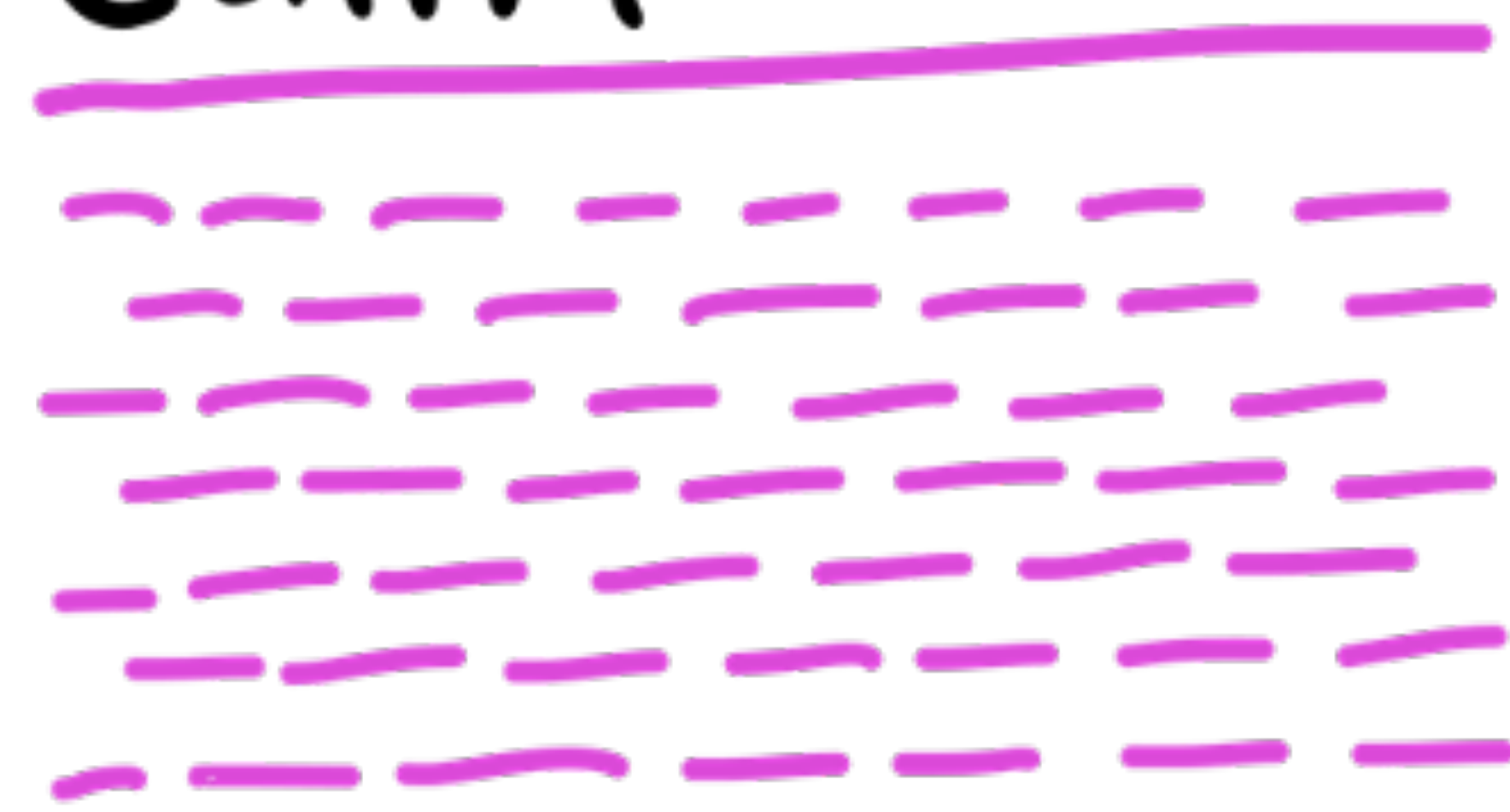


MAPPING

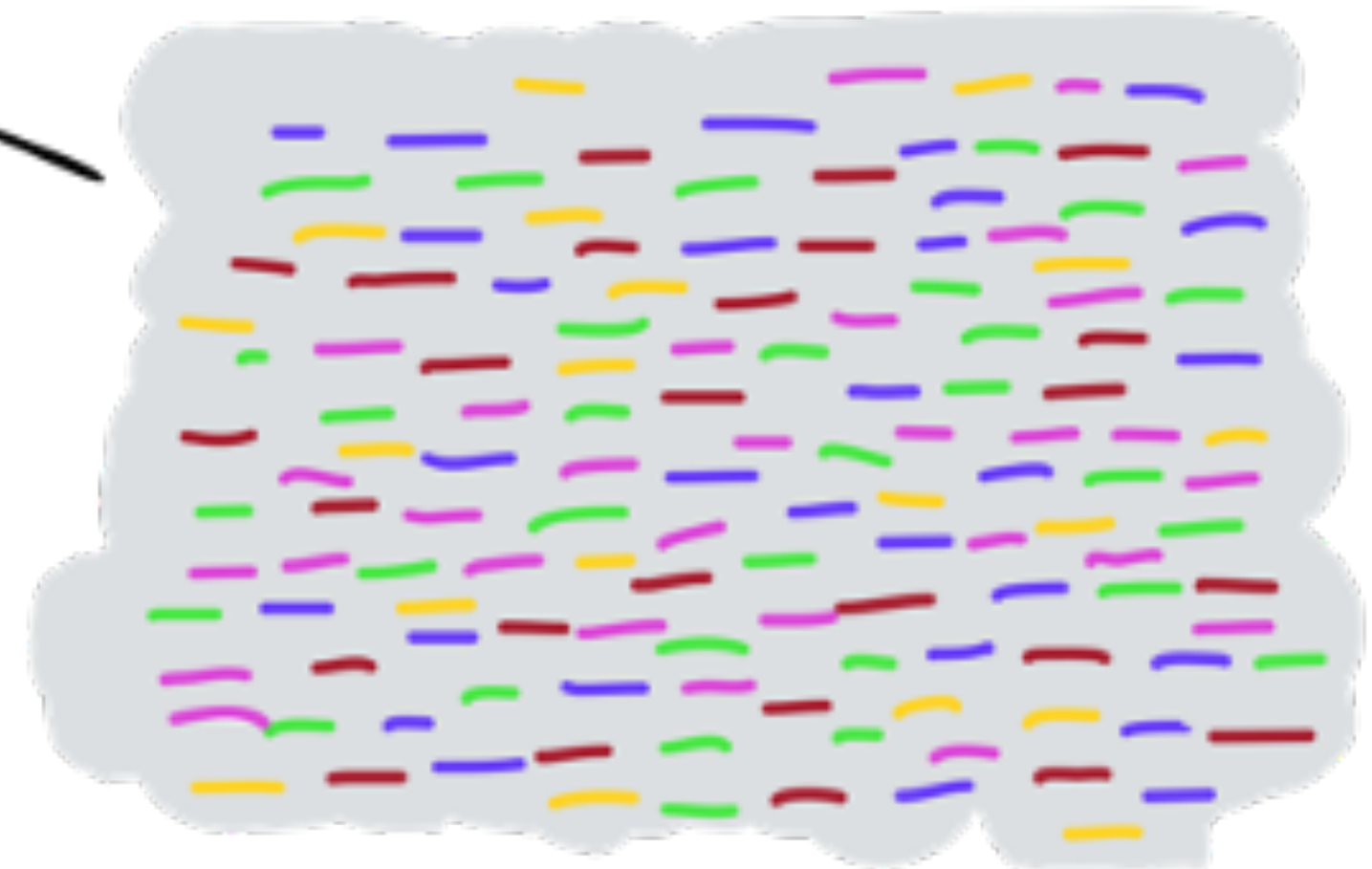


METAGENOMIC READS

CONTIG #1



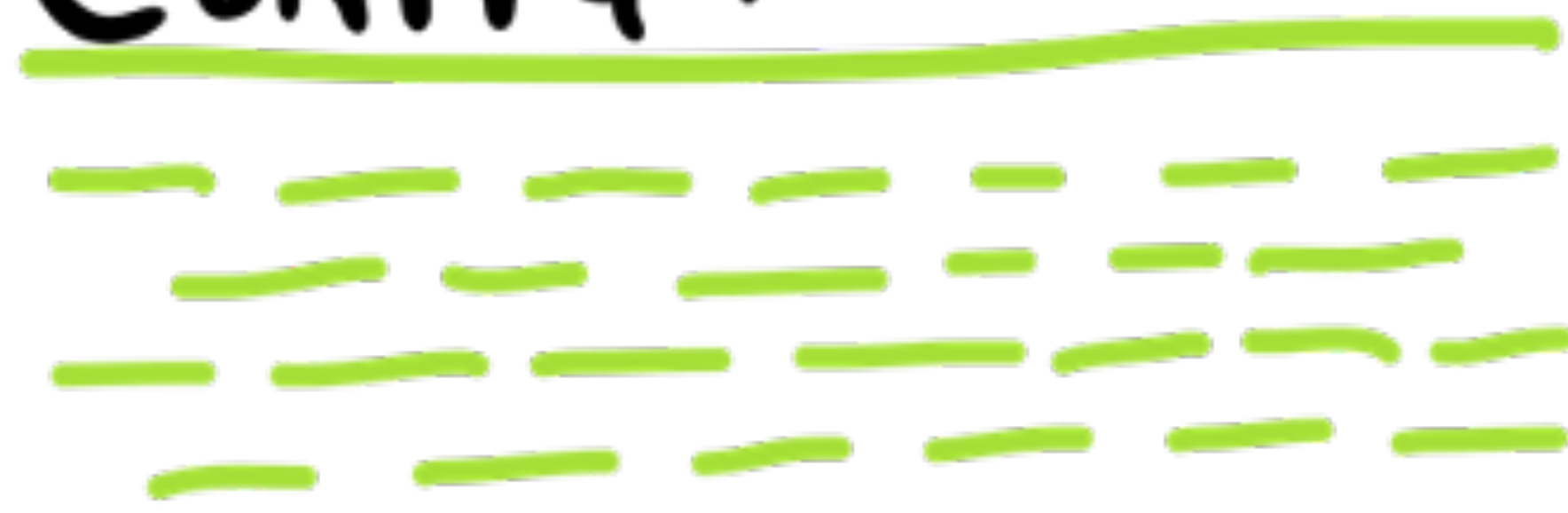
MAPPING



METAGENOMIC READS

COVERAGE: ~7X

CONTIG #2



COVERAGE: ~4X

A \_\_\_\_\_ B \_\_\_\_\_

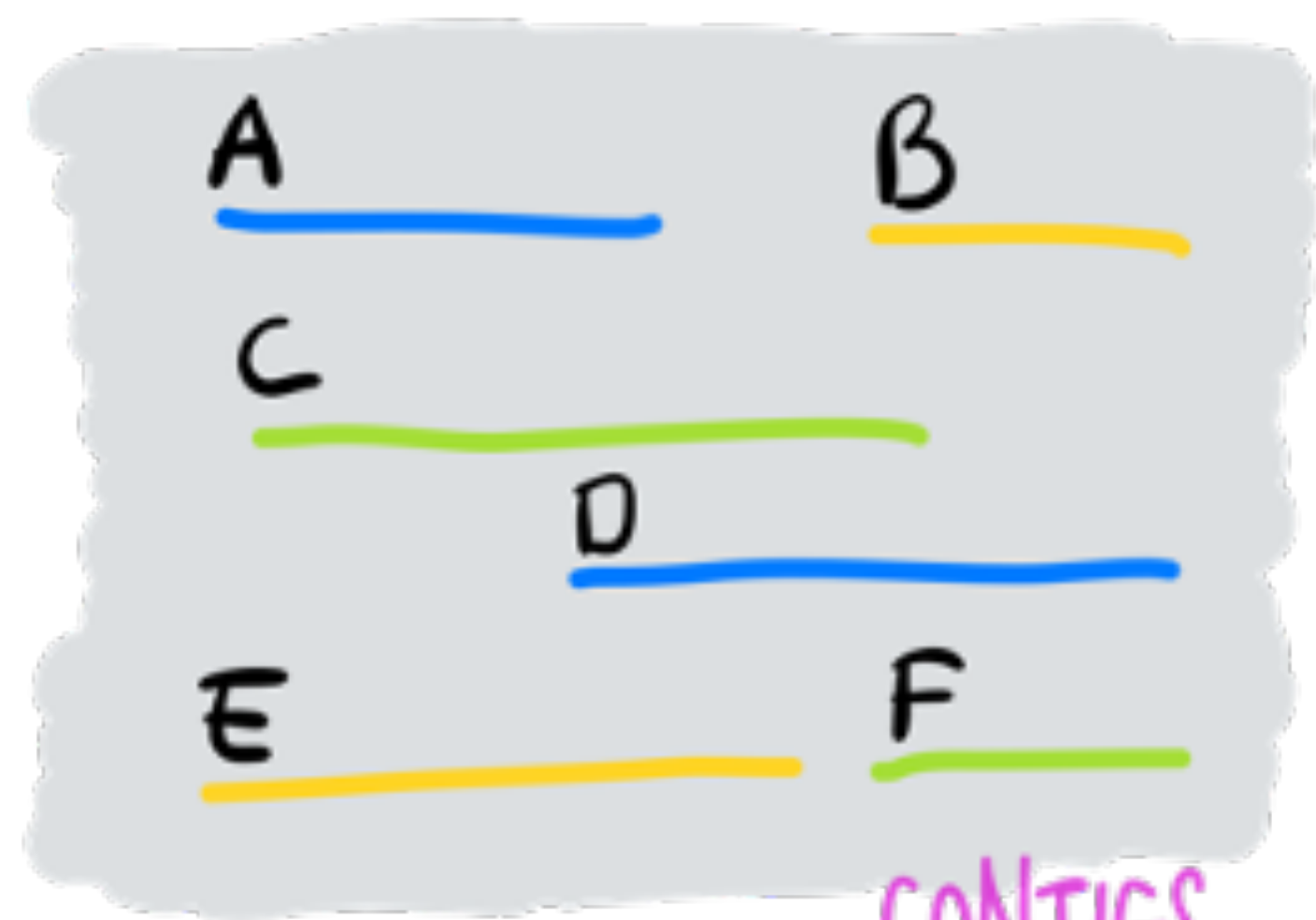
C \_\_\_\_\_

\_\_\_\_\_ D \_\_\_\_\_

E \_\_\_\_\_ F \_\_\_\_\_

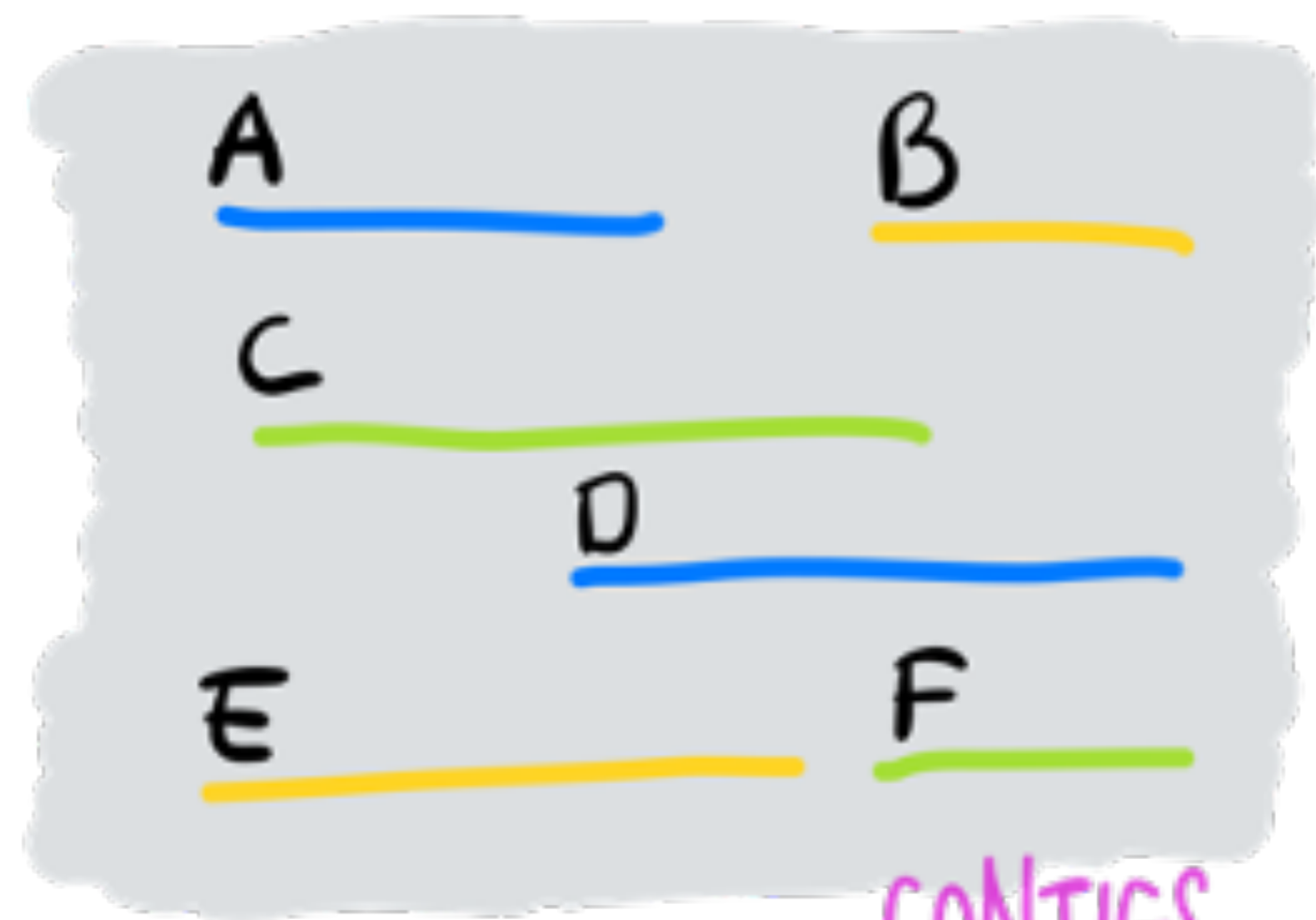
CONTIGS



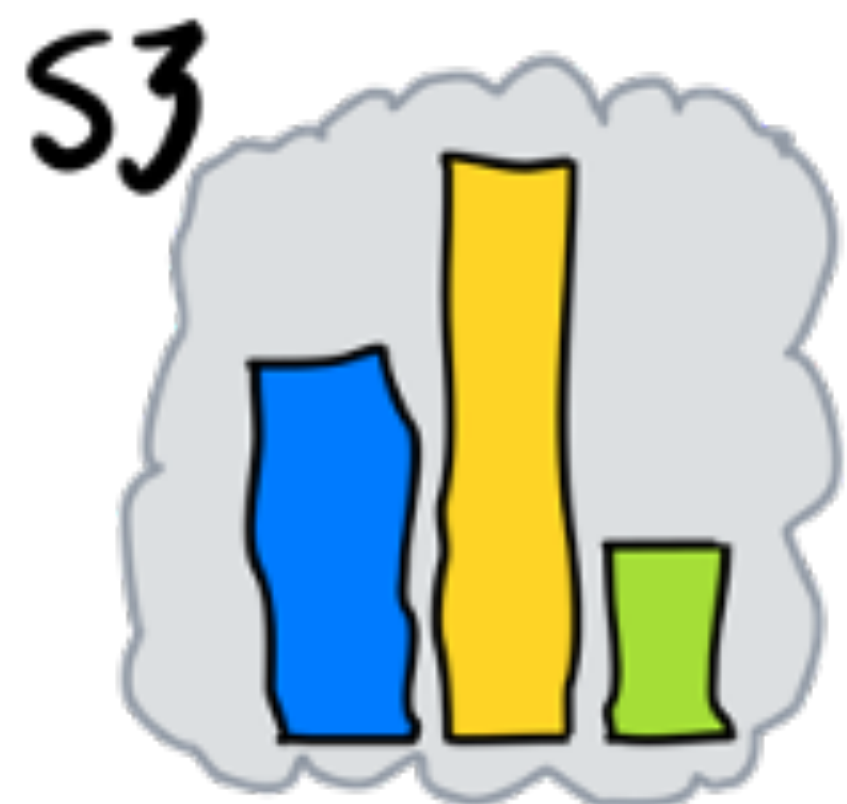


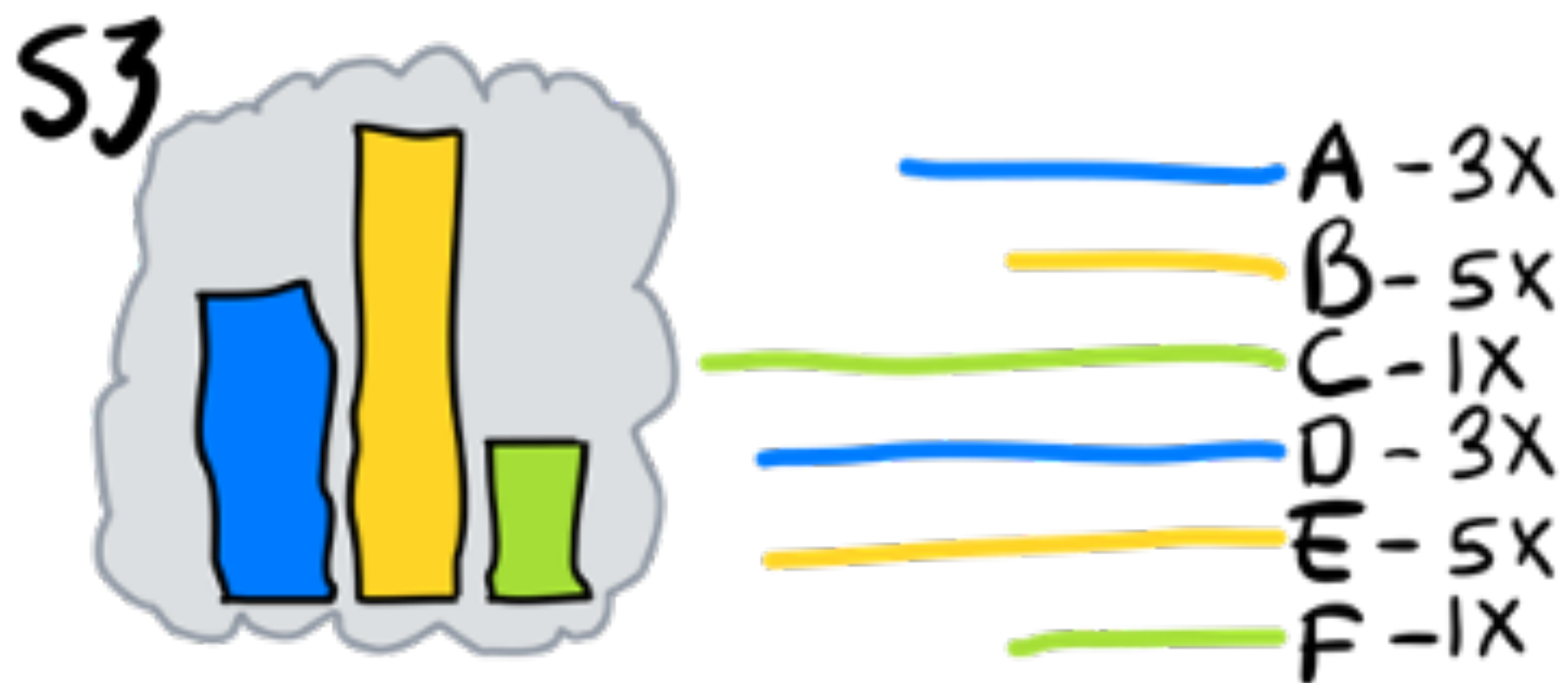
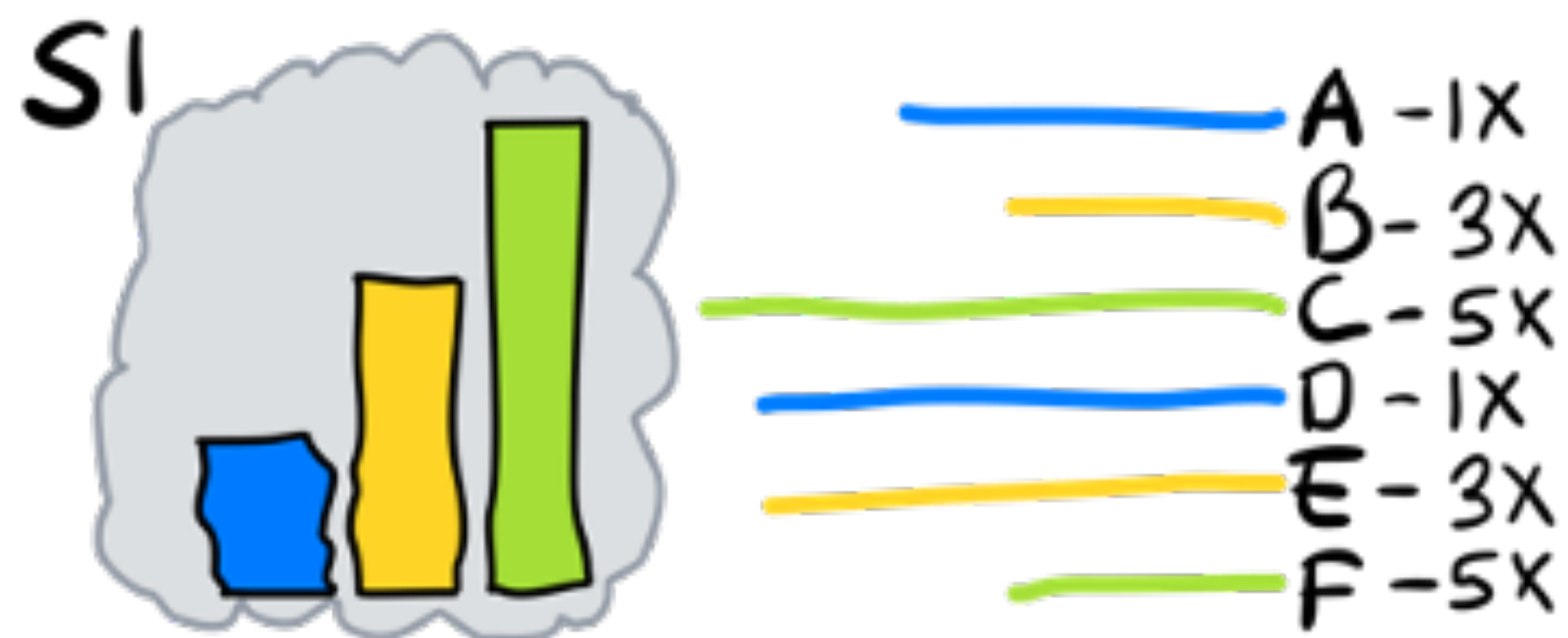
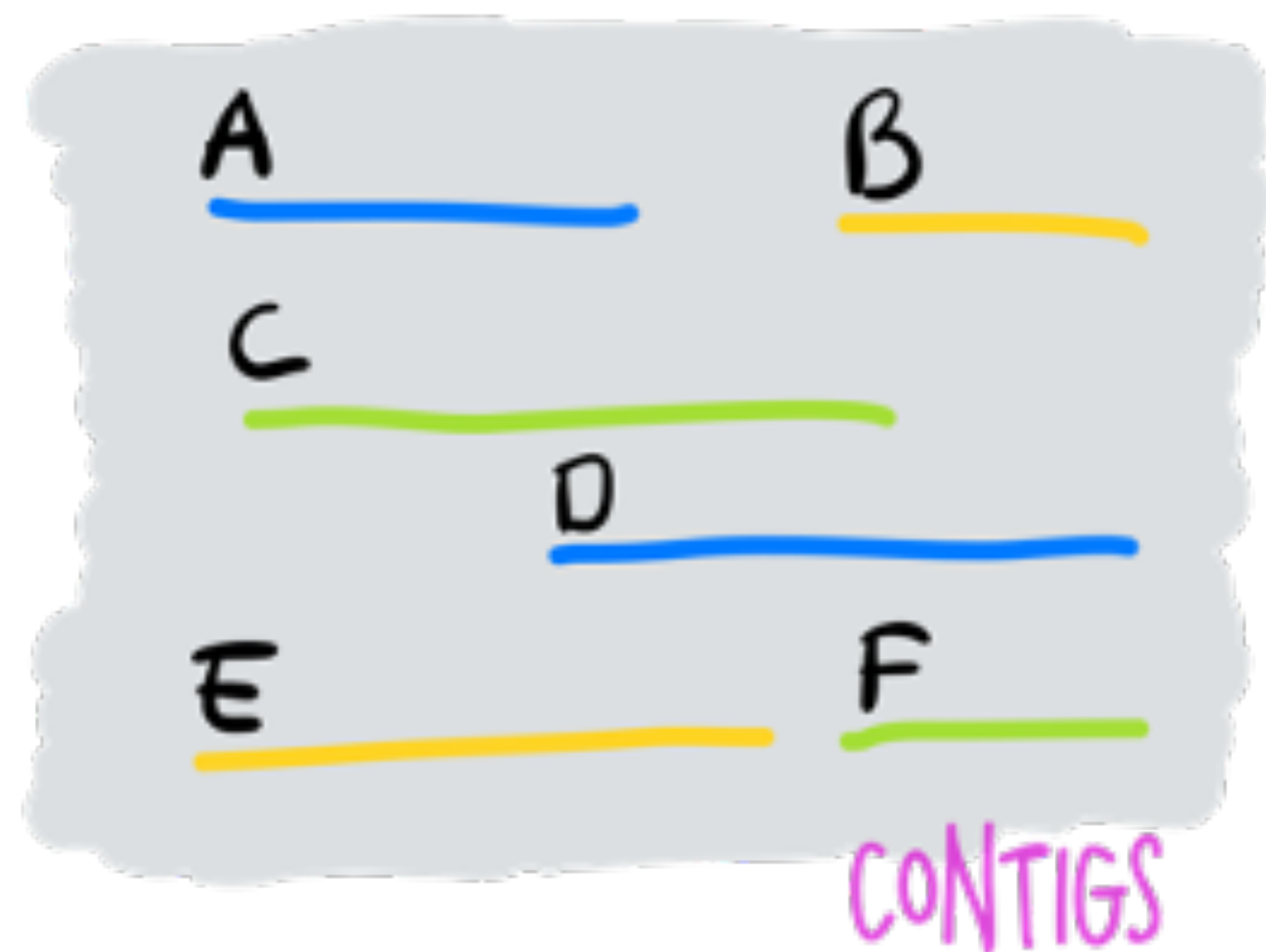
CONTIGS

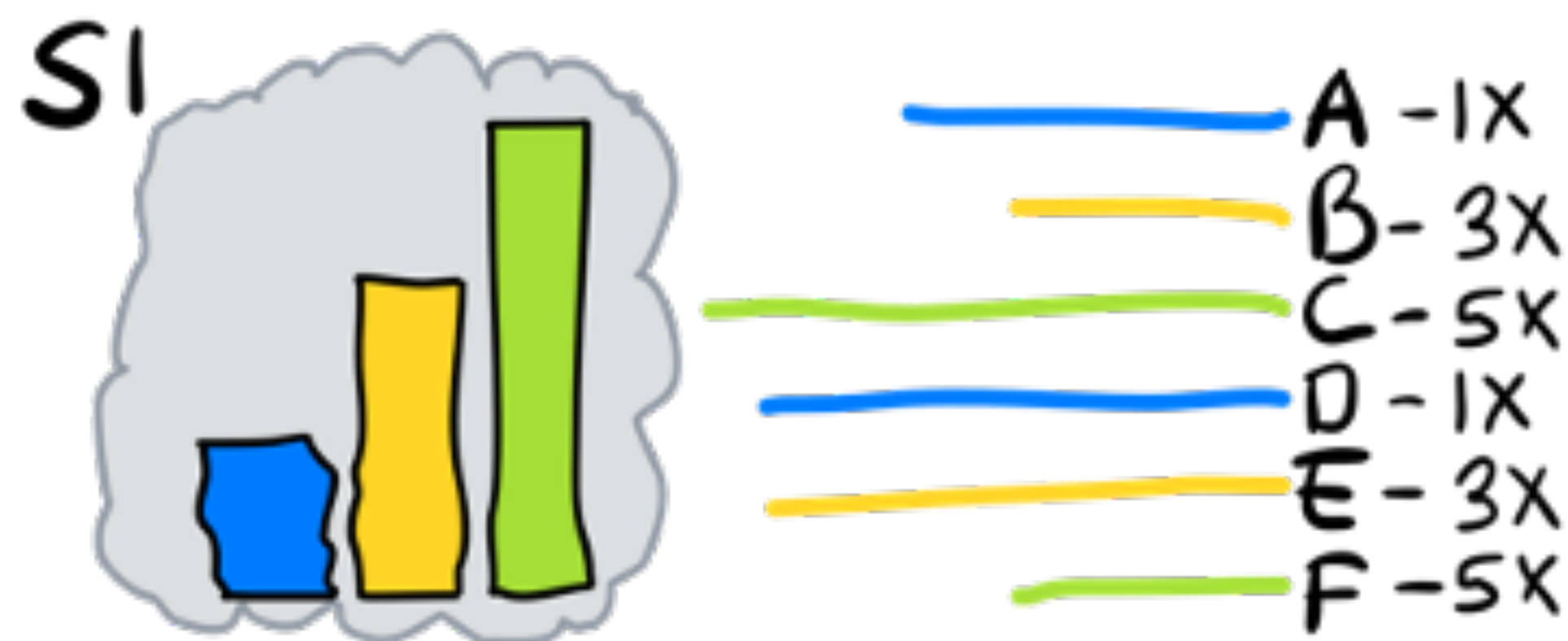
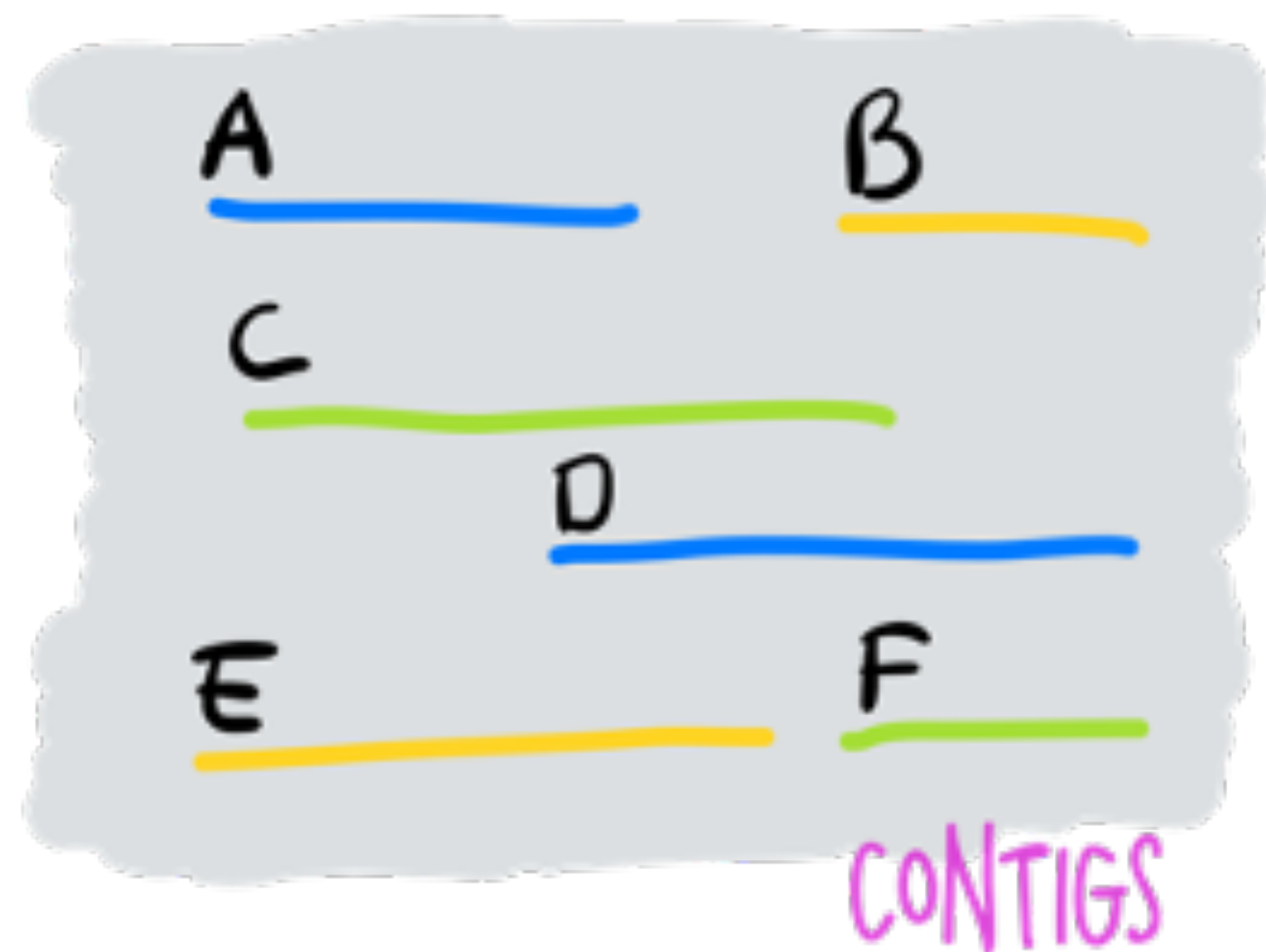




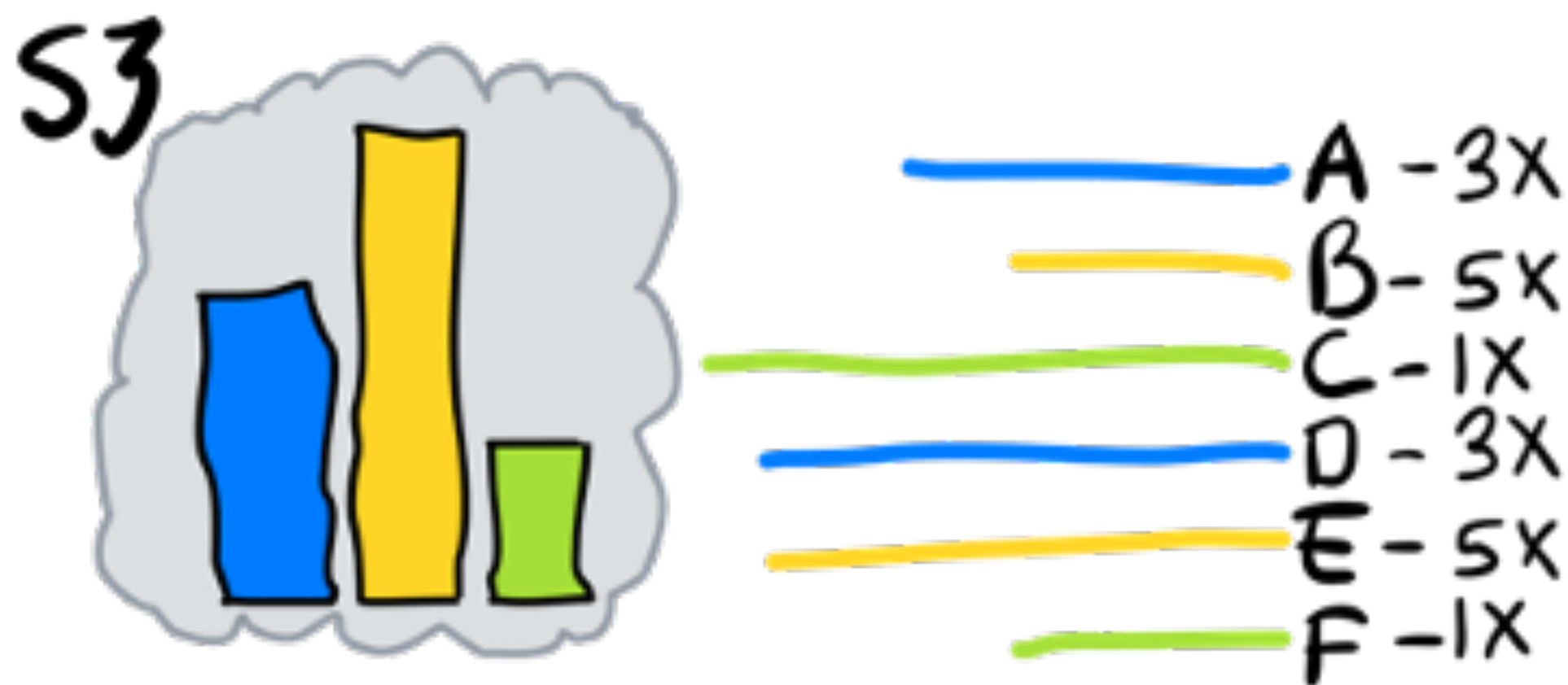
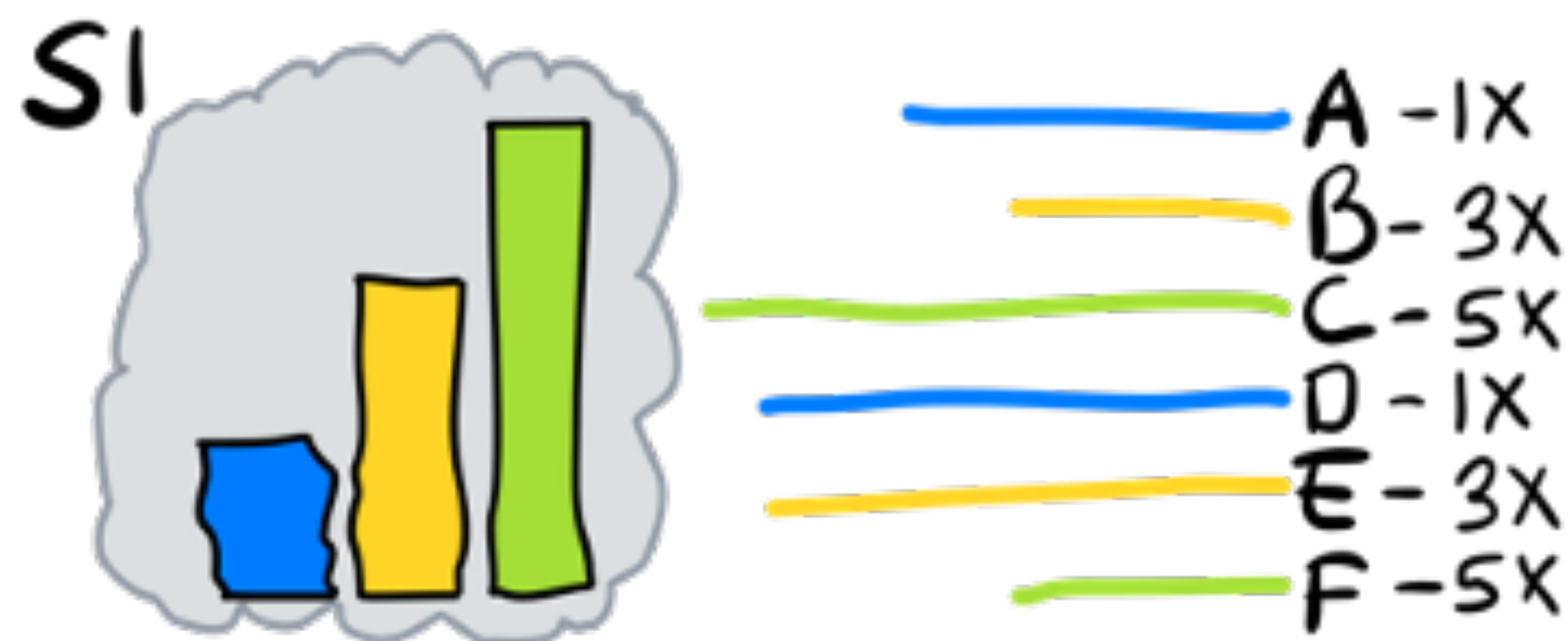
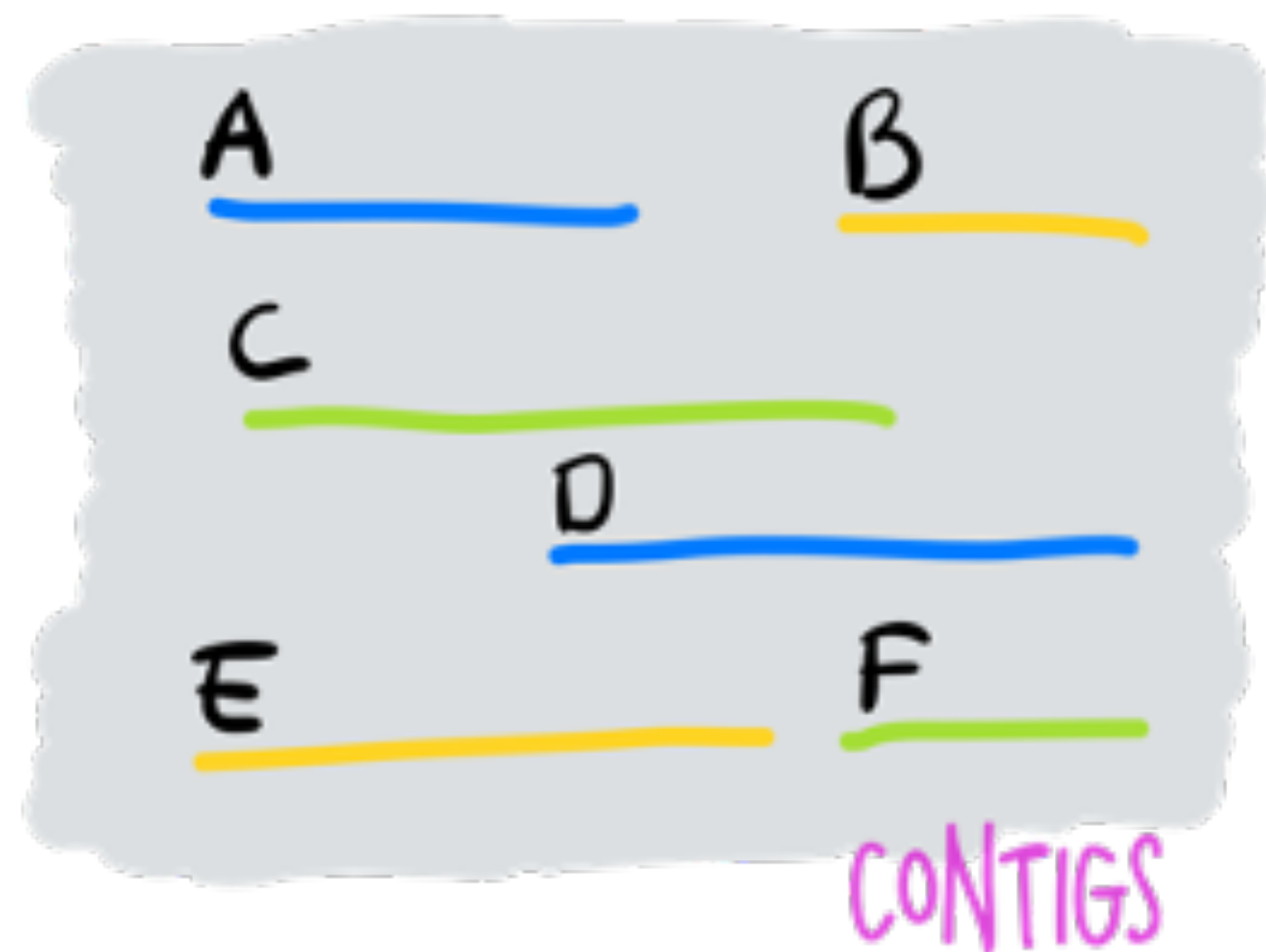
CONTIGS



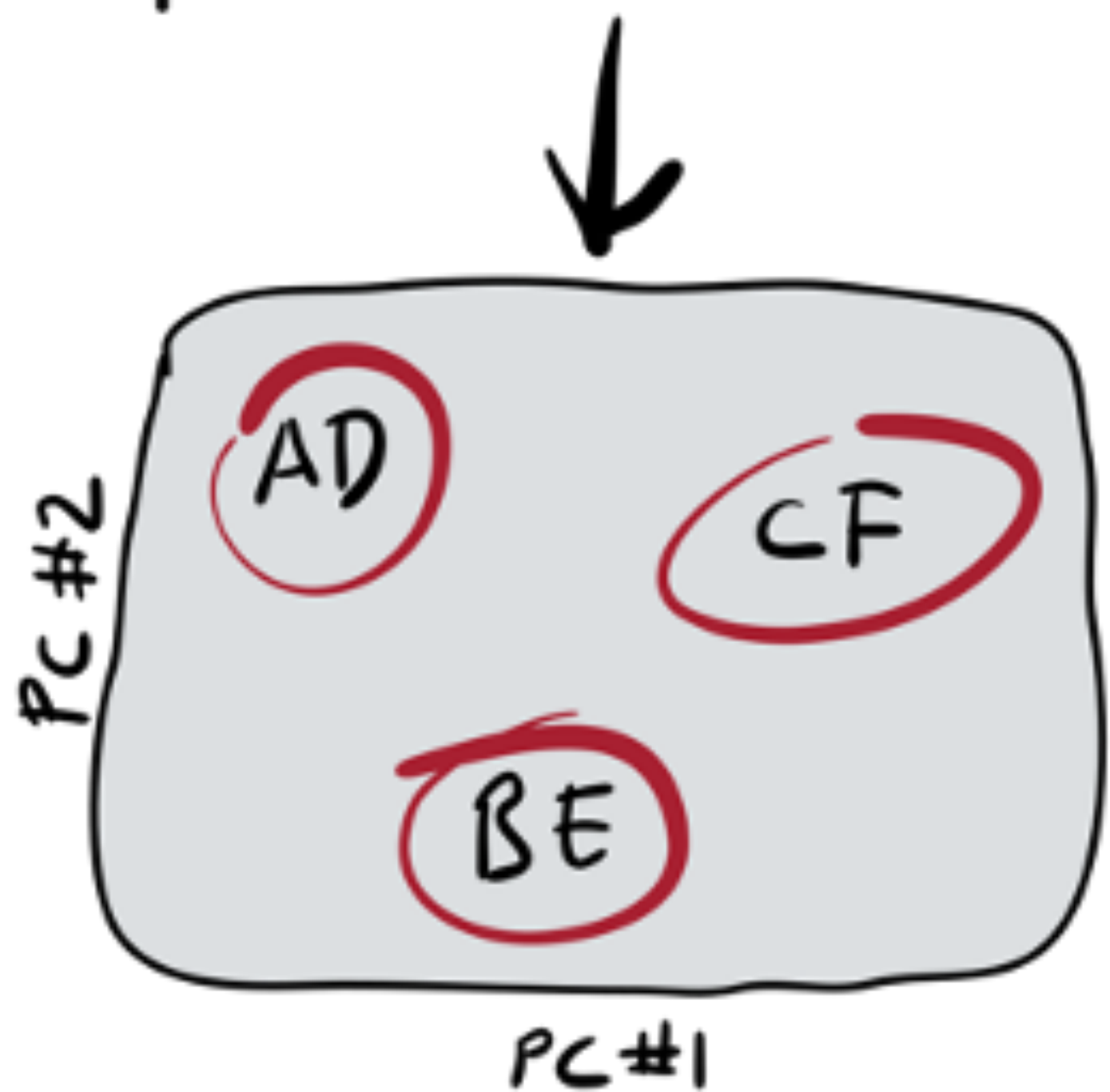




	A	B	C	D	E	F
S1	1	3	5	1	3	5
S2	5	1	3	5	1	3
S3	3	5	1	3	5	1



	A	B	C	D	E	F
S1	1	3	5	1	3	5
S2	5	1	3	5	1	3
S3	3	5	1	3	5	1



SEQUENCE COMPOSITION

CONTIGS



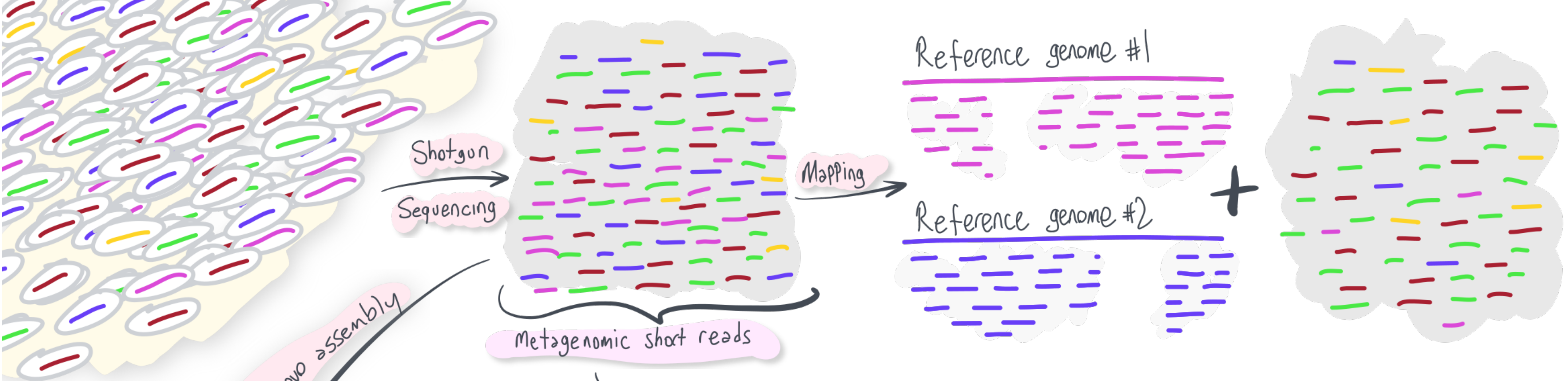
MAGs



DIFFERENTIAL COVERAGE

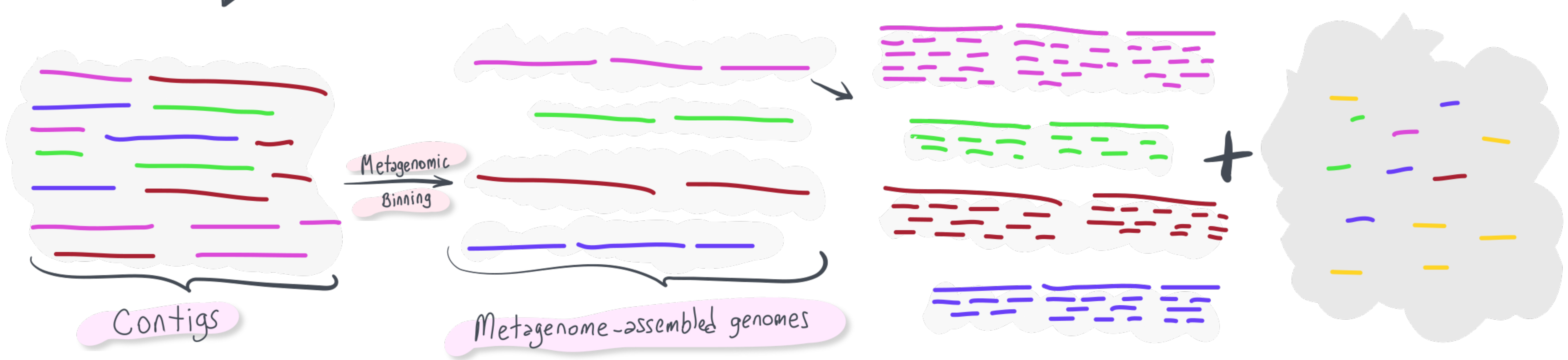
# Cultivation-independent genome-resolved metagenomics

A summary



De novo assembly

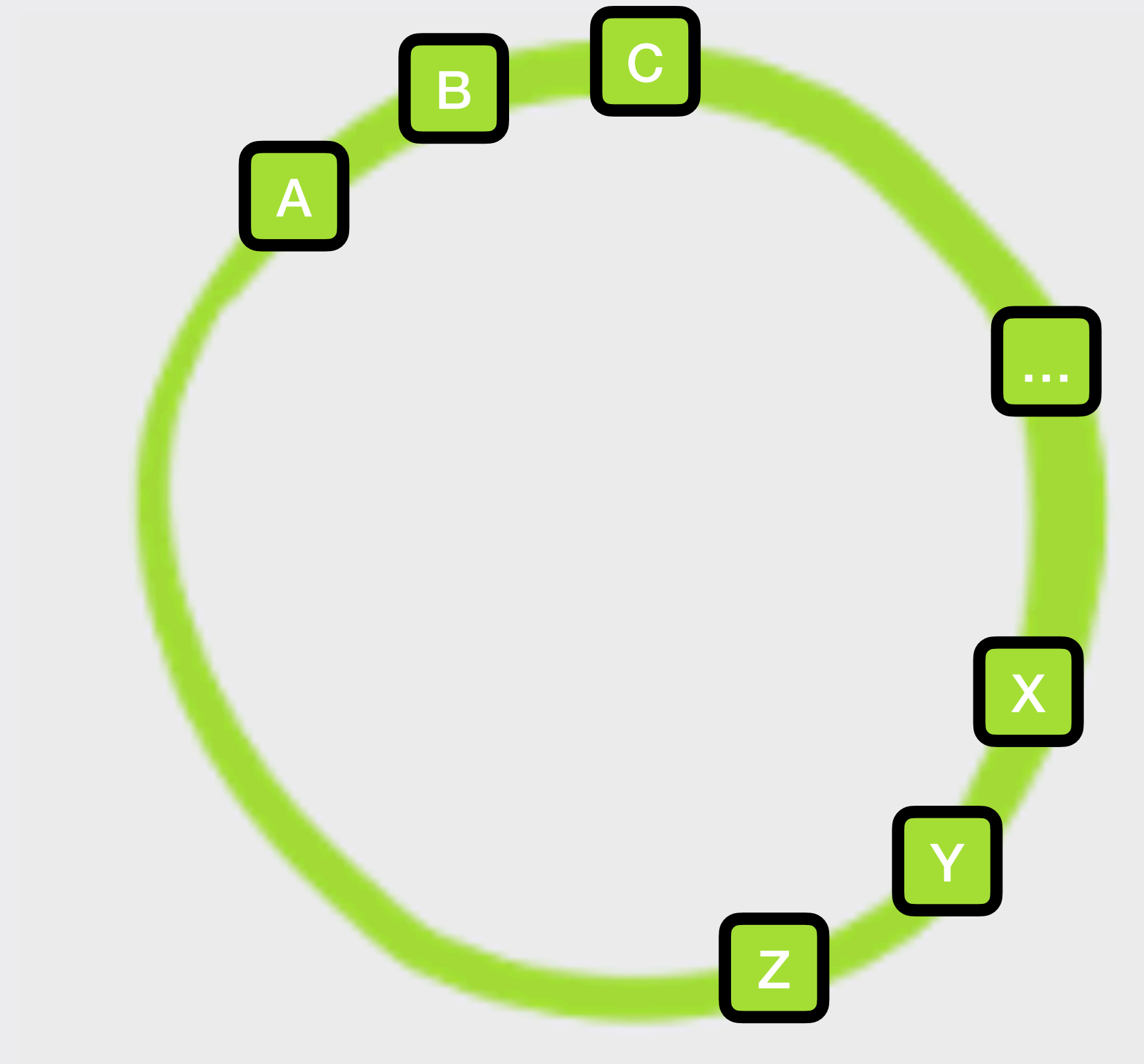
# GENOME RESOLVED METAGENOMICS



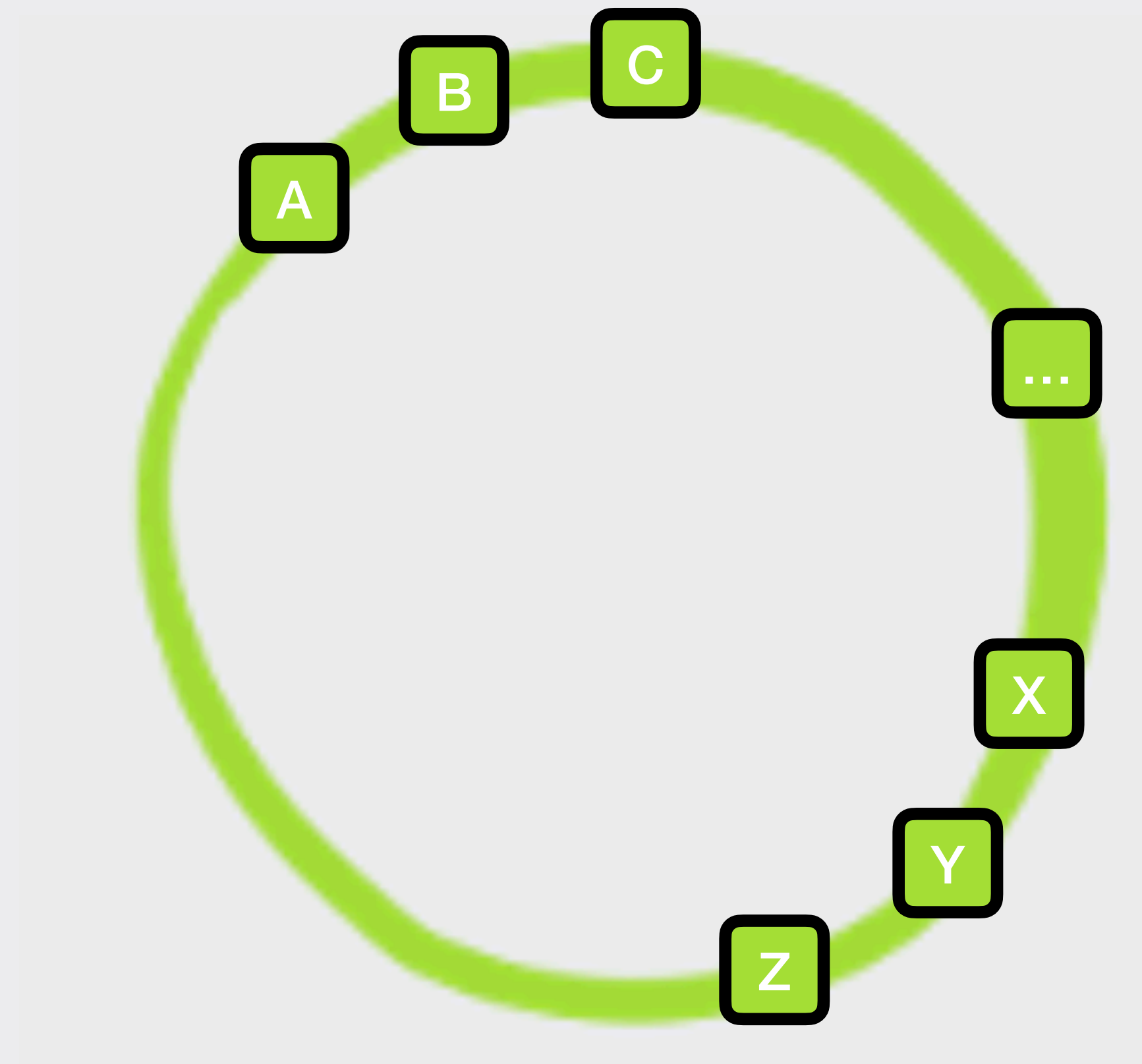


Nothing's perfect  
Evaluation and limits of MAGs

# Universal single-copy marker genes

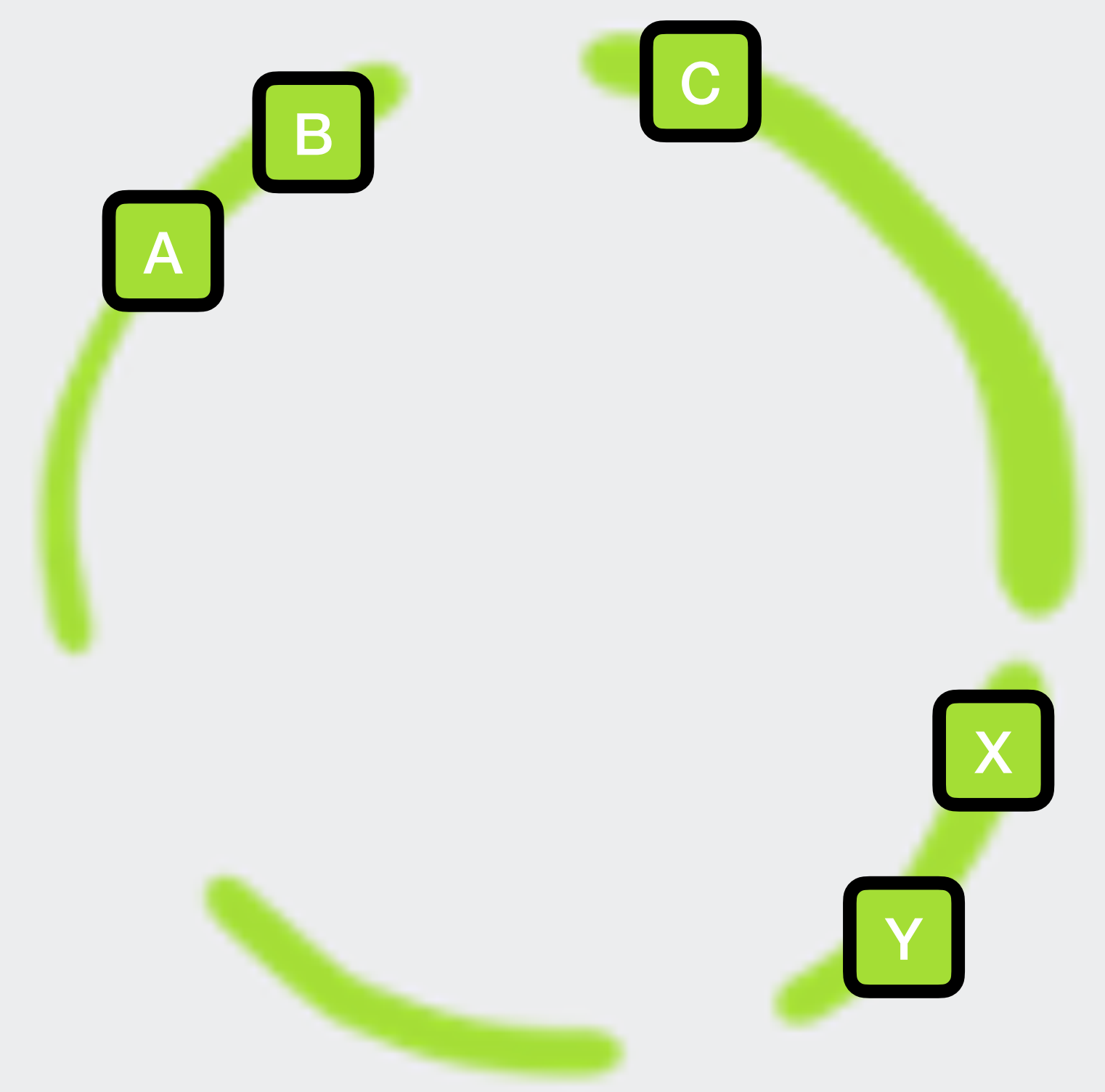
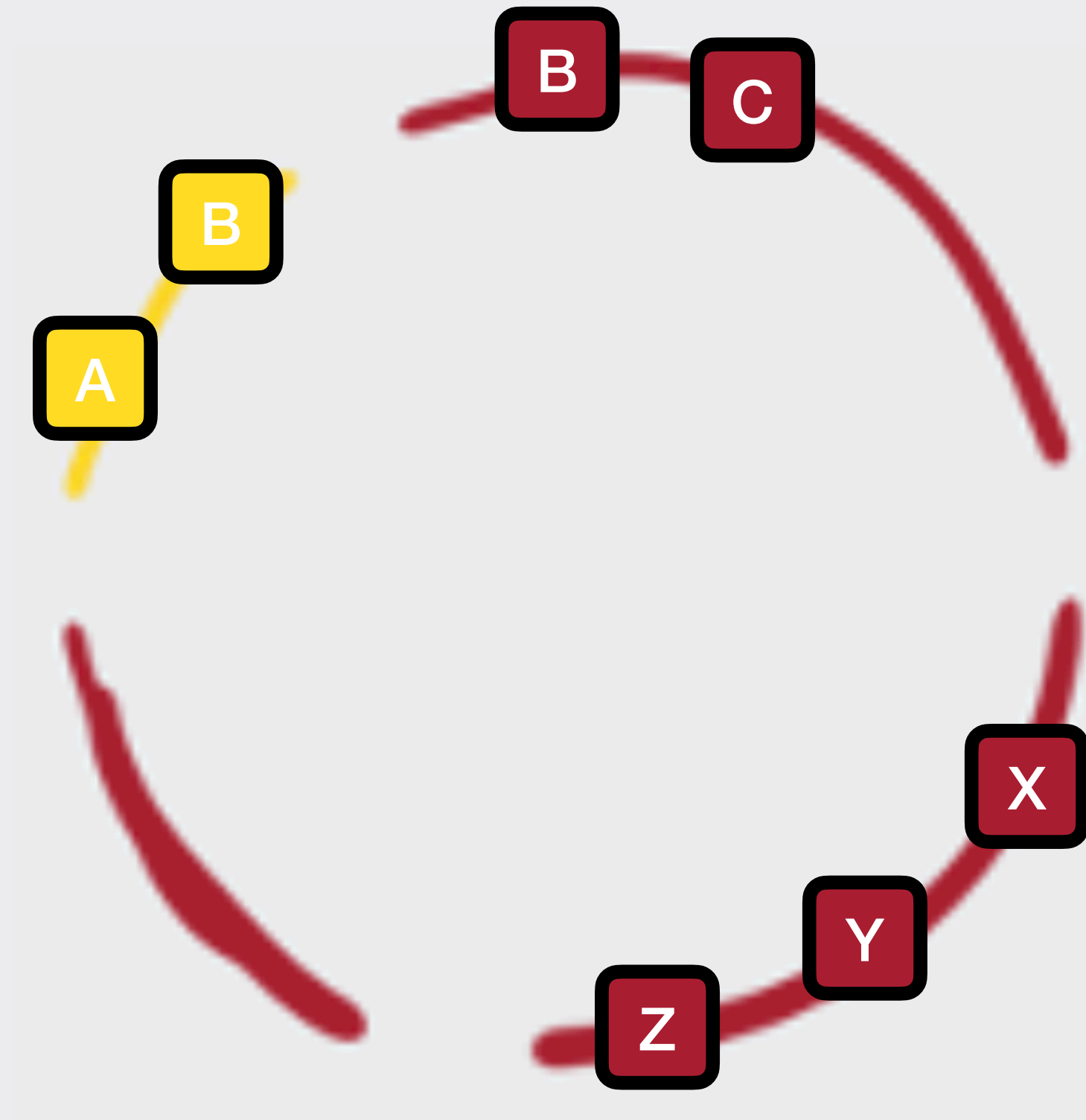
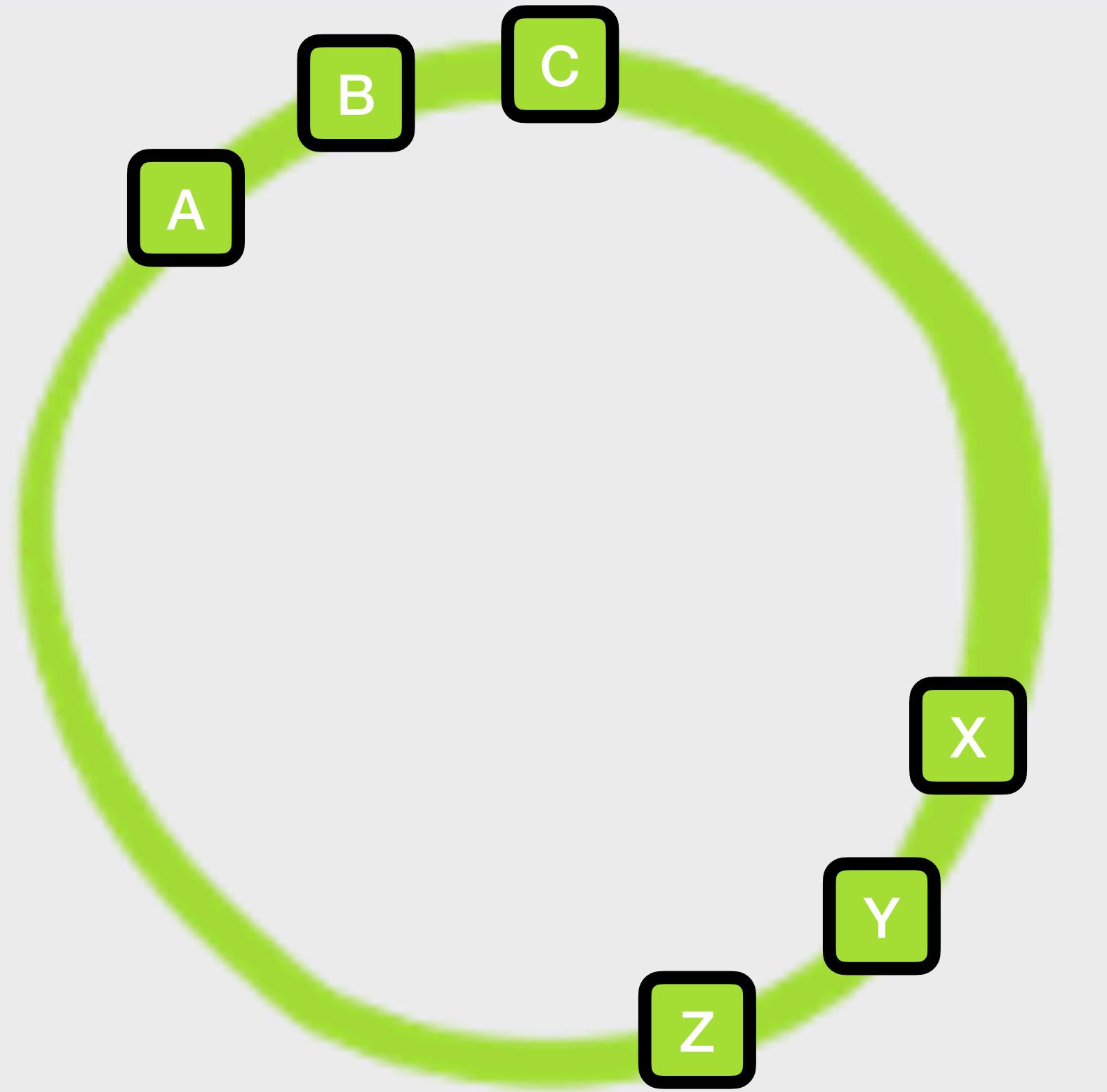


# Universal single-copy marker genes

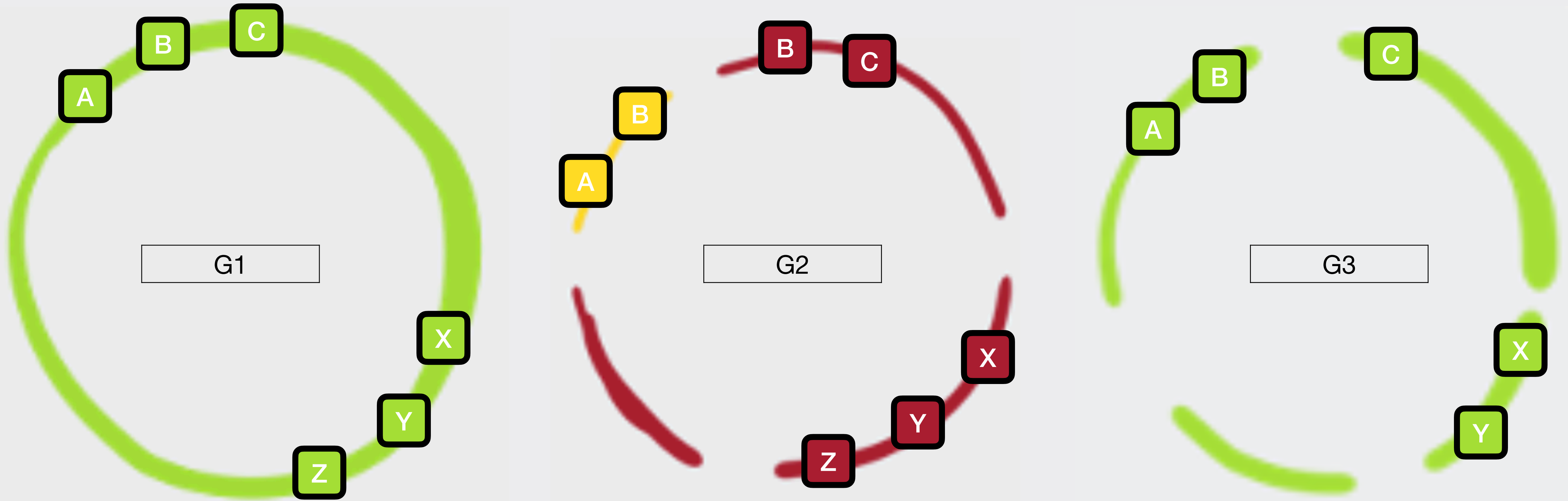


Between 40 and 120 for Bacteria/Archaea depending on cutoffs

# Universal single-copy marker genes

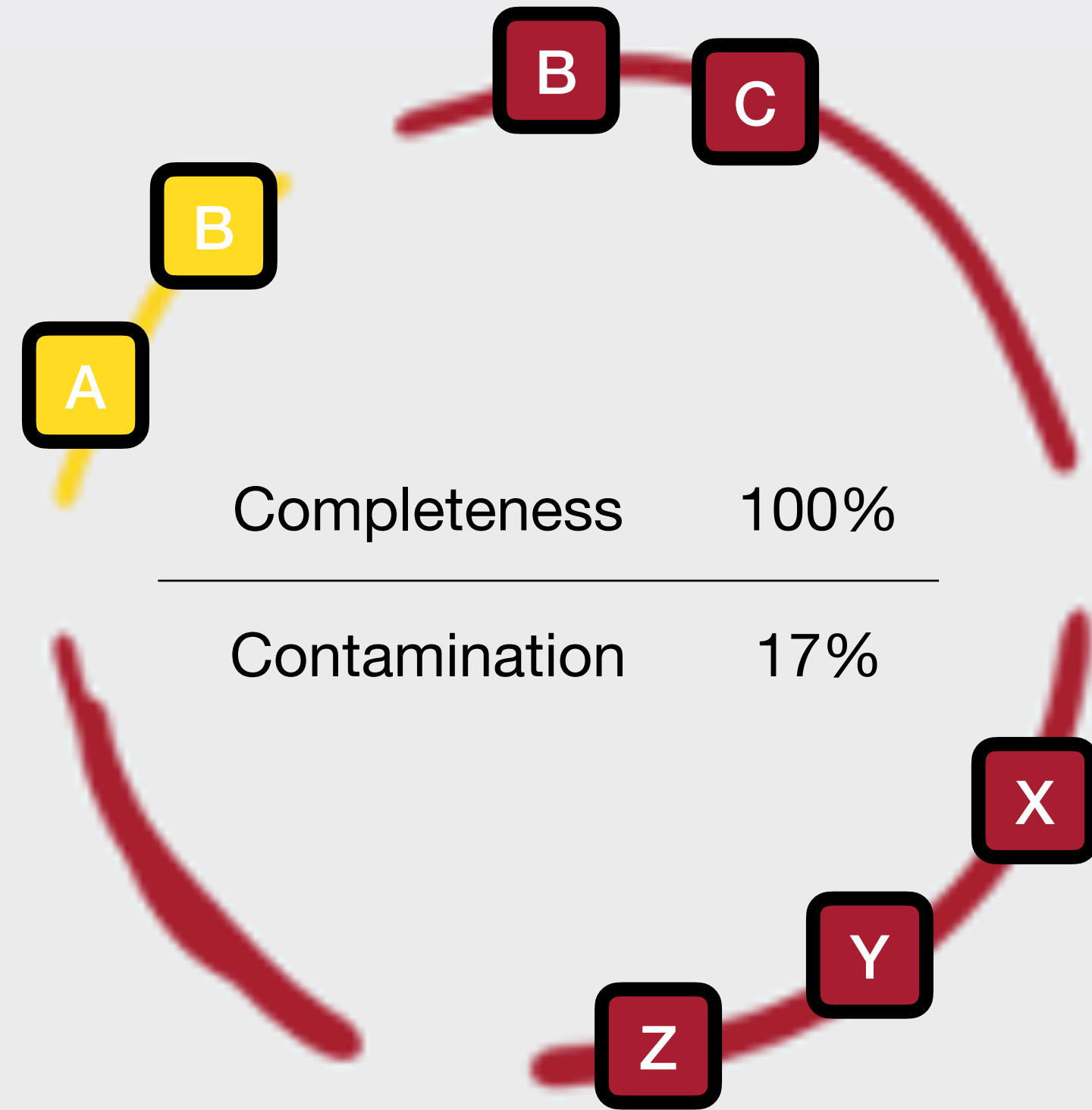


# Universal single-copy marker genes



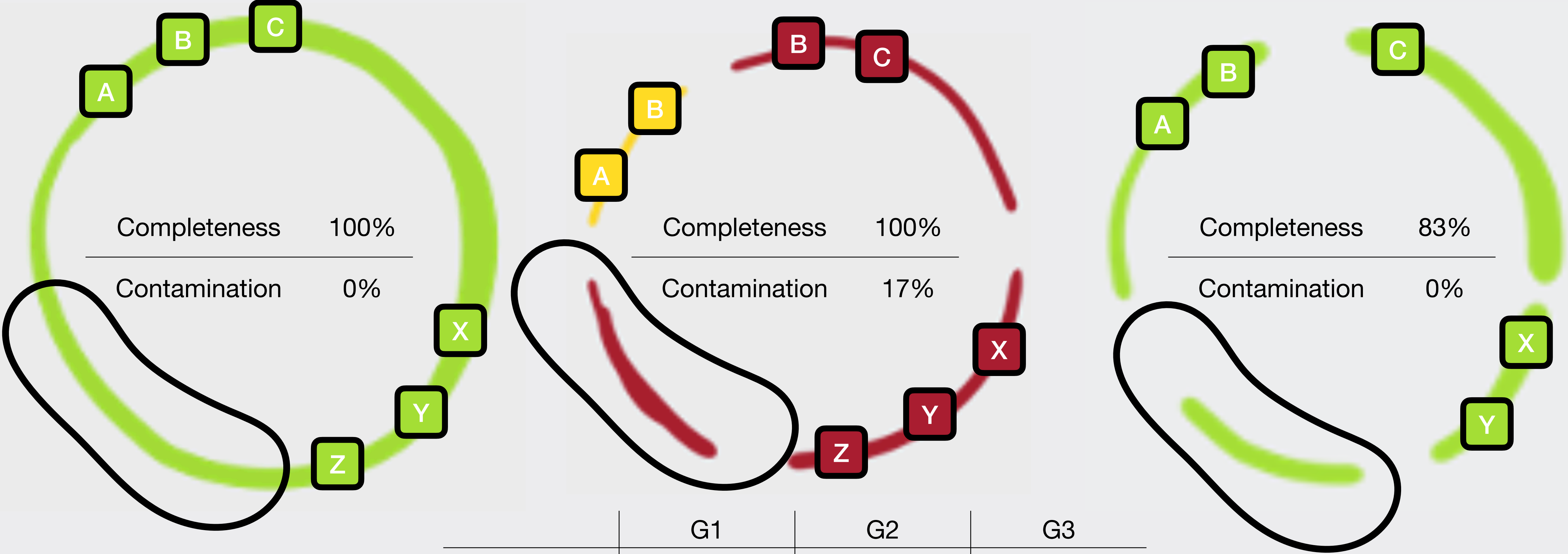
	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# Universal single-copy marker genes



	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# Universal single-copy marker genes



	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# Material

This course uses a lot of material from <https://merenlab.org/momics/>, I invite you to have a look.

If you want details on the bioinformatics behind you can start by having a look here:

[https://astrobiomike.github.io/genomics/metagen\\_anvio](https://astrobiomike.github.io/genomics/metagen_anvio)



## A. Murat Eren (Meren) (PI)

Web  Email  Twitter  LinkedIn  Github  ORCID

Address: Knapp Center for Biomedical Discovery, 900 E. 57th St., MB 9, RM 9118, Chicago, IL 60637 USA

Phone: +1-773-702-5935  Fax: +1-773-702-2281

*I am a computer scientist with a deep appreciation for the complexity of life. I design algorithms and experiments to better understand microbes and their ecology. [photos: 1, 2, 3].*

- » MBL Fellow, [Marine Biological Laboratory](#).
- » Assistant Professor, [The Department of Medicine at the University of Chicago](#).
- » Committee on Microbiology, [The Biomedical Sciences Cluster at the University of Chicago](#).



## Mike Lee

Web  Email  Twitter  LinkedIn  Github

» NASA Space Biology Fellow, [NASA Ames Research Center](#).

» JCVI Research Fellow, [J. Craig Venter Institute](#).

- 👉 [Combining reference genome annotations with your own in pangenomes](#) (Sat, Dec 01, 2018)
- 👉 [Anvi'o 'views' demystified](#) (Mon, May 08, 2017)
- 👉 [Making anvi'o use your own HMM collection](#) (Sat, May 21, 2016)

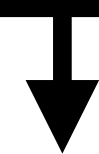


Application to the ocean microbiome  
Uncharted biosynthetic potential

# Why explore environmental microbiomes?

- Ubiquitous across earth's ecosystems
- Support global food webs
- Underpin biogeochemical cycles
- Determine Host's health and disease
- ...
- **Untapped metabolic diversity**
  - **Biosynthetic potential**

- **New Enzymes**
- **New Natural products**



Applications

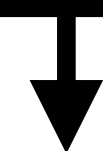
# Why explore environmental microbiomes?

- Ubiquitous across earth's ecosystems
- Support global food webs
- Underpin biogeochemical cycles
- Determine Host's health and disease
- ...
- **Untapped metabolic diversity**
- **Biosynthetic potential**

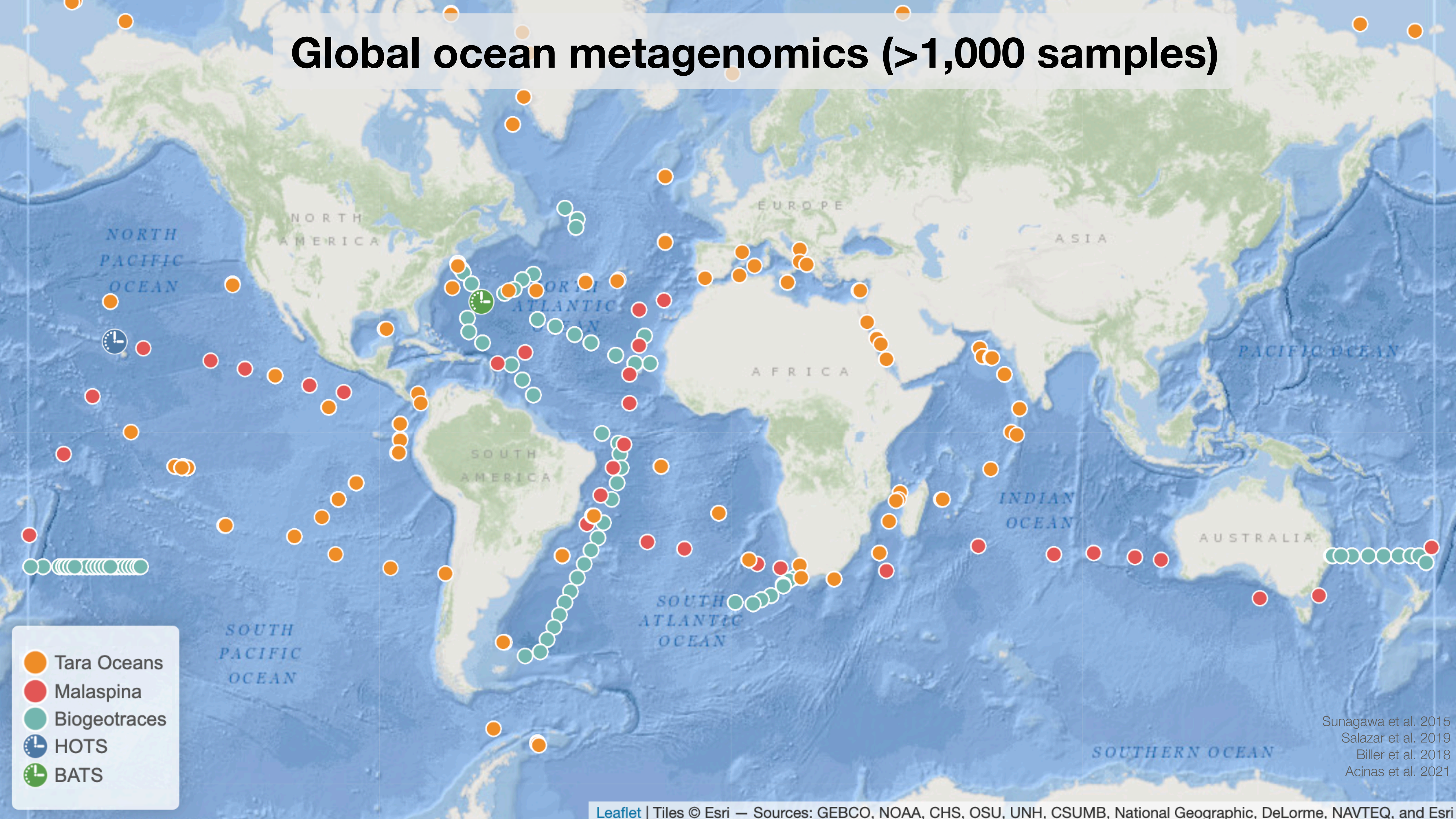
- **New Enzymes**
- **New Natural products**

Applications

Microbial  
Interactions



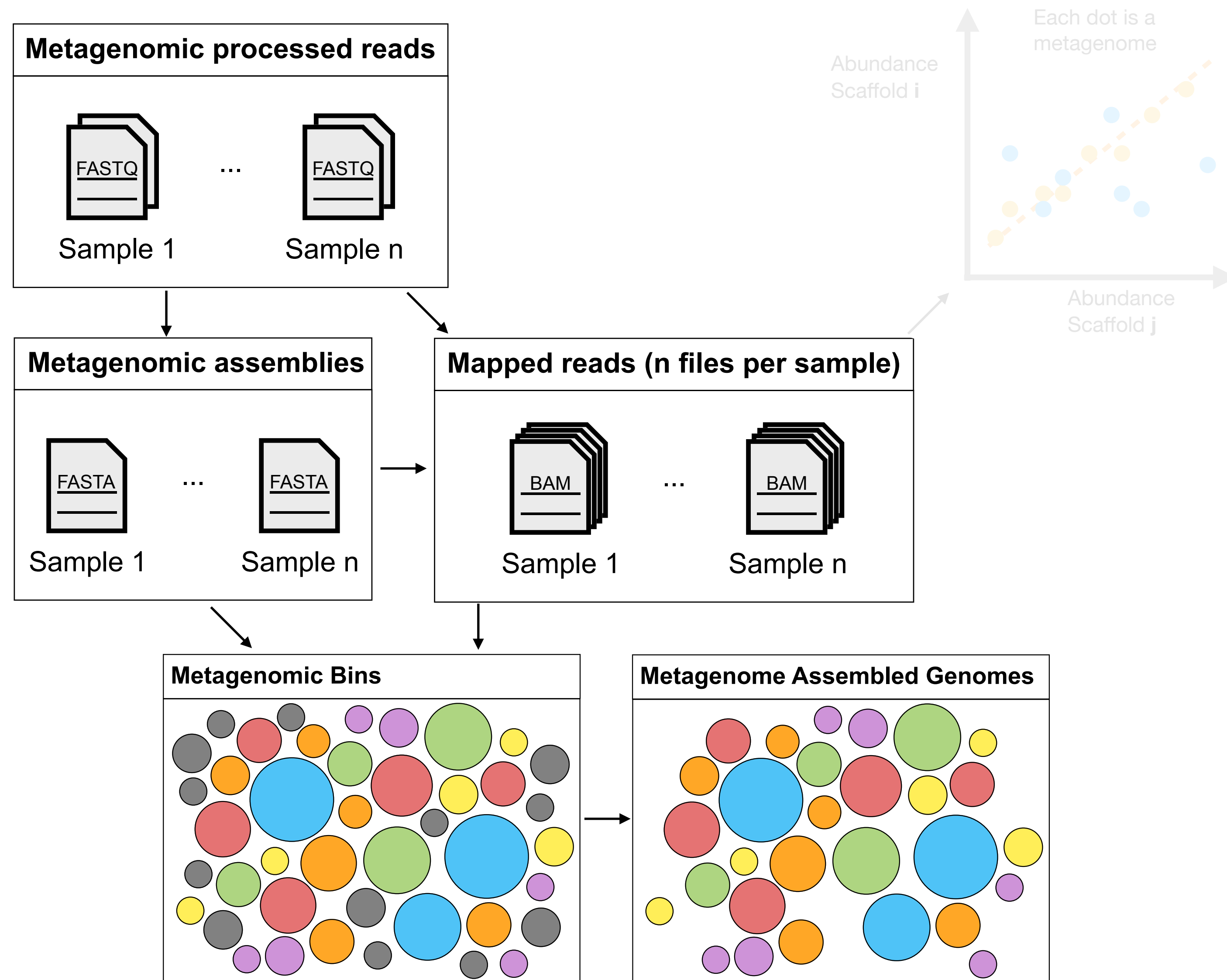
# Global ocean metagenomics (>1,000 samples)



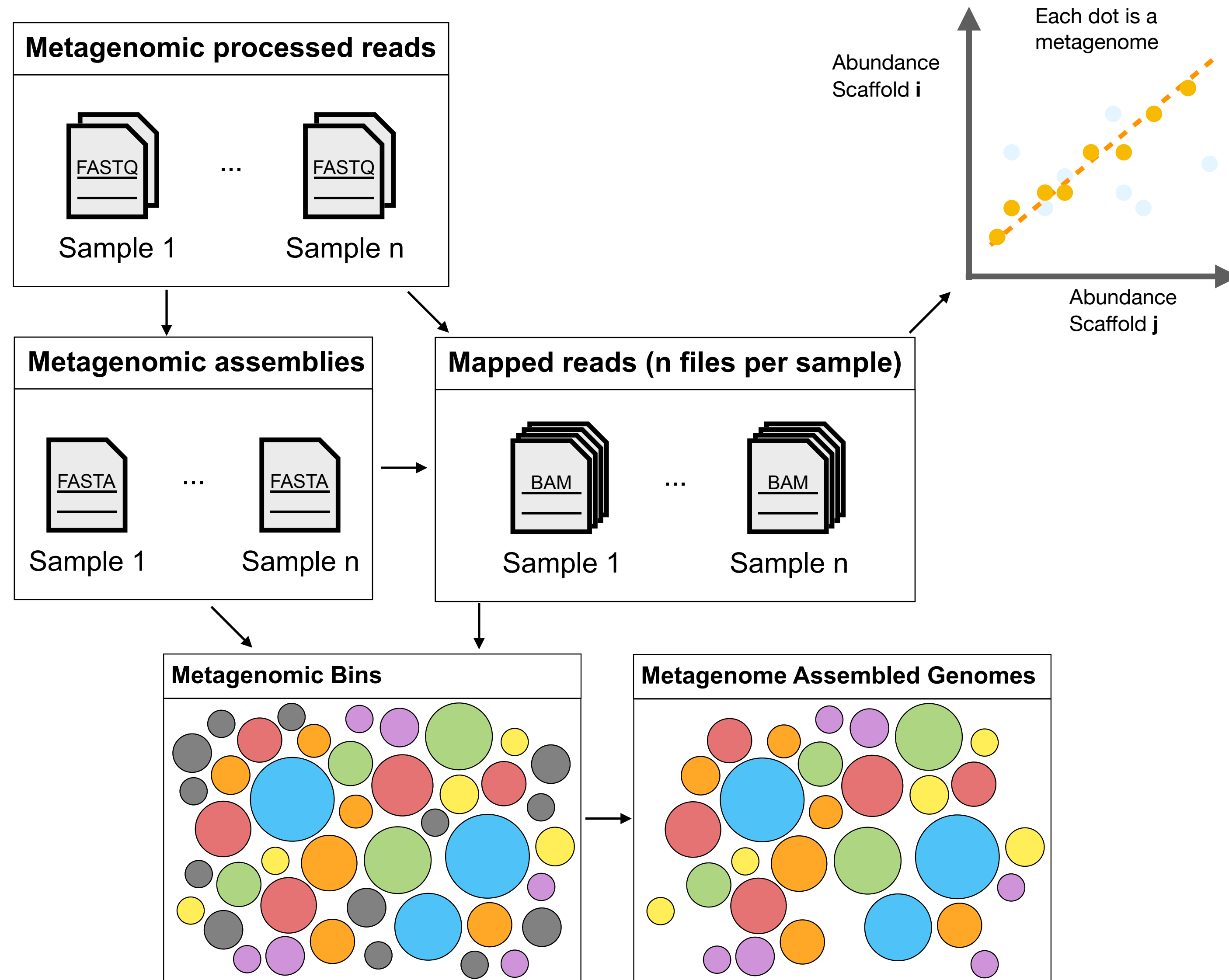
- Tara Oceans
- Malaspina
- Biogeotraces
- HOTS
- BATS

Sunagawa et al. 2015  
Salazar et al. 2019  
Biller et al. 2018  
Acinas et al. 2021

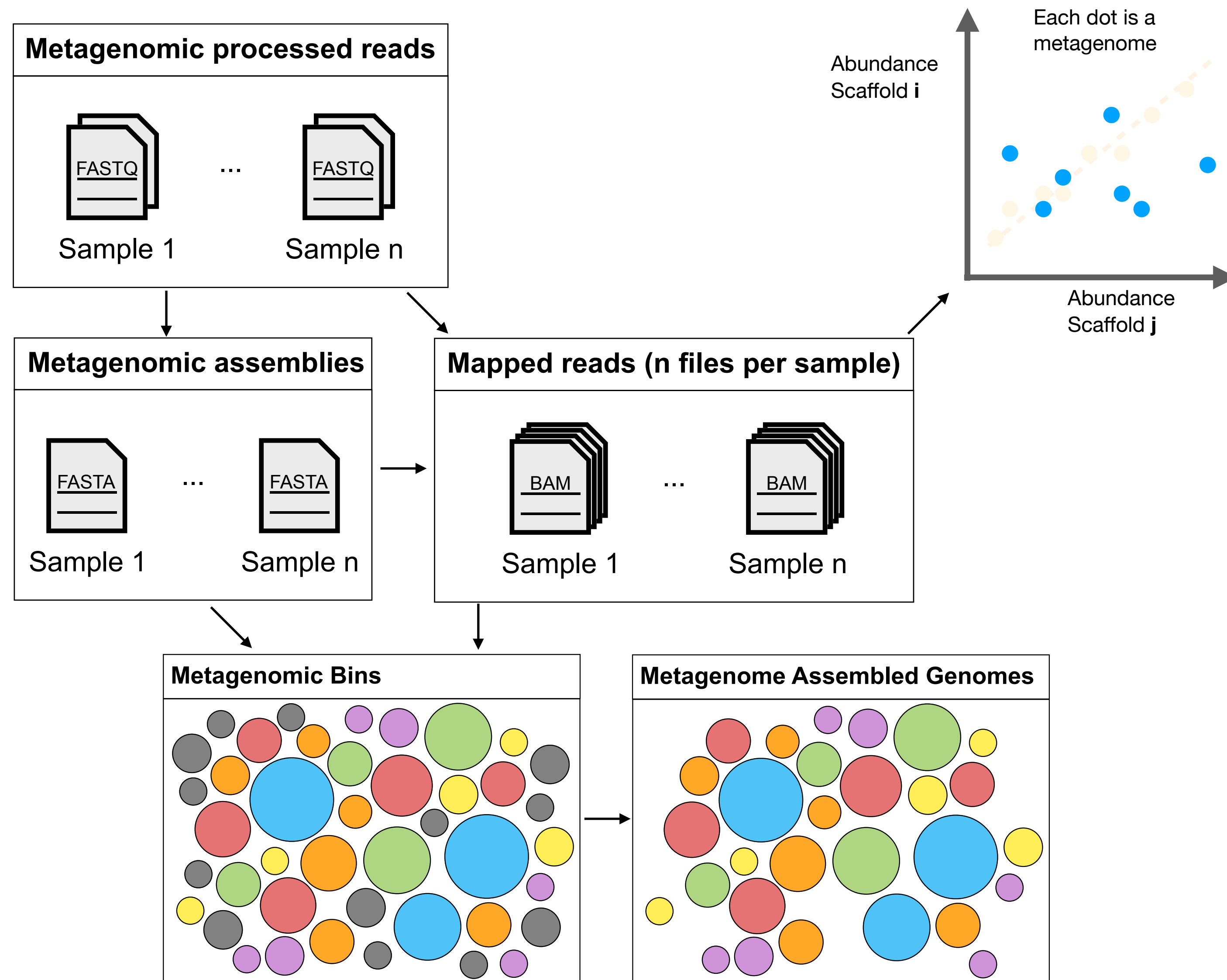
# Reconstructing genomes from metagenomes (MAGs)



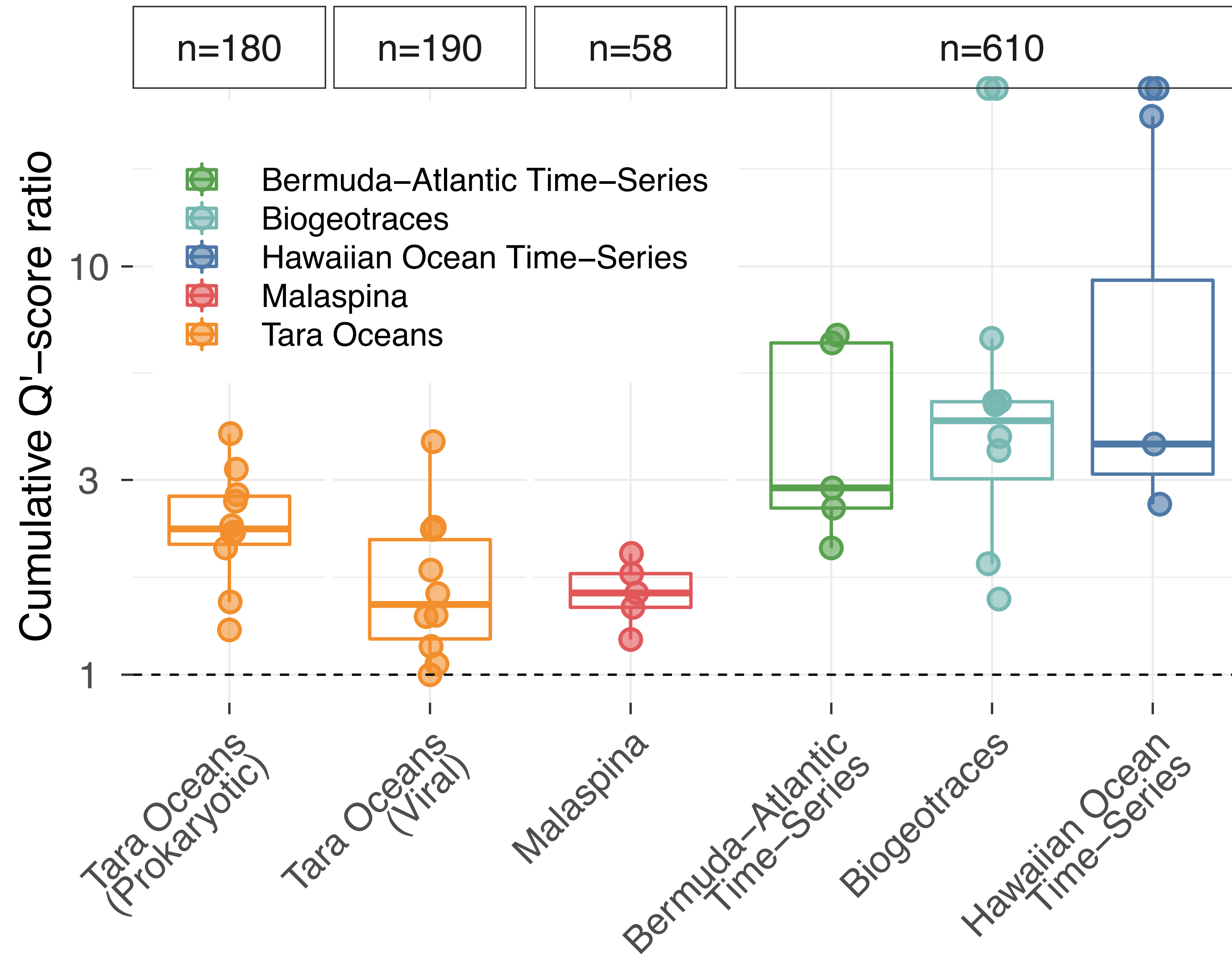
# Reconstructing genomes from metagenomes (MAGs)



# Reconstructing genomes from metagenomes (MAGs)



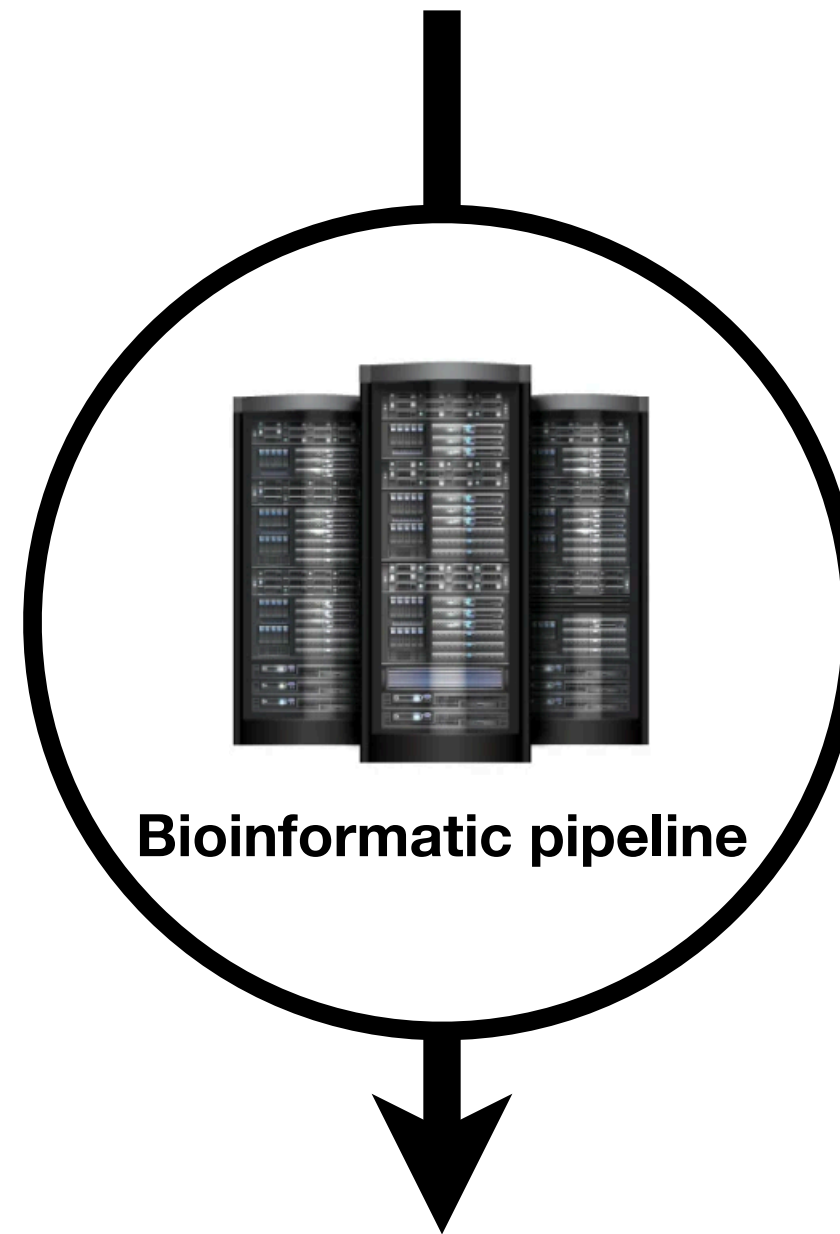
# Abundance correlation improves binning results three-folds





# Reconstructing microbial genomes from metagenomes

>1,000 Metagenomes



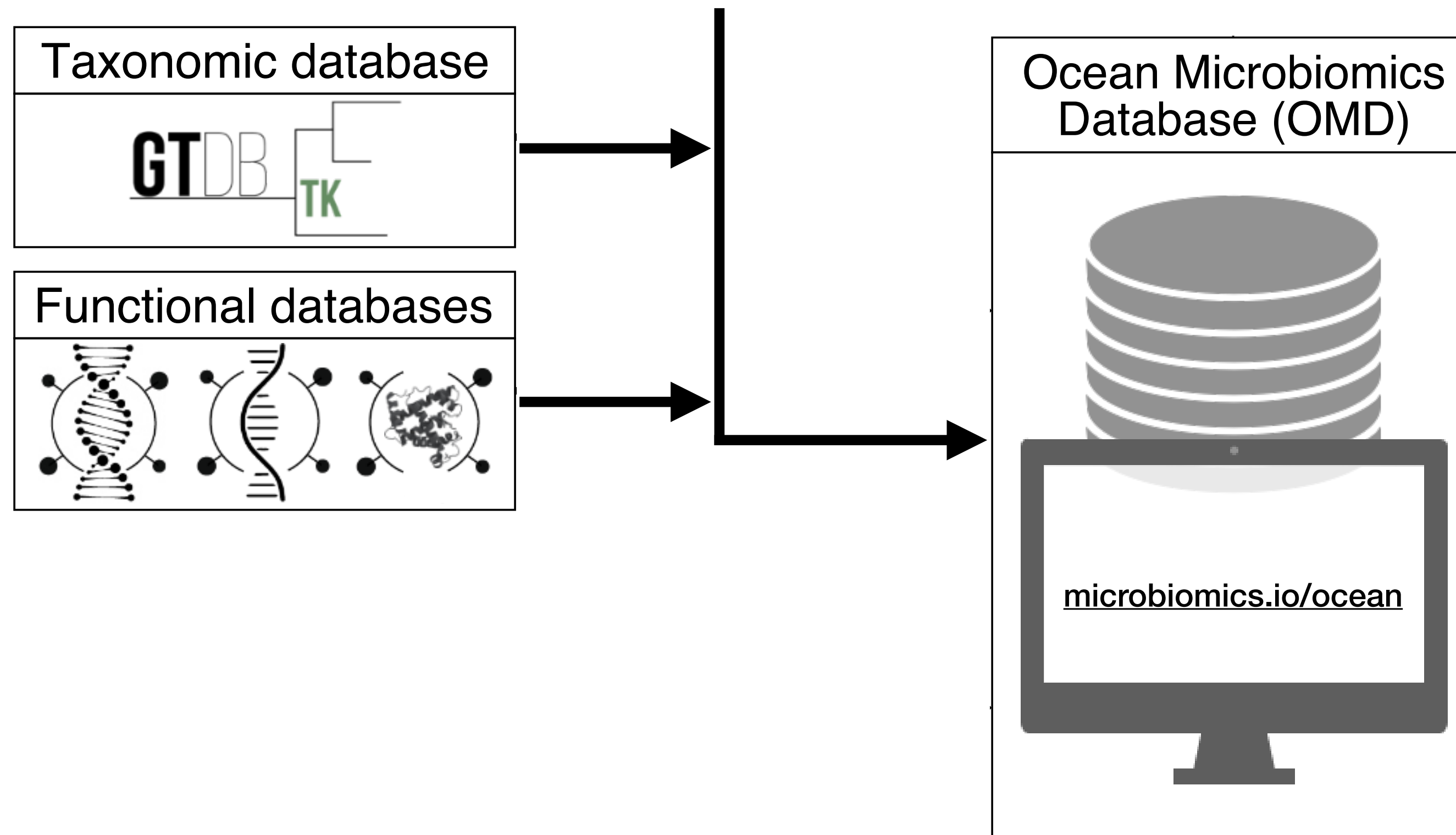
~26,000 Metagenome-Assembled Genomes (MAGs)

# Integrating cultivation dependent and independent methods

- ➔ Newly reconstructed MAGs (~26,000)
- ➔ Manually curated MAGs (~1,000)
- ➔ Single amplified genomes (SAGs) (~6,000)
- ➔ Reference genomes from isolates (~2,000)

# Establishing a rich ocean microbiome resource

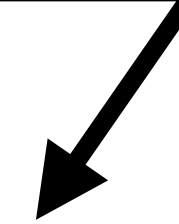
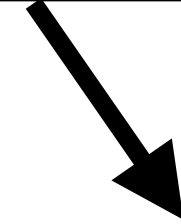
- ➔ Newly reconstructed MAGs (~26,000)
- ➔ Manually curated MAGs (~1,000)
- ➔ Single amplified genomes (SAGs) (~6,000)
- ➔ Reference genomes from isolates (~2,000)



# Improved representation of the ocean microbiome

Ocean Microbiomics DB

Metagenomes

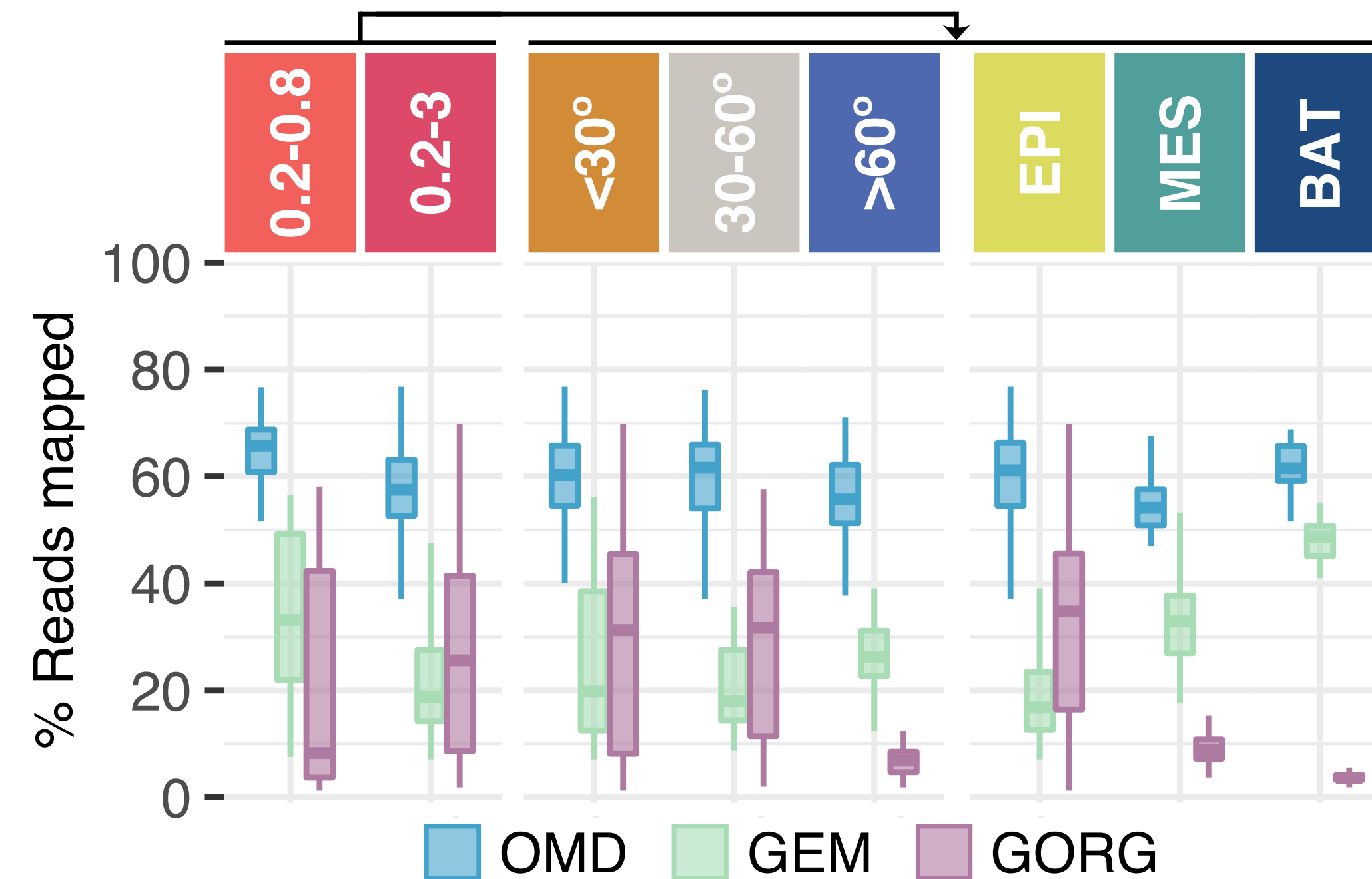
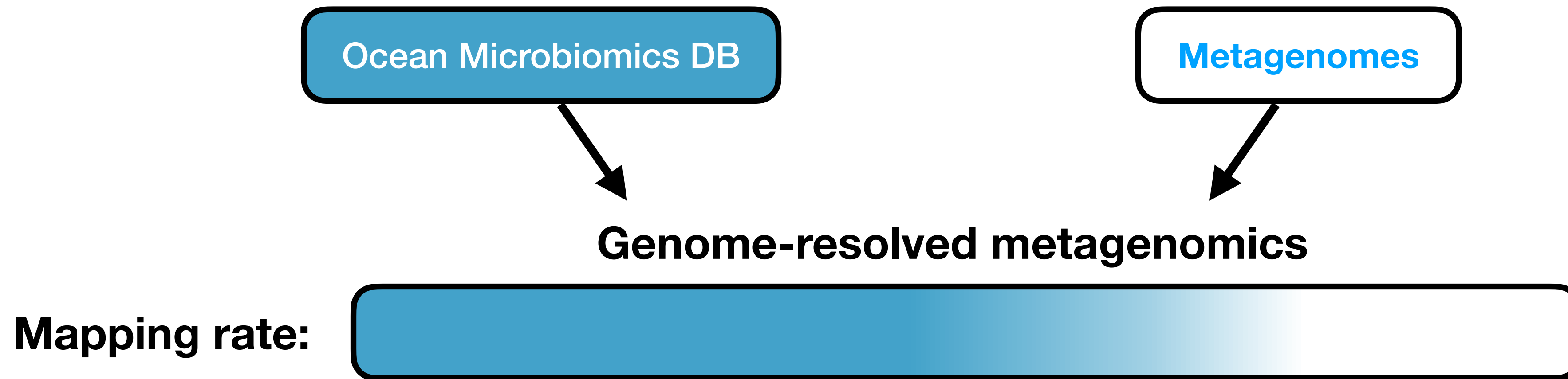


Genome-resolved metagenomics

Mapping rate:

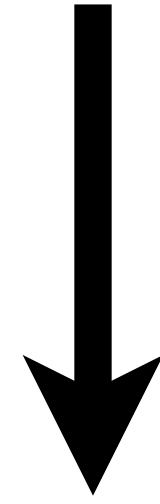


# Improved representation of the ocean microbiome



# **Providing access to its biosynthetic potential**

**~35k genomes**



**~40k Biosynthetic Gene Clusters (BGCs)**

**➔ How unique is it compared to currently sequenced microbes?**

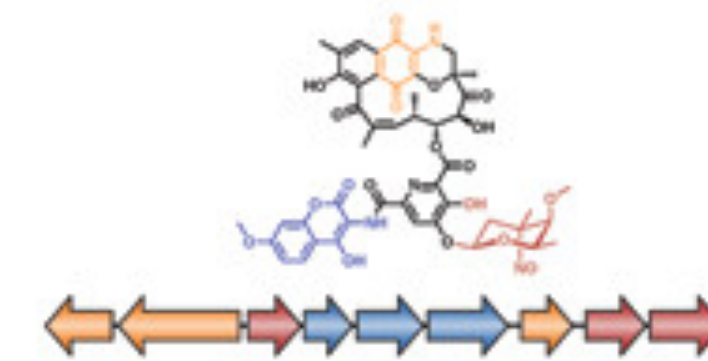
# Compared to sequenced and characterised biosynthetic potential



RefSeq  
~200,000 genomes



1.2 M BGCs



Characterised  
biosynthetic pathways



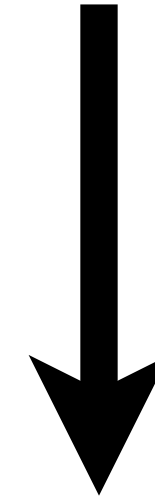
2,000 BGCs



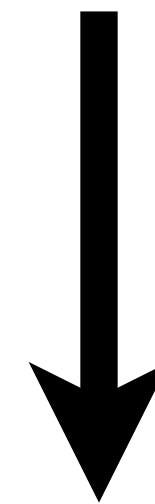


# Grouping BGCs into relevant units

~35k genomes

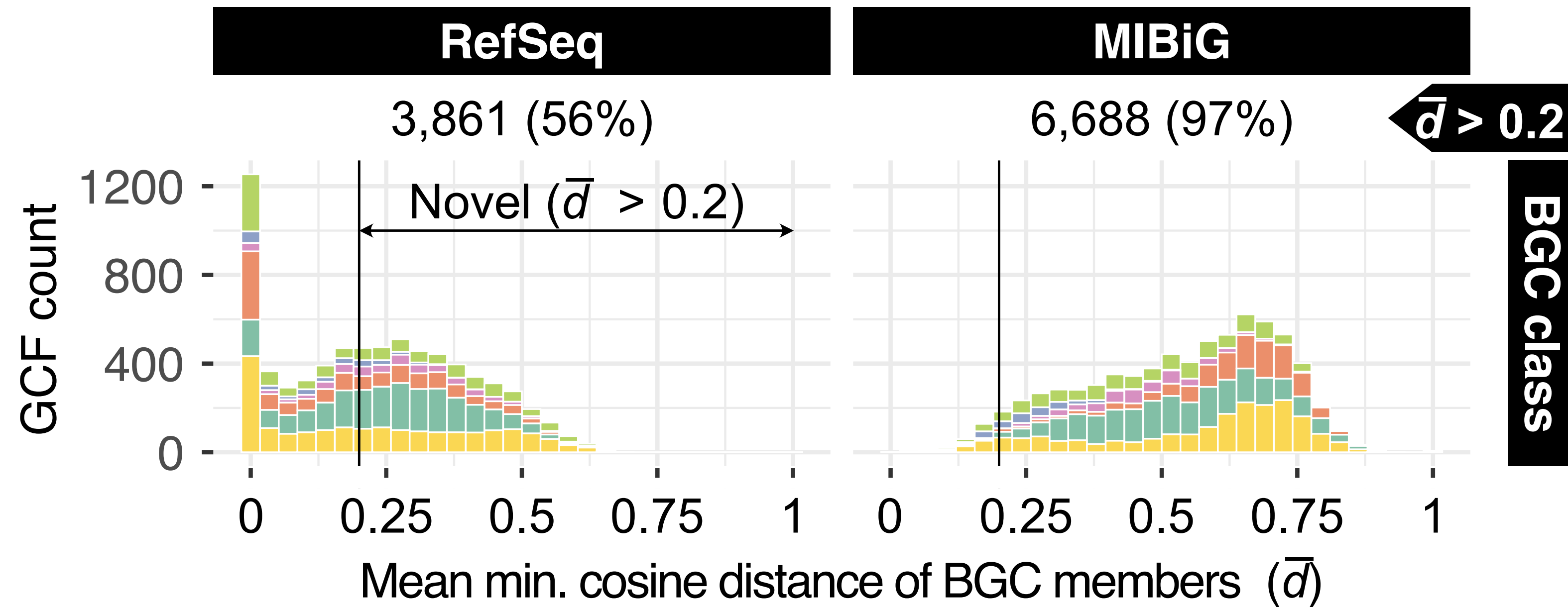


~40k Biosynthetic Gene Clusters (BGCs)



~7k Gene Cluster Families (GCFs)

# With large potential for new compounds

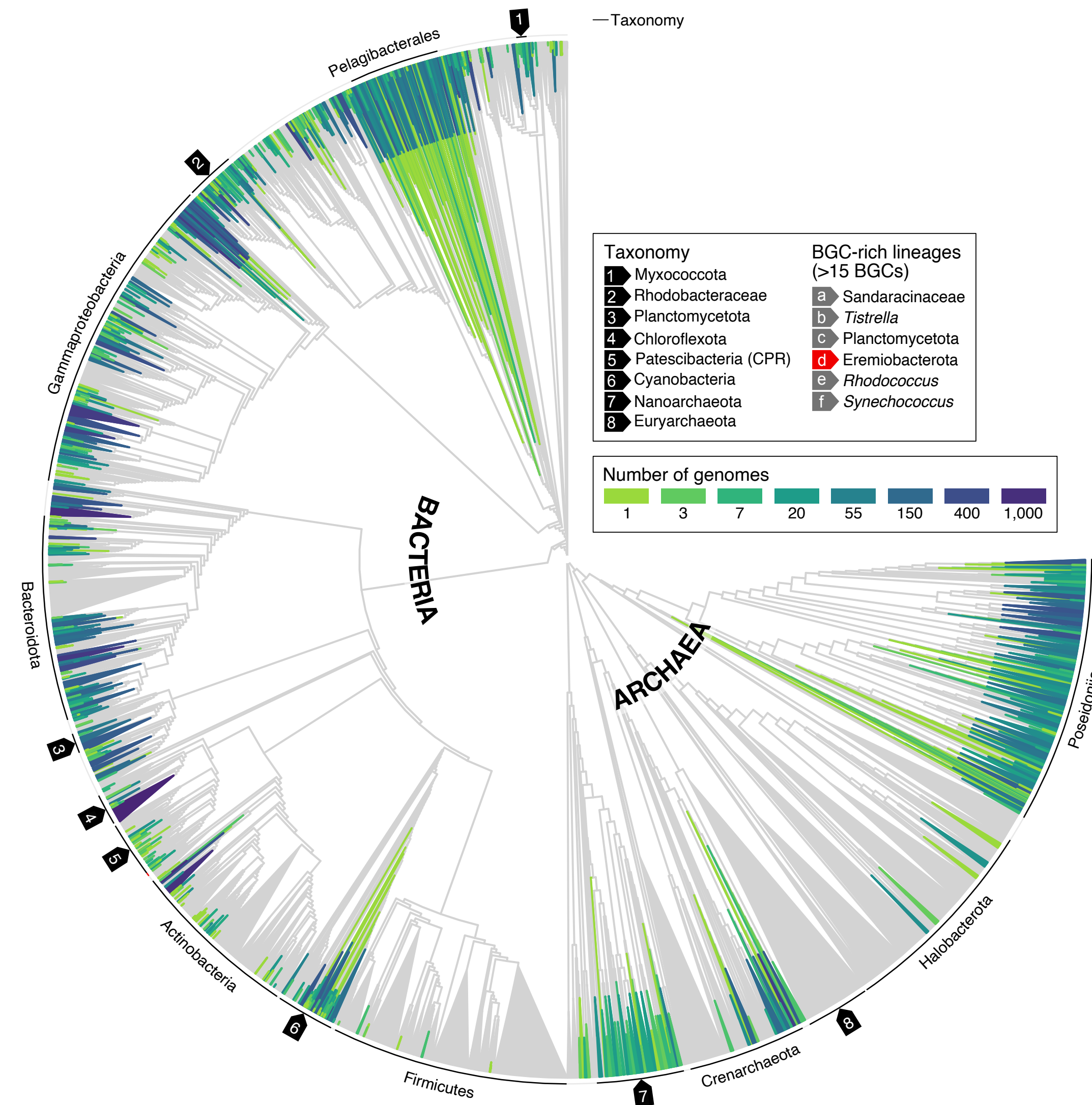


## BGC class

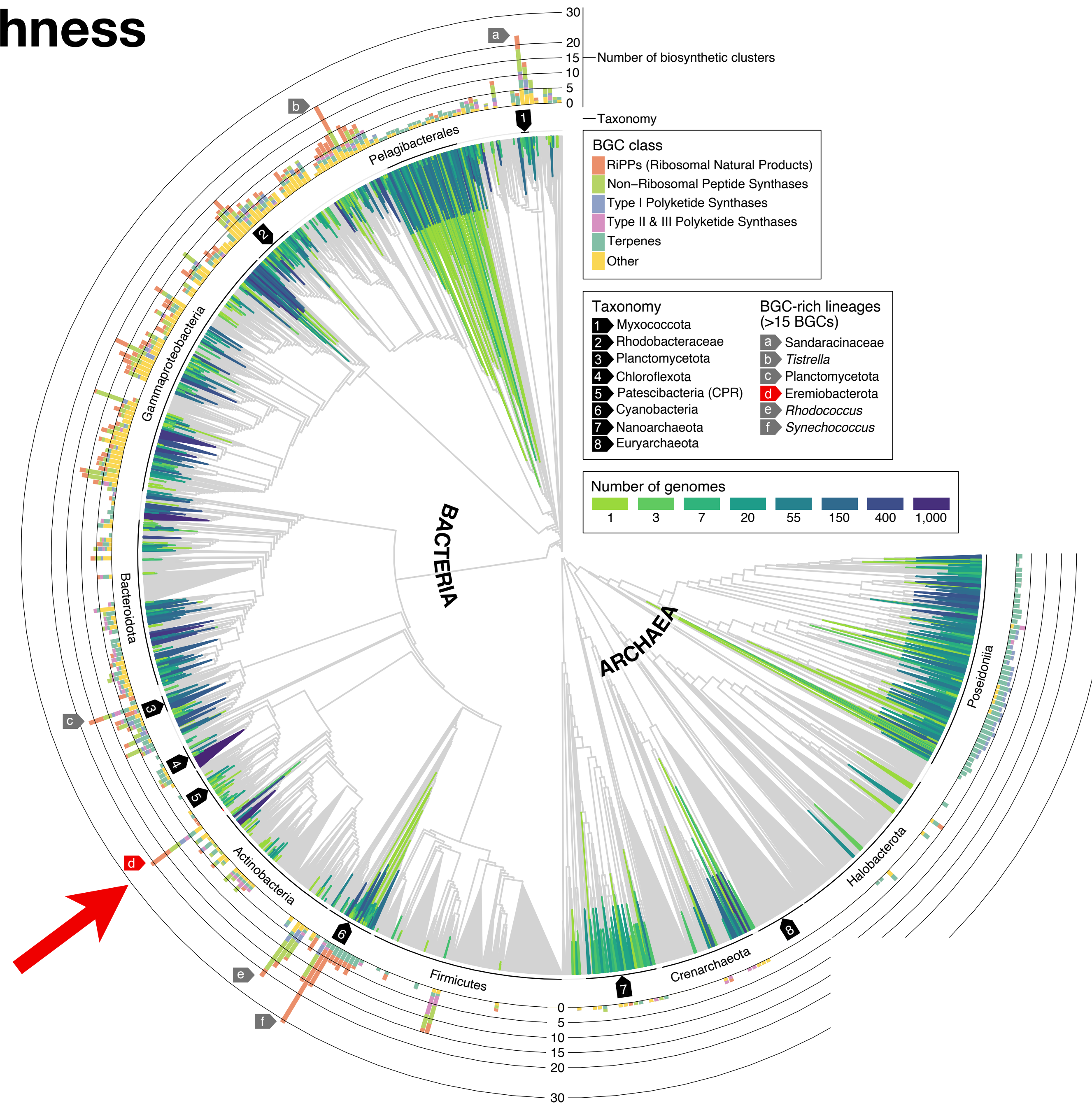
- RiPPs (Ribosomal Natural Products)
- Non-Ribosomal Peptide Synthases
- Type I Polyketide Synthases
- Type II & III Polyketide Synthases
- Terpenes
- Other

**➔ Are there BGC-rich microbial lineages to be discovered in the ocean?**

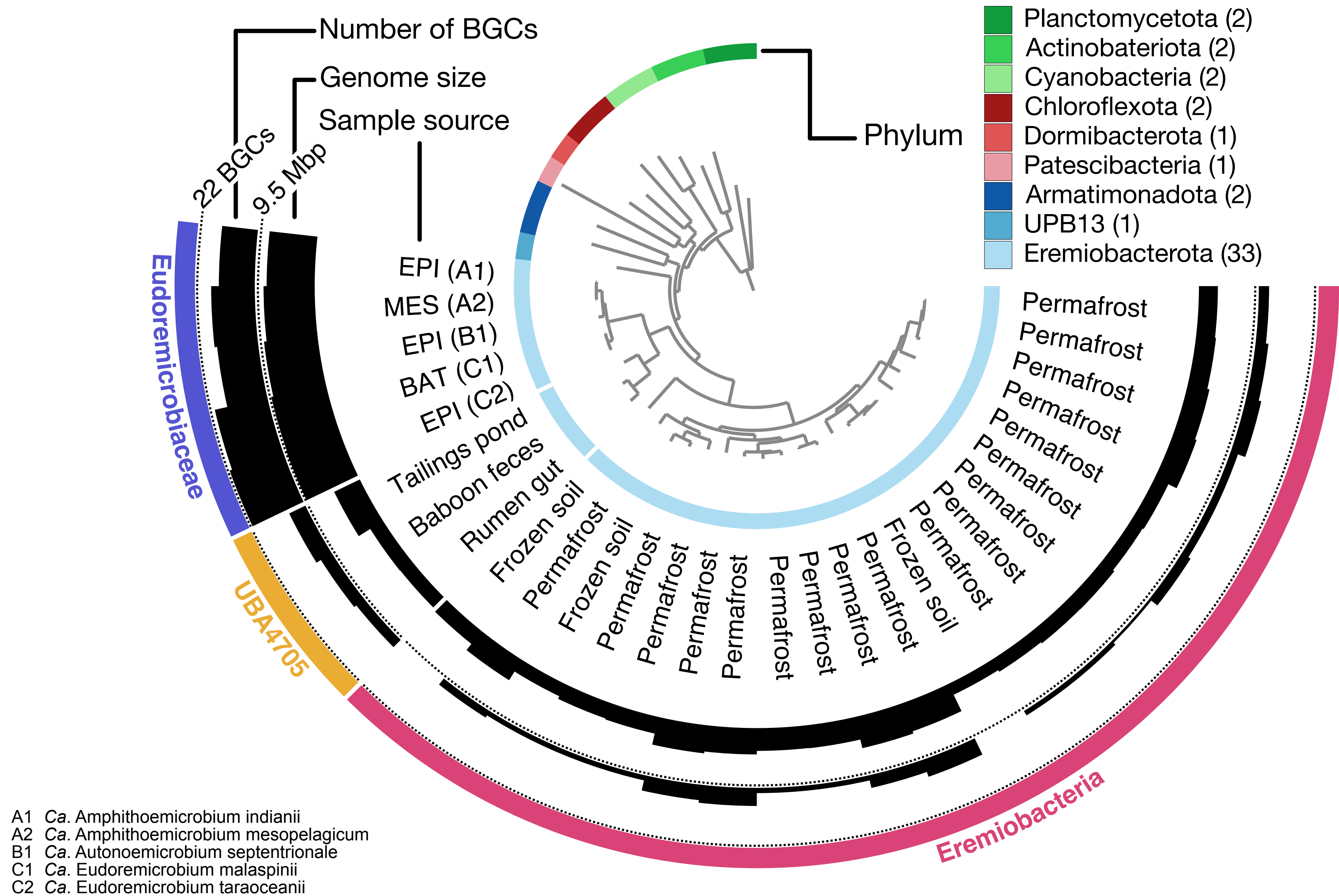
# Phylogenomic distribution of the ocean biosynthetic potential



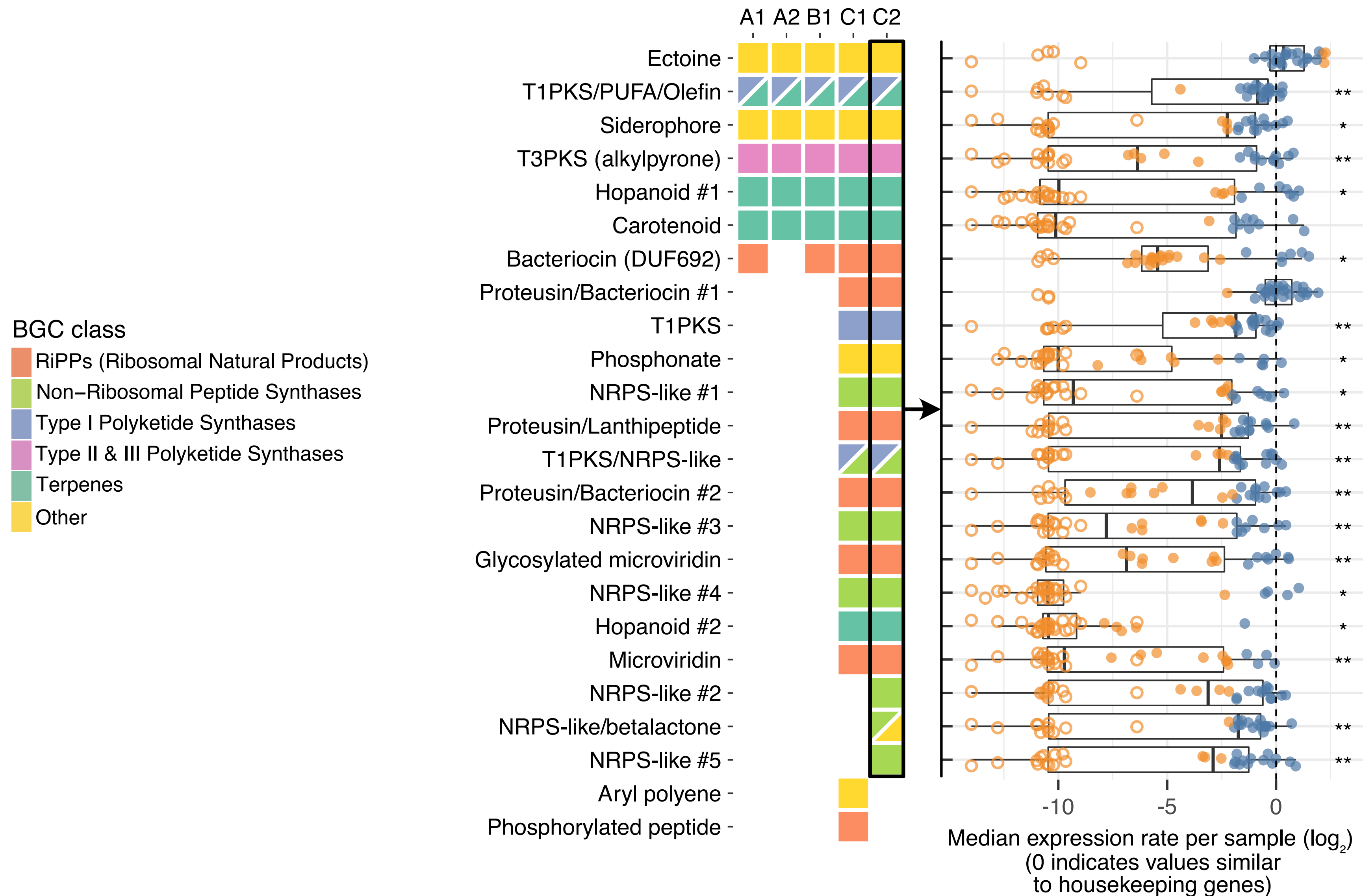
# Eremiobacterota, uncultivated phylum with unsuspected BGC richness



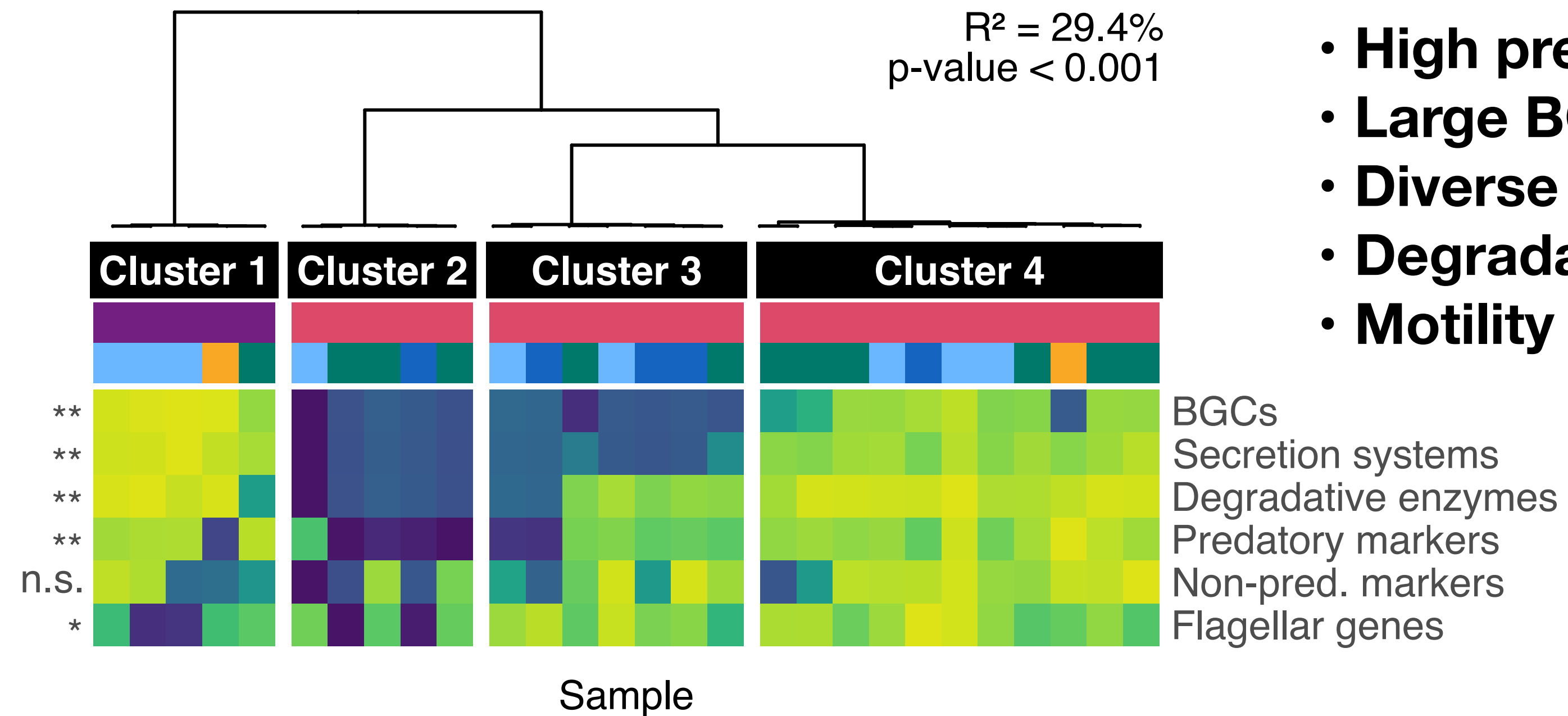
# A novel marine Eremiobacterota lineage: The candidate family Eudoremiobiaceae



# The newly identified bacterial family has a diverse and actively expressed BGC repertoire



# A biosynthetic potential that may supports a putative predatory lifestyle



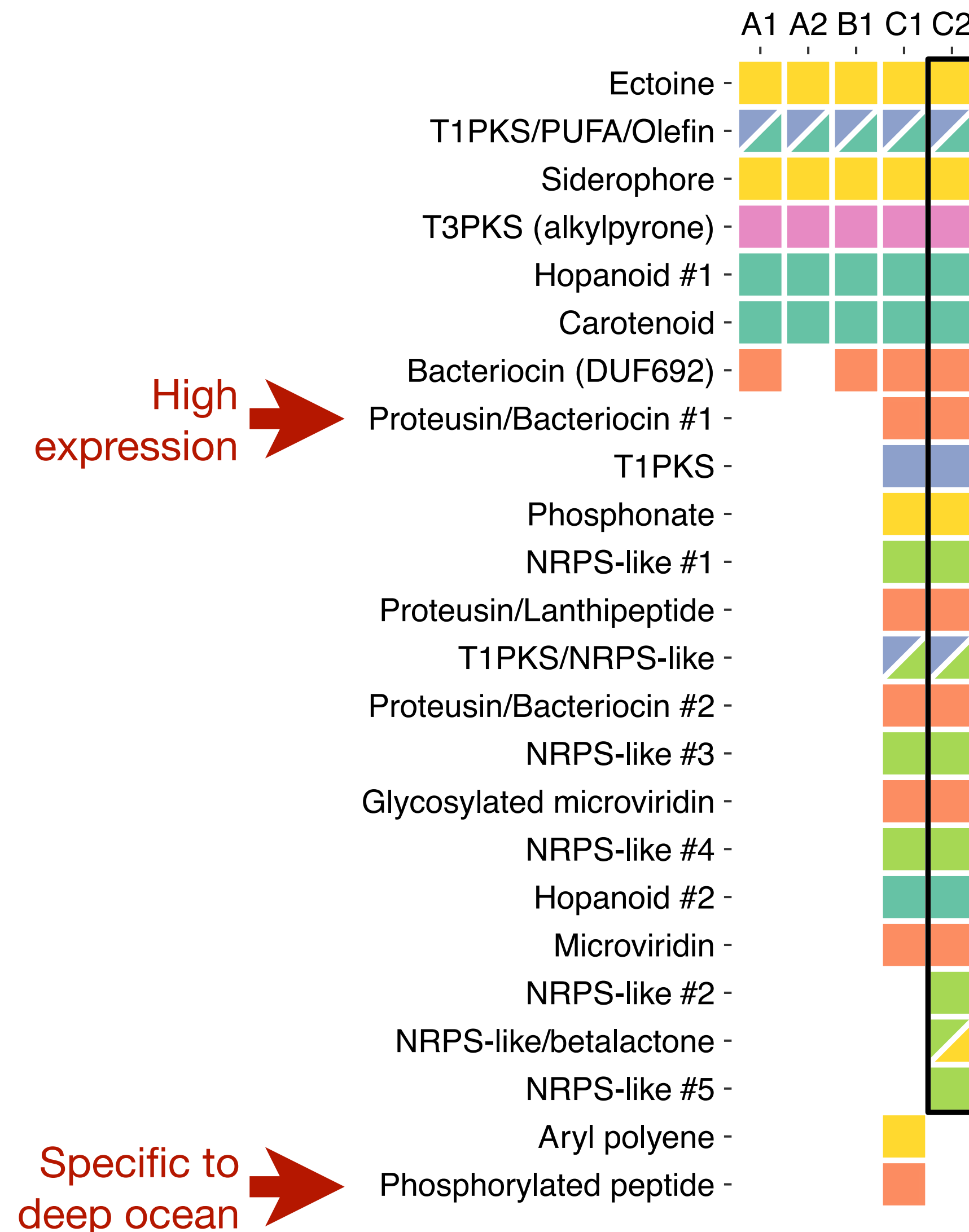
- **High predatory scores**
- **Large BGC diversity**
- **Diverse secretion systems**
- **Degradative enzymes**
- **Motility**





**➔ Is this computational approach sufficiently powerful to predict new enzymology and natural products?**

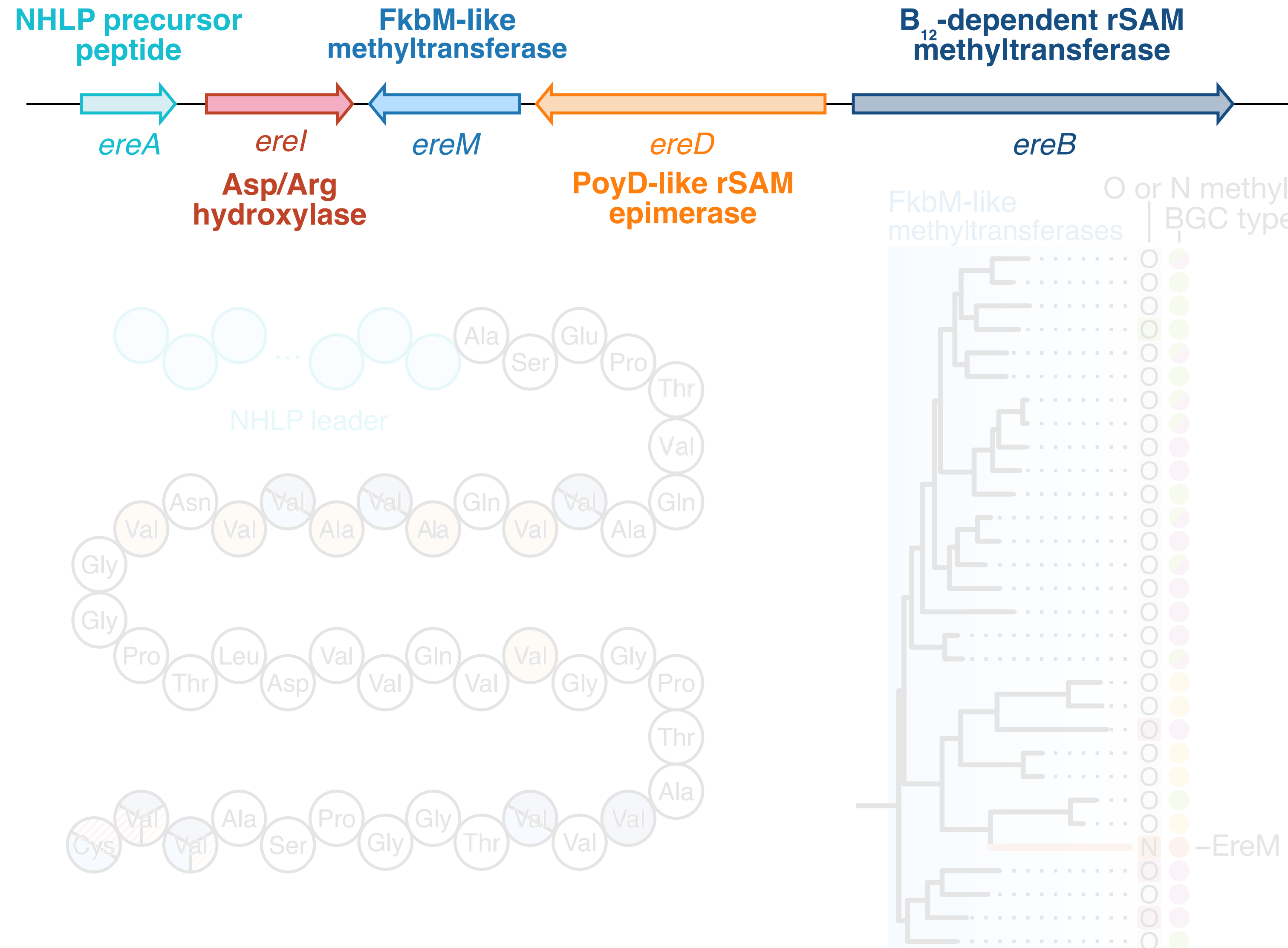
# Probing Eremiobacterota's biosynthetic potential for new enzymes and natural products



Characterising two **predicted novel** Ribosomal BGCs using:

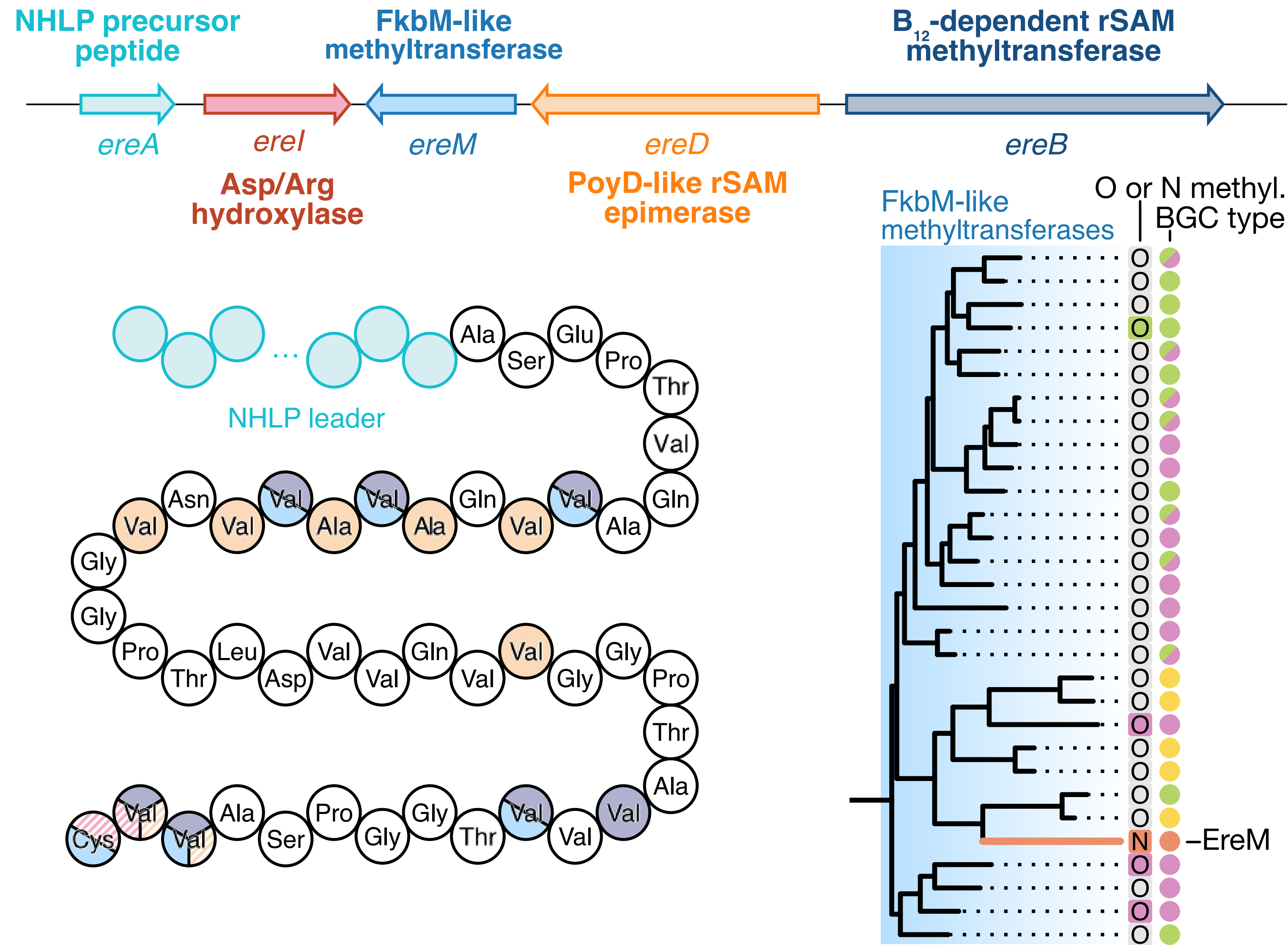
- Non-standard heterologous expression
- Tandem Mass Spectrometry (MS/MS)
- Isotope labelling
- Nuclear Magnetic Resonance (NMR)

# Intricate proteusin cluster reveals new enzymology

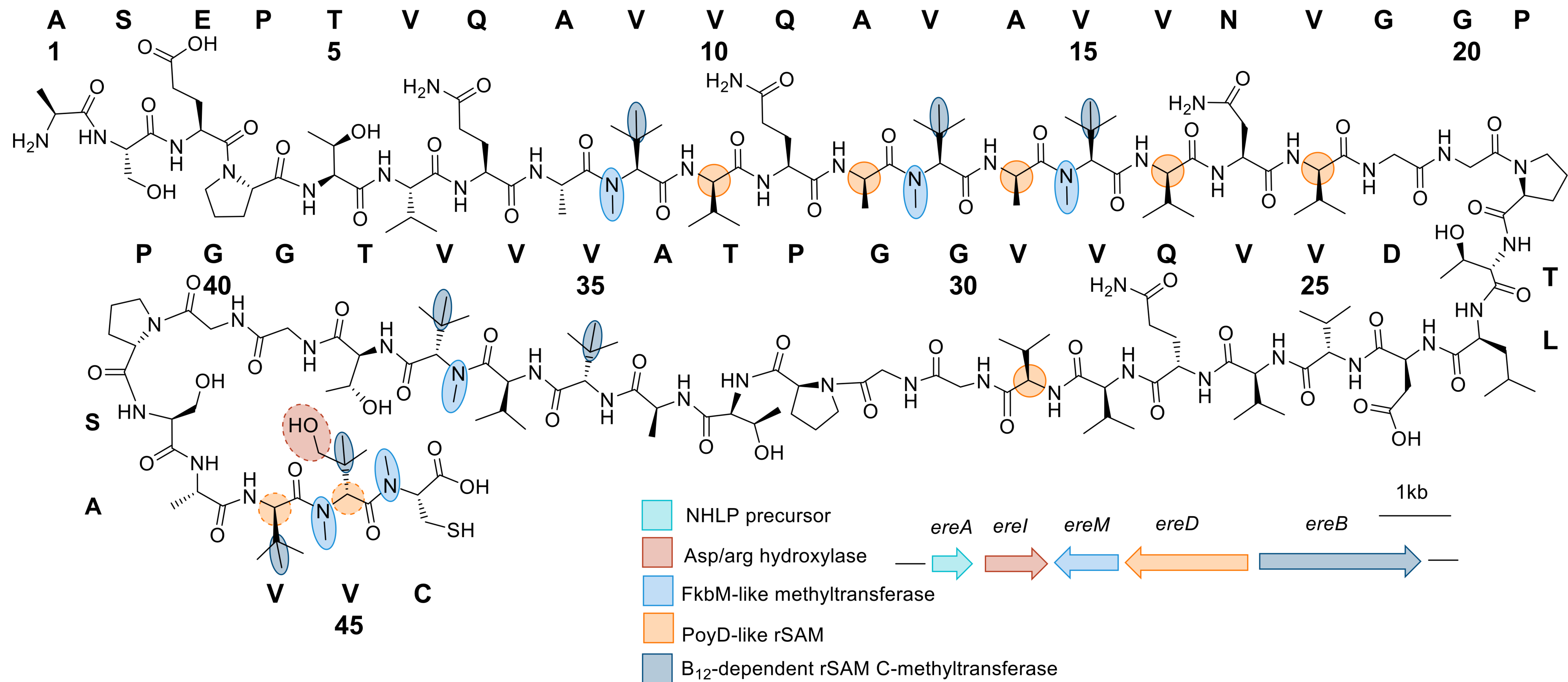




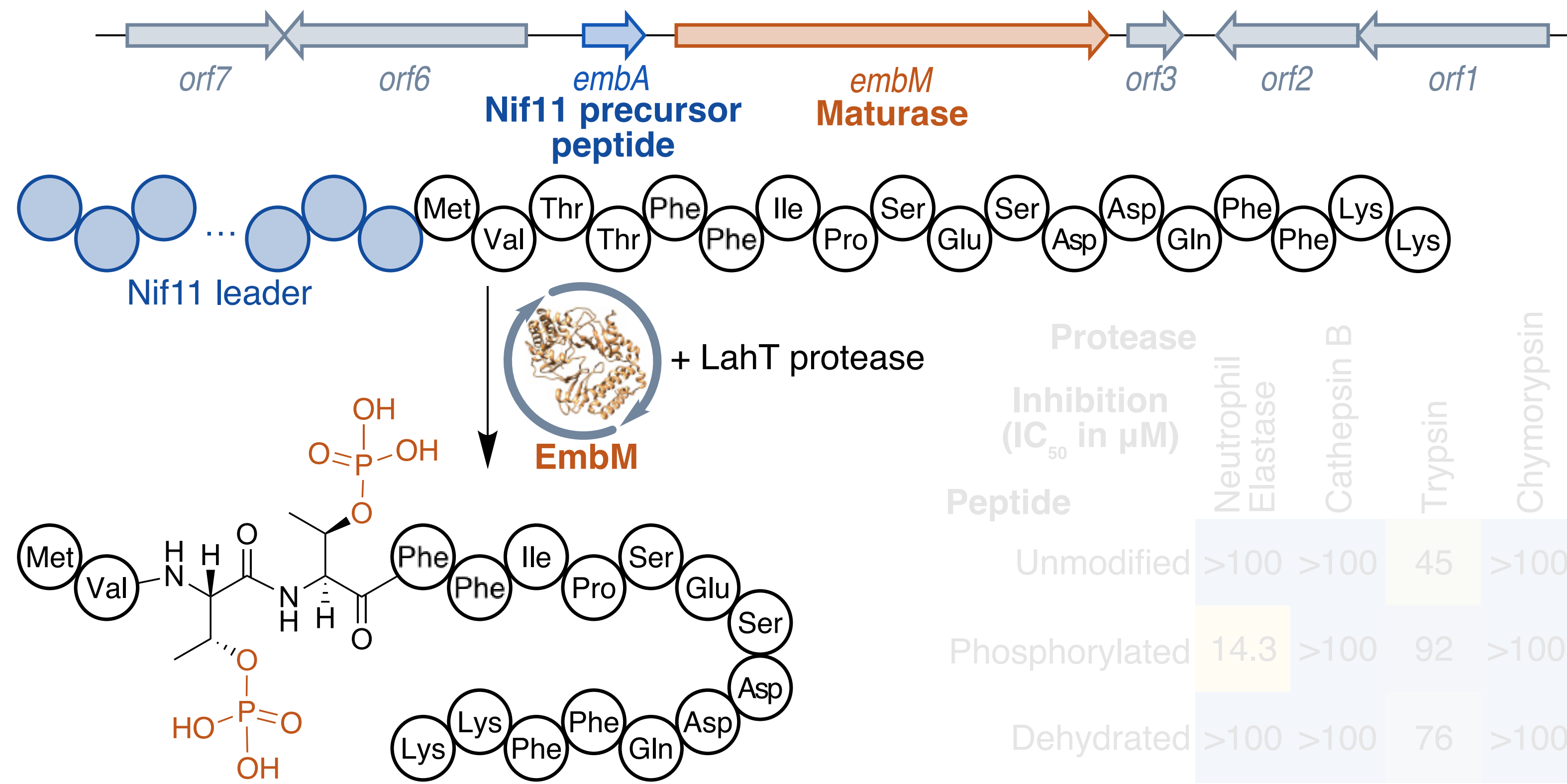
# Intricate proteusin cluster reveals new enzymology



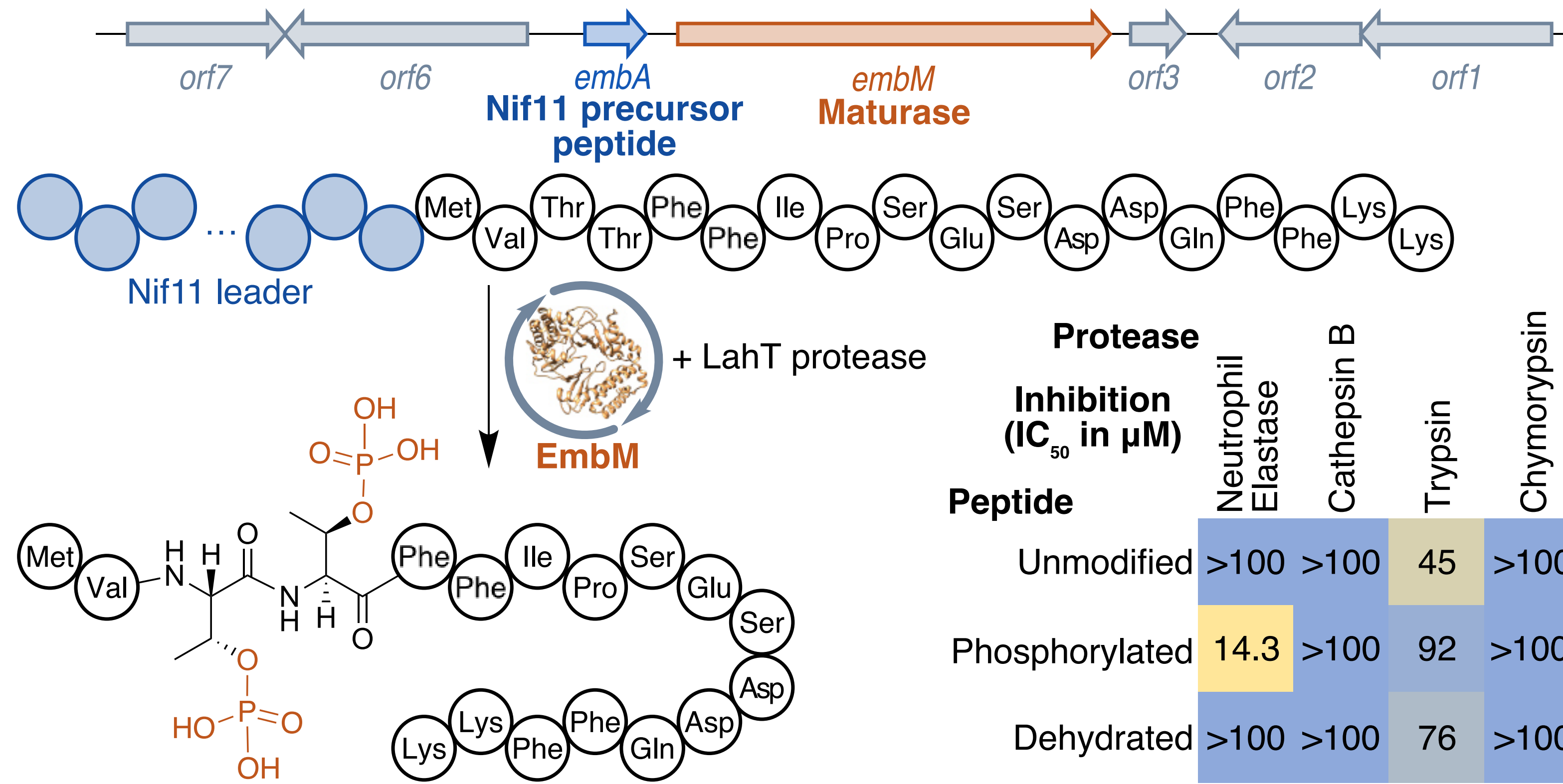
# A predicted O-methyltransferase with amide-N-methylation activity



# A new RiPP cluster with phosphorylation as sole modification



# With Human Elastase inhibition



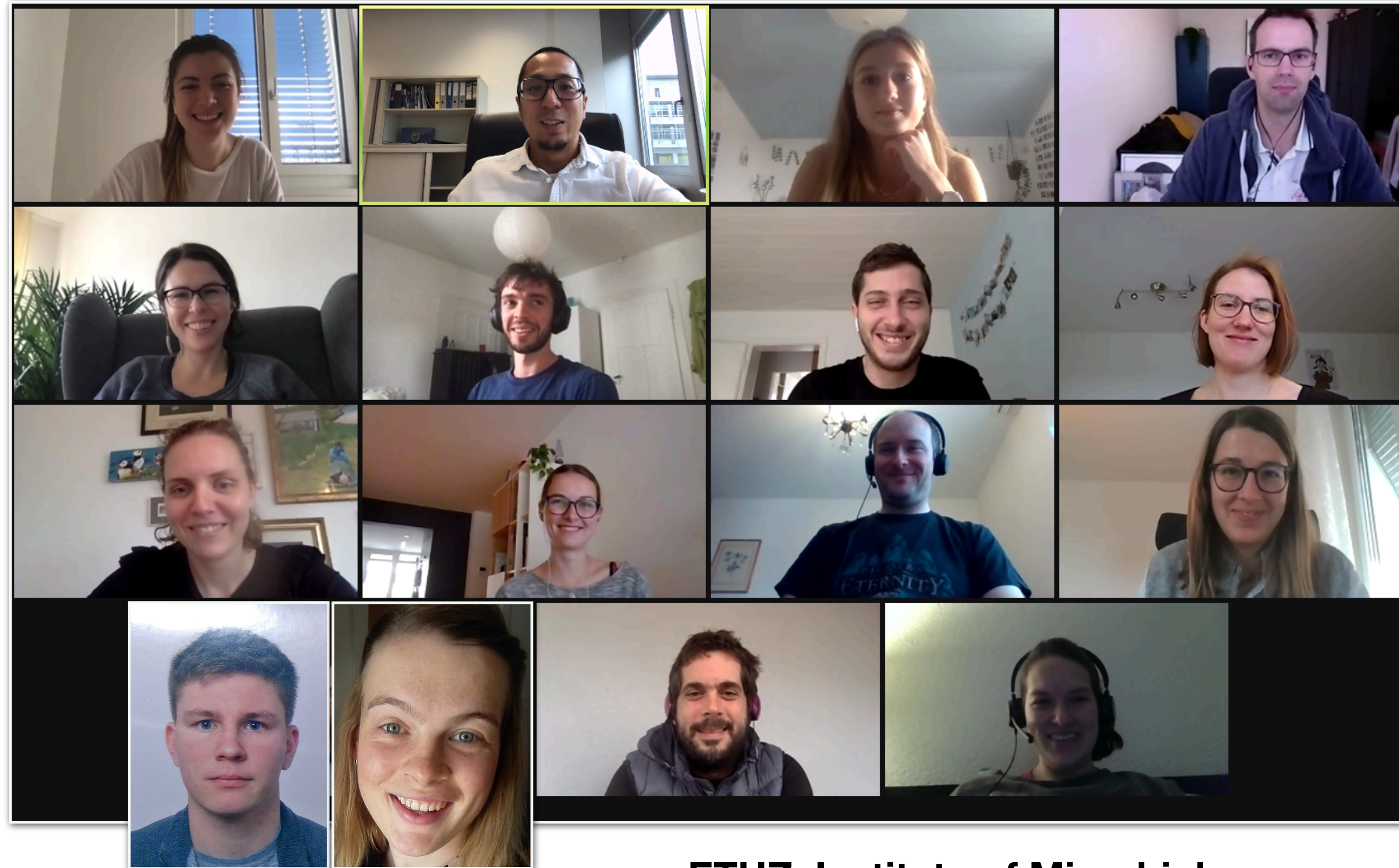


# Conclusions

- **Genome-resolved microbiomics** as a mean to **explore environmental microbiomes** and **discover novel enzymology** and **natural products**
- This approach provides **evolutionary and ecological context** to the **biosynthetic potential**
- Bioinformatics-guided **experimental characterisation** is **necessary** and can still lead to **unpredicted discoveries**

# Thank you for your attention

Sunagawa Lab



**ETH** zürich



Swiss Institute of  
Bioinformatics



FONDS NATIONAL SUISSE  
SCHWEIZERISCHER NATIONALFONDS  
FONDO NAZIONALE SVIZZERO  
SWISS NATIONAL SCIENCE FOUNDATION



Helmut Horten Stiftung

**TARA**  
**OCEANS**

Fondation  
**taraocéan**  
explorer et partager

**ETHZ, Institute of Microbiology:**



Jörn Piel



Serina Robinson



Clarissa Forneris



Florian Hubrich



**Questions?**

Image: François Aurat