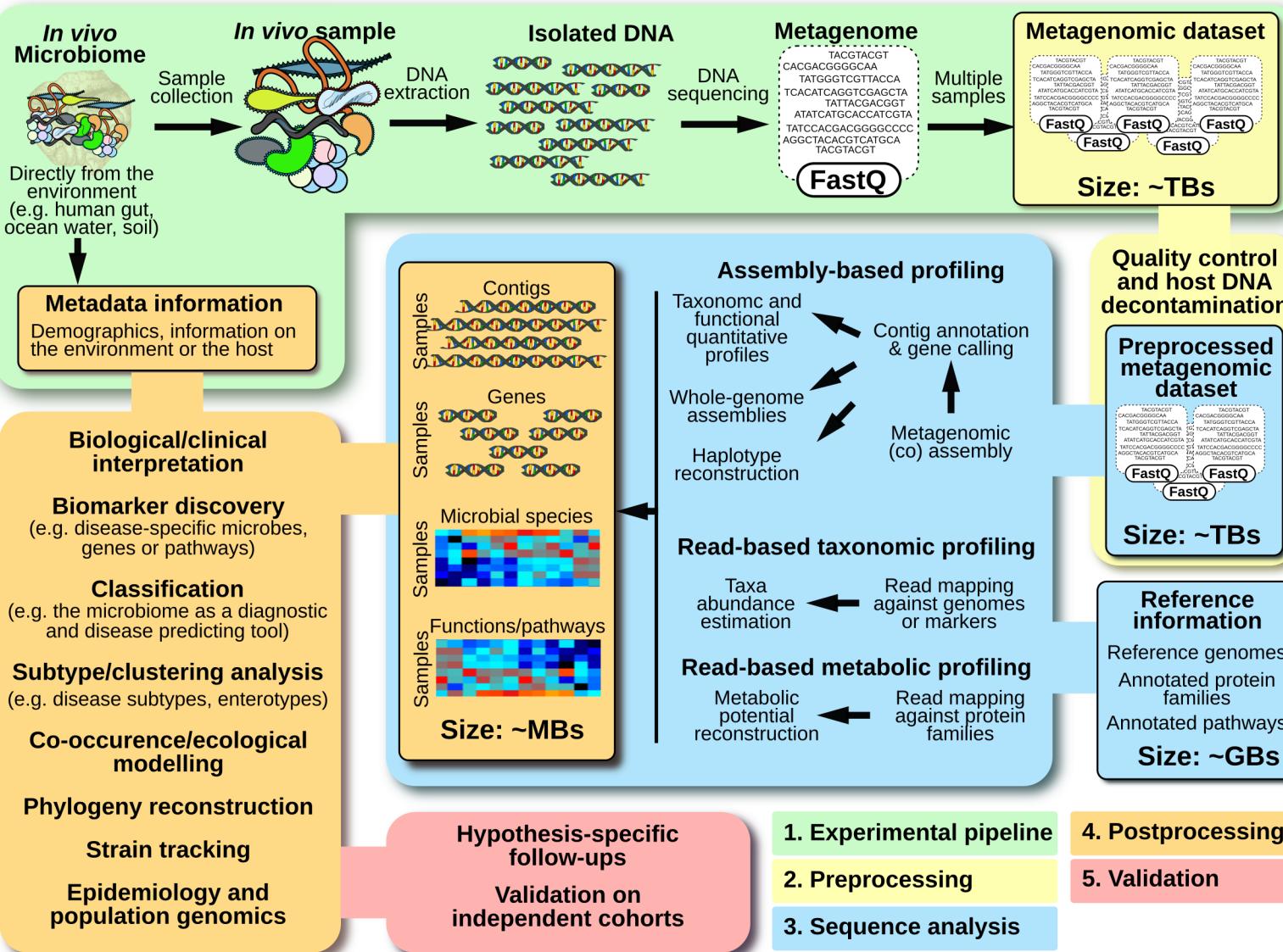


Species and strain resolution from metagenomic coassemblies

Dr Christopher Quince
University of Warwick

Introduction

- Metagenomics provides a culture free method to survey an entire microbial community
- Contig binning from coassemblies has allowed *de novo* species extraction:
 - Candidate Phyla radiation ([Brown et al. Nature 2015](#))
 - Comammox ([Maartje et al. 2015](#))
- Reference based methods exist for high resolution strain profiling ([Pathoscope – Hong et al. 2014](#); [SPARSE - Zhou et al. 2018](#))
- Still cannot obtain ‘population genomics’ data *de novo* from metagenomes

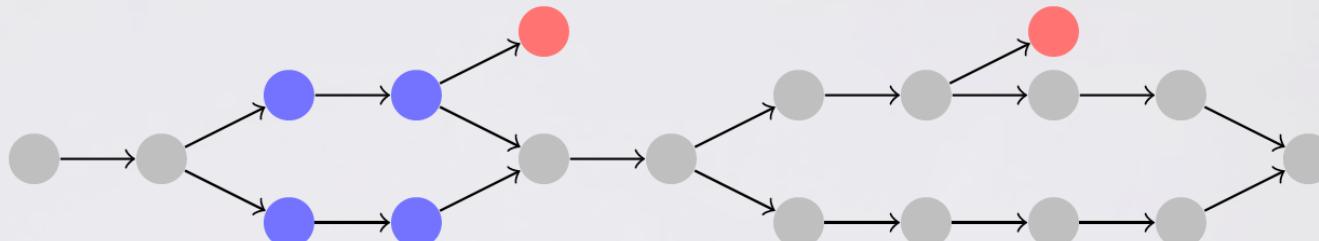


Overview

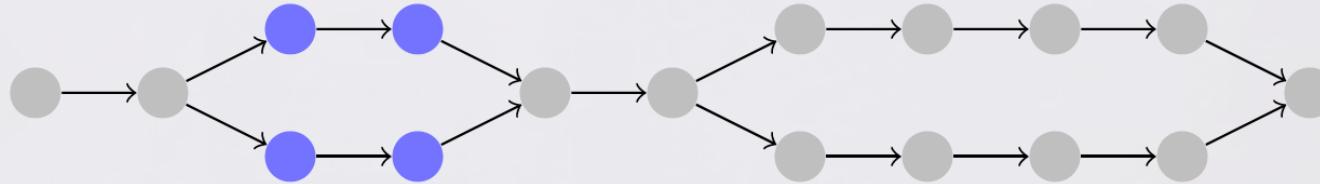
- Introduction: Why is metagenomics coassembly hard?
- Part I: CONCOCT, contig binning
- Part II: DESMAN, resolving intra-MAG diversity

How a metagenome assembler generally works

- 1) de Bruijn **graph** construction



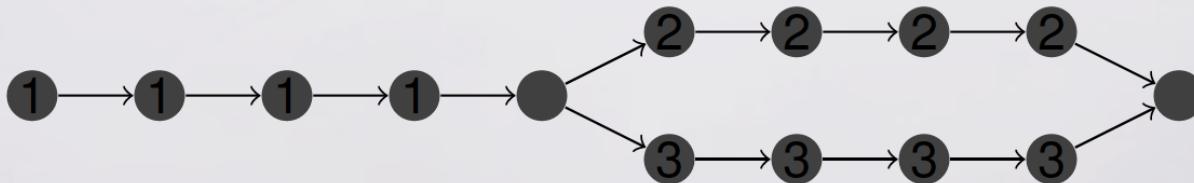
- 2) Likely **sequencing errors** are removed.



- 3) Variations (e.g. SNPs, similar repetitions) are removed.

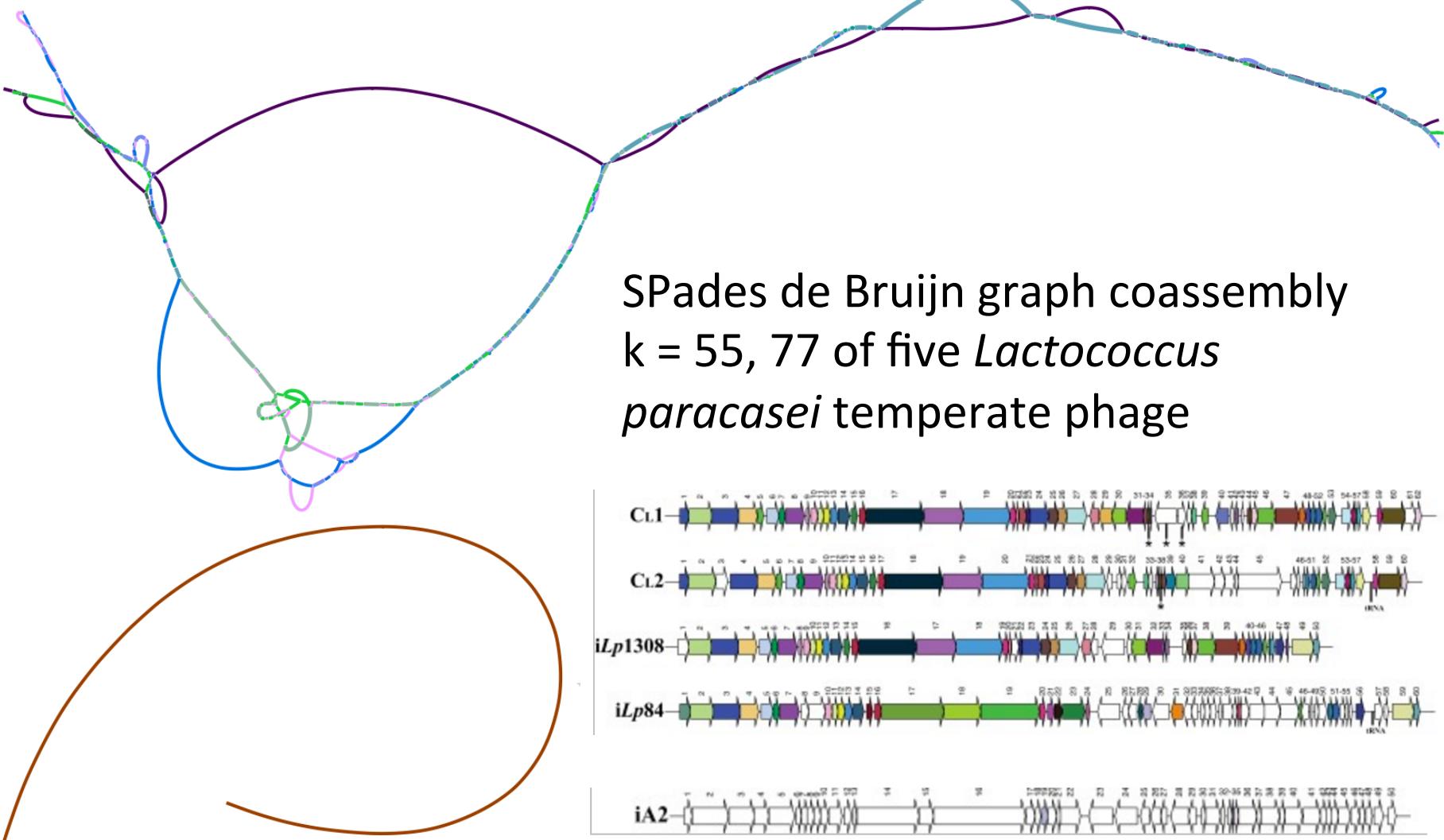
→ **Skipped in Strains #1**

- 4) **Simple paths** (i.e. contigs) are returned.

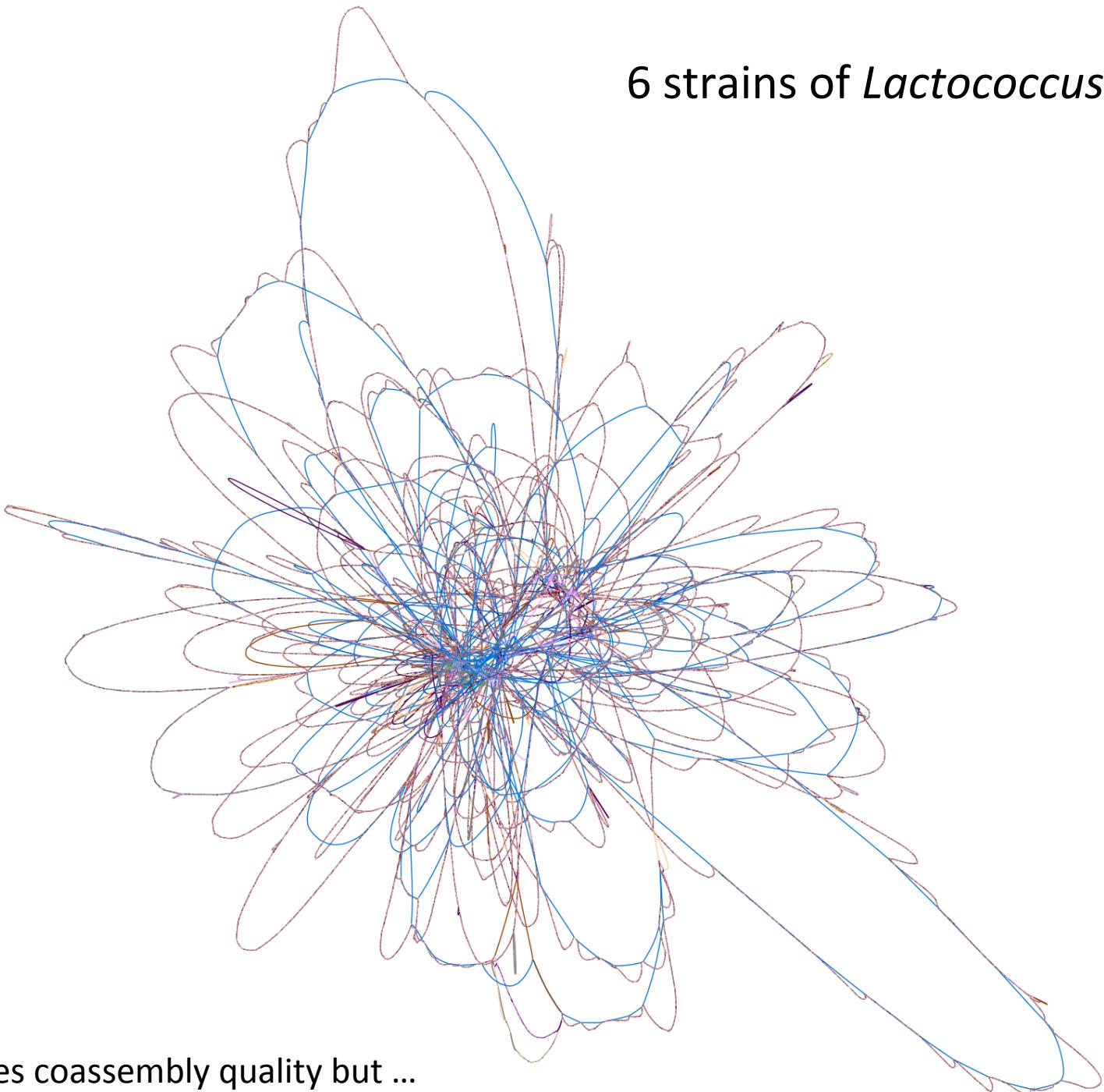


- 5) Extra steps: repeat-resolving, scaffolding (**not done in Minia**)

Strain variation leads to complex assembly graphs...



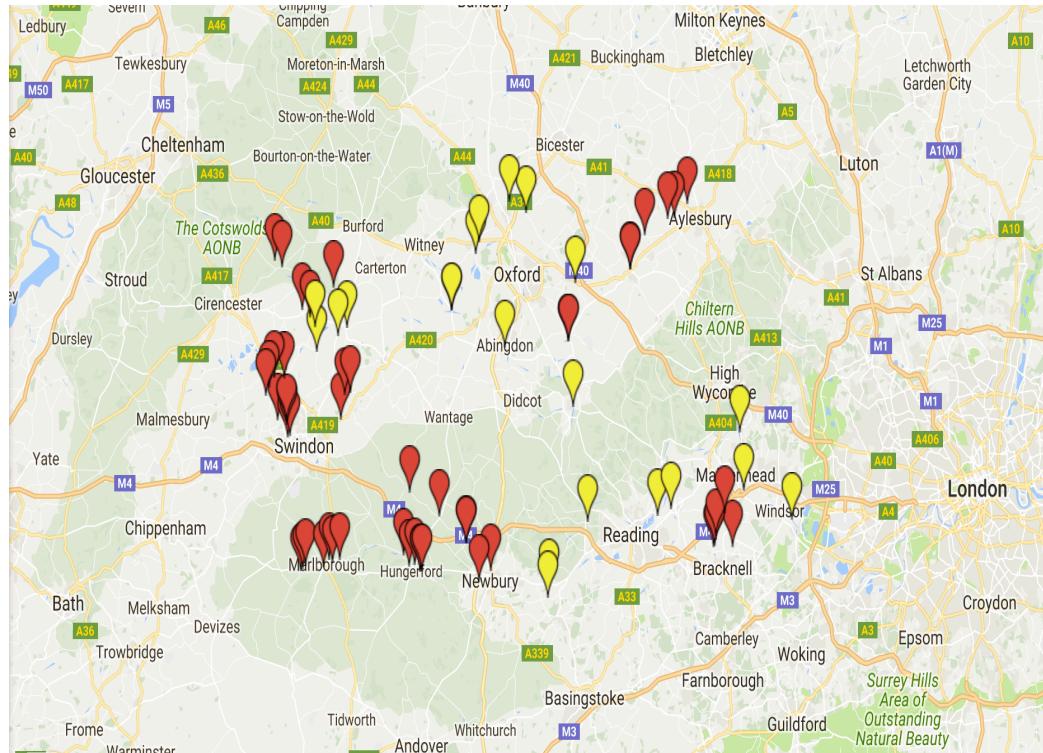
6 strains of *Lactococcus lactis*



This reduces coassembly quality but ...

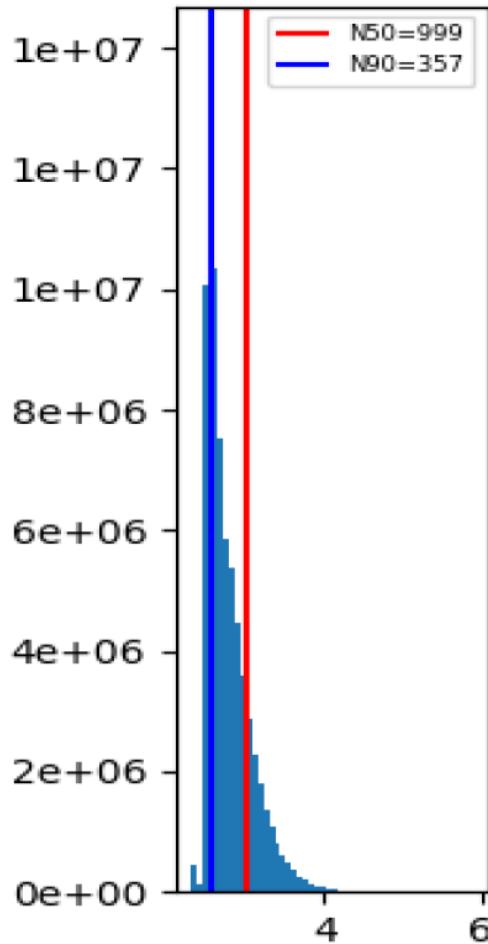
Thames AMRG survey

- 20 sites located along 7 tributaries of Thames catchment were chosen for sampling across three seasons and with up to three biological replicates per site

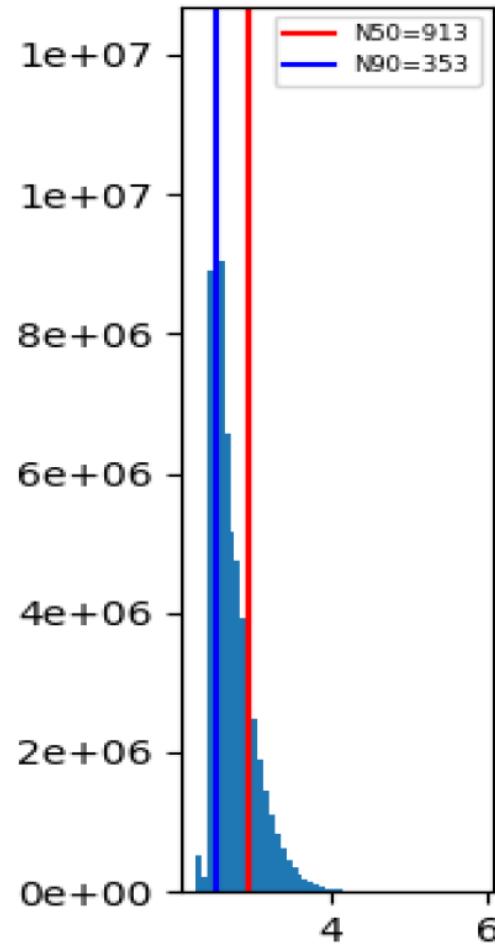


- Total 171 samples
- 30 – 114.7 million 2X150 bp reads per sample
- 2.3 Tbp

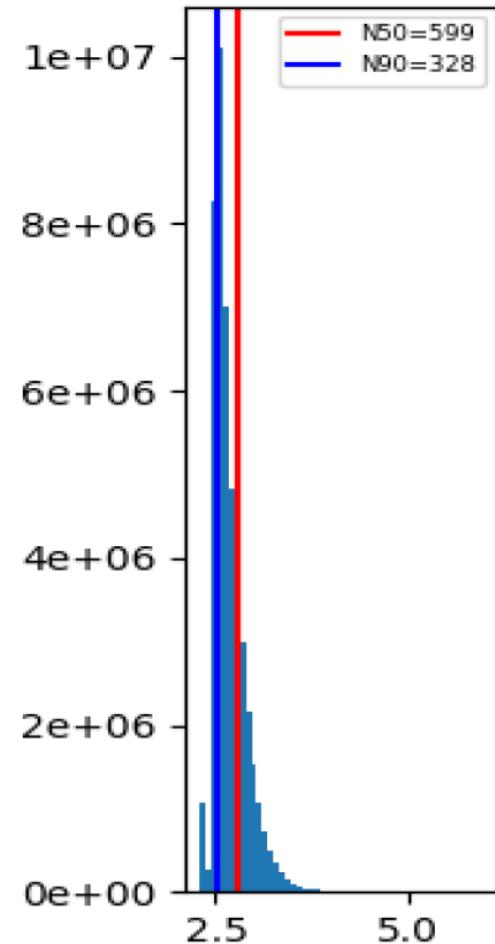
Global Assembly
nb contig=74,377,142
Total length=59,341,453,108



All Rivers
nb contig=64,287,670
Total length=48,753,041,062



All Samples
nb contig=45,205,443
Total length=26,107,237,118



Performed assemblies with megahit

PART I: CONCOCT

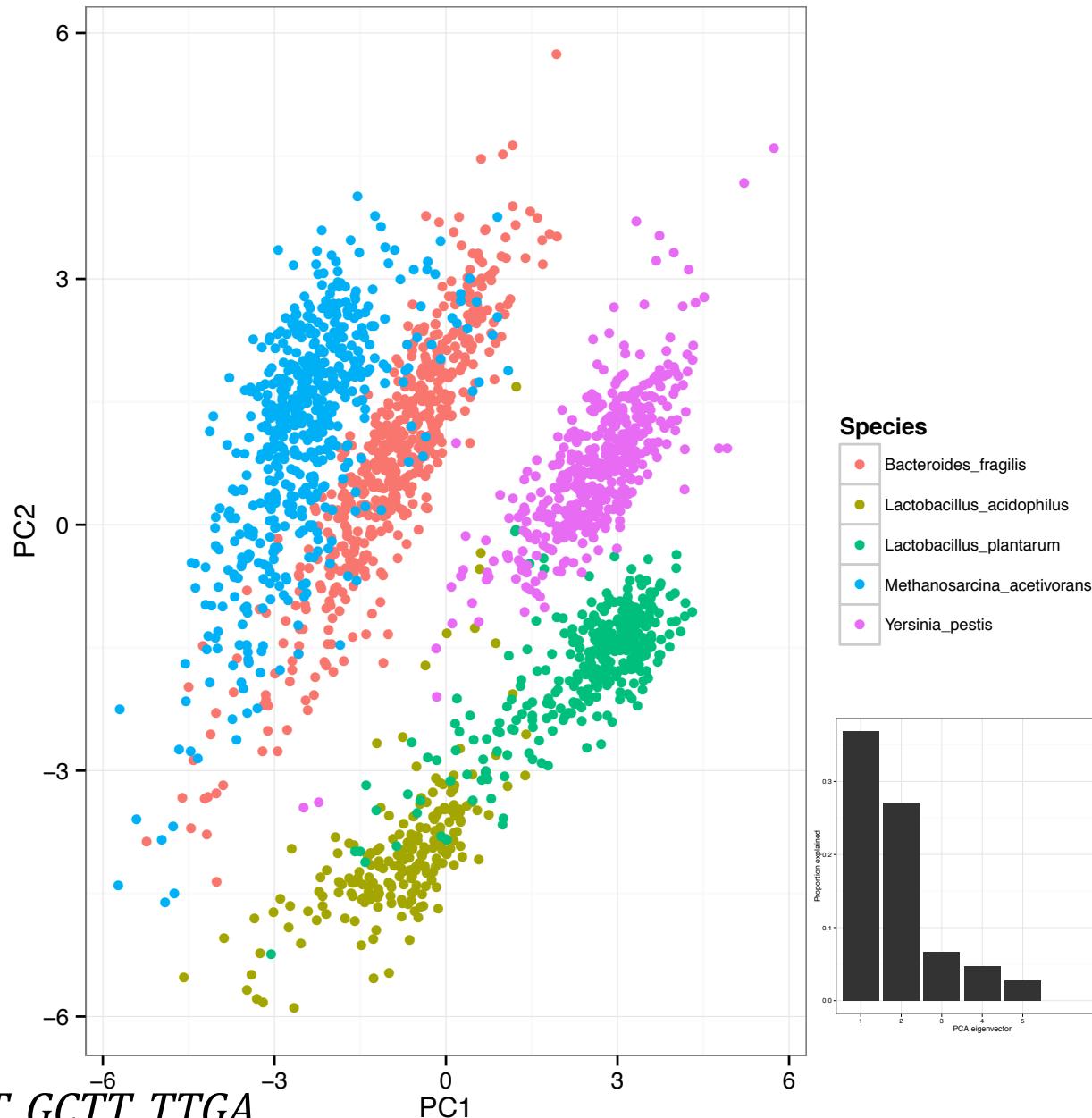
Contig binning

Contig clustering by sequence composition

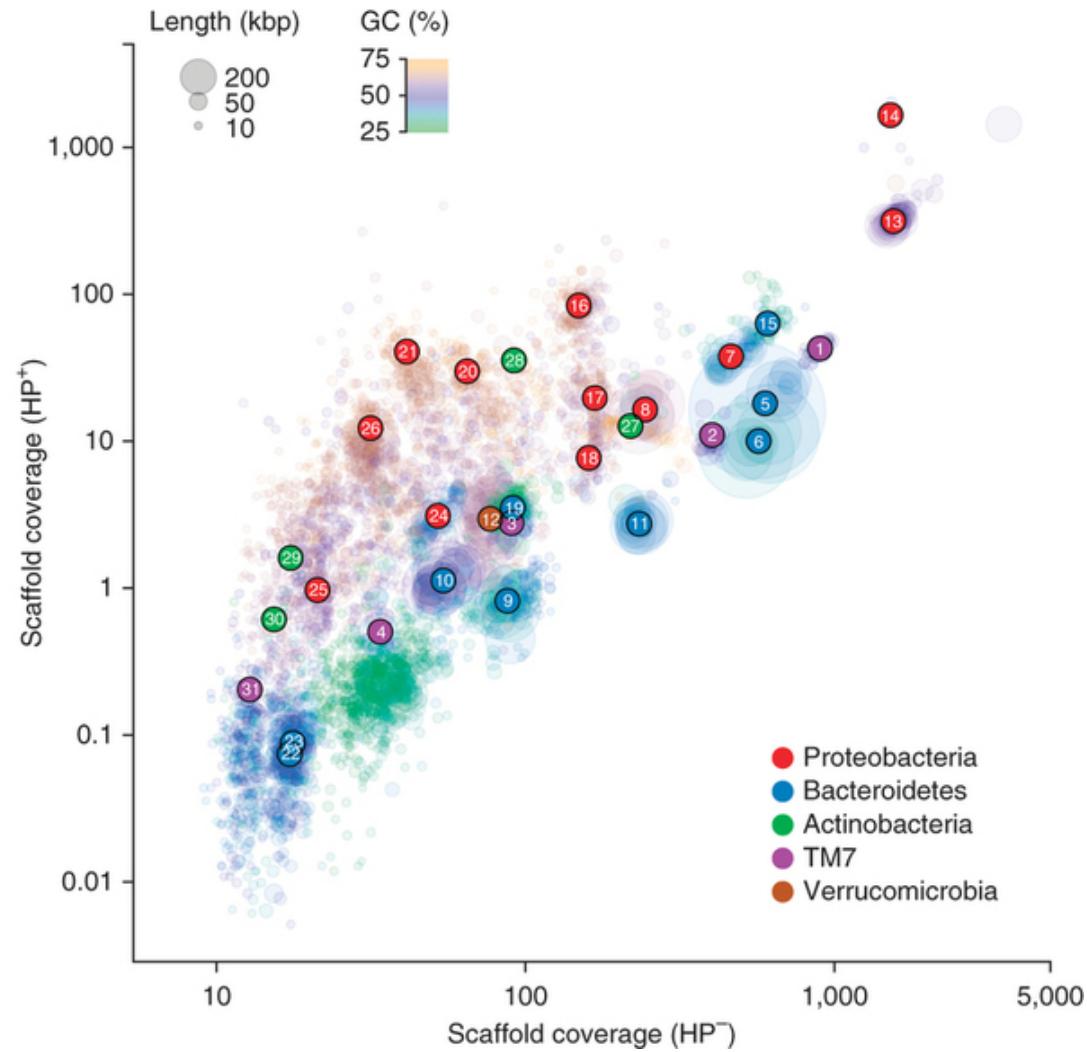
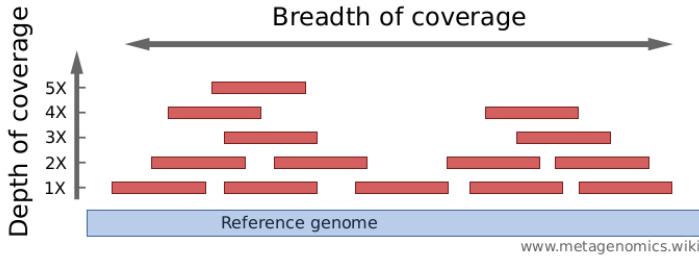
Key to contig binning is to find features that are shared by sequence from the same genome...

sequence = $(\overline{CTGG}, \overline{TGGC})$

4-mers: 2×CTGG, TGGC, GGCT, GCTT, TTGA



Contig clustering by differential coverage (Albertsen et al. Nature Biotech. 2014)



CONCOCT: Clustering cONTigs on COverage and Composition

(Alneberg et al. Nat. Methods 2014)

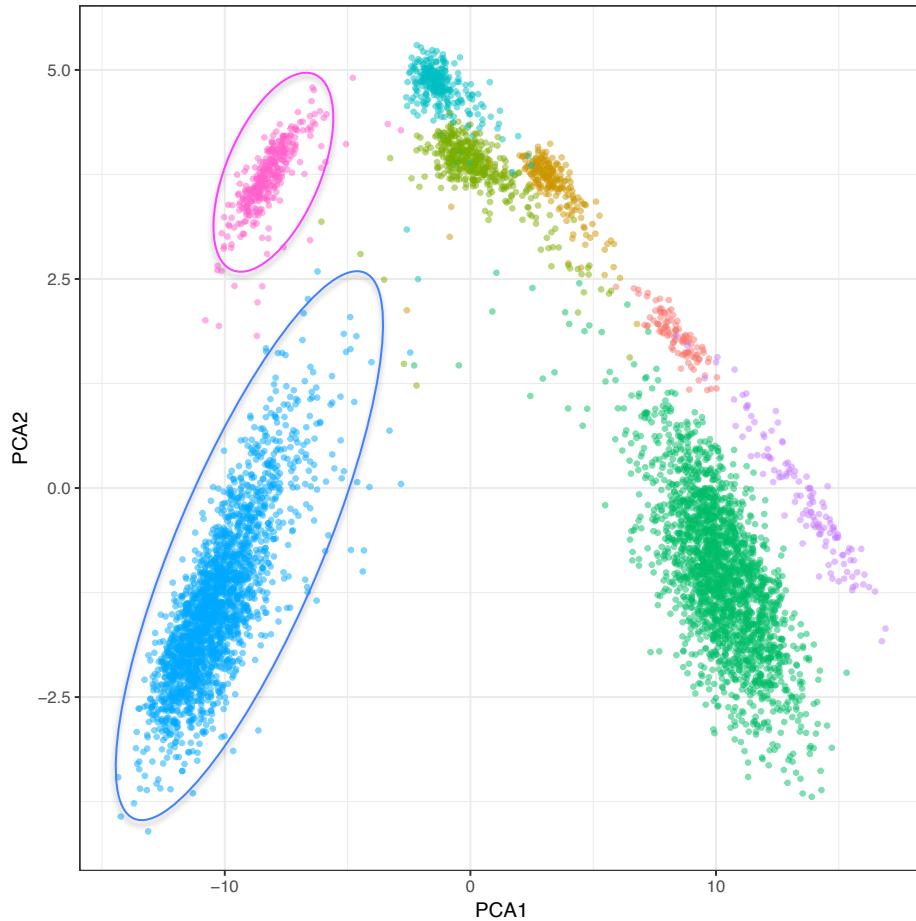


- Albertsen's strategy and ESOMs ([Sharon et al. 2013](#)) required human supervision
- Created CONCOCT as the first? automatic binner to combine coverage and composition
- Crowded field now e.g. GroopM ([Imelfort et al. PeerJ 2014](#)), MetaBAT ([Kang et al. PeerJ 2015](#)), MaxBin2 ([Wu et al. Bioinf. 2015](#)), COCACOLA ([Lu et al. 2017](#))
- CONCOCT still used and is highly cited, first step in the second generation human supervised clustering platform Anvi'o ([Eren et al. PeerJ 2015](#))



CONCOCT algorithm

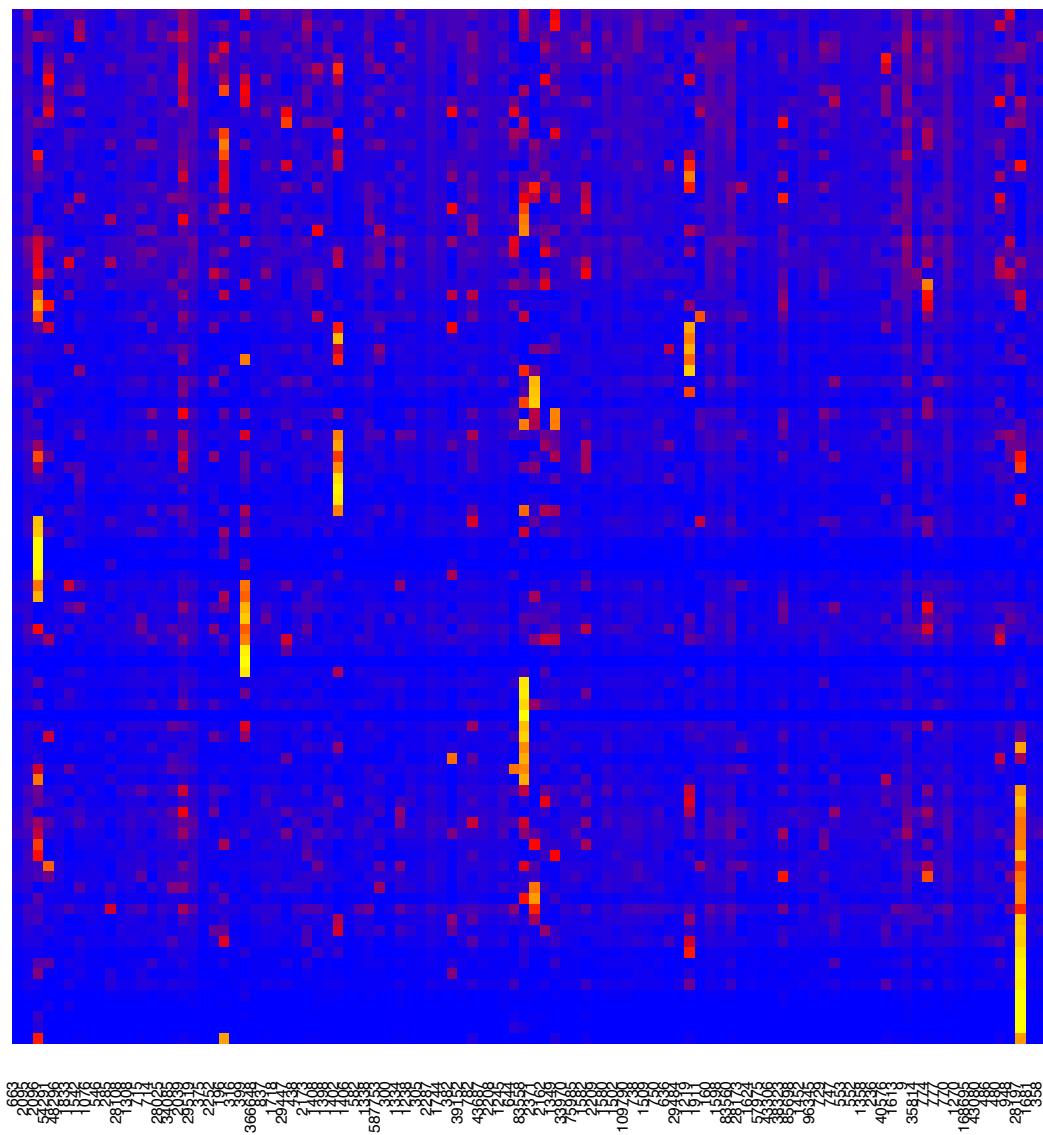
- Perform **coassembly** across all samples
- Map reads back to get coverage depth in **each** sample
- Use coverage and composition but perform PCA to reduce dimensionality
- Cluster with Gaussian Mixture Model
- Variational Bayes to select number of components through automatic relevance determination (ARD)
- Added multithreading and refinement step in development branch



Metagenomic binner	Description	Original publication
CONCOCT	Binner using differential coverage, tetranucleotide frequencies, paired-end linkage	Alneberg, J., Bjarnason, B. S., de Bruijn, I., Schirmer, M., Quick, J., Ijaz, U. Z., ... Quince, C. (2014). Binning metagenomic contigs by coverage and composition. <i>Nature Methods</i> , 11(11), 1144–1146.
CONCOCT with refinement step	Binner using differential coverage, tetranucleotide frequencies, paired-end linkage	-
MaxBin 2.0	Binner using multi-sample coverage, tentranucleotide frequencies	Wu, Y. W., Simmons, B. A., & Singer, S. W. (2015). MaxBin 2.0: An automated binning algorithm to recover genomes from multiple metagenomic datasets. <i>Bioinformatics</i> , 32(4), 605–607.
MetaBAT	Binner using multi-sample coverage, tetranucleotide frequencies, paired-end	Kang, D. D., Froula, J., Egan, R., & Wang, Z. (2015). MetaBAT, an efficient tool for accurately

Strain diverse synthetic community

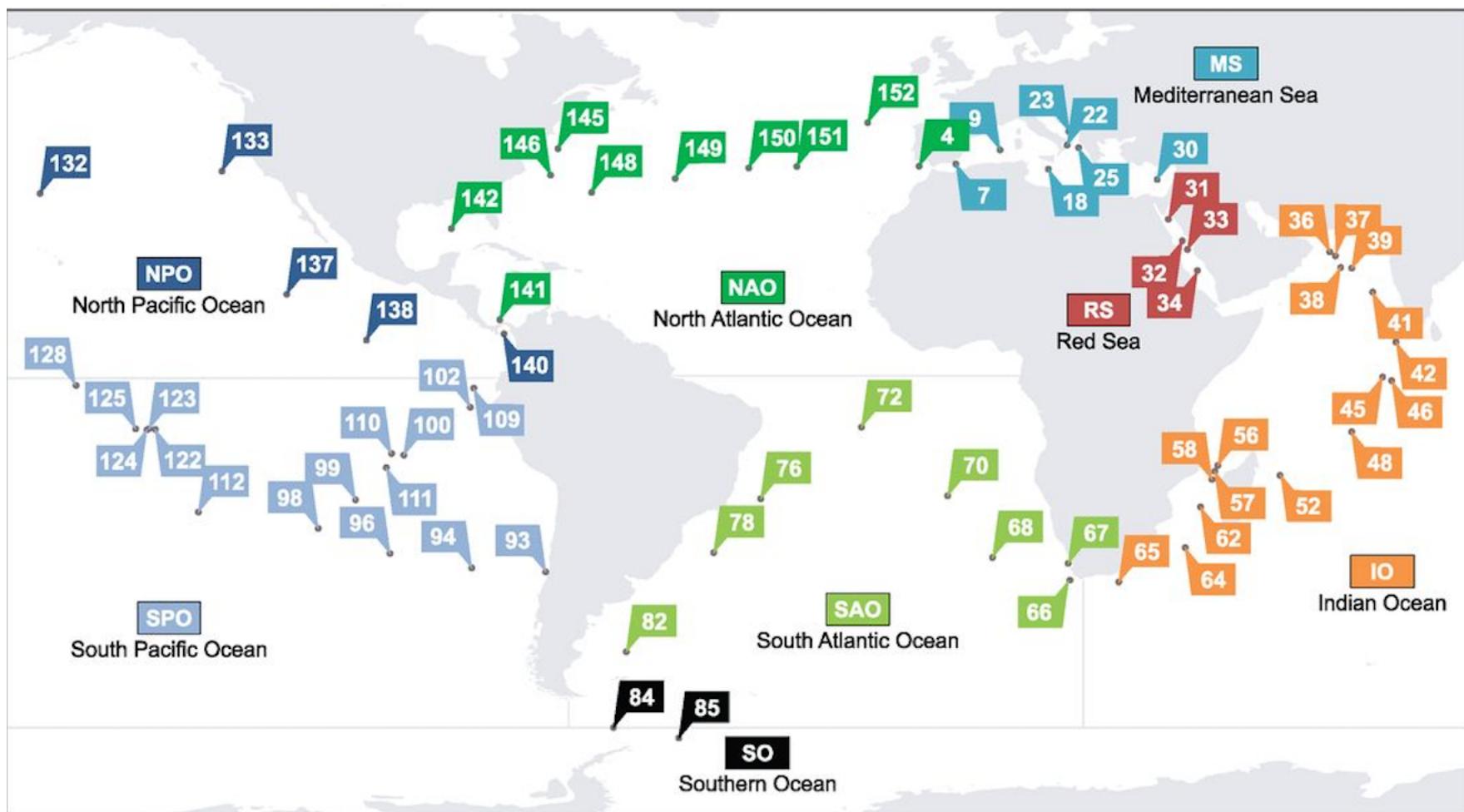
- 100 different species and 210 genomes (1:50, 2:20, 3:10, 4:10, 5:10)
- Simulated 96 samples of 6.25 million 2X150 bp HiSeq paired end reads (<https://github.com/chrisquince/StrainMetaSim>)
- Realistic species and strain coverages across samples
- Coassembly with megahit
- Randomly select 10, 20, 30 and 60* samples
- Evaluate with adjusted Rand index and number of complete genomes



Results

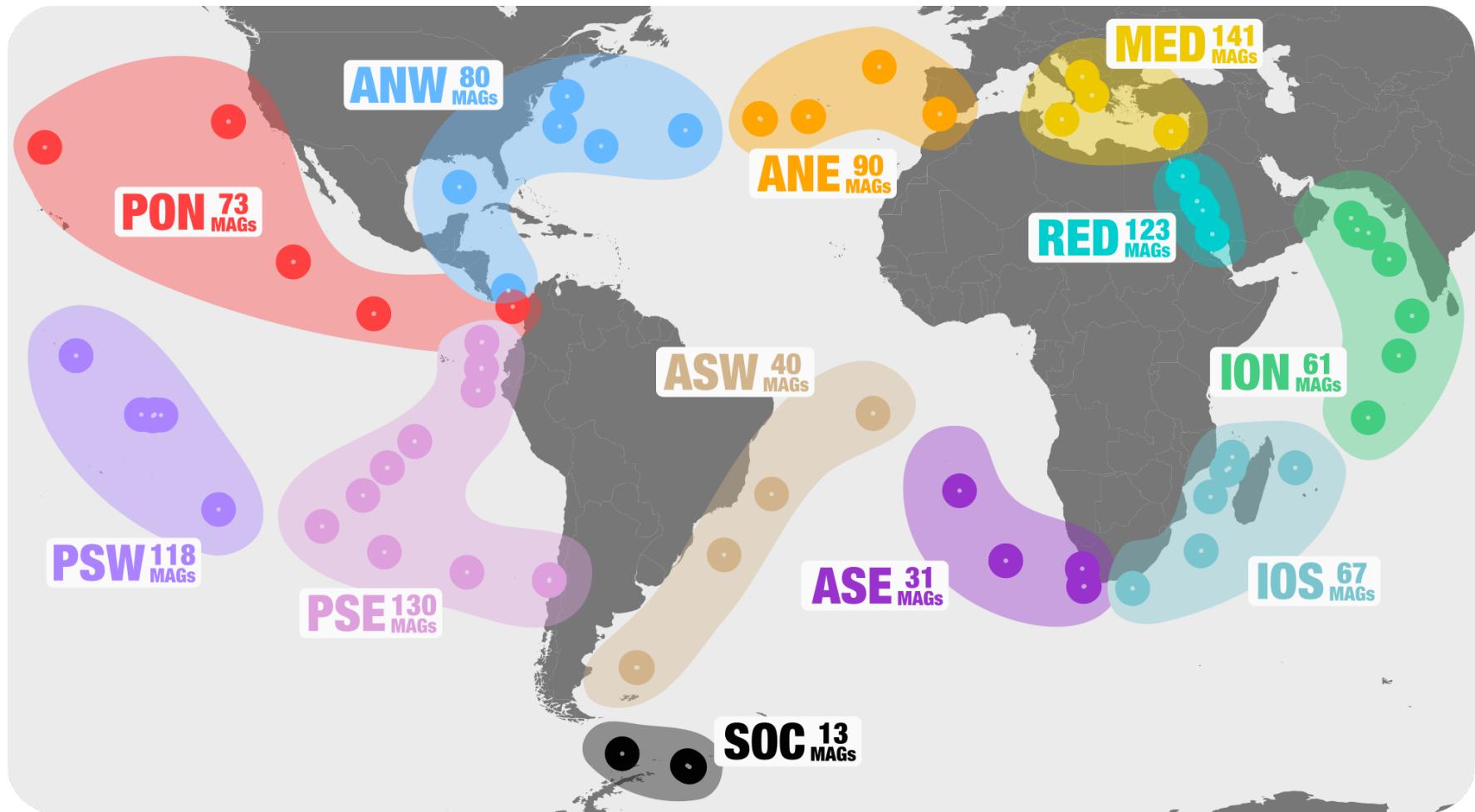
Sample size	Metagenomic binner	Contigs binned	Number of bins (75% COGs)	ARI
N=10	CONCOCT	73458	170 (59)	0.8086
	CONCOCT refinement	73458	130 (66)	0.8197
	MAXBIN 2.0	72512	110 (48)	0.7105
	MetaBAT	71603	123 (48)	0.2895
	GATTAGA	73379	95 (39)	0.4881
N=20	CONCOCT	79863	130 (52)	0.7800
	CONCOCT refinement	79863	165 (70)	0.8061
	MAXBIN 2.0	79377	116 (64)	0.7398
	MetaBAT	76875	137 (45)	0.2768
	GATTAGA	79778	87 (42)	0.3638
N=30	CONCOCT	78382	115 (53)	0.6456
	CONCOCT refinement	78382	165 (75)	0.7848
	MAXBIN 2.0	77983	117 (54)	0.7579
	MetaBAT	74936	136 (46)	0.2836
	GATTAGA	78431	89 (44)	0.50988
N=60	CONCOCT	75486	104 (58)	0.6806
	CONCOCT refinement	75486	152 (79)	0.7944
	MAXBIN 2.0	74936	117(70)	0.7688
	MetaBAT	72625	136 (43)	0.2967

TARA Oceans sampling sites: Sunagawa et al. Science 2015



The 93 TARA Oceans metagenomes we analyzed represent the planktonic size fraction (0.2-3 μ m) of 61 surface samples and 32 samples from the deep chlorophyll maximum layer of the water column

TARA Oceans: Sunagawa et al. Science 2015



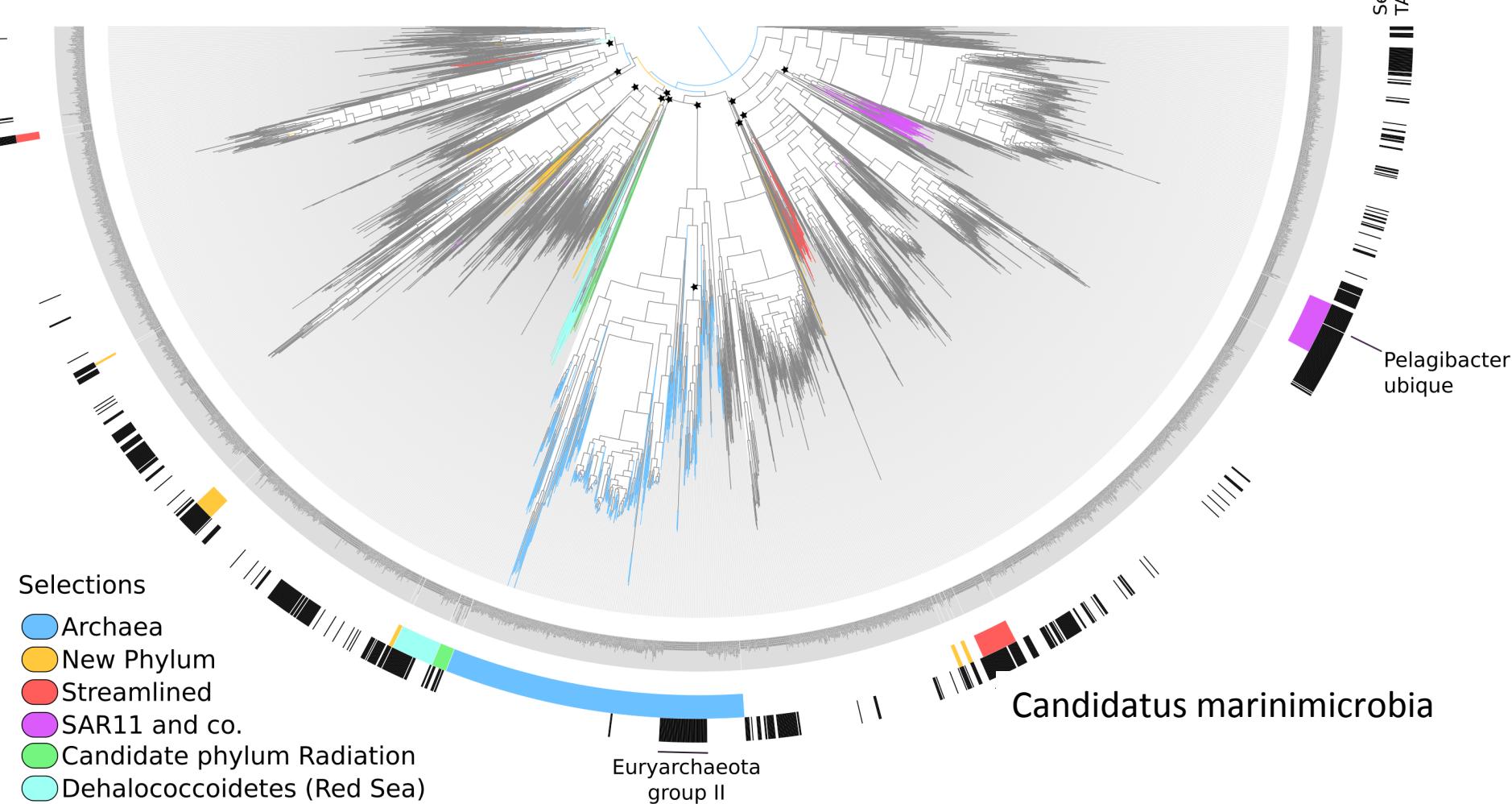
- 93 metagenomes from the planktonic size fraction for which we performed 12 metagenomic co-assemblies
- Generated 957 non-redundant MAGs encompassing the three domains of life
Delmont, Quince, ..., Eren “Nitrogen-fixing populations of Planctomycetes and Proteobacteria are abundant in the surface ocean” Nature Microbiology 2018

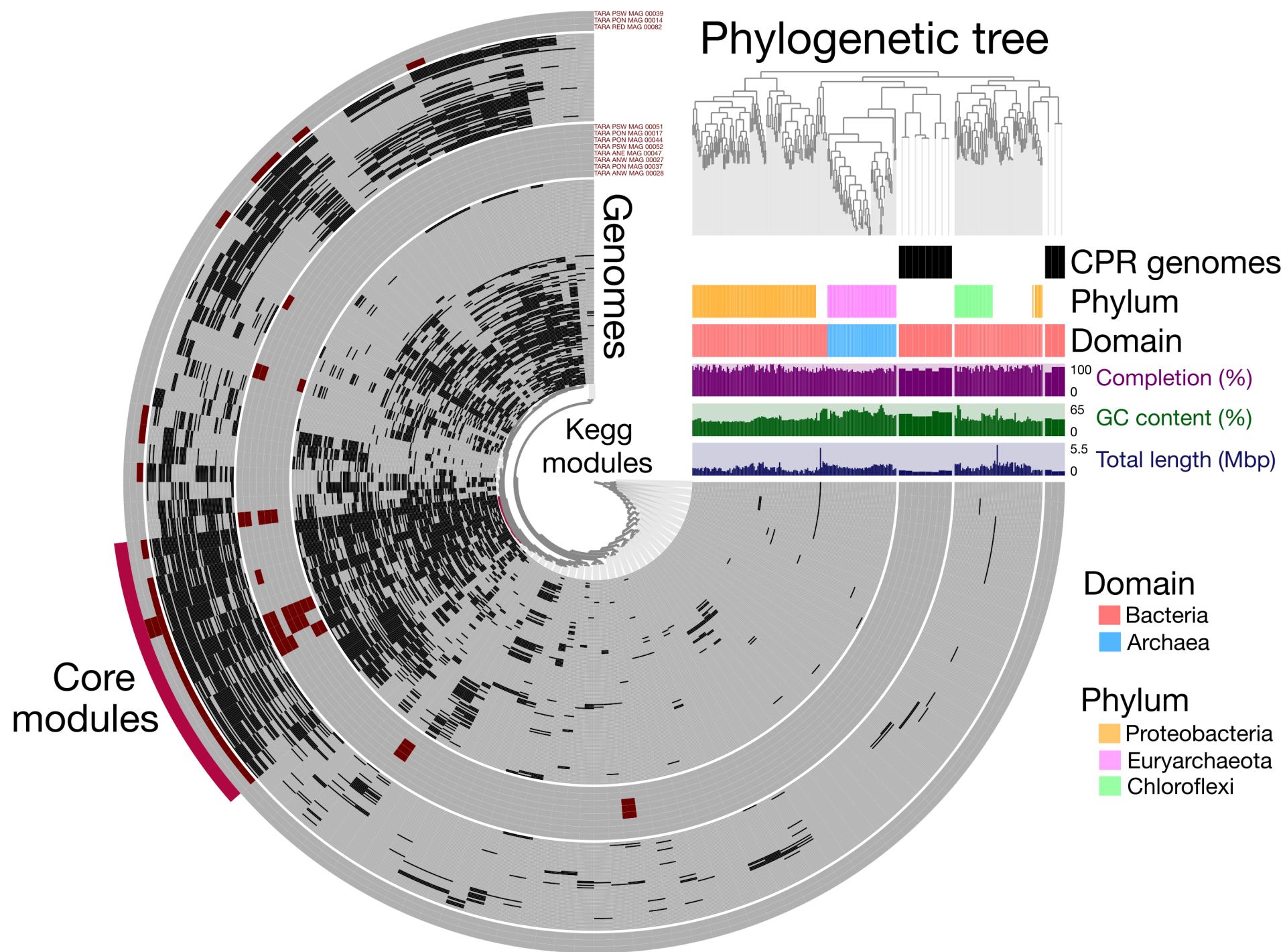


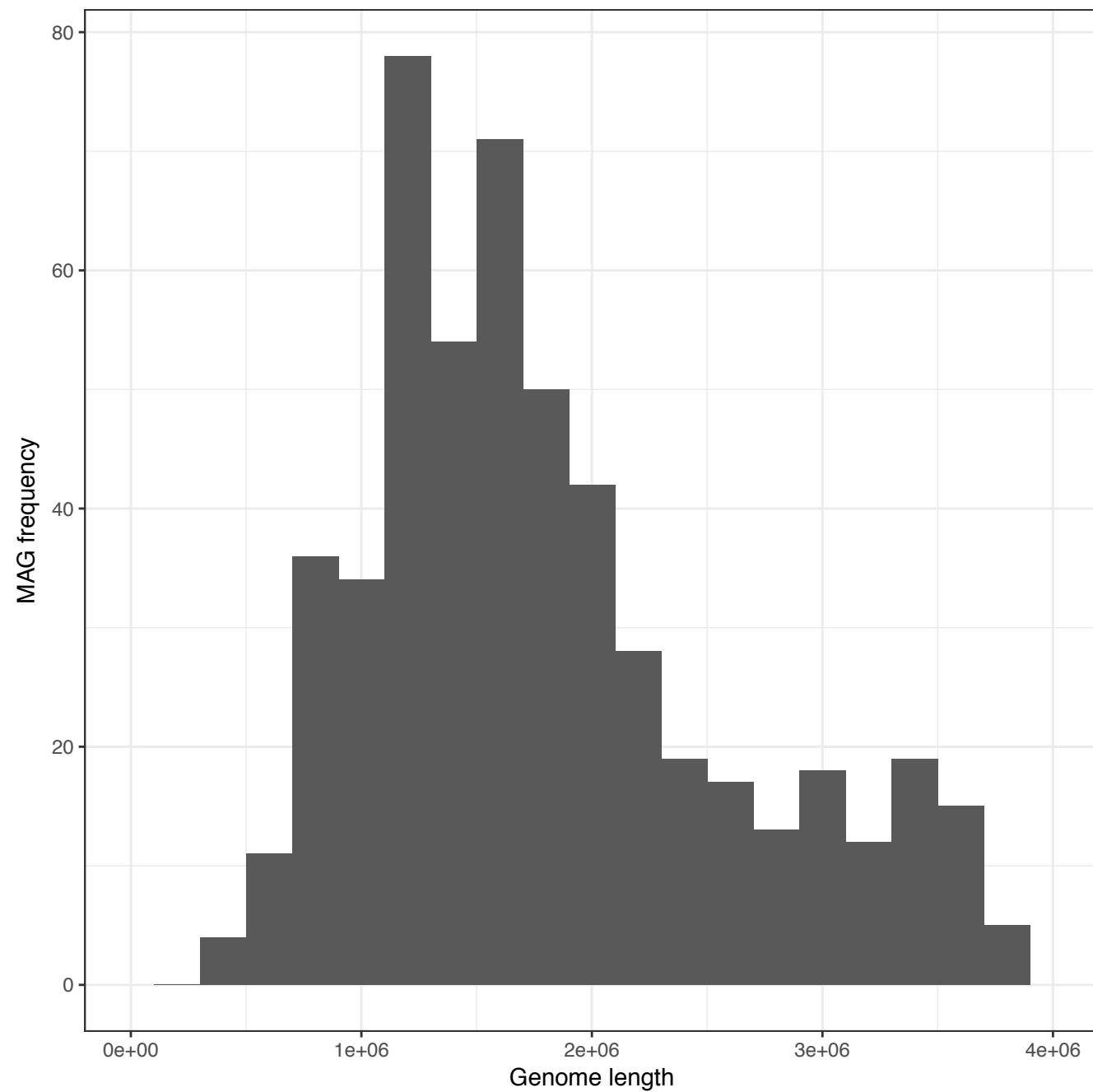
Phylogenetic analysis of 660 MAGs

plus 1,755 ref genomes (1 genome per genus for bacteria)

→ Concatenated single copy genes







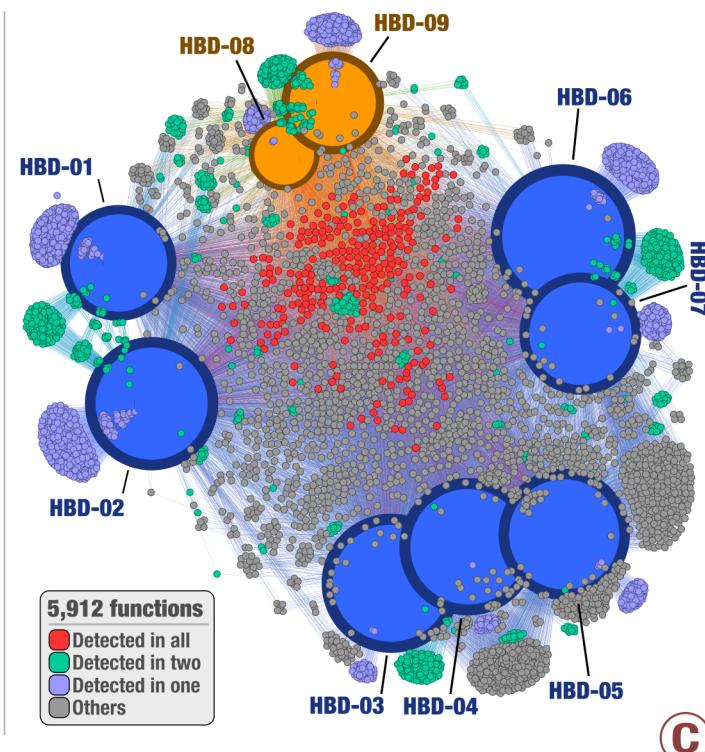
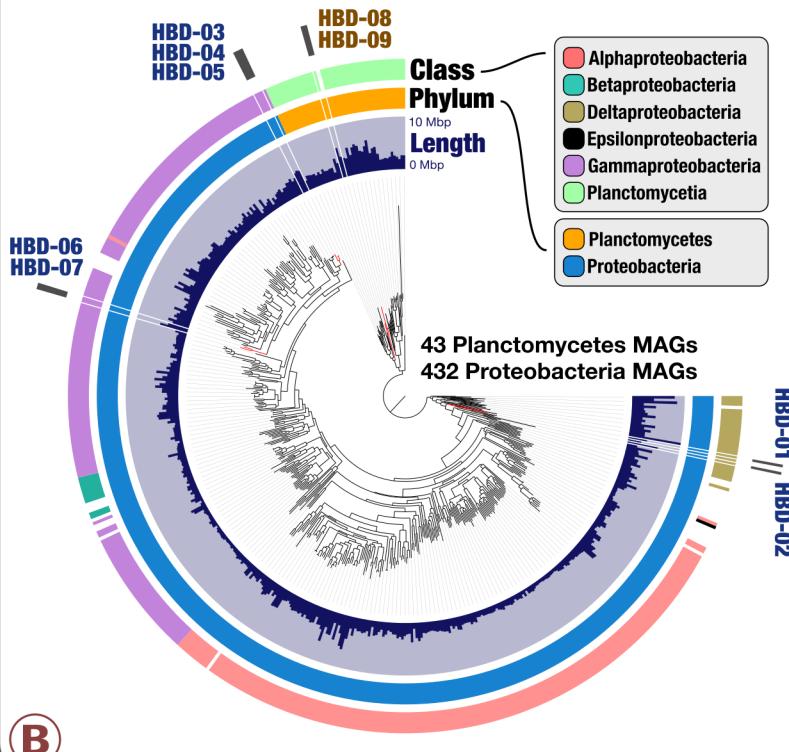
Oceanic heterotrophic bacterial diazotrophs (HBDs)

- Nitrogen fixation in the surface regulates microbial primary productivity and the sequestration of carbon through the biological pump
- Cyanobacteria populations have long been thought to represent the main suppliers of the bio-available nitrogen
- Nitrogenase reductase gene surveys revealed the existence of non-cyanobacterial populations that can also fix nitrogen - heterotrophic bacterial diazotrophs (HBDs)

- Search for MAGs containing the catalytic (*nifH*, *nifD*, *nifK*) and the biosynthetic proteins (*nifE*, *nifN* and *nifB*) required for nitrogen fixation
- One of six Cyanobacteria
- Nine other MAGs possessed all six genes

A

	Region of Recovery	Length (Mbp)	N50	Num Contigs	Num Genes	GC Content	Percent Compl.	Percent Redund.	Phylum	Class	Order	Family	Genus
HBD-01	PSW	3.67	48,153	118	3,592	52.6%	97.7%	4.4%	Proteobacteria	Deltaproteo.	Desulfovibionales	Desulfovibrionaceae	Desulfovibrio
HBD-02	PSW	6.00	20,964	405	5,385	53.1%	97.1%	5.9%	Proteobacteria	Deltaproteo.	Desulfobacterales	Desulfobacteraceae	-
HBD-03	ION	4.47	57,949	110	4,155	52.4%	97.5%	8.1%	Proteobacteria	Gammaproteo.	Oceanospirillales	Oceanospirillaceae	-
HBD-04	PON	4.29	48,897	138	3,957	52.4%	89.7%	6.1%	Proteobacteria	Gammaproteo.	Oceanospirillales	Oceanospirillaceae	-
HBD-05	PSE	4.15	65,098	94	3,867	53.3%	47.7%	5.8%	Proteobacteria	Gammaproteo.	Oceanospirillales	Oceanospirillaceae	-
HBD-06	ANW	5.49	76,792	112	4,842	54.2%	98.1%	5.6%	Proteobacteria	Gammaproteo.	Pseudomonadales	-	-
HBD-07	ANW	3.99	10,488	487	3,585	58.7%	91.2%	4.3%	Proteobacteria	Gammaproteo.	Pseudomonadales	-	-
HBD-08	PSW	4.03	10,413	480	3,327	52.6%	33.5%	6.0%	Planctomycetes	-	-	-	-
HBD-09	PSW	5.86	79,495	113	4,872	50.0%	97.3%	4.6%	Planctomycetes	Planctomycetia	Planctomycetales	Planctomycetaceae	-

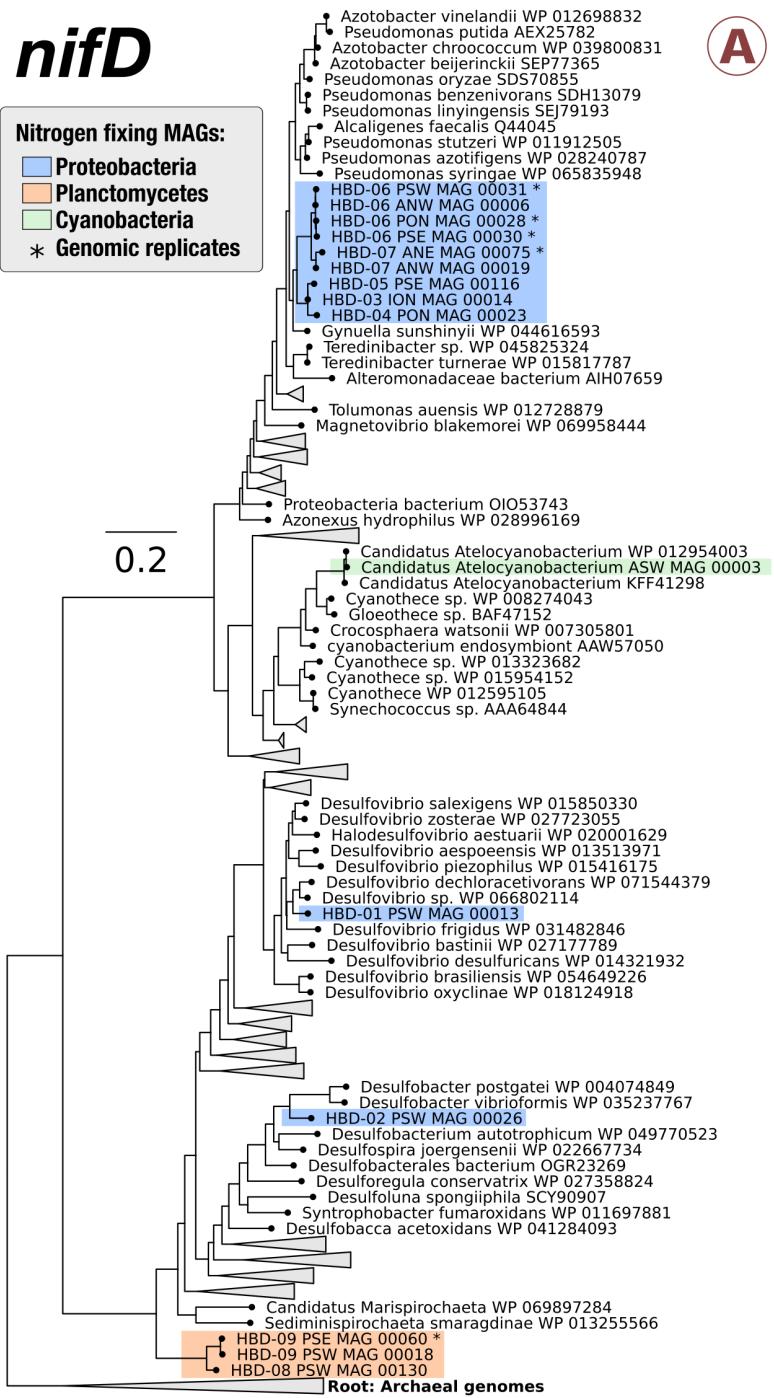


nifD

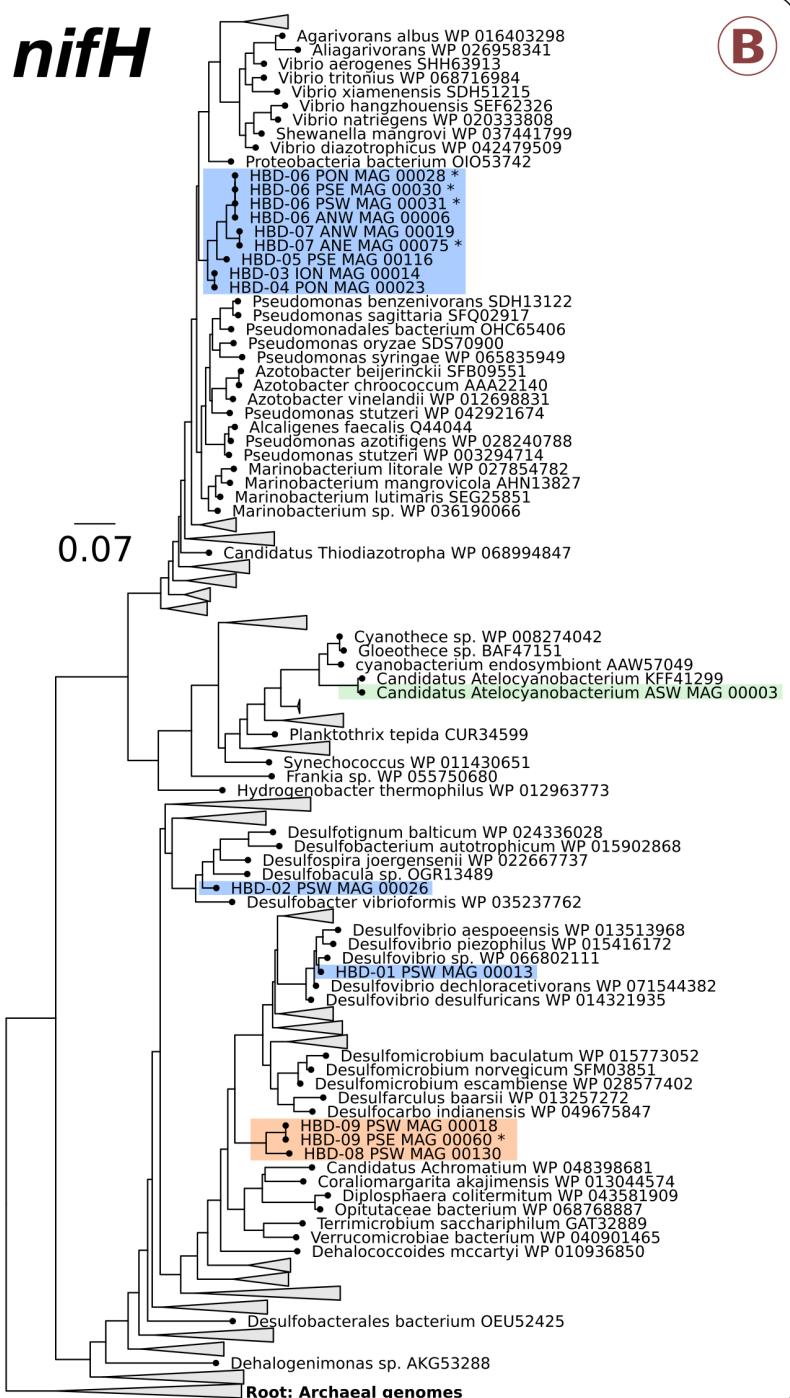
A

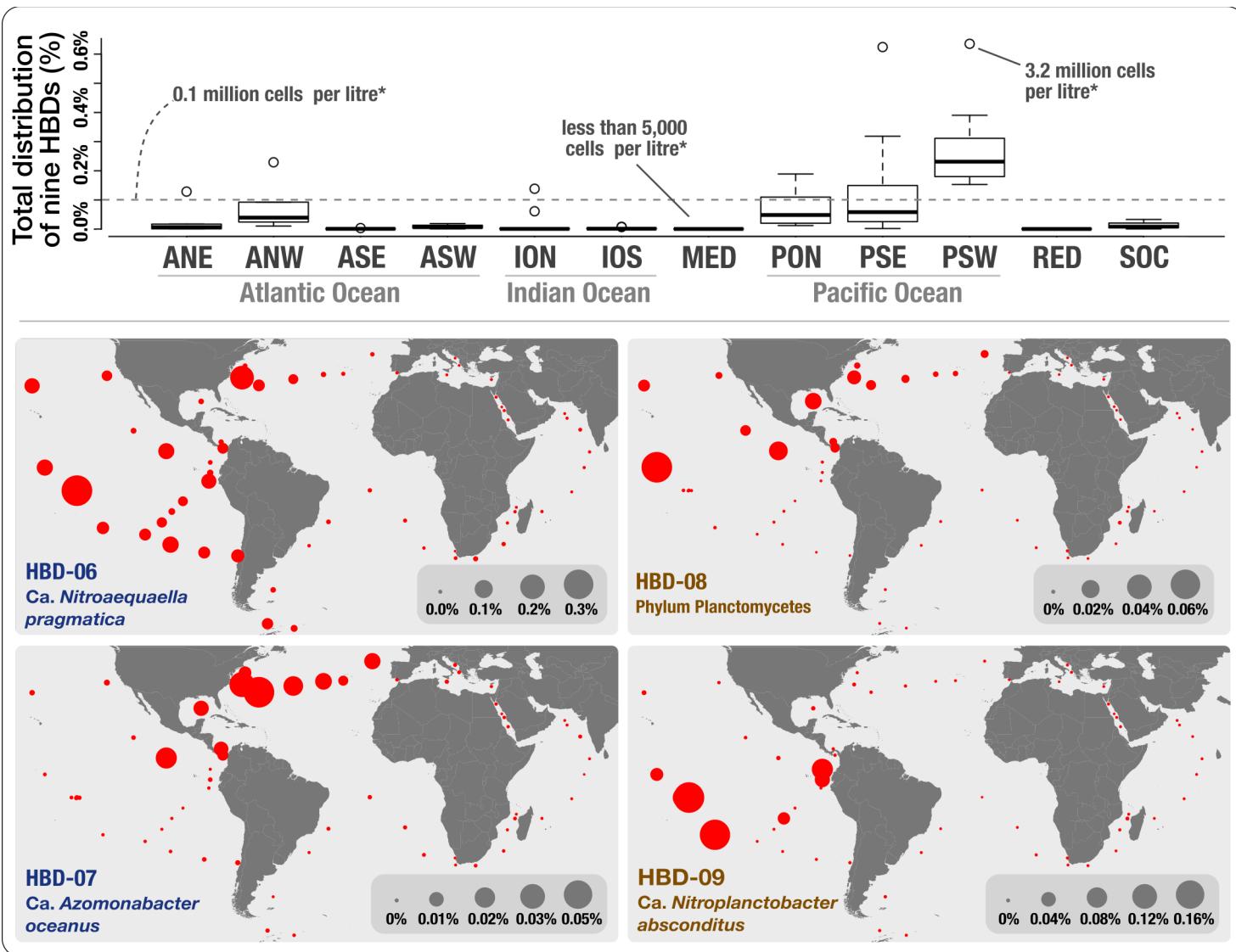
Nitrogen fixing MAGs:

- Proteobacteria
- Planctomycetes
- Cyanobacteria
- * Genomic replicates



nifH

B



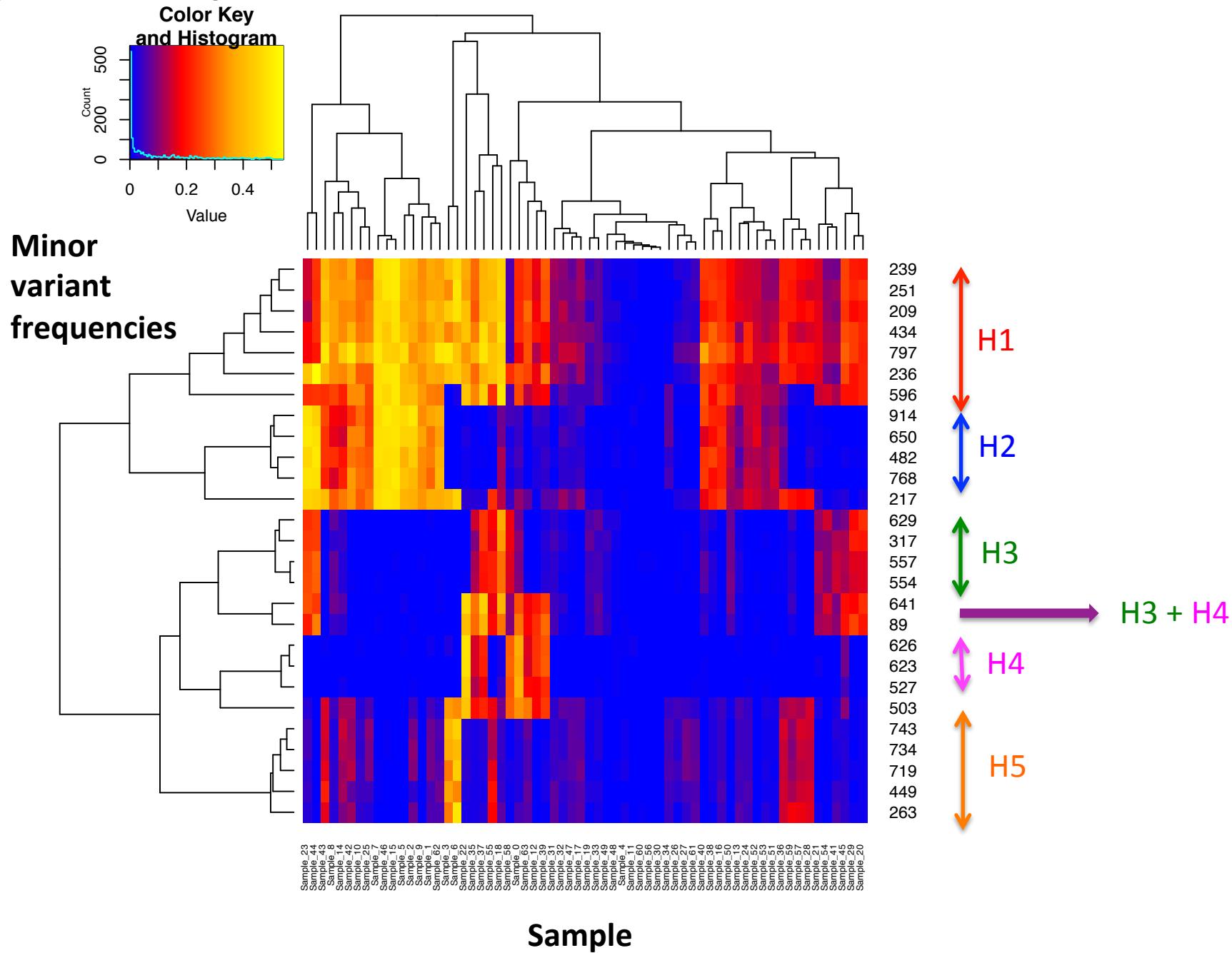
PART II: **DESMAN**

Resolving Intra-MAG diversity

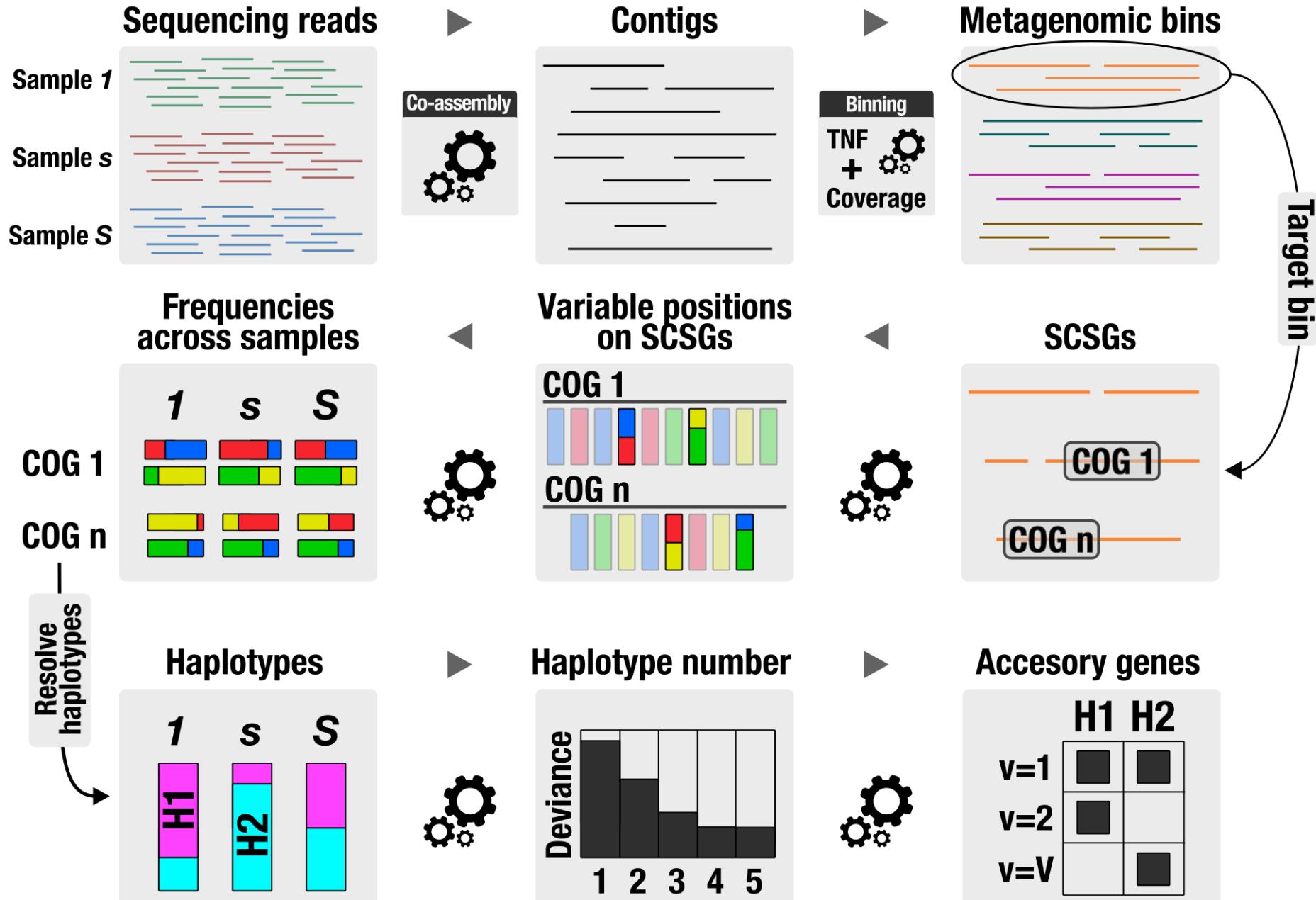
Resolving intra-MAG diversity

- Additional variation exists within MAGs both nucleotide variants on shared genes and variation in the accessory genome
- Read based haplotype resolution cannot span contigs or regions of low variation
- Methods exist using co-occurrence across multiple samples to resolve strain mixtures *de novo* after mapping to references:
 - Constrains (Luo et al. Nature Biotech. 2015)
 - Lineage (O'Brien et al. Genetics 2014)

Linking variants by co-occurrence

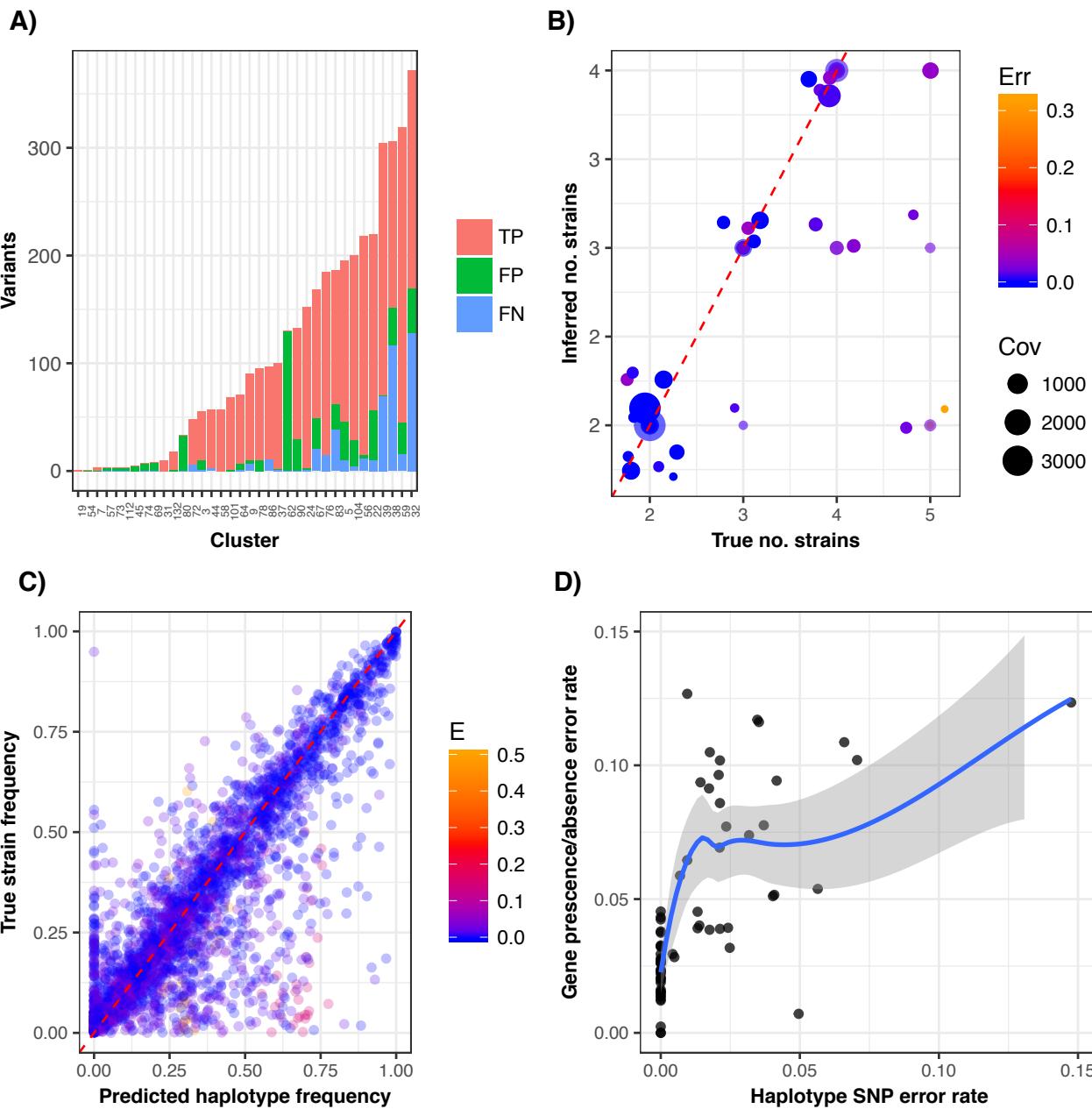


DESMAN <https://github.com/chrisquince/DESMAN>



Complex synthetic community

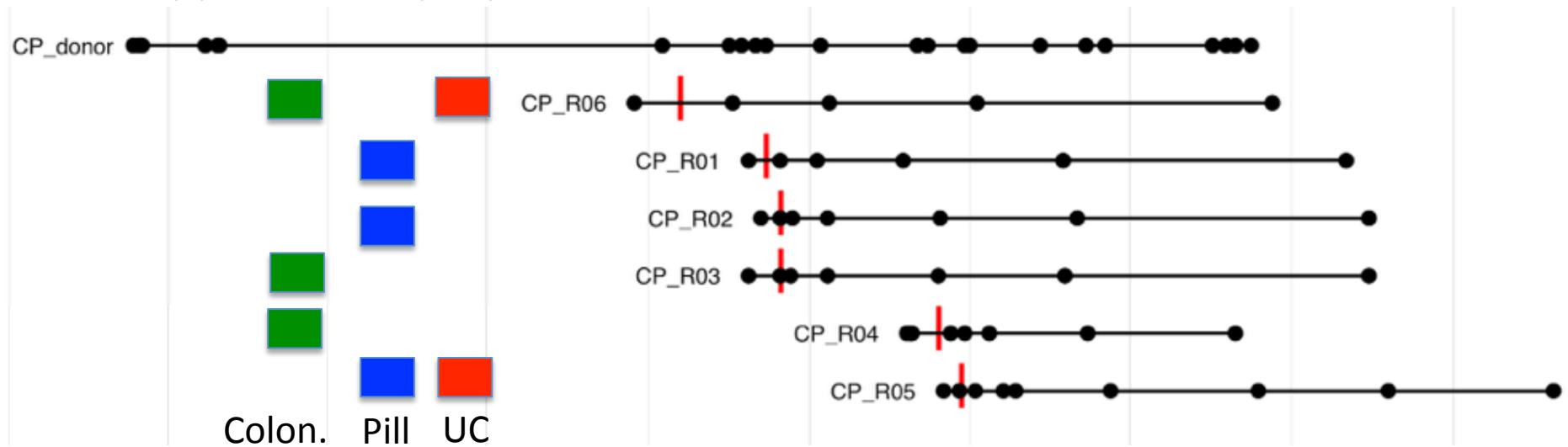
- Detected variants with a mean precision of 92.32% and a mean recall of 91.85%
- Predicted the correct haplotype number for 18/25 (72%) of the clusters
- SNV error rate: median of 0.25% and a mean of 2.38%
- True frequency against predicted gave a slope of 0.820 (R-squared 0.741, $p < 2.2e-16$)
- Overall accessory gene prediction accuracy was 95.7%.



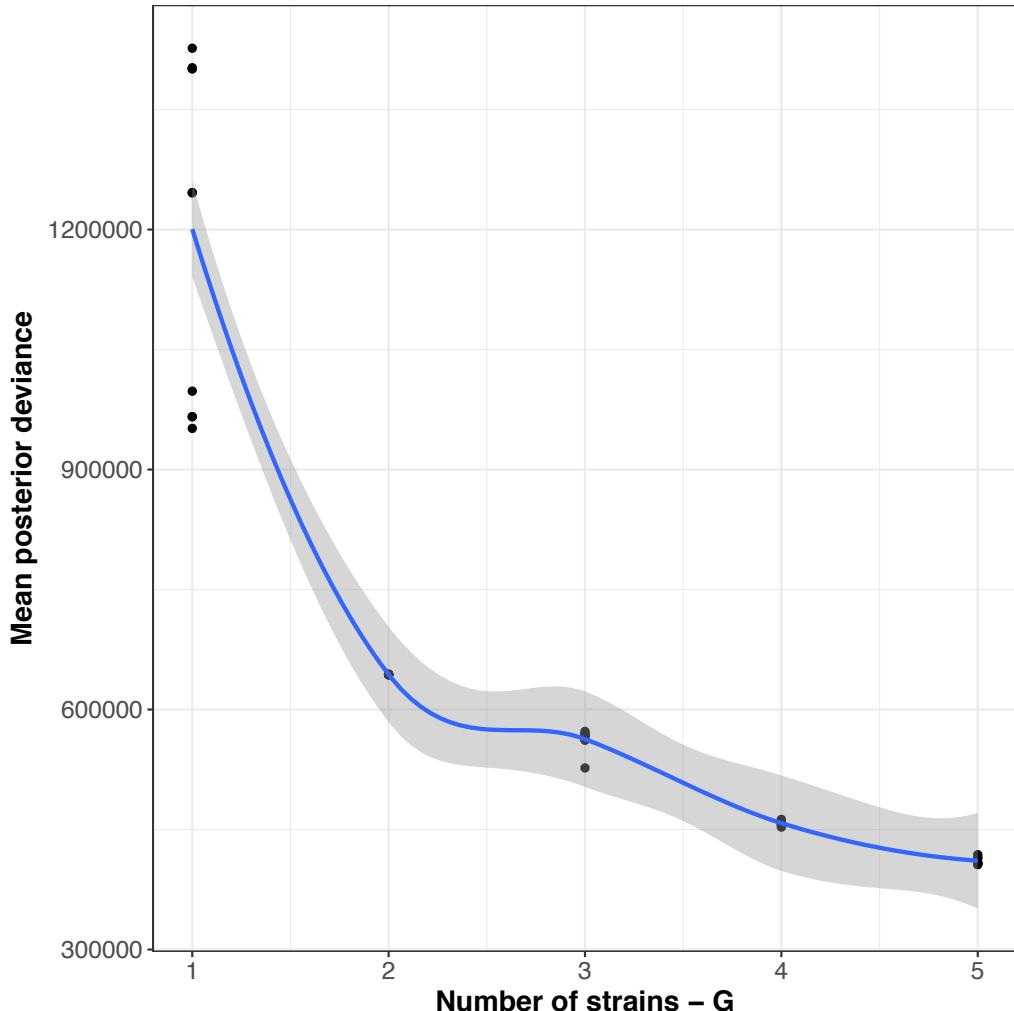
Fecal Microbiome Transplants

(Andrea Watson, A. Murat Eren)

- Healthy donors, CP, 24 samples 2 years
- 6 recipients, with *C. difficile* infection, and 2 also diagnosed with ulcerative colitis.
- Half received FMT through pill, and the other half received FMT through colonoscopy.
- For each of the recipients we have at least one sample from pre-FMT, and a median of 5 samples taken post-FMT over the course of approximately a year.

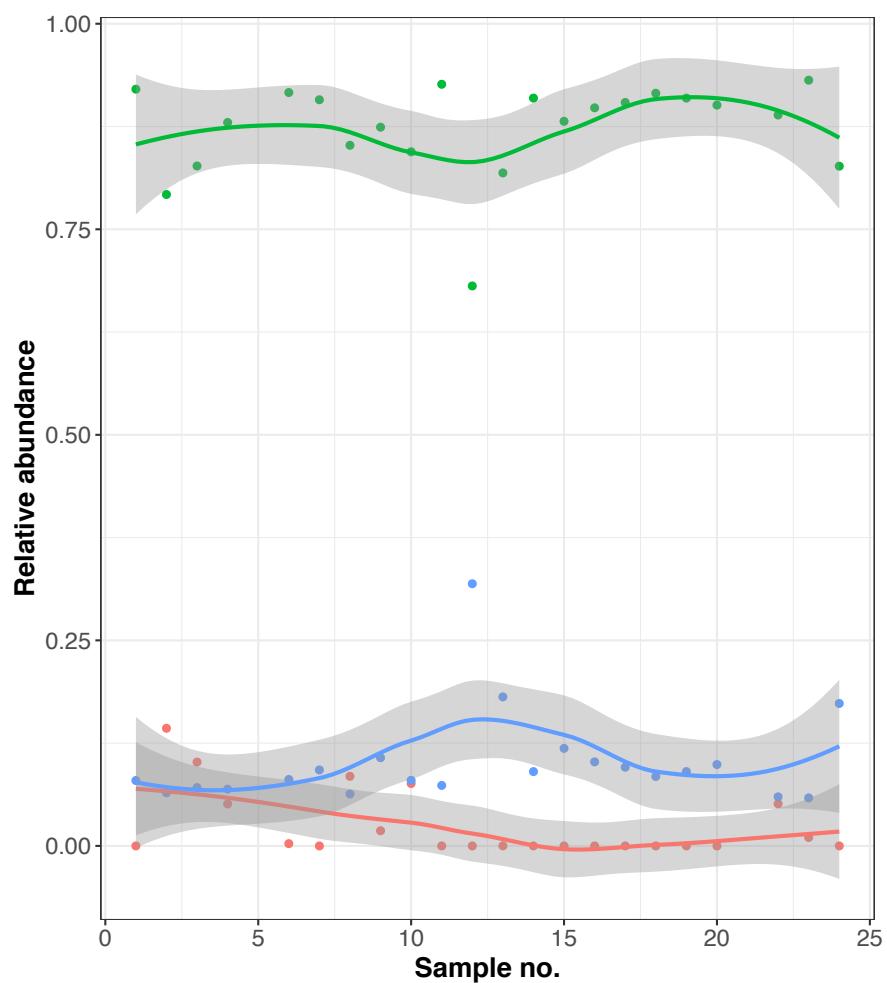


FMT DESMAN strain analysis

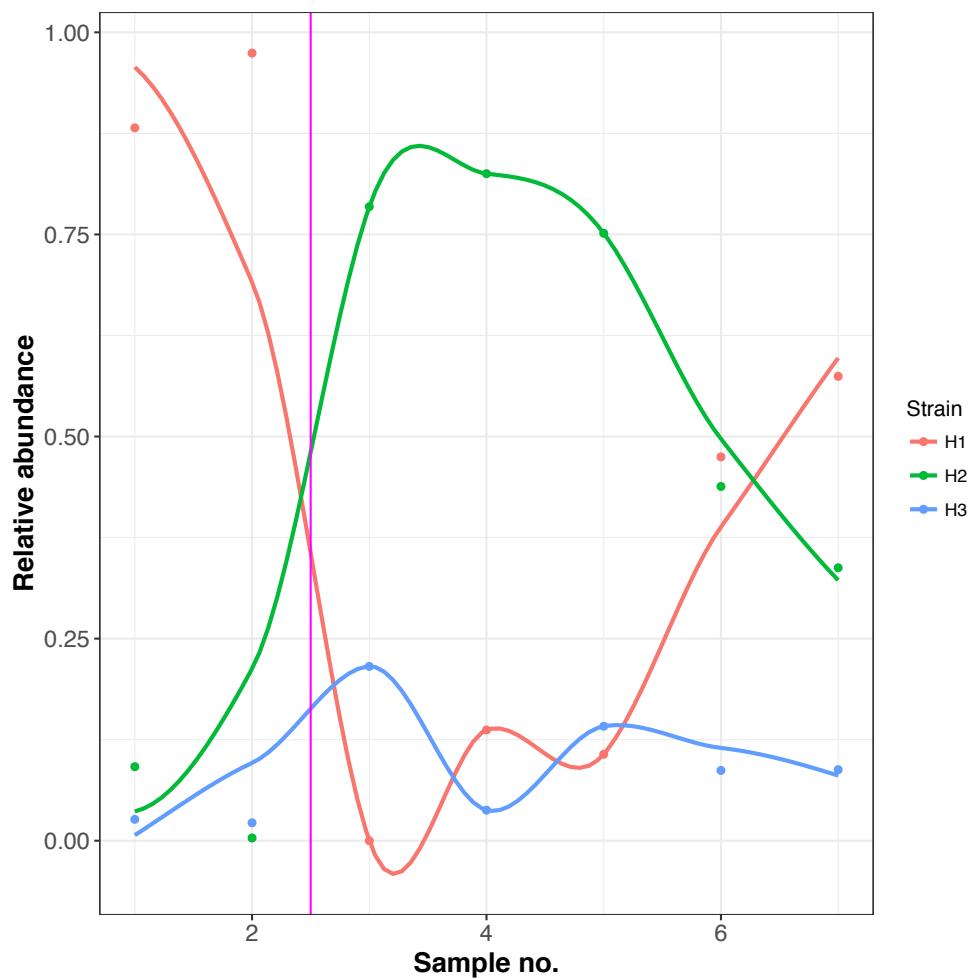


- Detected variants on one Anvi'o MAG assigned to *Allistipes Finegoldii*
- Detected 41,544 variants on 2.8 Mbp
- Applied DESMAN to 1,000 variant positions increasing strain numbers from 1 to 5
- 4 strains selected as best fit, three of which were ‘reproducible’
- Varying between 35 – 50% of SNV positions

CP Donor



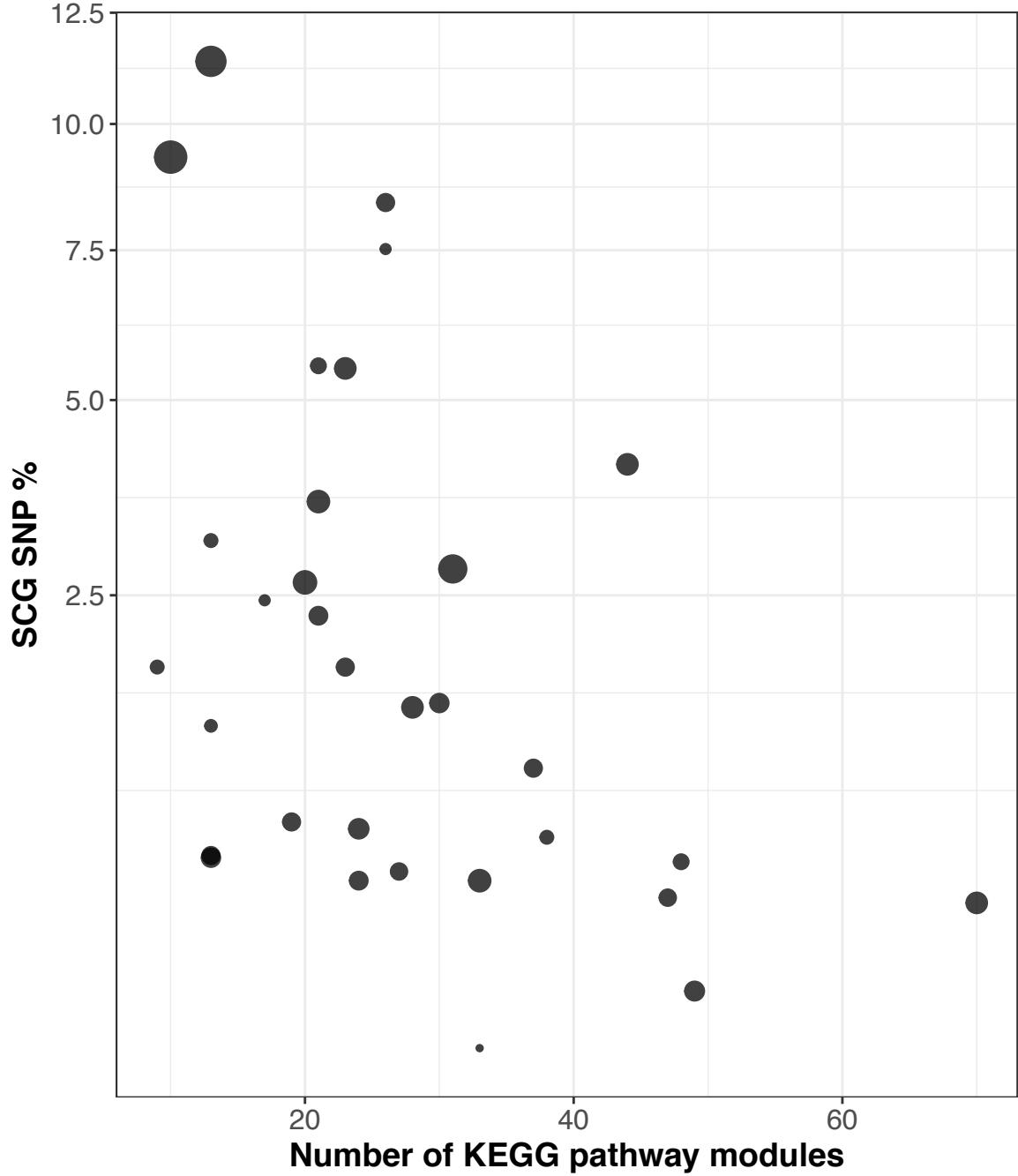
Recipient R03



Pre FMT

Post FMT

Tara MAG variants



Consider 32 MAGs with cov. > 100.

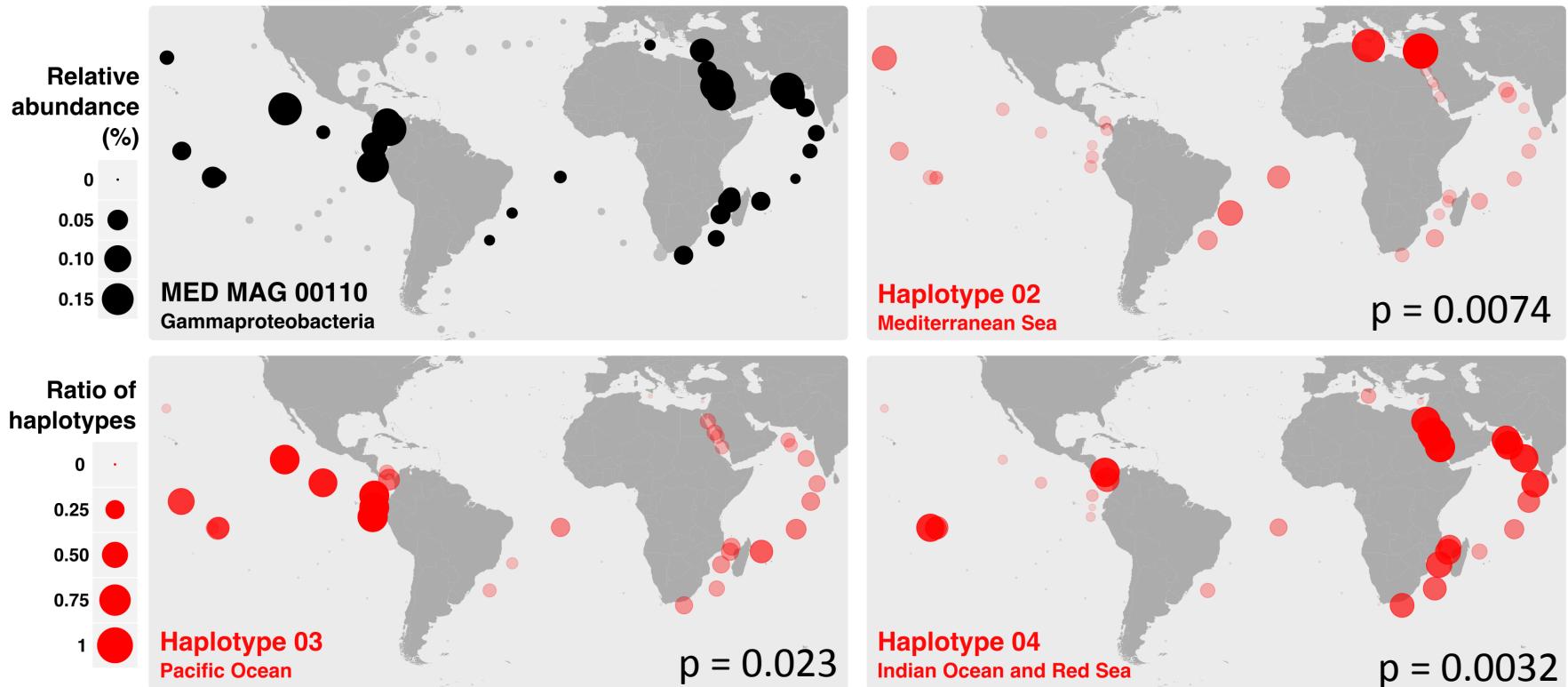
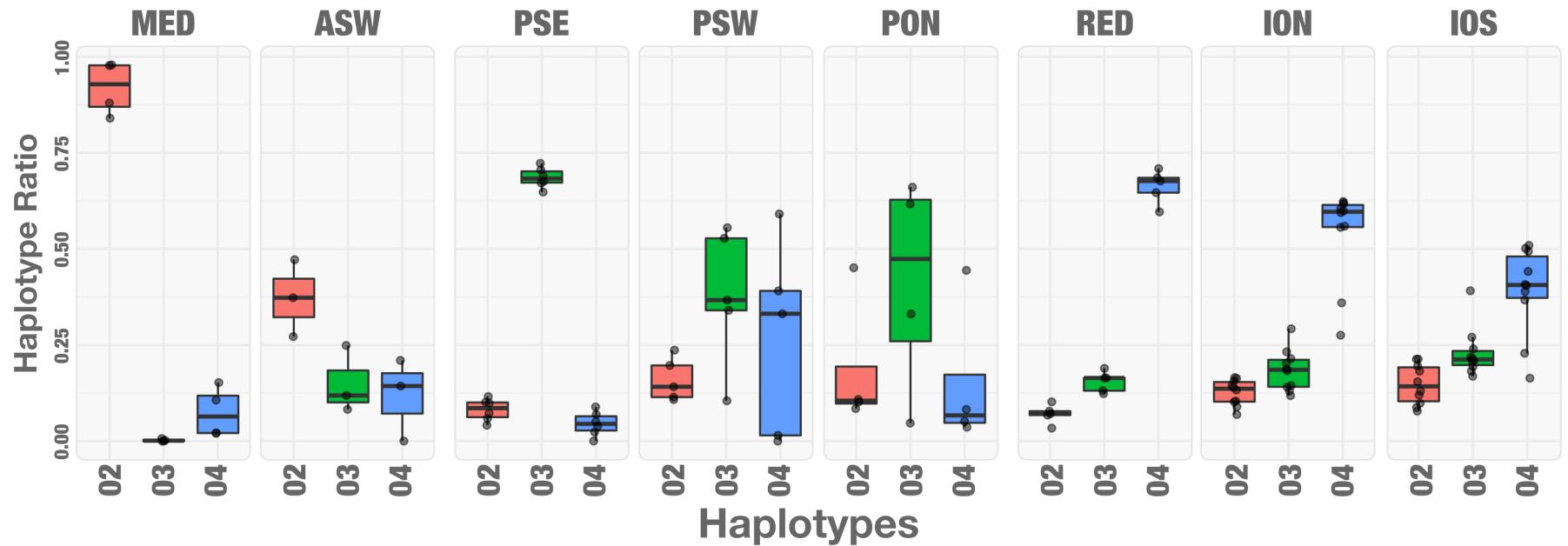
Observe a negative correlation with genome length (Spearman's $p = 0.016$)

Sites > 1
● 20
● 30
● 40
● 50

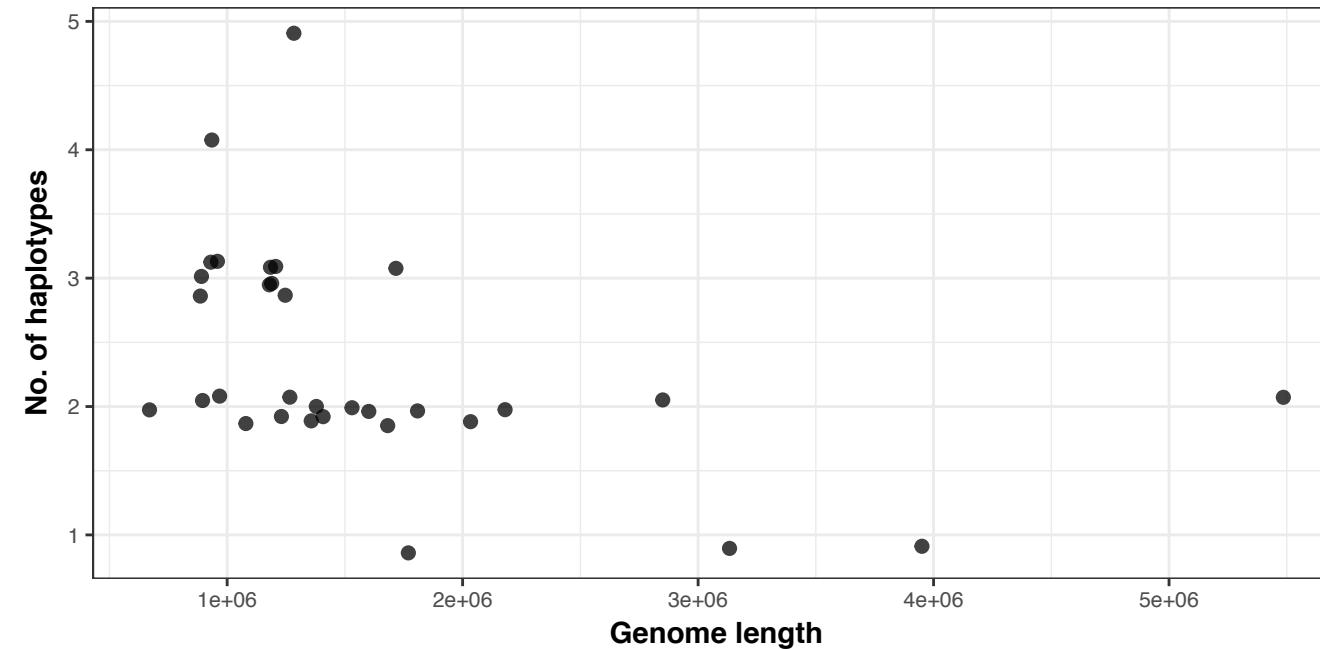
Stronger negative relationship with number of KEGG Pathway modules (Spearman's $p = 0.0045$)

TARA DESMAN analysis

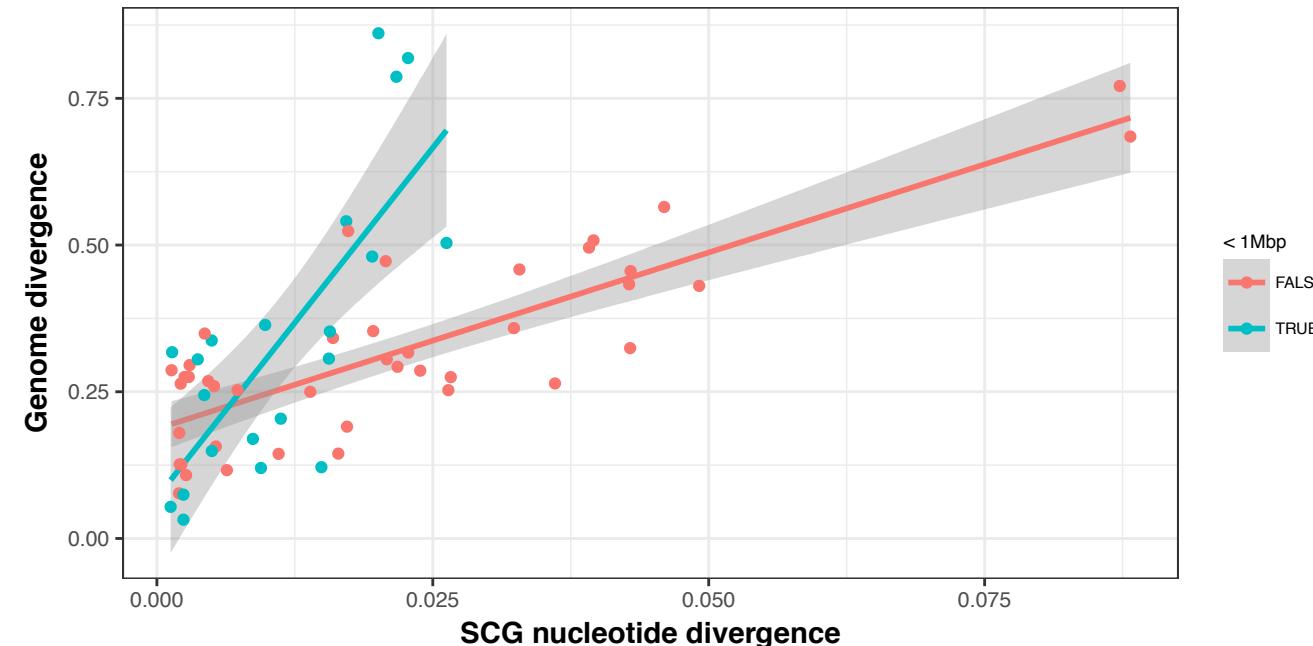
- Out of the 32 MAGs tested for haplotypes 29/32 had strain variation (1->3 2->17 3->10 4->1 5->1)
- The haplotypes were geographically localised e.g. TARA MAG 00110 a Proteobacteria with a highly streamlined 890,789 bp genome
 - Large group of uncultured organisms (relatives *Candidatus Evansia muelleri* and *Riesia pediculicola*)
 - Three haplotypes that differed by around 2% ANI on core genes and between 79-86% of accessory genes at 5% ANI clusters



TARA DESMAN results



- Significant negative correlation between strain number in MAG and genome length ($p = 0.000068$)
- $42/73 = 57\%$ of inferred strains had a significant correlation with geographic region
- Significant interaction ($p = 3.51e-06$) between rate of whole genome divergence relative to core gene divergence and highly streamlined genome (<1Mb)



Summary

- Strain diversity is endemic in environmental and host associated communities
- Now have the tools to generate hundreds or even thousands of genomes from metagenomes
- CONCOCT followed by DESMAN can resolve both species and strain diversity *de novo* but with limitations
- The question is then what will we do with all this data!
 - Machine learning based trait inference from genomes
[https://www.biorxiv.org/content/early/
2018/04/25/307157](https://www.biorxiv.org/content/early/2018/04/25/307157)

Acknowledgements



Unilever

- Medical Research Council funding and CLIMB, BBSRC/Unilever BBSRC
bioscience for the future
MRC
Cloud Infrastructure
for Microbial
Bioinformatics
- Tom Delmont and A. Murat Eren (U. of Chicago)
- Leonidas Souliotis and Sebastian Raguideau (Warwick)

