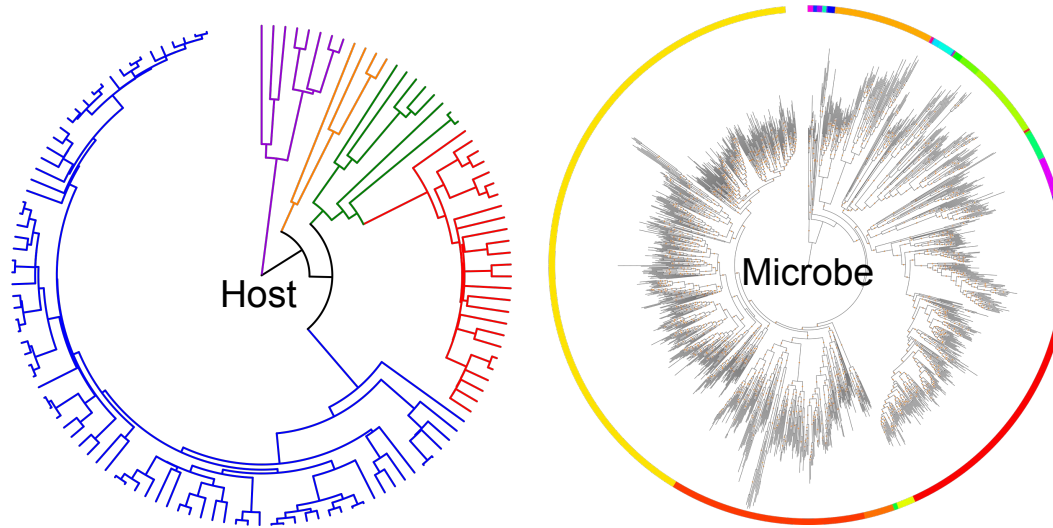


Bioinformatics in environmental and host-associated microbiome research



Nick Youngblut
Group Leader
Department of Microbiome Science
Max Planck Institute for Developmental Biology

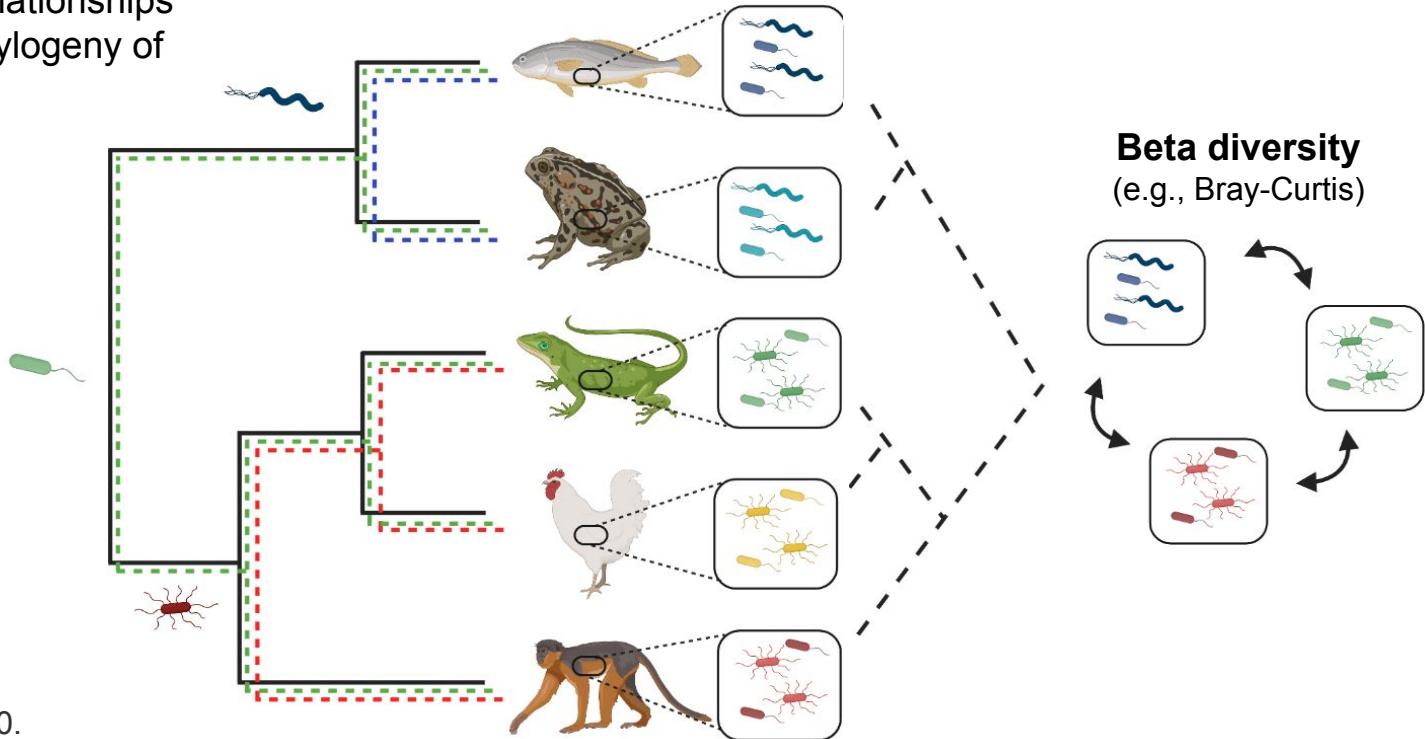




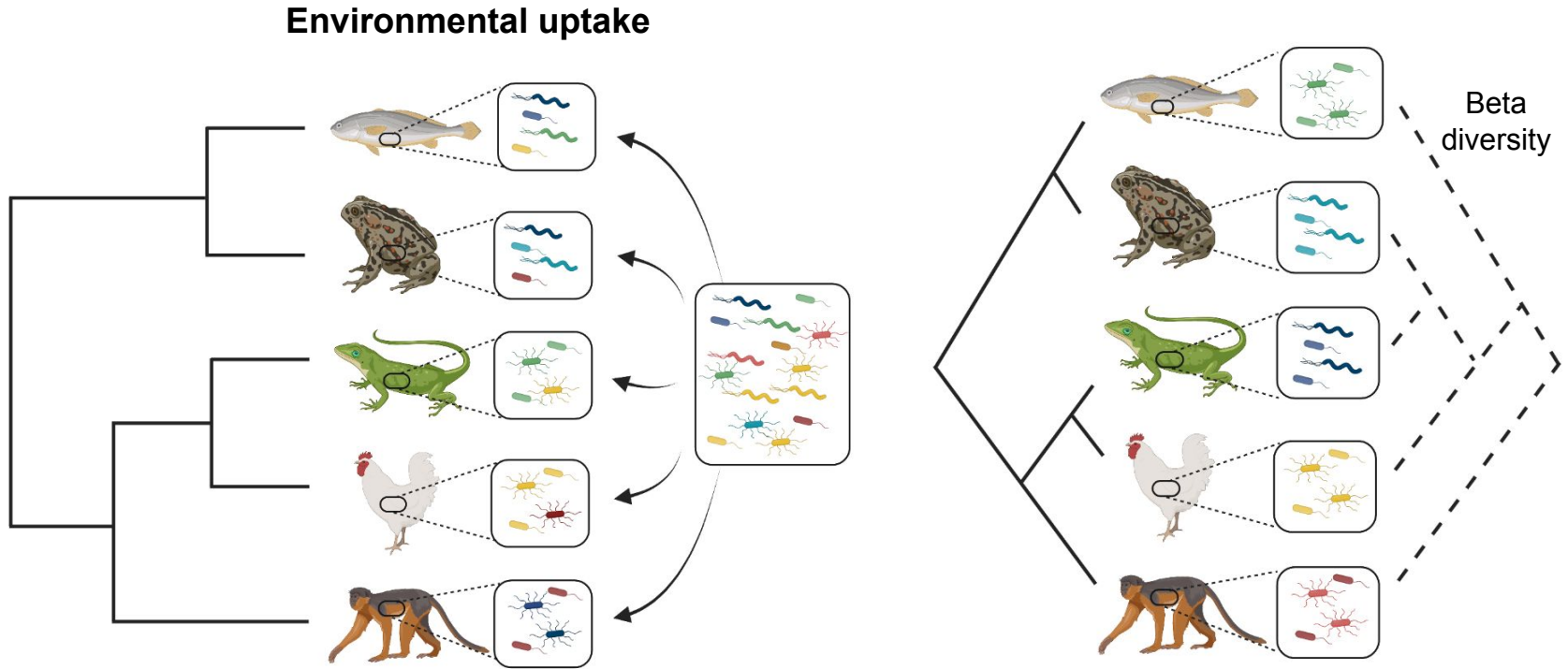
Microbial community assembly in the animal gut

Phylosymbiosis:

microbial community relationships that recapitulate the phylogeny of their host

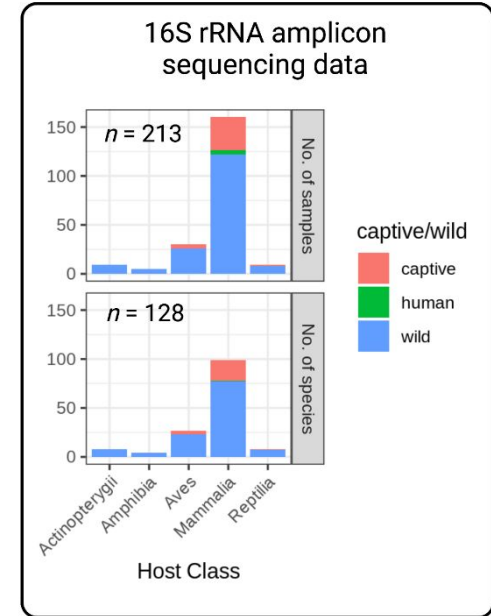
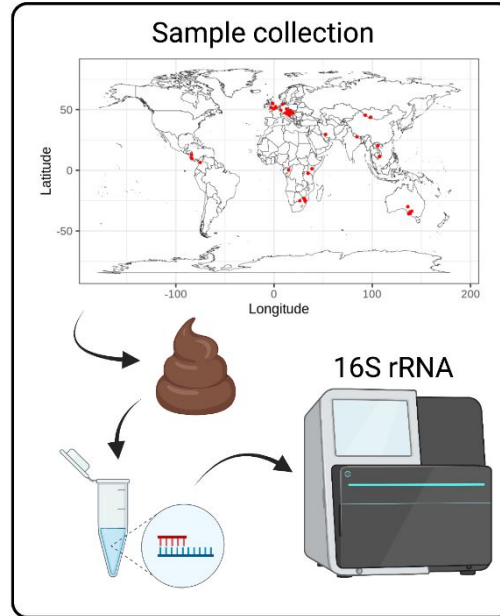
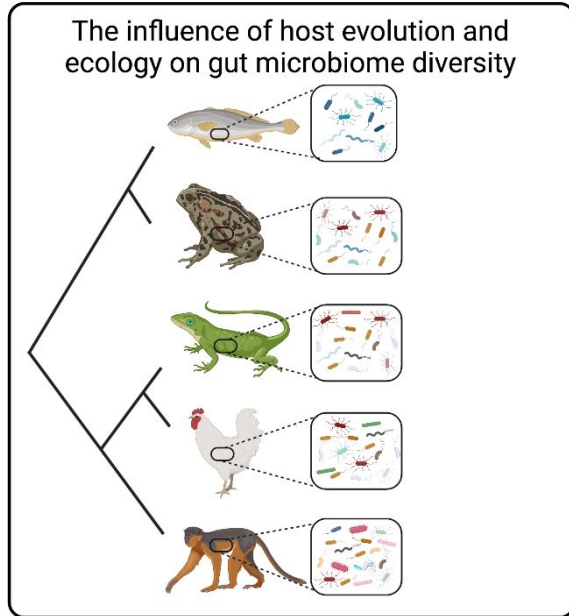


Microbial community assembly in the animal gut





The vertebrate gut microbiome: influence of host evolution & ecology



Georg Reischer



Andreas Farnleitner



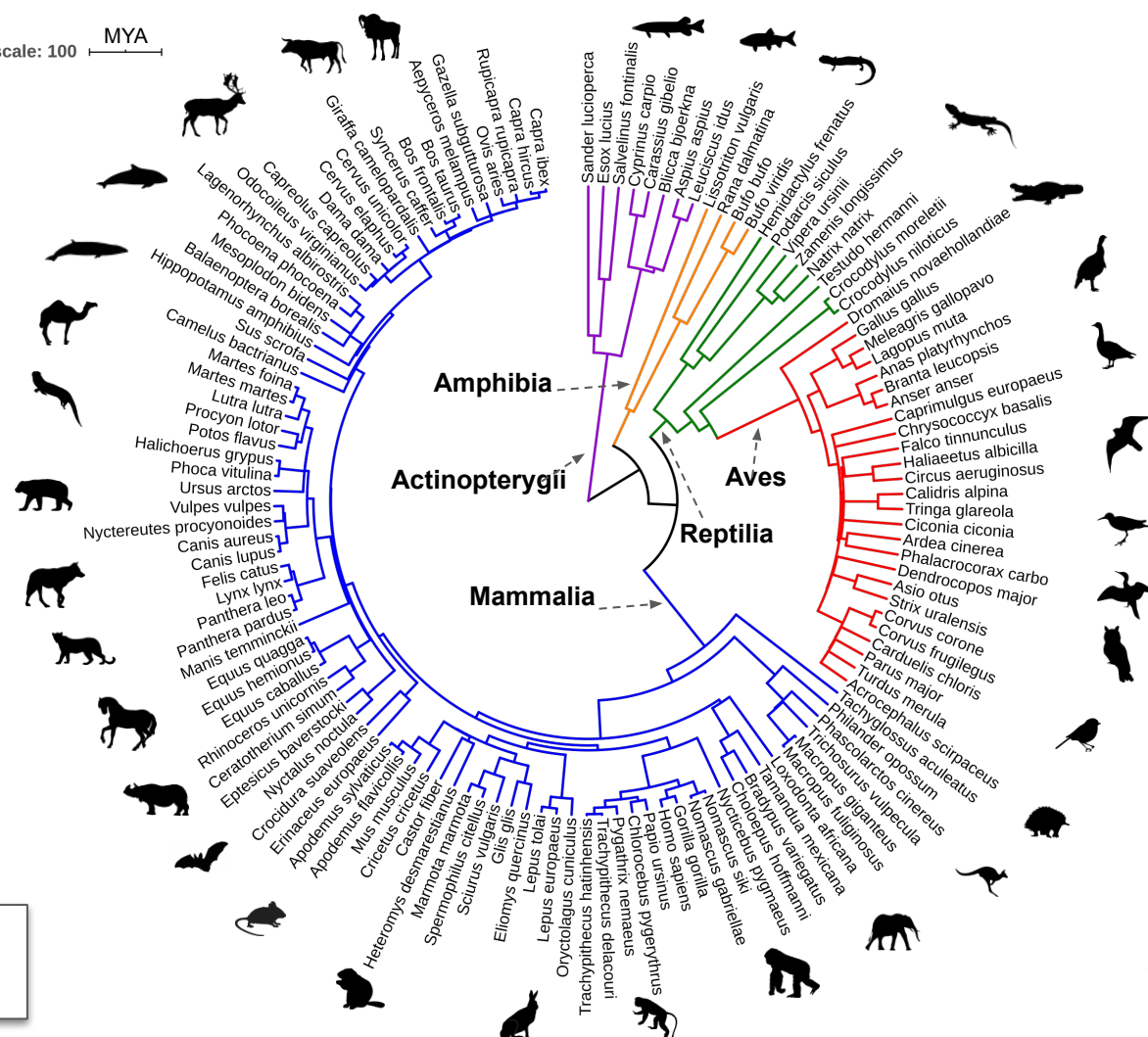
Gabriella Stalder



Chris Walzer

Youngblut et al., 2019.
Nature Communications

Tree scale: 100 MYA

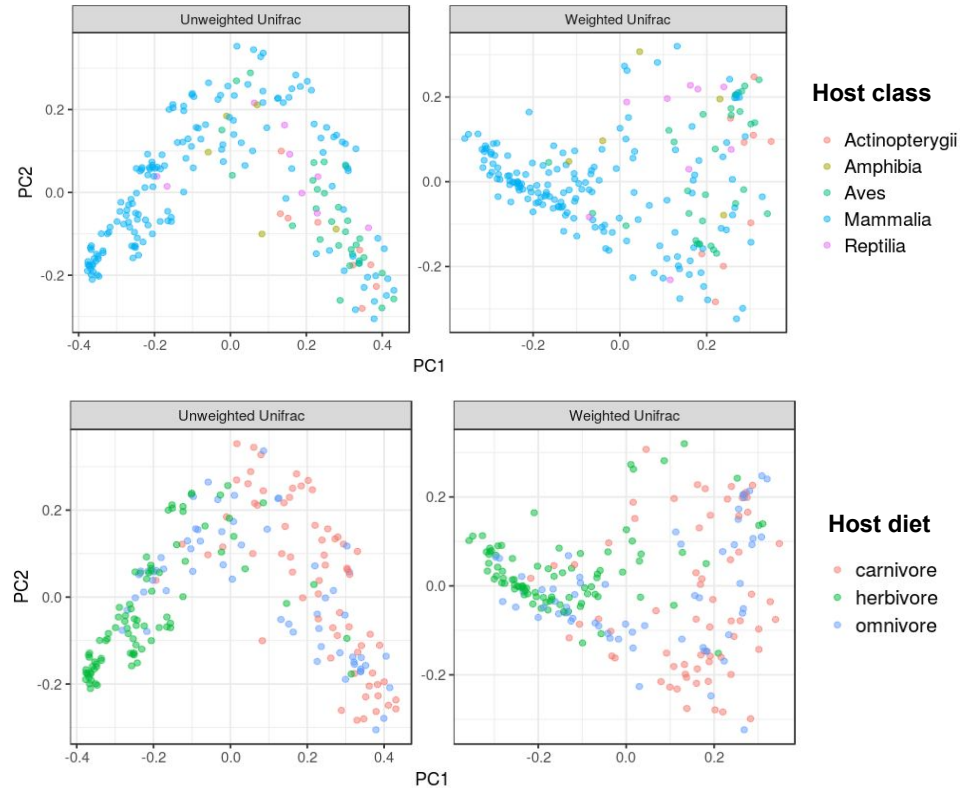
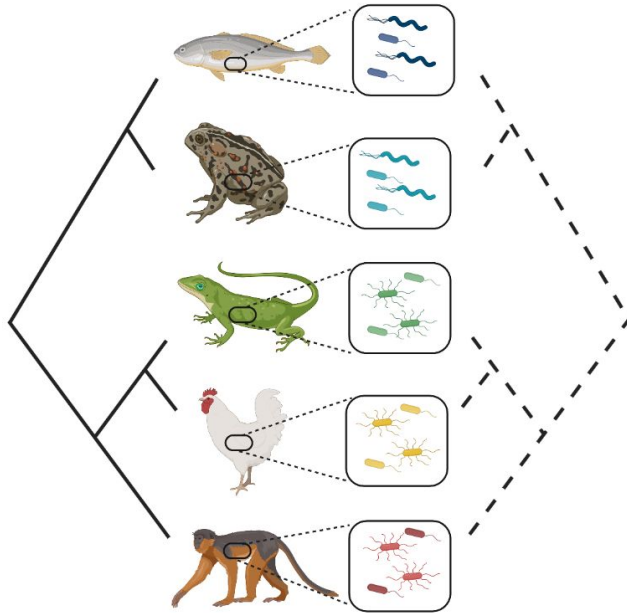


➤ 213 samples
➤ 128 species



Pattern of phyllosymbiosis?

Phyllosymbiosis?





Phylosymbiosis versus other factors

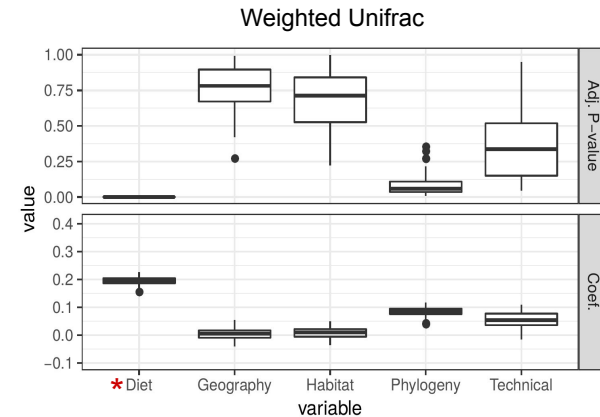
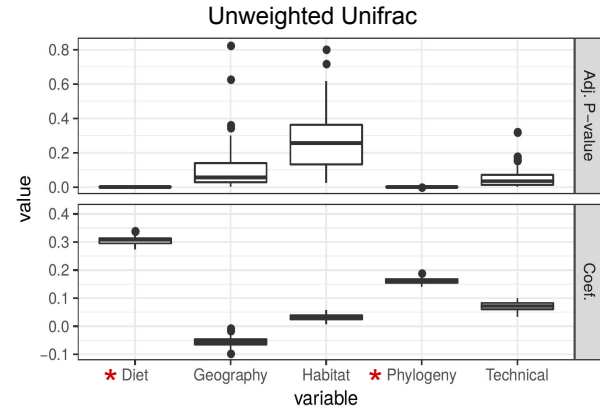
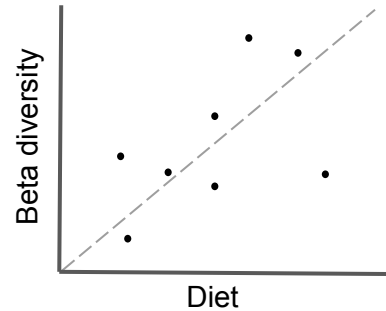
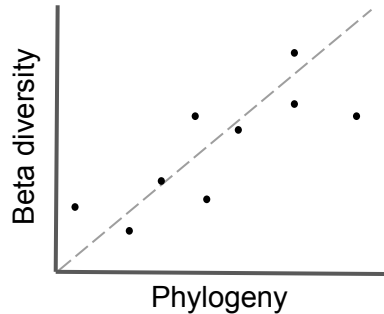
Multiple regression on matrices (MRM)

Phylogeny
(patristic distance)

| | S1 | S2 | S3 |
|----|-----|-----|-----|
| S1 | 0 | ... | ... |
| S2 | 0.1 | 0 | ... |
| S3 | 0.3 | 0.2 | 0 |

Diet
(Gower's distance)

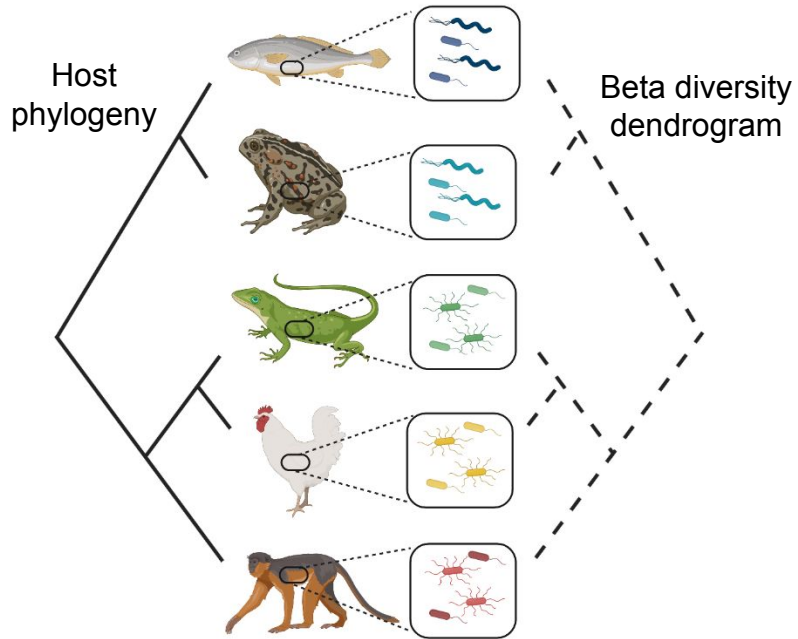
| | S1 | S2 | S3 |
|----|-----|-----|-----|
| S1 | 0 | ... | ... |
| S2 | 0.3 | 0 | ... |
| S3 | 0.6 | 0.1 | 0 |



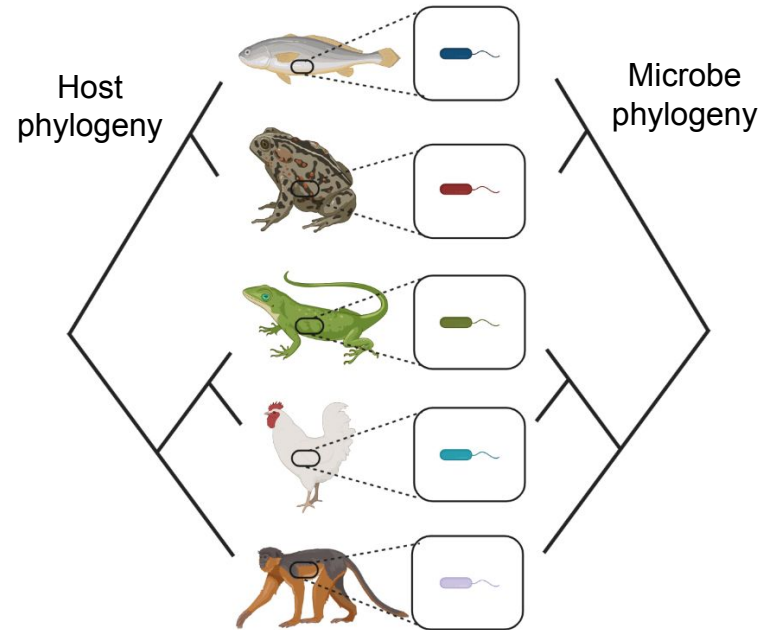


Phylosymbiosis versus cophylogeny

Phylobiosis: community similarity



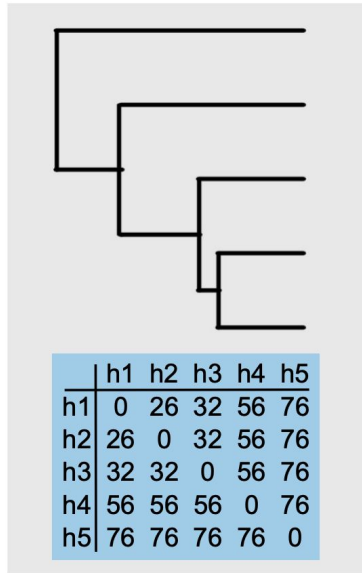
Cophylogeny: congruent evolutionary histories



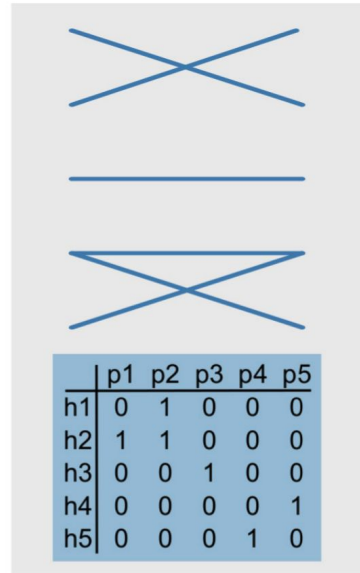


Phylosymbiosis versus cophylogeny

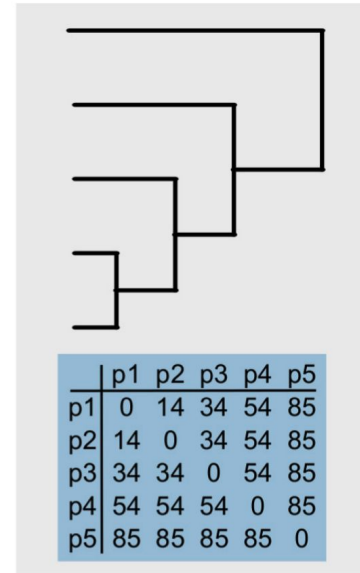
Cophylogeny for the entire microbial community Procrustean Approach to Cophylogeny (PACo)



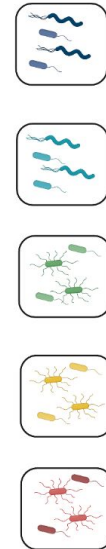
Host phylogeny



Host-microbe
associations



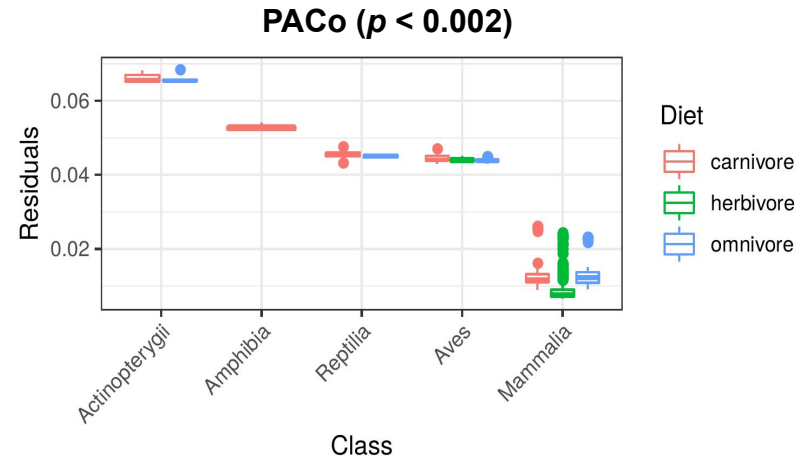
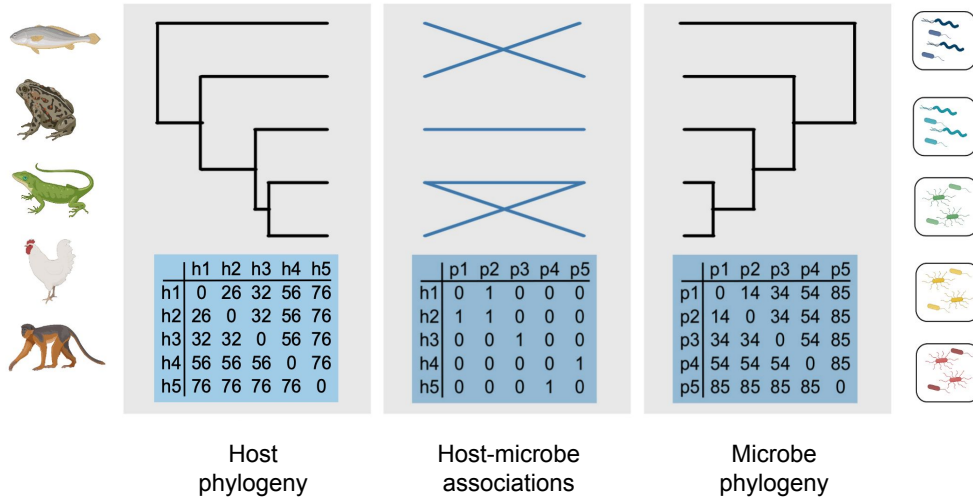
Microbe phylogeny





Pattern of cophylogeny?

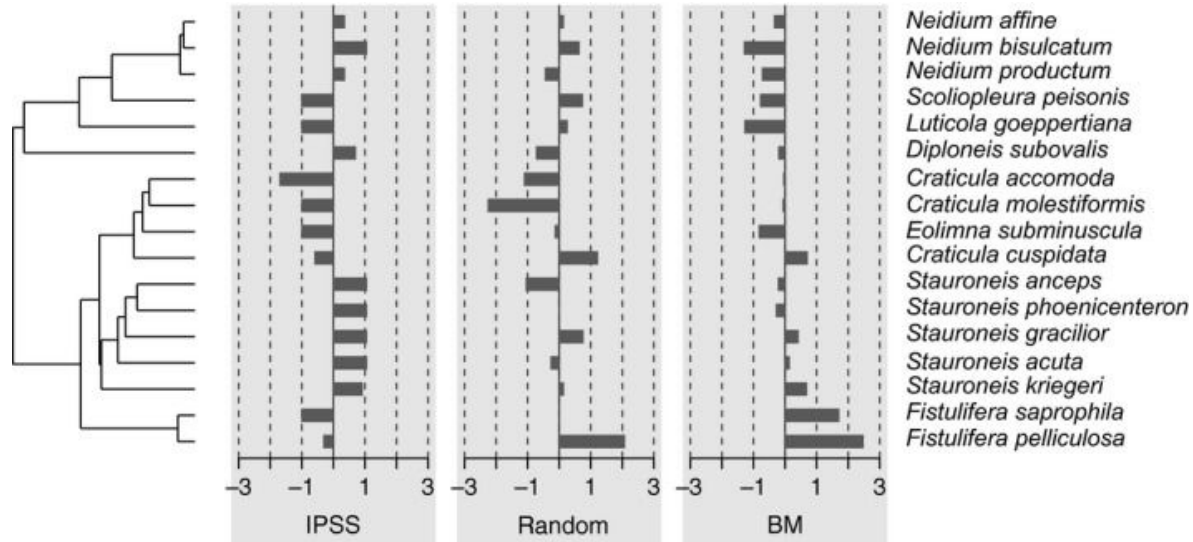
Cophylogeny for the entire microbial community Procrustean Approach to Cophylogeny (PACo)





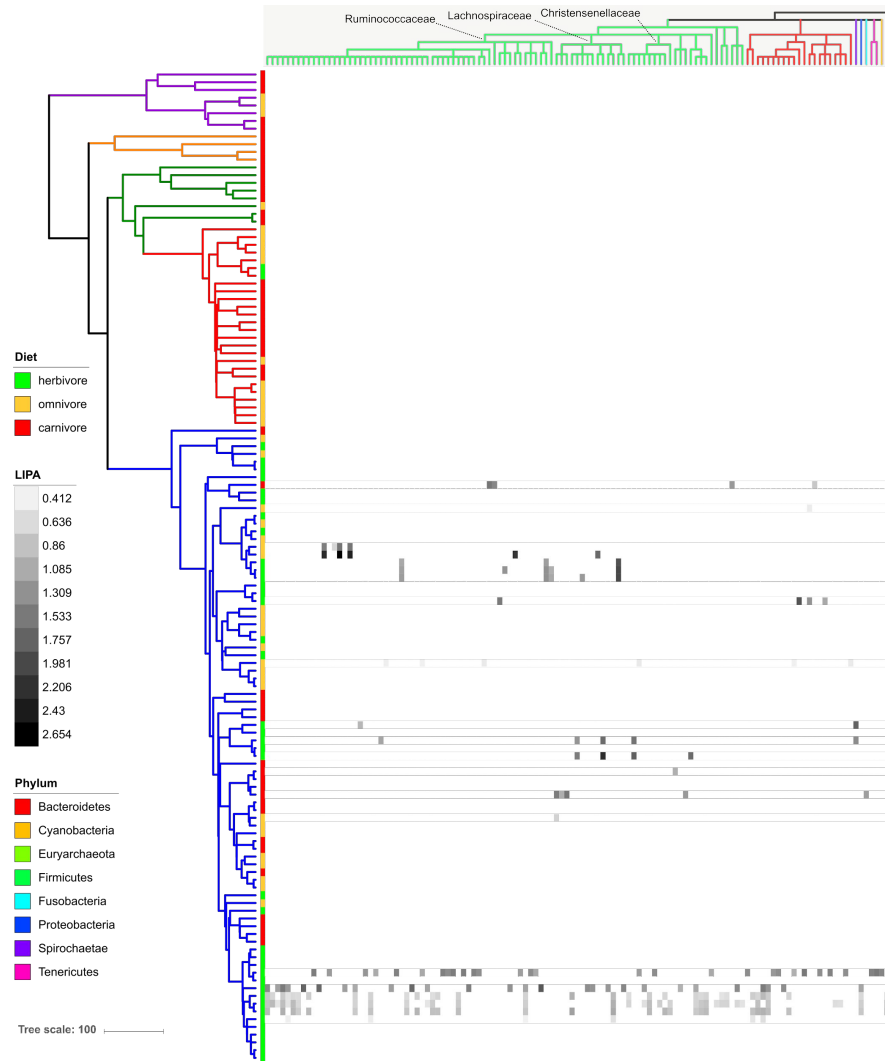
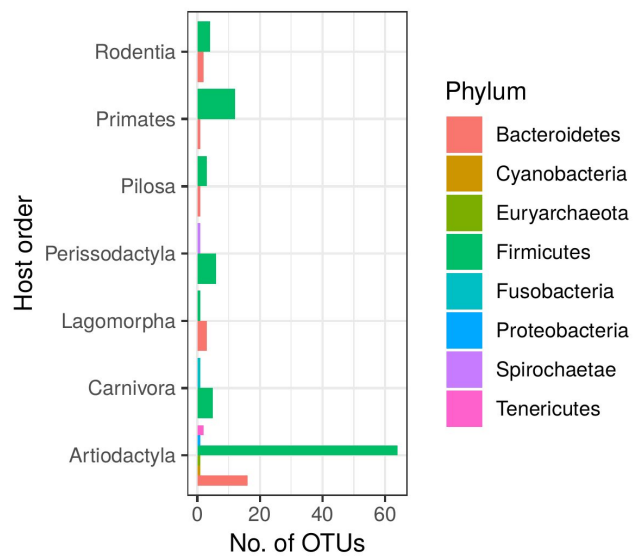
Specific microbe-host clade associations?

Phylogenetic signal: autocorrelation of ≥ 1 trait across the phylogeny



Specific host-microbe associations

LIPA: Local Indicator of Phylogenetic Association

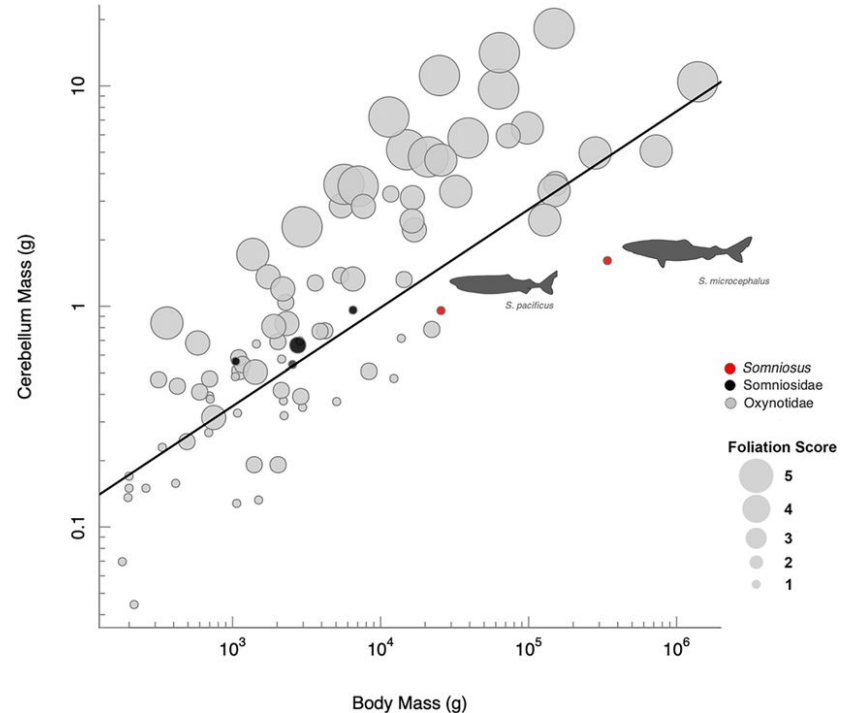




Trait-trait associations, while controlling for diet

- Phylogenetic Generalised Least Squares (PGLS)
- Trait₁ => host diet
- Trait₂ => microbial diversity

Example of PGLS on body morphology traits

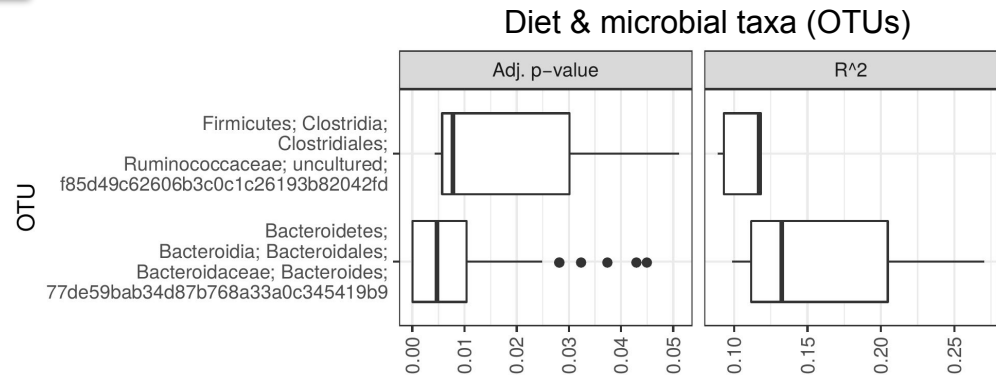
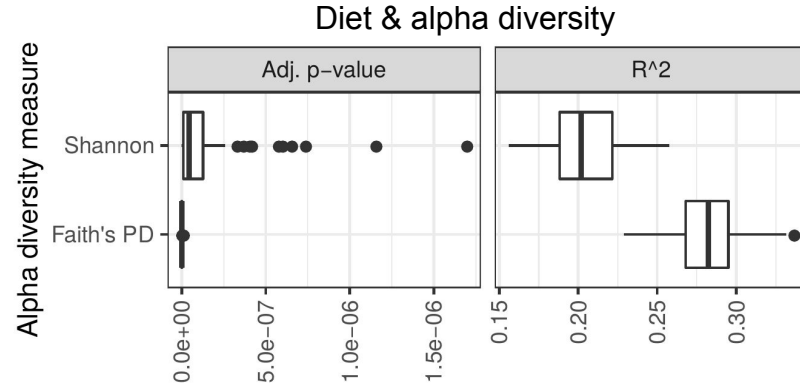




Associations with host diet

Association with diet while controlling for host phylogeny?

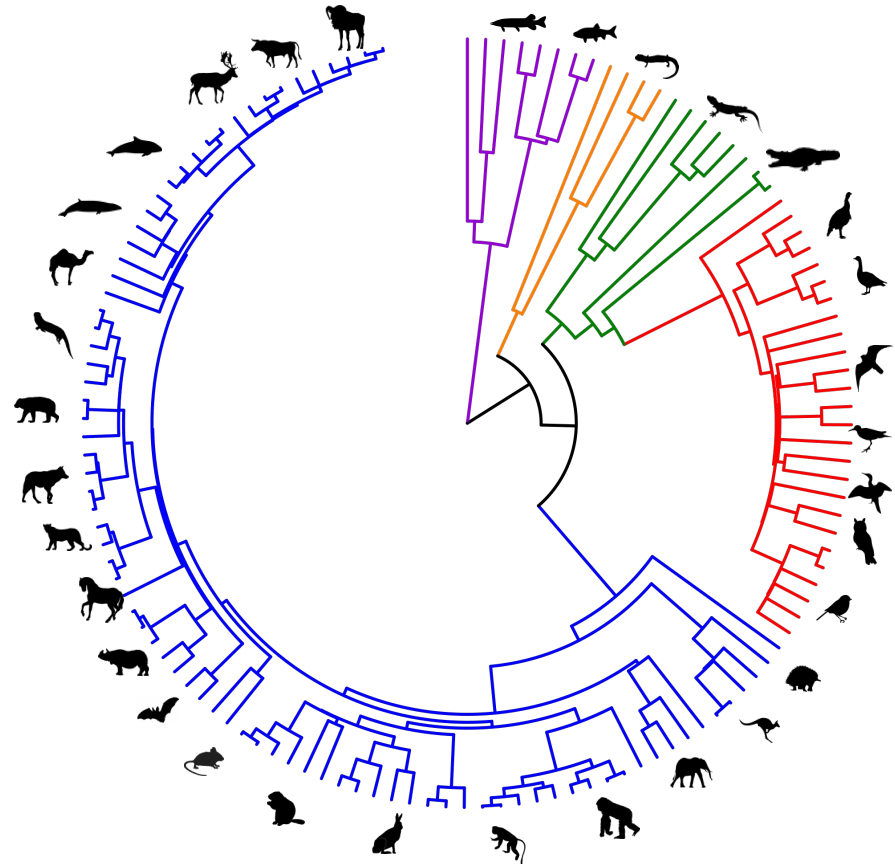
- Phylogenetic Generalised Least Squares (PGLS)



Summary

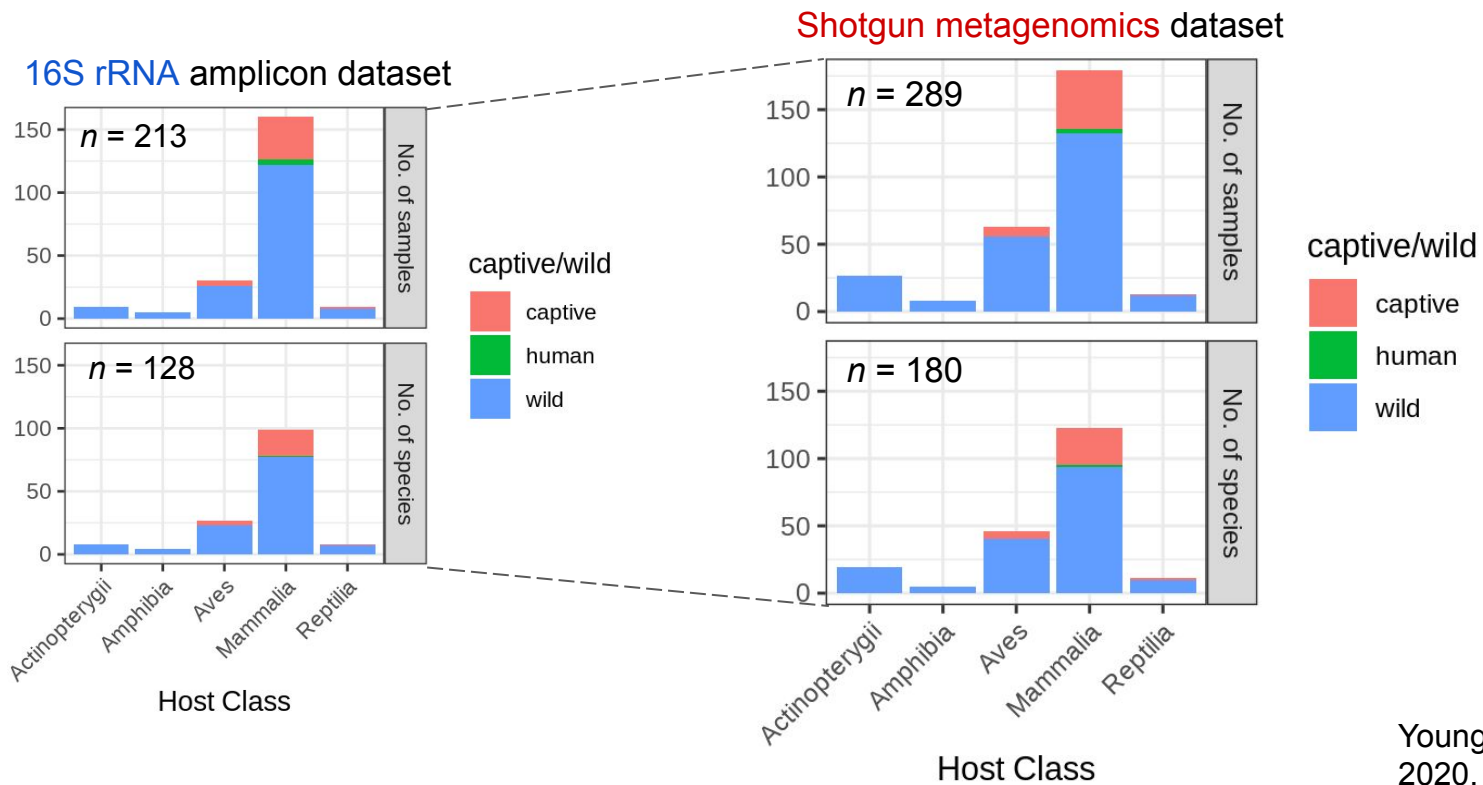


- **Phylosymbiosis pattern detected**
 - Multiple regression on matrices (MRM)
 - Accounting for other factors
- **Significant cophylogeny pattern**
 - PACo
 - Strongest for mammals
- **Diet stronger factor than host phylogeny**
 - Selecting for total diversity versus particular microbes (PGLS)

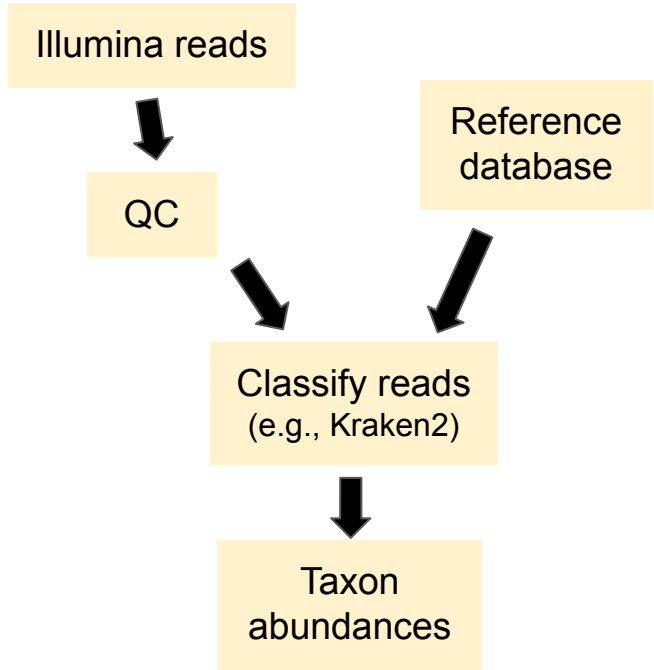




Beyond 16S rRNA: utilizing metagenomics



Metagenome profiling: improving the reference database



Genome Taxonomy Database (GTDB)



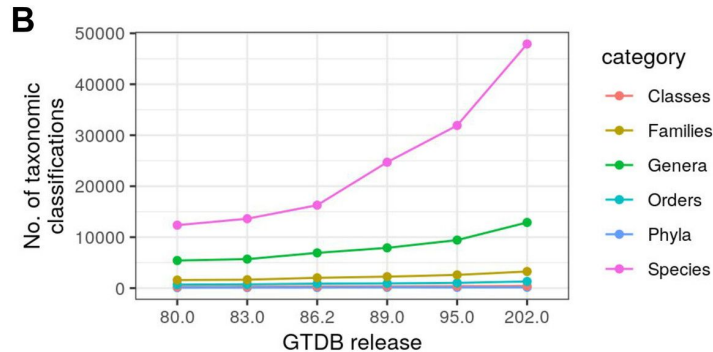
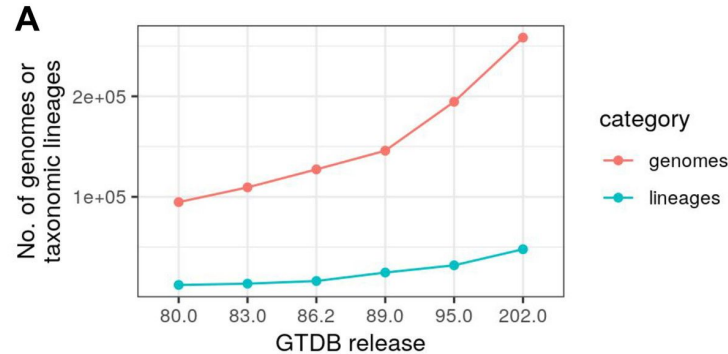
Phil Hugenholtz



Metagenome profiling: improving the reference database

Problem: rapidly expanding size of the GTDB database.
How to scale?

The GTDB is expanding rapidly



Solution

Struo2: efficient metagenome profiling database construction for ever-expanding microbial genome datasets

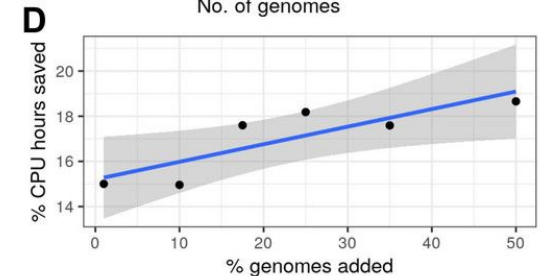
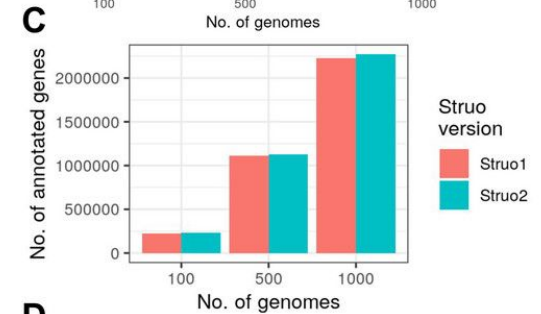
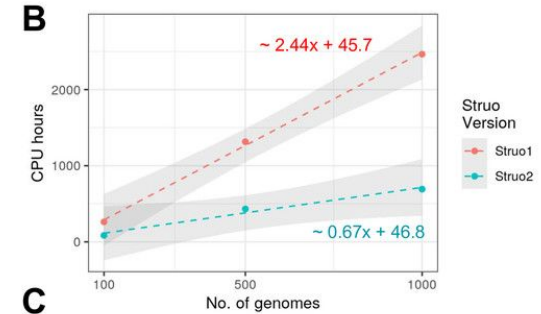
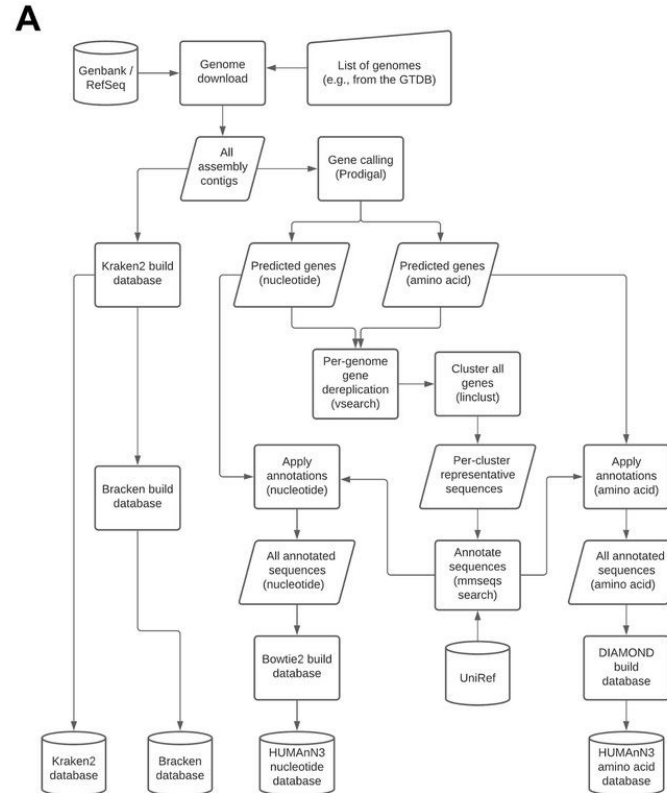


Metagenome profiling: improving the reference database

Struo2: efficient metagenome profiling database construction for ever-expanding microbial genome datasets

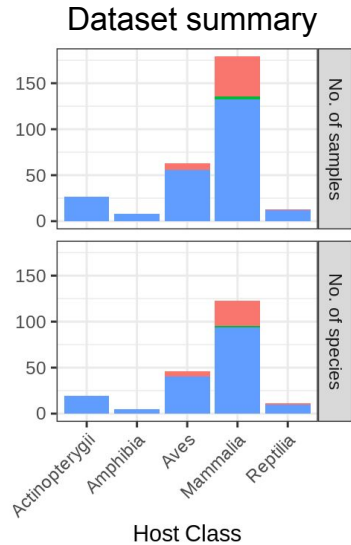


MMseqs2



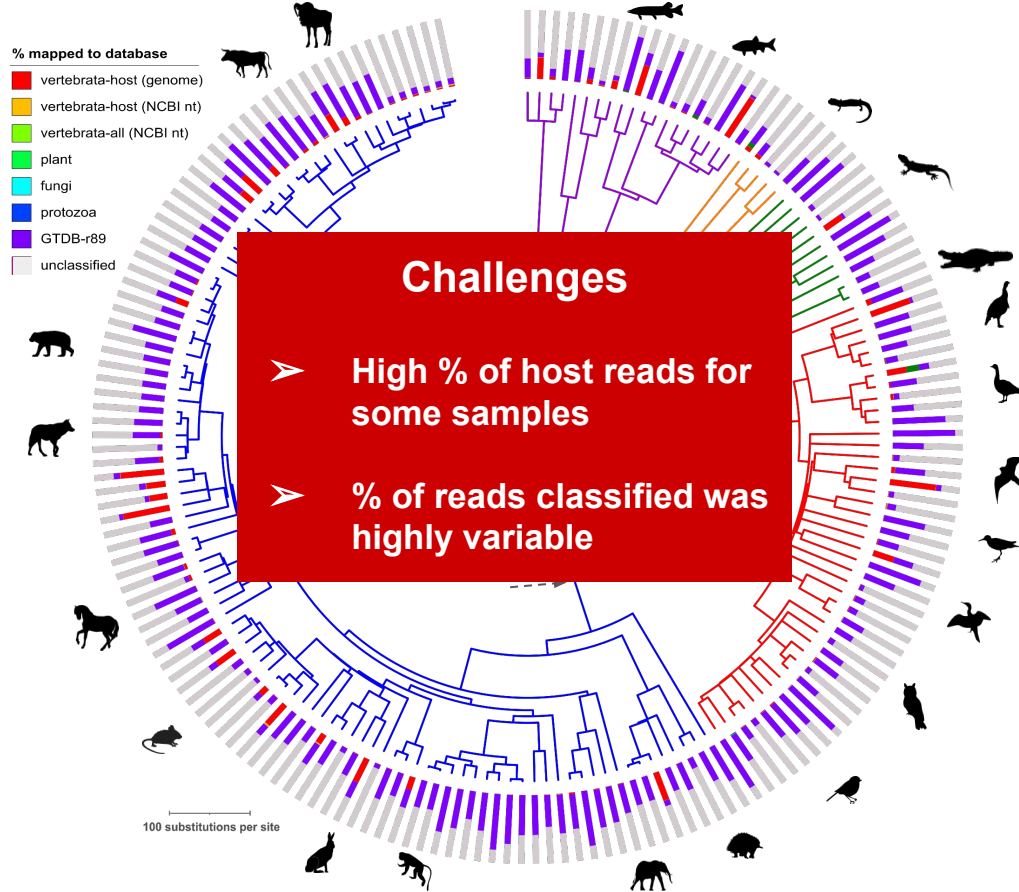


Beyond 16S rRNA: utilizing metagenomics



% mapped to database

- vertebrata-host (genome)
- vertebrata-host (NCBI nt)
- vertebrata-all (NCBI nt)
- plant
- fungi
- protozoa
- GTDB-r89
- unclassified



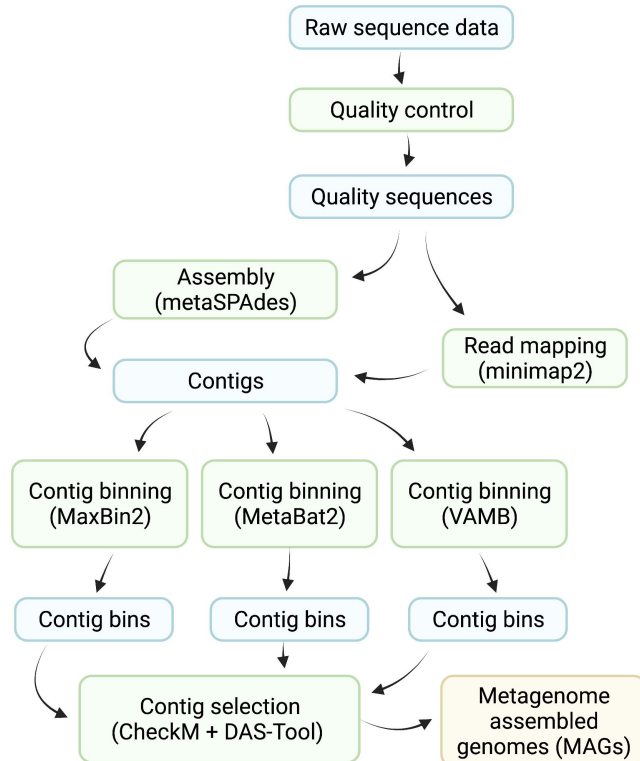
Challenges

- High % of host reads for some samples
- % of reads classified was highly variable

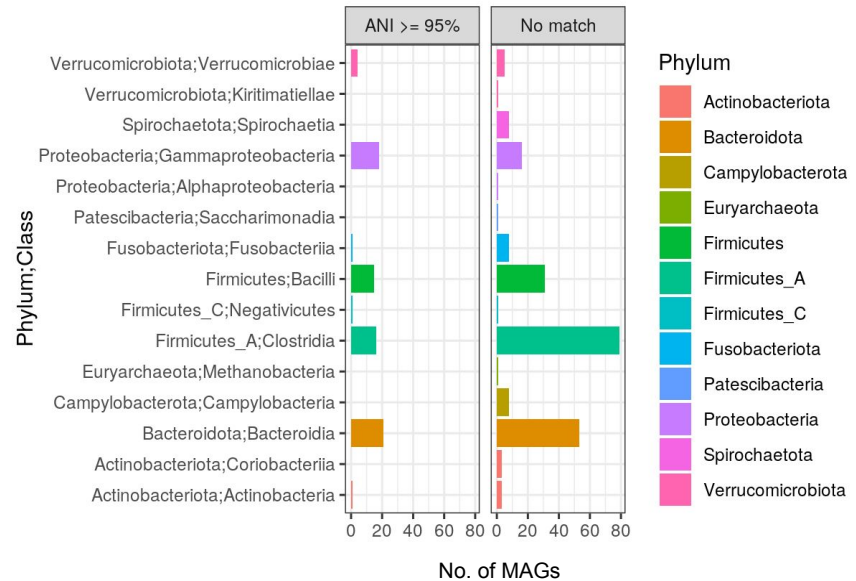
Expanding known microbial genomic diversity via metagenome assembly



Metagenome assembly pipeline (simplified)



Metagenome-assembled genomes (MAGs)

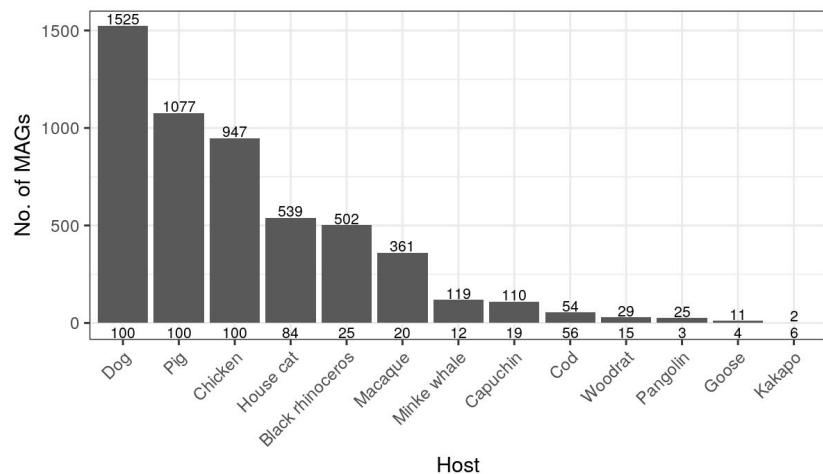


- 296 non-redundant, quality MAGs
- 248 species

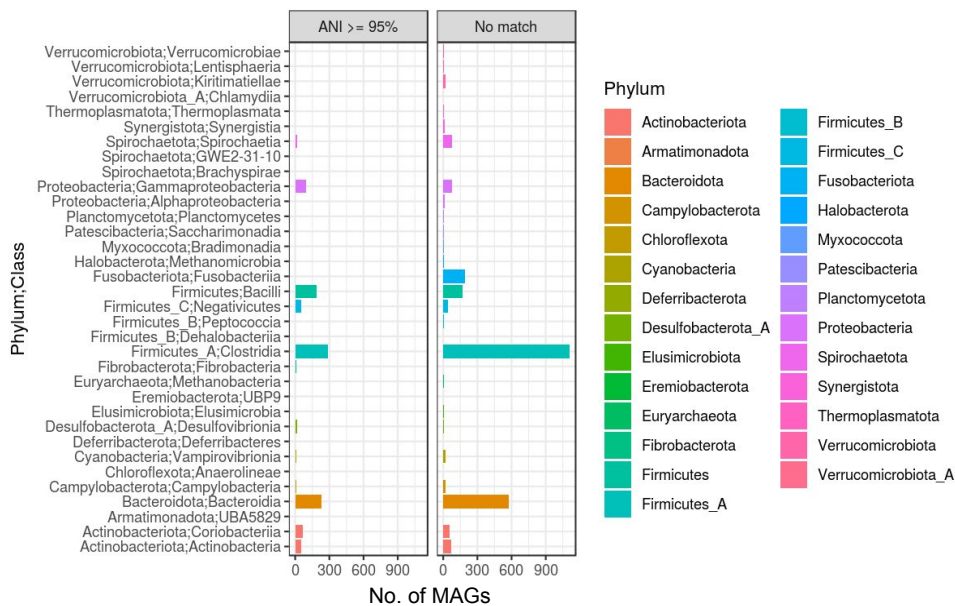


Expanding known microbial genomic diversity via metagenome assembly

Extending the assembly pipeline to 14 other vertebrate gut datasets



Metagenome-assembled genomes (MAGs)



- 5596 non-redundant, quality MAGs
- 1522 species

Using deep learning to improve metagenome assemblies



Bioinformatics, 36(10), 2020, 3011–3017

doi: 10.1093/bioinformatics/btaa124

Advance Access Publication Date: 25 February 2020

Original Paper

OXFORD

Genome analysis

DeepMAsED: evaluating the quality of metagenomic assemblies

Olga Mineeva^{1,2,†}, Mateo Rojas-Carulla^{1,†}, Ruth E. Ley³, Bernhard Schölkopf¹ and Nicholas D. Youngblut ^{3,*}

Goal: identify errors in metagenome assemblies via deep learning



Olga Mineeva



Mateo Rojas-Carulla



Bernhard Schölkopf

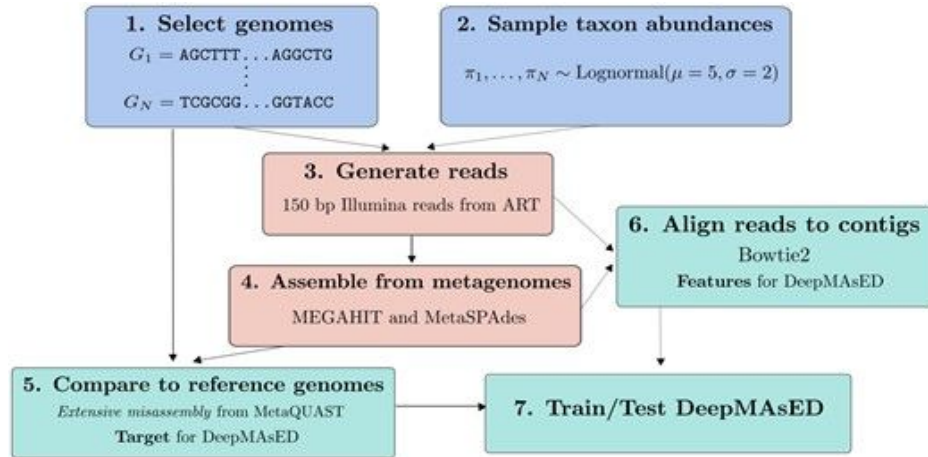


Gunnar Rätsch

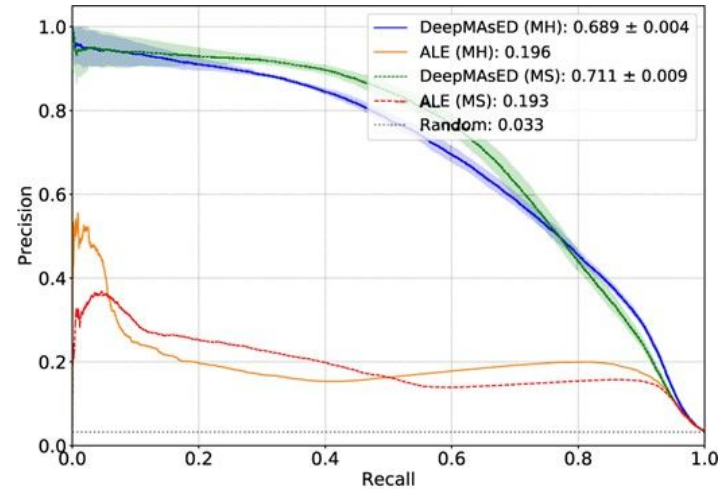


Using deep learning to improve metagenome assemblies

Task: map reads back to contigs & predict misassembled contigs

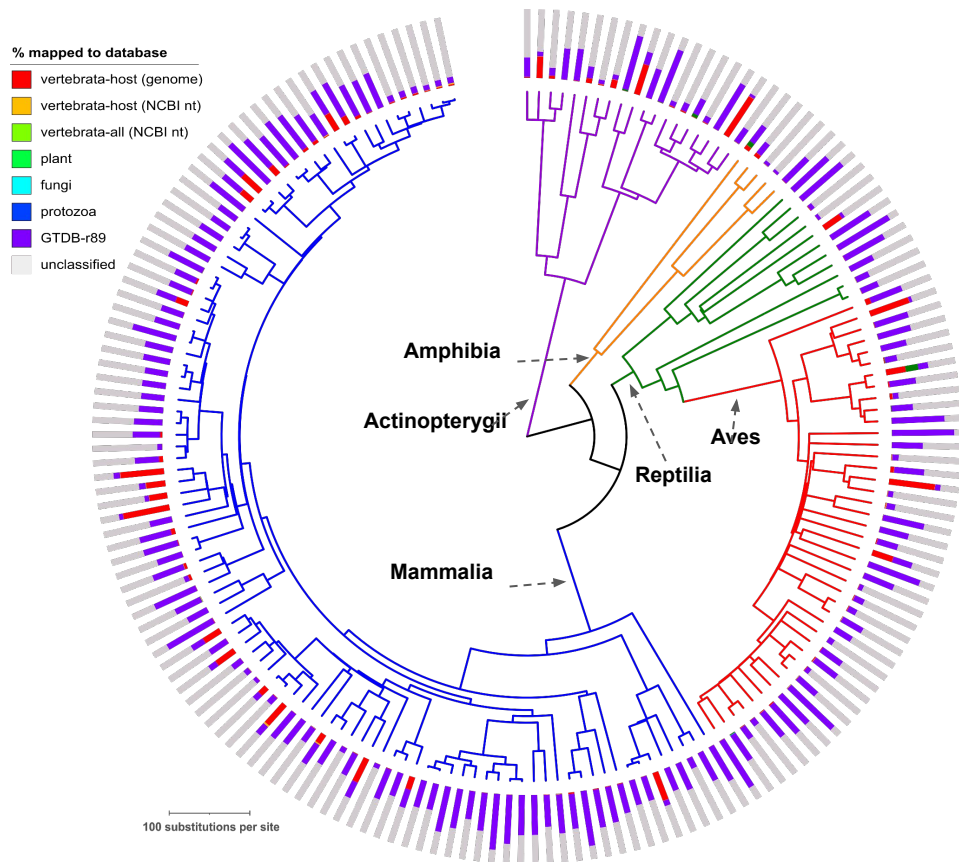


DeepMAS-ED is substantially better than the state-of-the-art

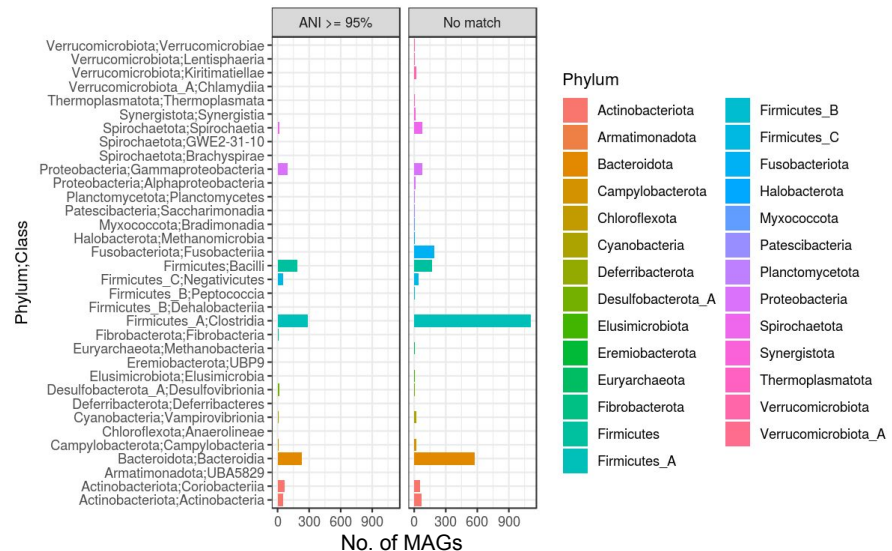




MAGs significantly increase the % of reads classified



Metagenome-assembled genomes (MAGs)

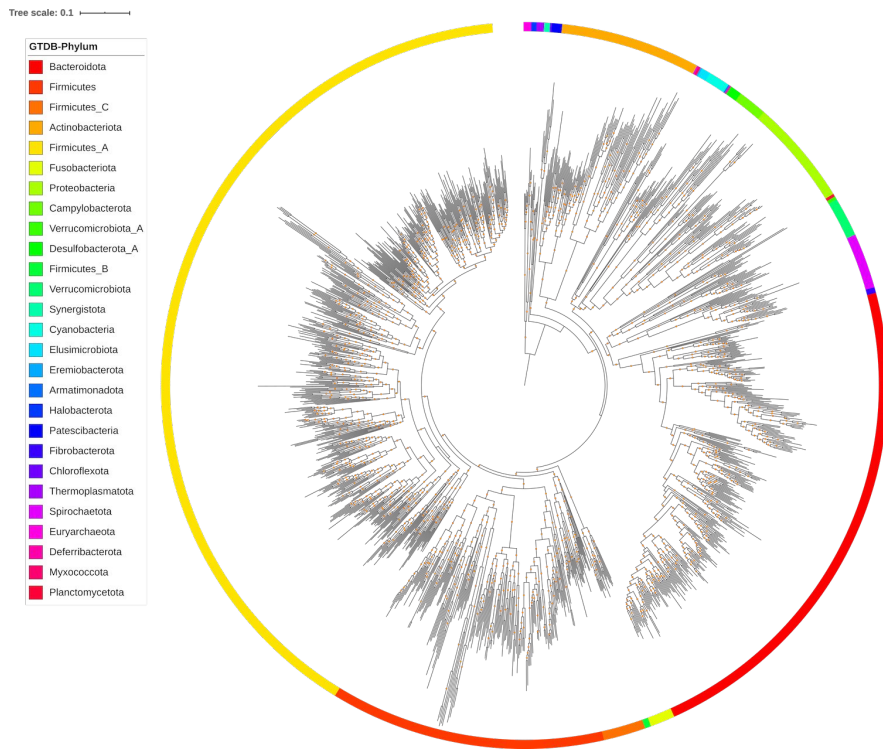


- 5596 non-redundant, quality MAGs
- 1522 species

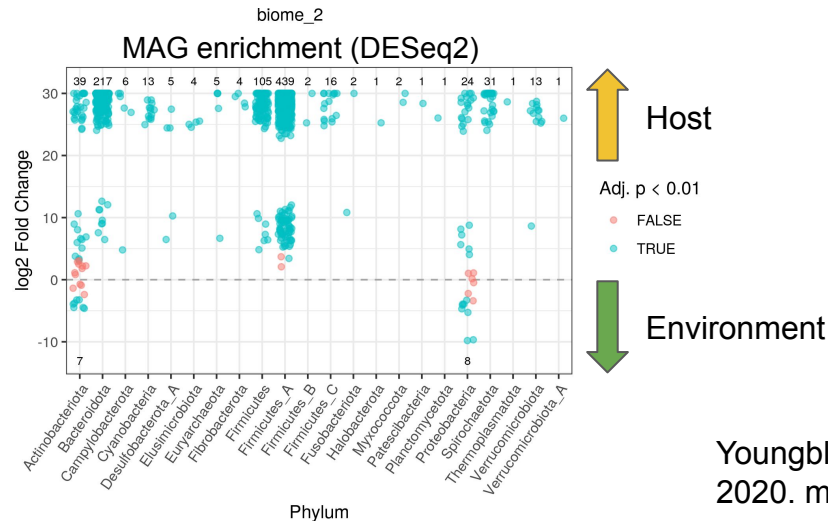
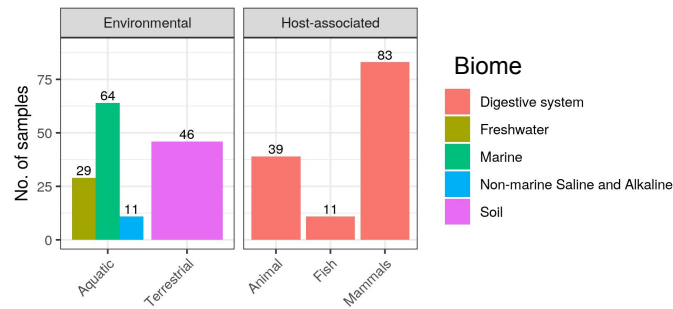


Most MAGs are host-associated

1522 species-level MAGs
(5596 at the strain-level)

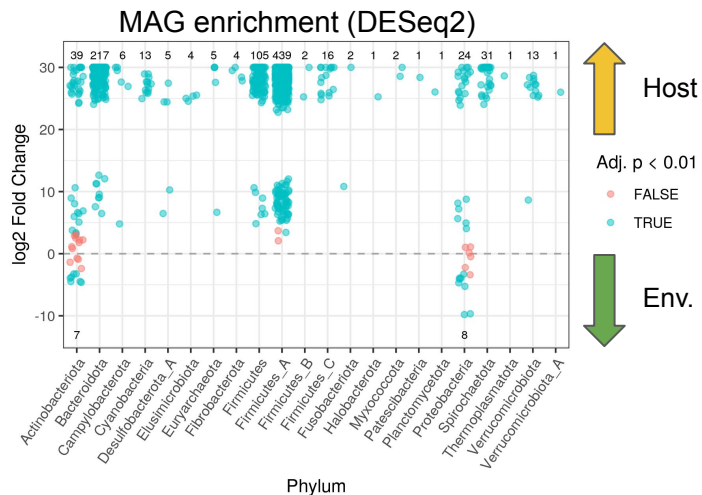


“host-environment” metagenome dataset

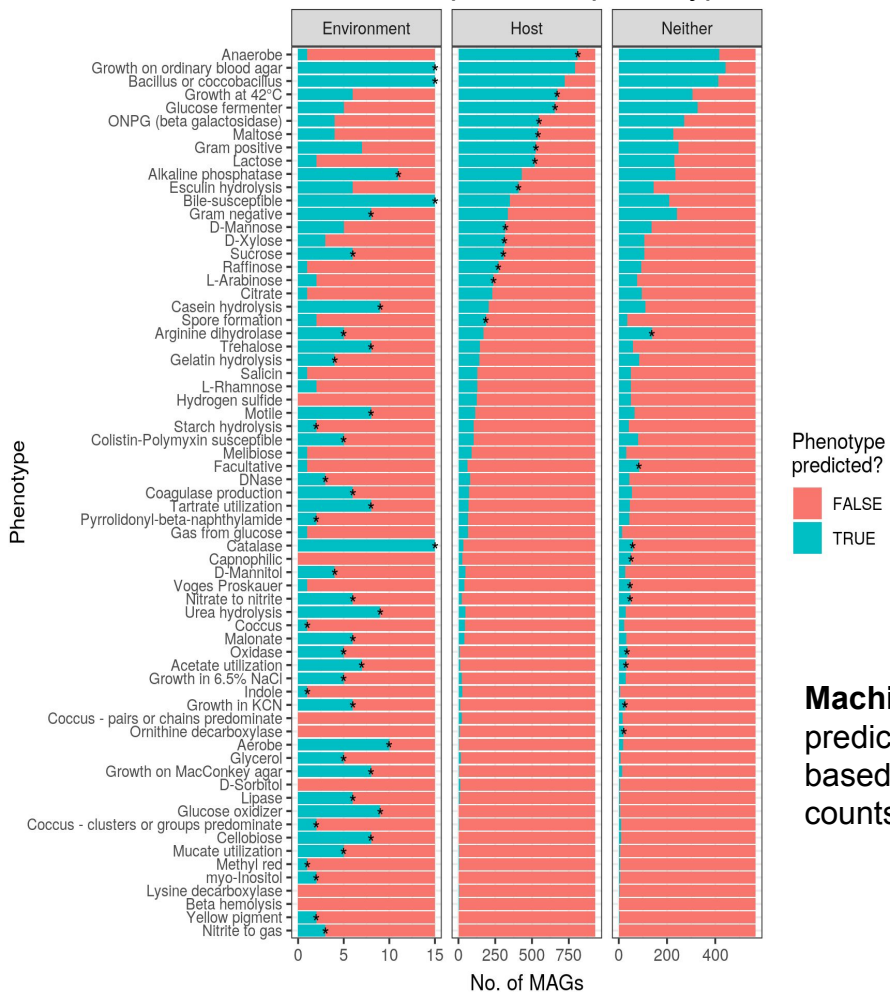




Host & env. associated traits



MAGs with predicted phenotypes

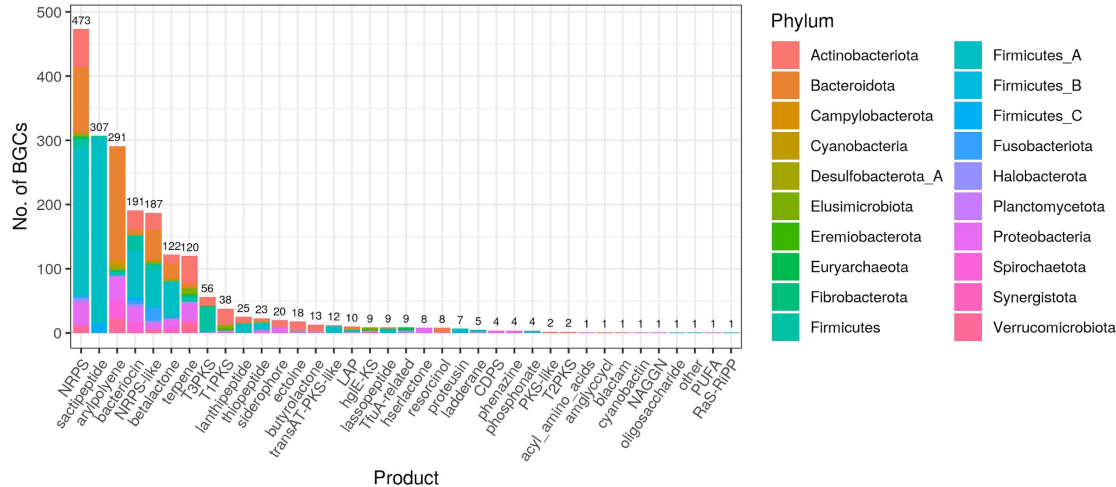


Machine learning:
predicting traits
based on Pfam
counts



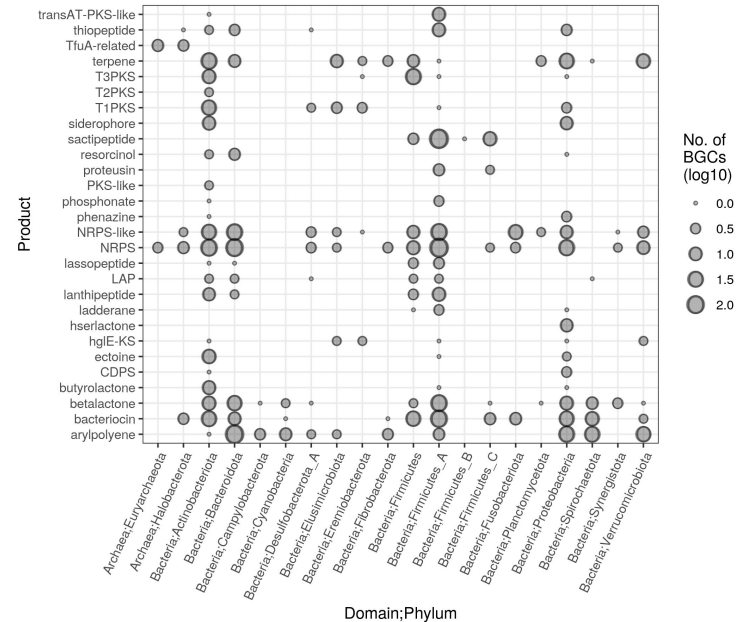
A large diversity of biosynthetic gene clusters (BGCs)

BGCs identified with antiSMASH



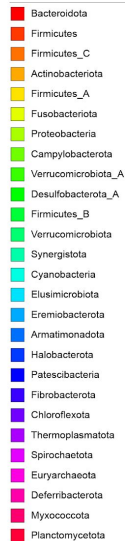
- 1986 BGCs
- High novelty based on MIBiG

BGCs identified with antiSMASH

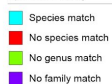


Species with ≥ 3 BGCs ($n = 233$)

1) Phylum



2) GTDB novelty



3) DESeq2 Adj. $P < 0.01$



5)

50
40
30
20
10

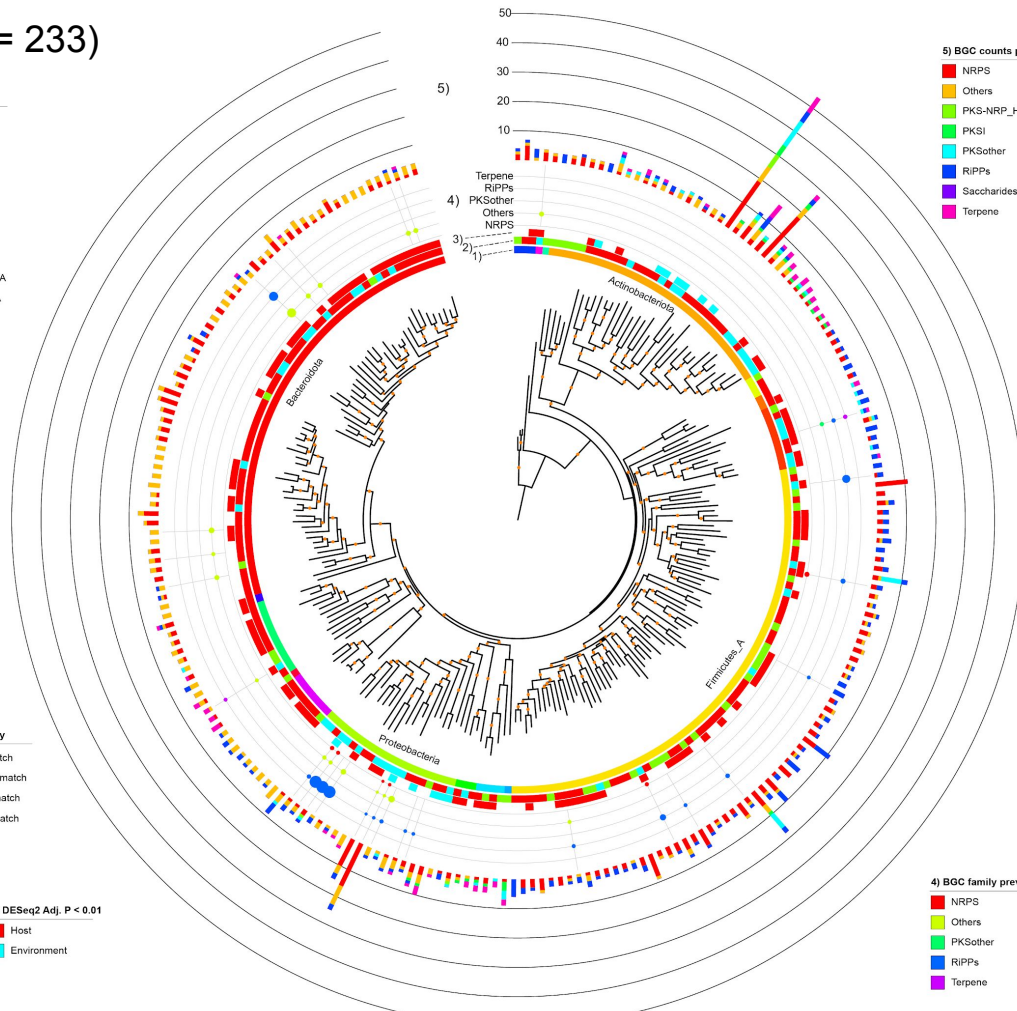
4)

Terpene
RiPPs
Others
PKSother
NRPS

5) BGC counts per MAG



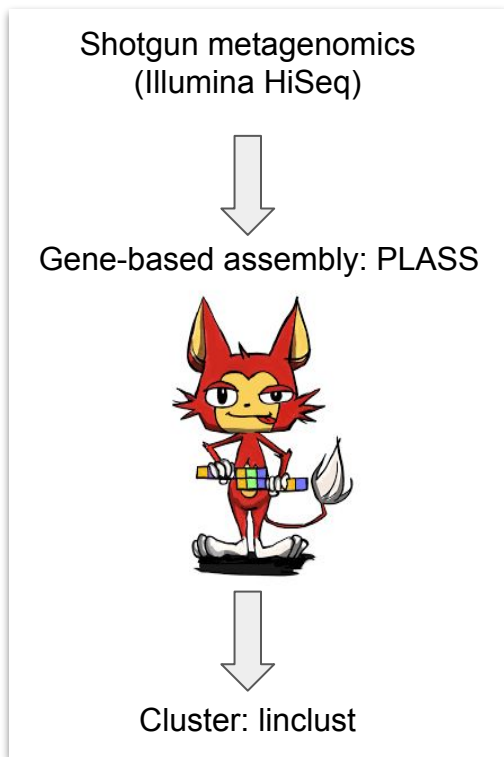
4) BGC family prevalence



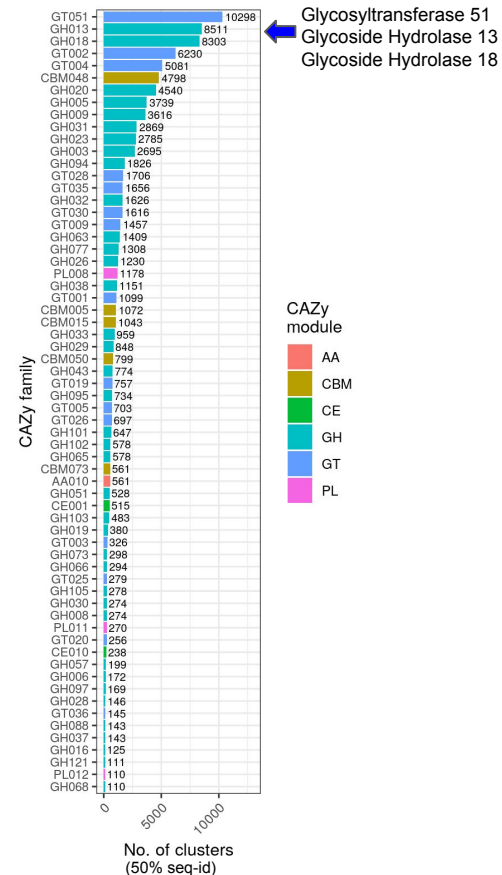
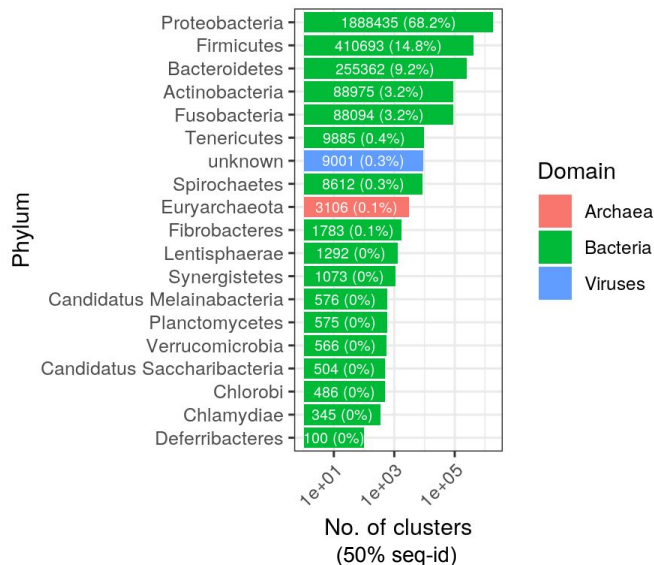
Tree scale: 1



A large amount of gene-level diversity

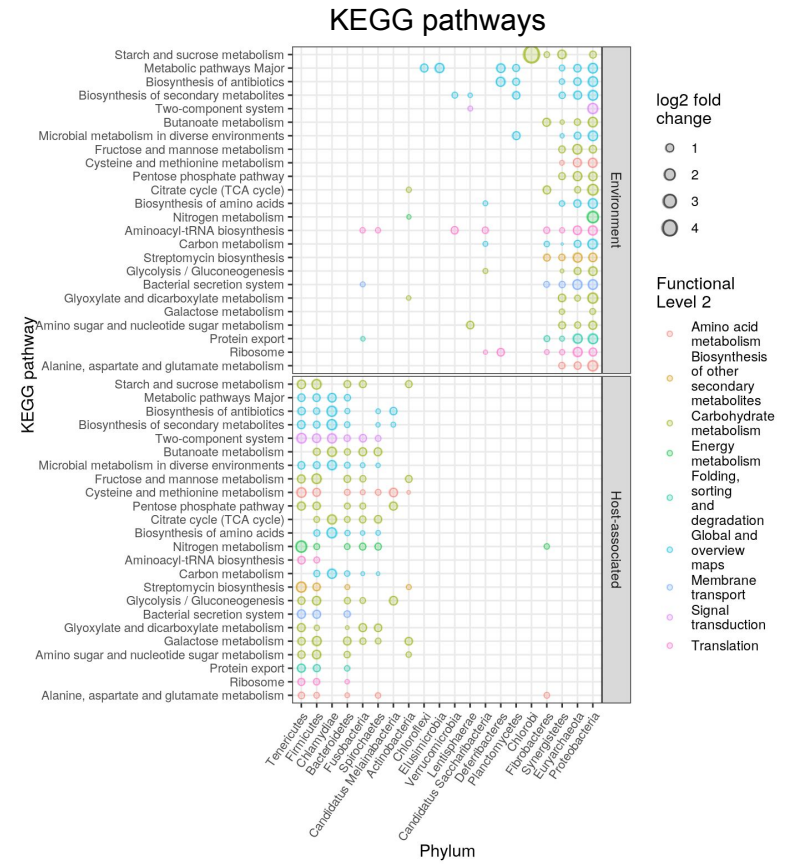
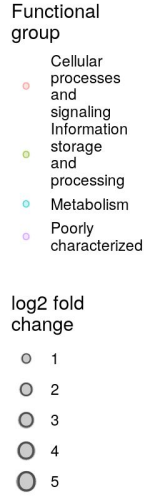
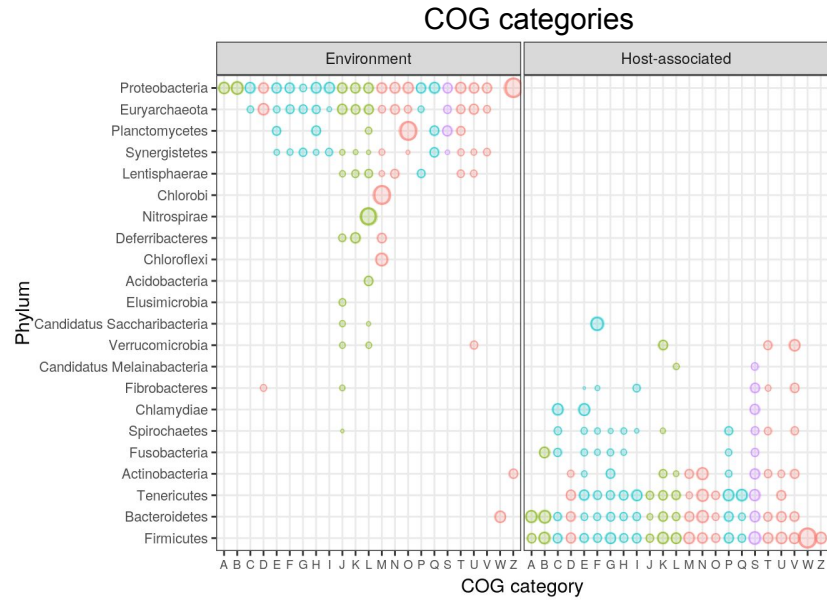


>150 million gene sequences





Functional redundancy for major gene categories



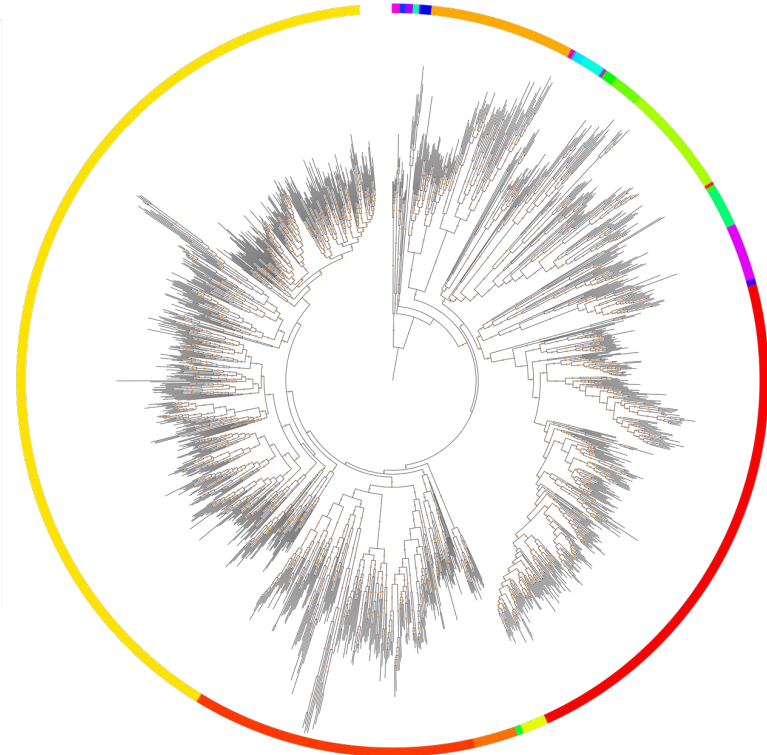
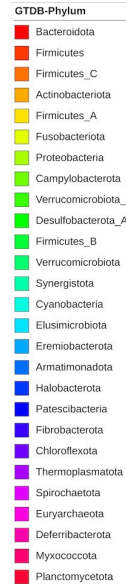
Summary



- Novel taxonomic diversity
- Most MAGs were host-associated
- Traits enriched in host- and environment-specific species
- Large BGC diversity
 - RiPPs & NRPS most prevalent
- Large gene diversity
 - Functional redundancy

1522 species-level MAGs
(5596 at the strain-level)

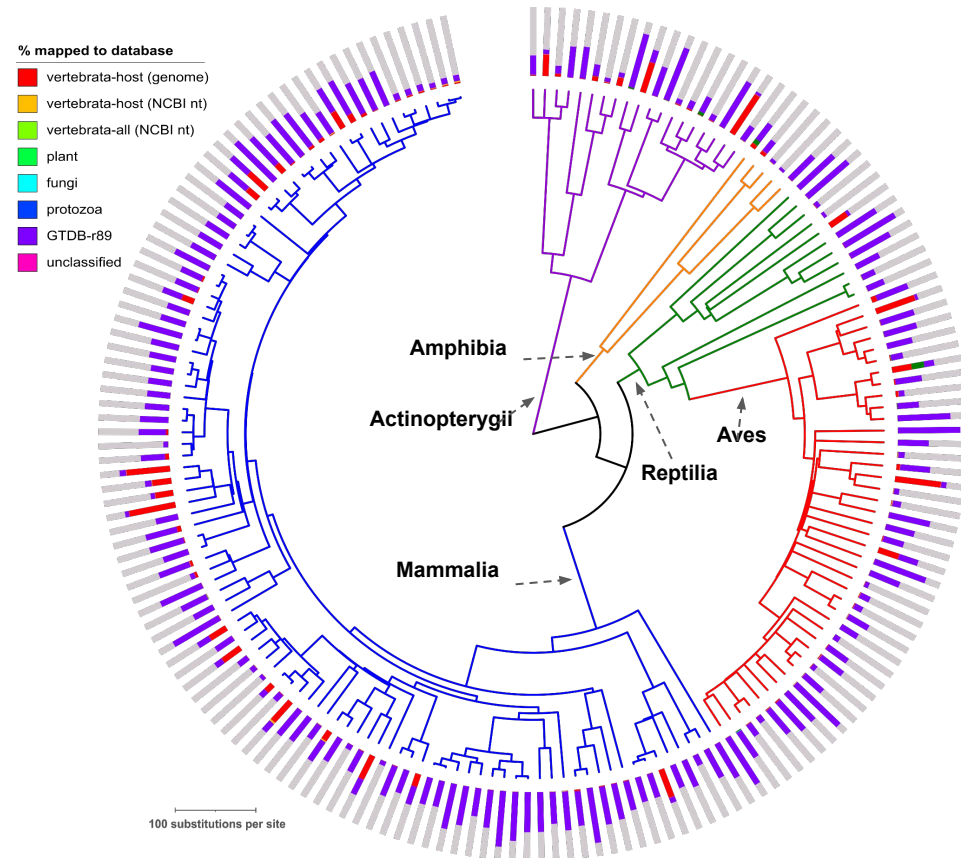
Tree scale: 0.1



Future directions



- Species-level phyllosymbiosis & cophylogeny
 - Genome phylogeny
- Phyllosymbiosis at the functional level
 - Genes/pathways
- Diet-function associations
- Adaptation to the gut environment & symbiosis



Department of Microbiome Science (Ley Lab)



LEY LAB



Georg Reischer



Andreas Farnleitner



Gabriella Stalder



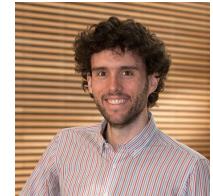
Chris Walzer



Silke Dauser



Tony Walters



Jacobo de la Cuesta



Ruth Ley

