



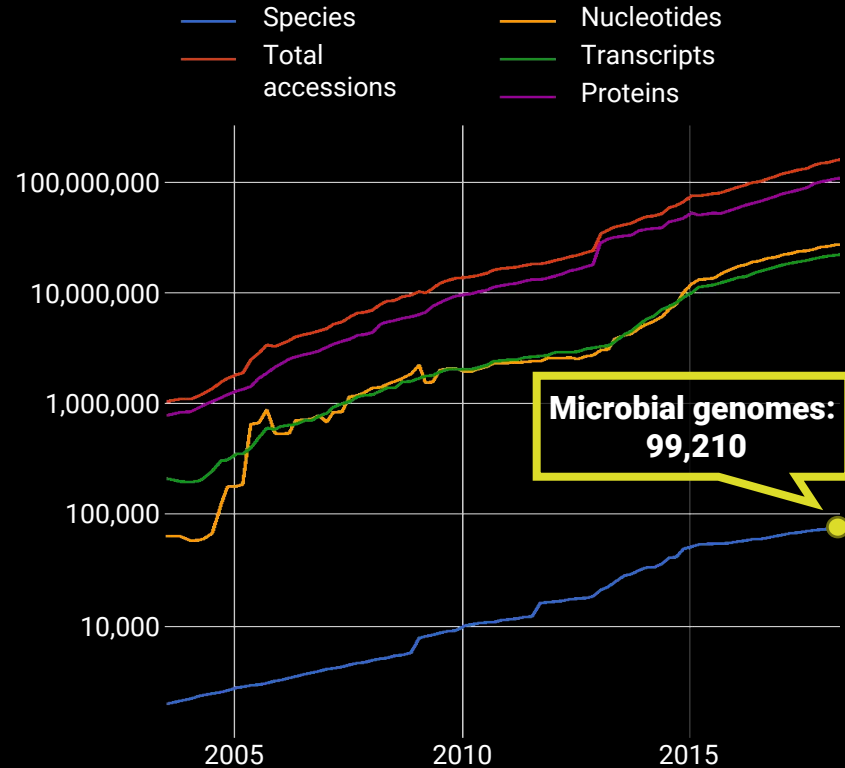
Multi-resolution phylogenetic analysis of known and unexplored members of the human microbiome with PhyloPhlAn 2

**3rd Workshop on Statistical and Algorithmic Challenges
in Microbiome Data Analysis**

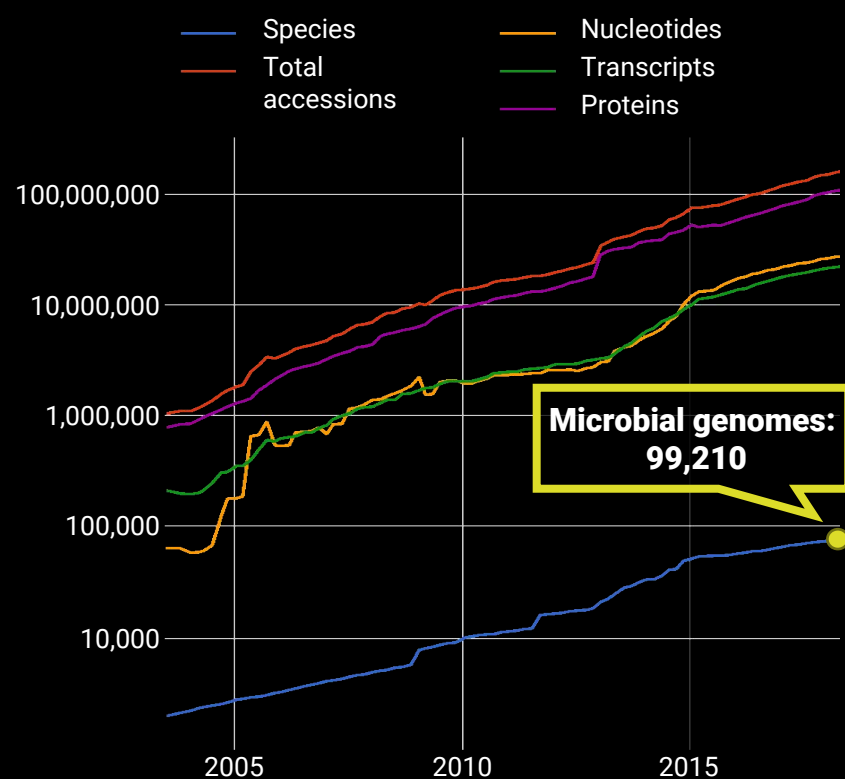
Francesco Asnicar - University of Trento
Computational Metagenomics lab - Prof. Nicola Segata



Increasing availability of genomes from isolates and MAGs

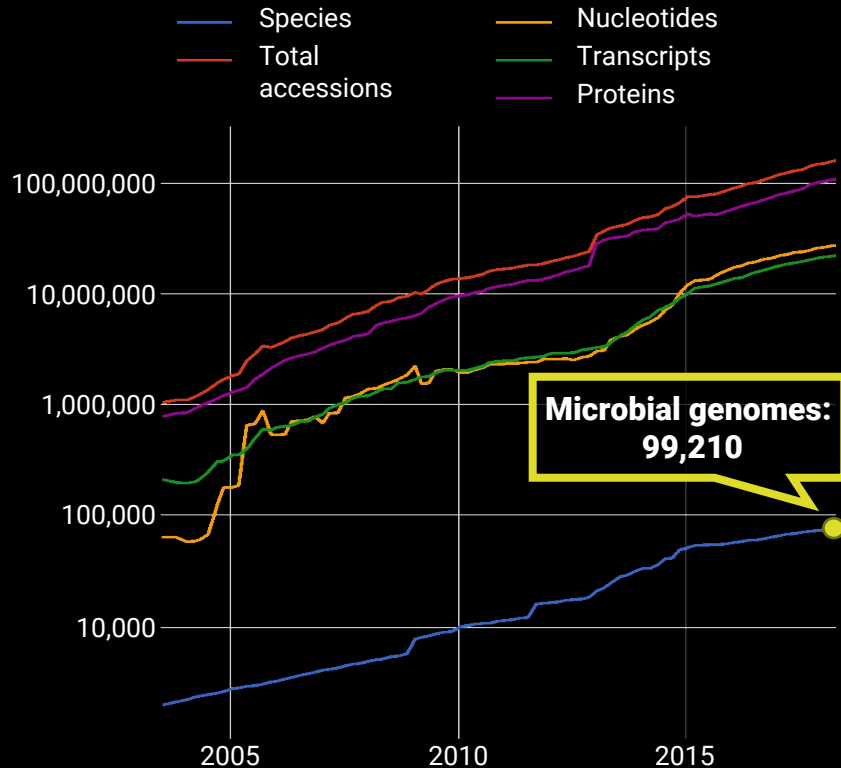


Increasing availability of genomes from isolates and MAGs



10 most sequenced species	# gen.
<i>Streptococcus pneumoniae</i>	>7500
<i>Staphylococcus aureus</i>	>6500
<i>Escherichia coli</i>	>5000
<i>Mycobacterium tuberculosis</i>	>3800
<i>Pseudomonas aeruginosa</i>	>2000
<i>Salmonella enterica</i>	>1700
<i>Acinetobacter baumannii</i>	>1400
<i>Klebsiella pneumoniae</i>	>1100
<i>Campylobacter coli</i>	>1000
<i>Neisseria meningitidis</i>	>1000

Increasing availability of genomes from isolates and MAGs



10 most sequenced

Streptococcus pneumoniae

Staphylococcus aureus

Escherichia coli

Mycobacterium tuberculosis

Pseudomonas aeruginosa

Salmonella enterica

Acinetobacter baumannii

Klebsiella pneumoniae

Campylobacter jejuni

Neisseria meningitidis

nature microbiology

ARTICLES

DOI: 10.1038/s41564-017-0012-7

OPEN

Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life

Donovan H. Parks¹, Christian Rinke², Maria Chuvpochina, Pierre-Alain Chaumeil, Ben J. Woodcroft, Paul N. Evans, Philip Hugenholtz^{1*} and Gene W. Tyson^{1*}

LETTER

doi:10.1038/nature14486

Unusual biology across a group comprising more than 15% of domain Bacteria

Christopher T. Brown¹, Laura A. Hug², Brian C. Thomas², Itai Sharon², Cindy J. Castelle², Andrea Singh², Michael J. Wilkins^{3,4}, Kelly C. Wrighton⁴, Kenneth H. Williams⁵ & Jillian F. Banfield^{2,5,6}

nature microbiology

ARTICLES

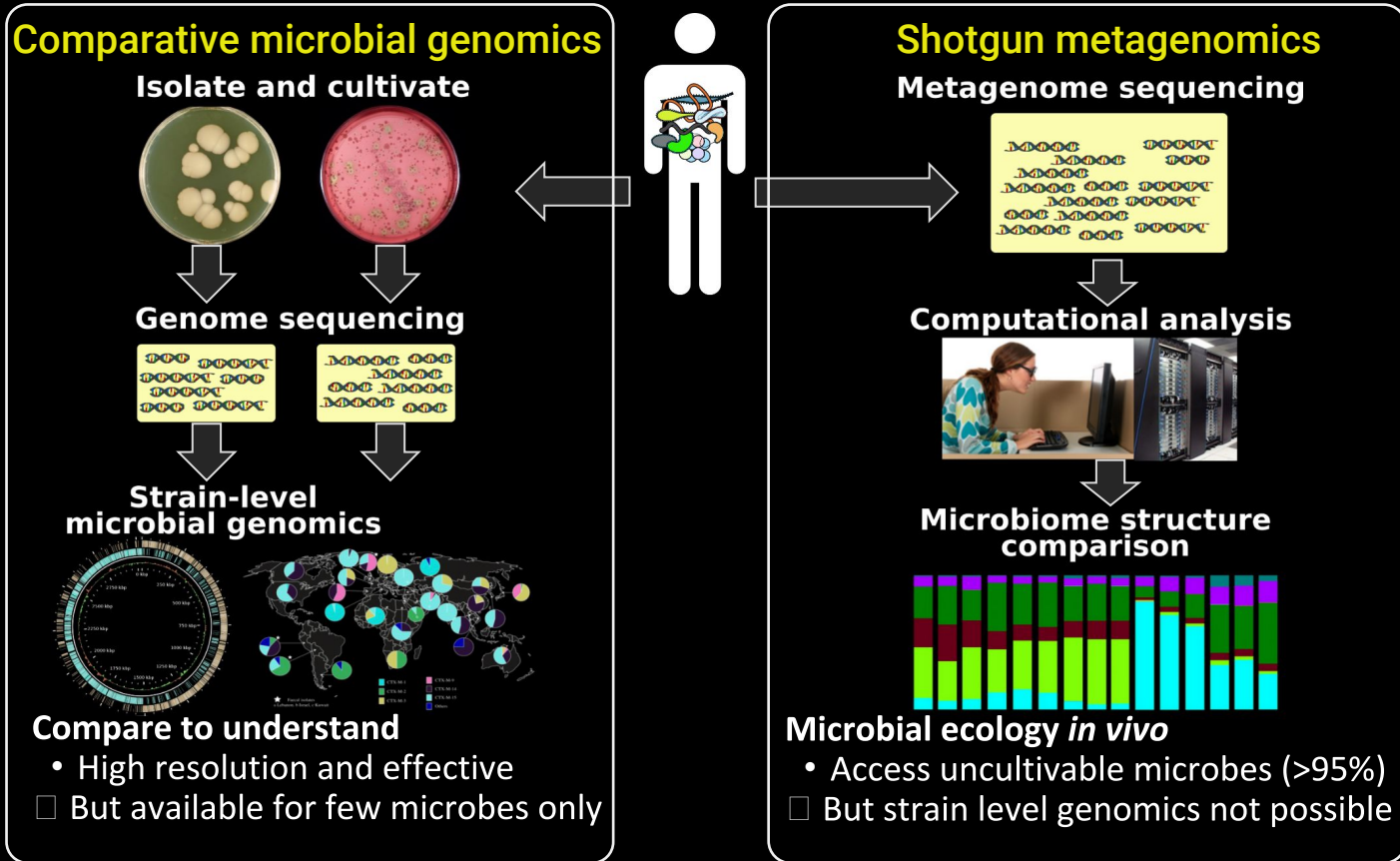
https://doi.org/10.1038/s41564-018-0171-1

OPEN

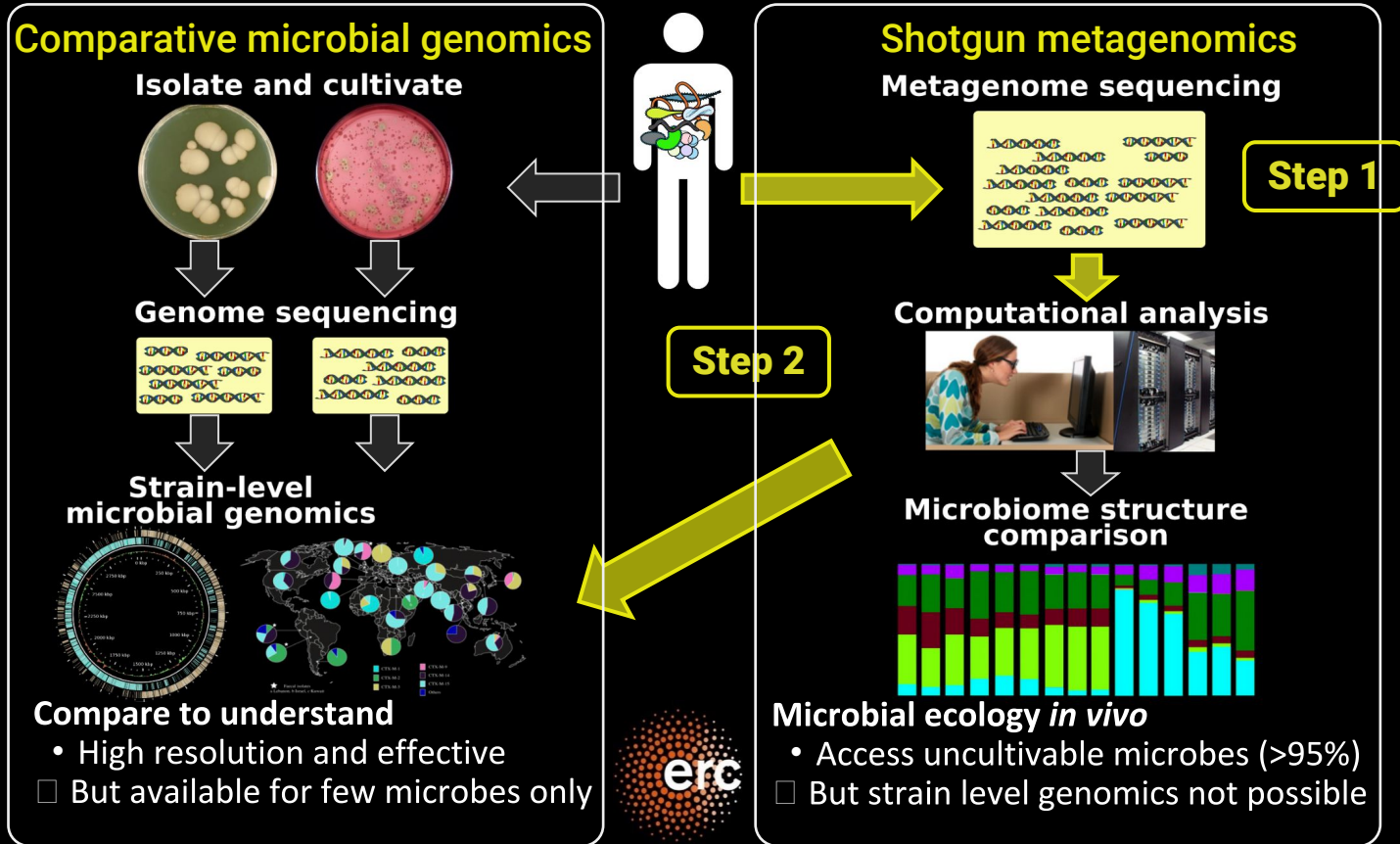
Recovery of genomes from metagenomes via a dereplication, aggregation and scoring strategy

Christian M. K. Sieber^{1,2}, Alexander J. Probst², Allison Sharrar², Brian C. Thomas², Matthias Hess³, Susannah G. Tringe^{1*} and Jillian F. Banfield^{2*}

Large-scale whole-genome comparative genomics



Large-scale whole-genome comparative genomics



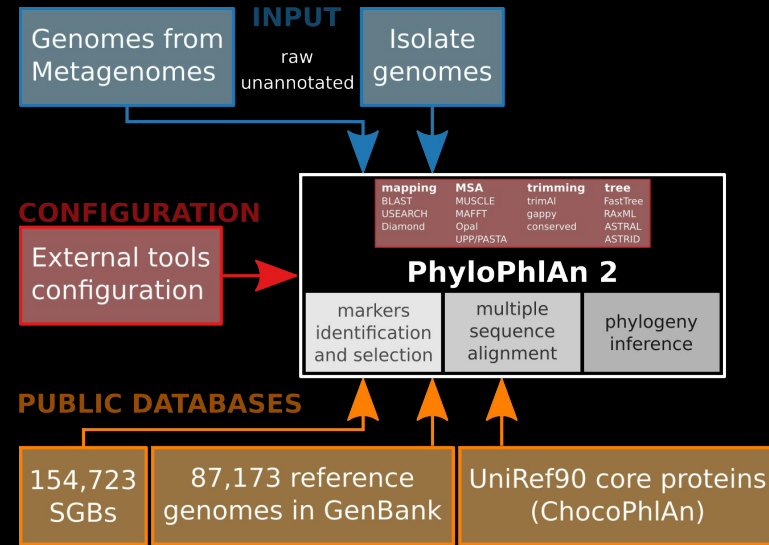
PhyloPhlAn (version 2)

An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

Main features:

- Scalable, flexible, automatic, modular, customizable



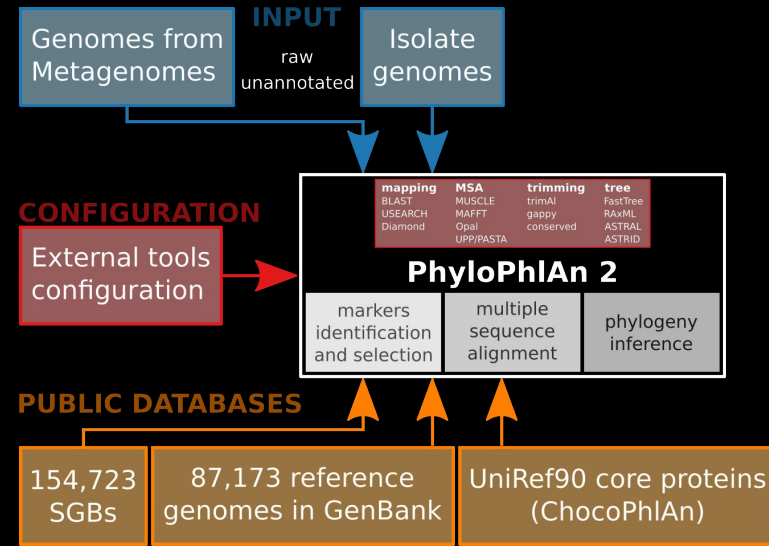
PhyloPhlAn (version 2)

An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

Main features:

- Scalable, flexible, automatic, modular, customizable



MSAs QC

- several scoring position functions implemented, e.g. MUSCLE and trident

MSAs subsample

- several subsampling approaches implemented

trimming

- trimAl
- gappy positions
- conserved positions
- all of the above

other QCs

- discarding low-quality inputs
- discarding low-quality markers
- discarding fragmentary sequences extracted from inputs

PhyloPhlAn (version 2)

An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

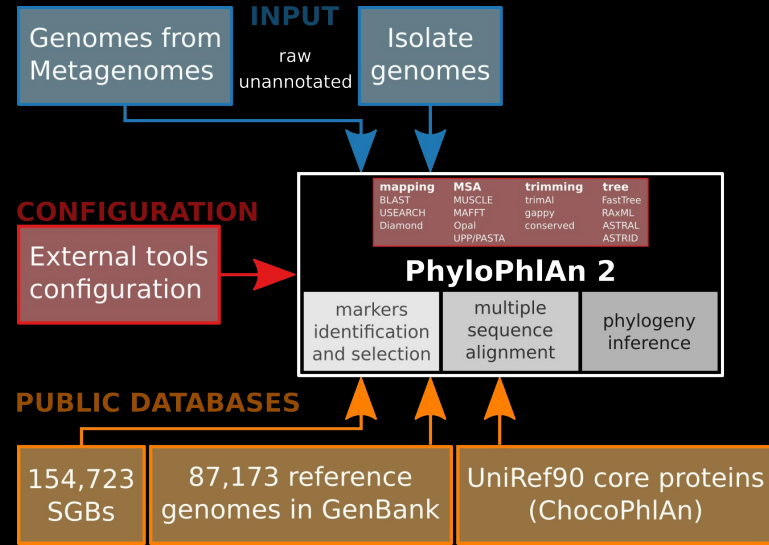
Main features:

- Scalable, flexible, automatic, modular, customizable

Examples of use-cases:

**Tree-of-life size
phylogenies**

High number of inputs
Universal markers



PhyloPhlAn (version 2)

An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

Main features:

- Scalable, flexible, automatic, modular, customizable

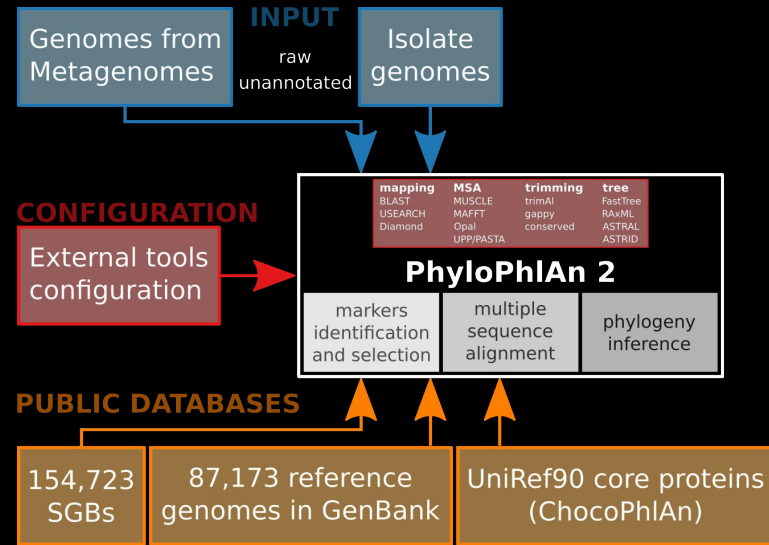
Examples of use-cases:

**Tree-of-life size
phylogenies**

High number of inputs
Universal markers

**Species-level
phylogenies**

Species-specific markers



PhyloPhlAn (version 2)

An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

Main features:

- Scalable, flexible, automatic, modular, customizable

Examples of use-cases:

**Tree-of-life size
phylogenies**

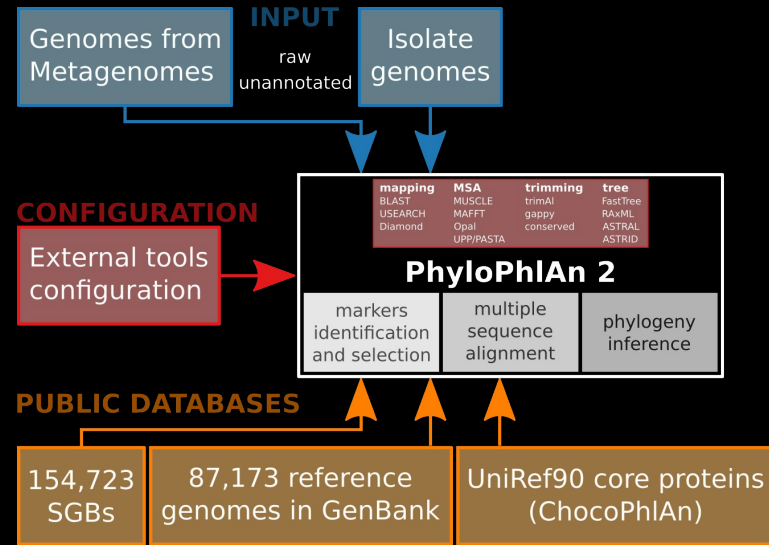
High number of inputs
Universal markers

**Species-level
phylogenies**

Species-specific markers

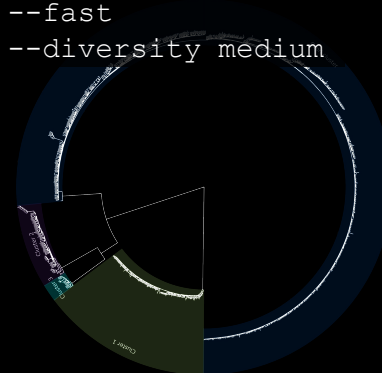
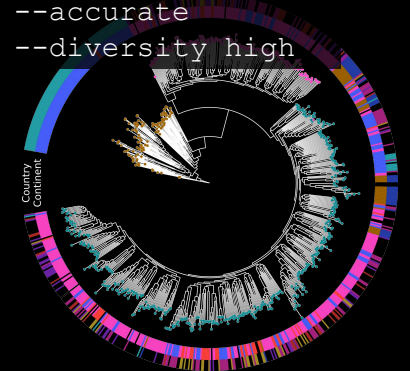
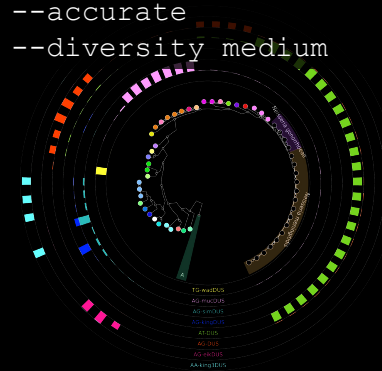
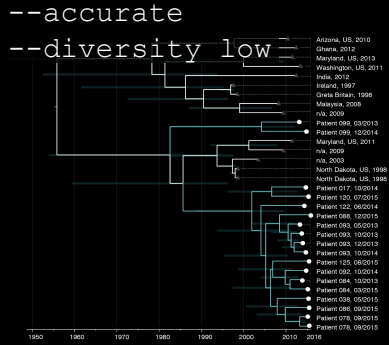
**Metagenomic
application**

Newly assembled MAGs



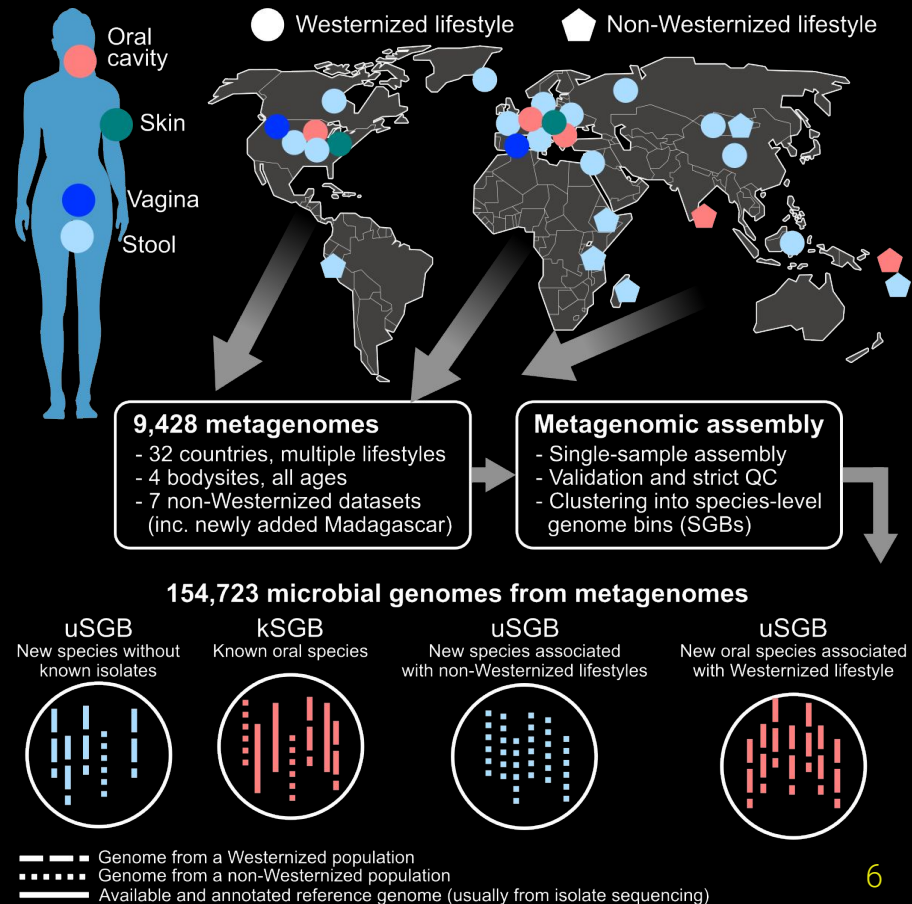
PhyloPhlAn (version 2)

INCREASING RESOLUTION



INCREASING DIVERSITY / SAMPLE SIZE

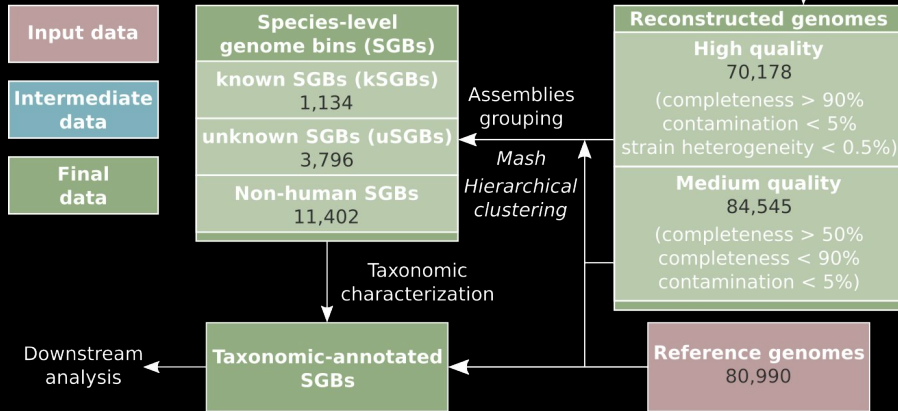
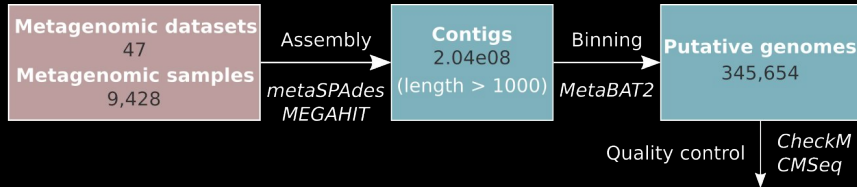
Large-scale metagenomic assembly



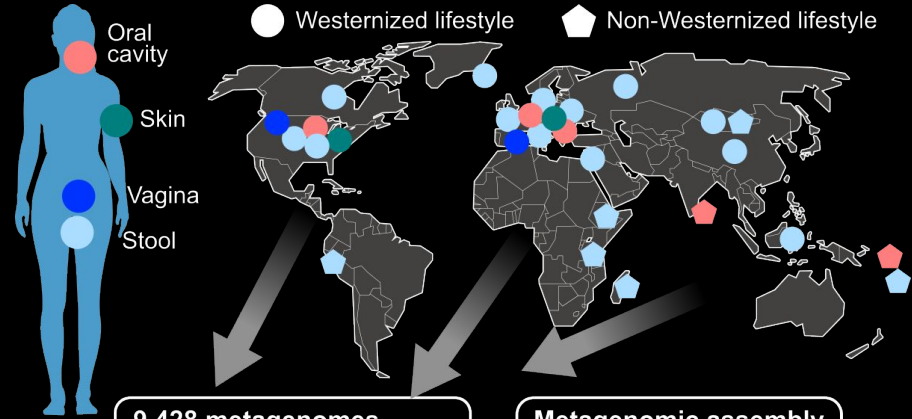
Pasolli et al., *Cell* (2019)

Large-scale metagenomic assembly

- Single-sample assembly-based
- Strict QC on reconstructed genomes
- Implemented strain-heterogeneity filtering [new]
- Species-level genome bin (SGB) clustering [new]



Pasolli et al., *Cell* (2019)



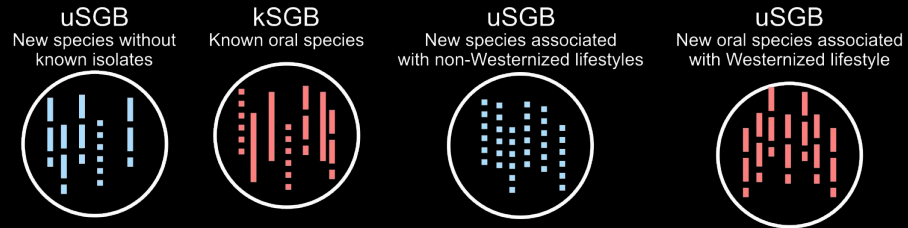
9,428 metagenomes

- 32 countries, multiple lifestyles
- 4 bodysites, all ages
- 7 non-Westernized datasets (inc. newly added Madagascar)

Metagenomic assembly

- Single-sample assembly
- Validation and strict QC
- Clustering into species-level genome bins (SGBs)

154,723 microbial genomes from metagenomes

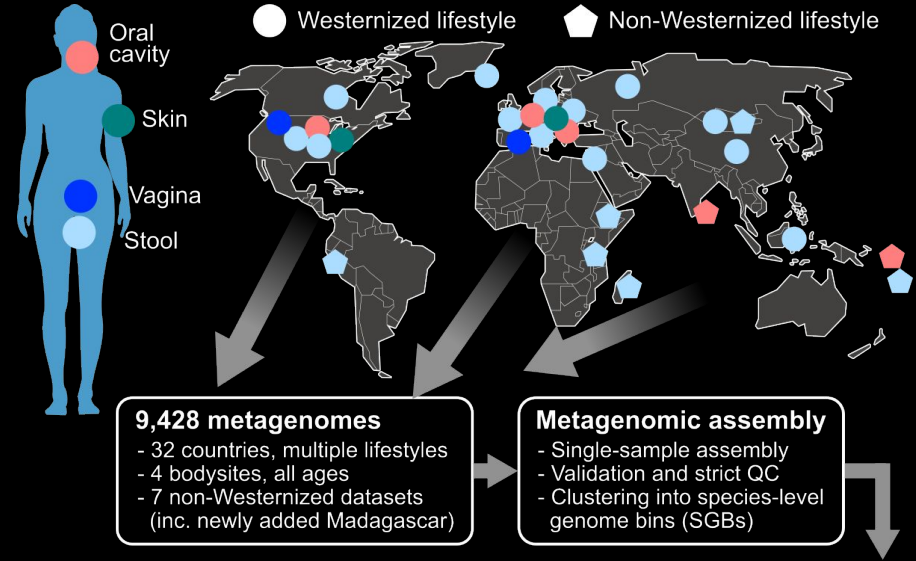


Legend for genome types:

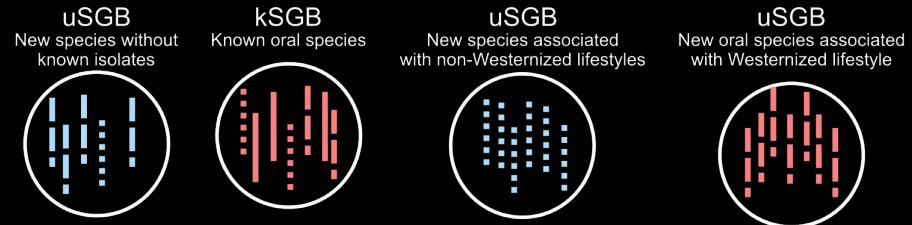
- Genome from a Westernized population
- ⋯ Genome from a non-Westernized population
- ▬ Available and annotated reference genome (usually from isolate sequencing)

Large-scale metagenomic assembly

9,428 metagenomes



154,723 microbial genomes from metagenomes



Pasolli et al., *Cell* (2019)

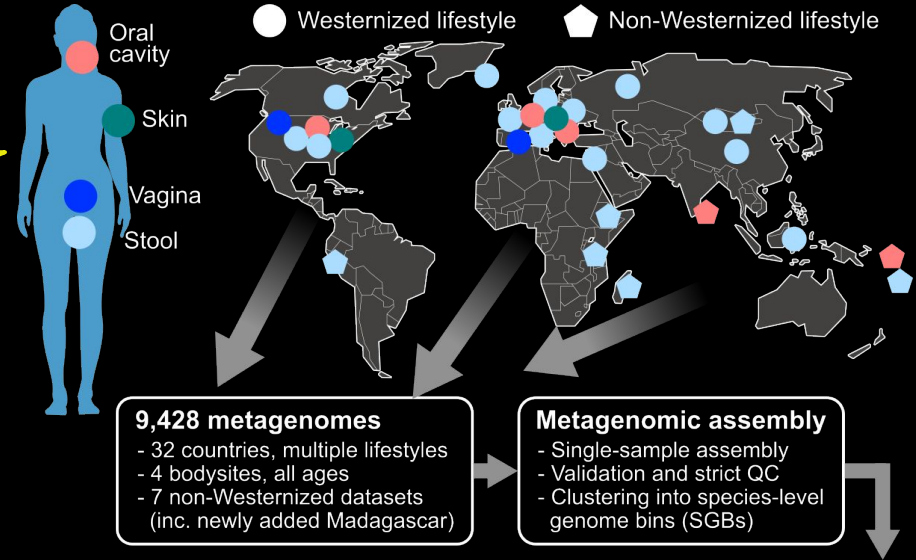
--- Genome from a Westernized population
... Genome from a non-Westernized population
— Available and annotated reference genome (usually from isolate sequencing)

Large-scale metagenomic assembly

9,428 metagenomes

single-sample
assembly

154,723 reconstructed genomes



154,723 microbial genomes from metagenomes



Large-scale metagenomic assembly

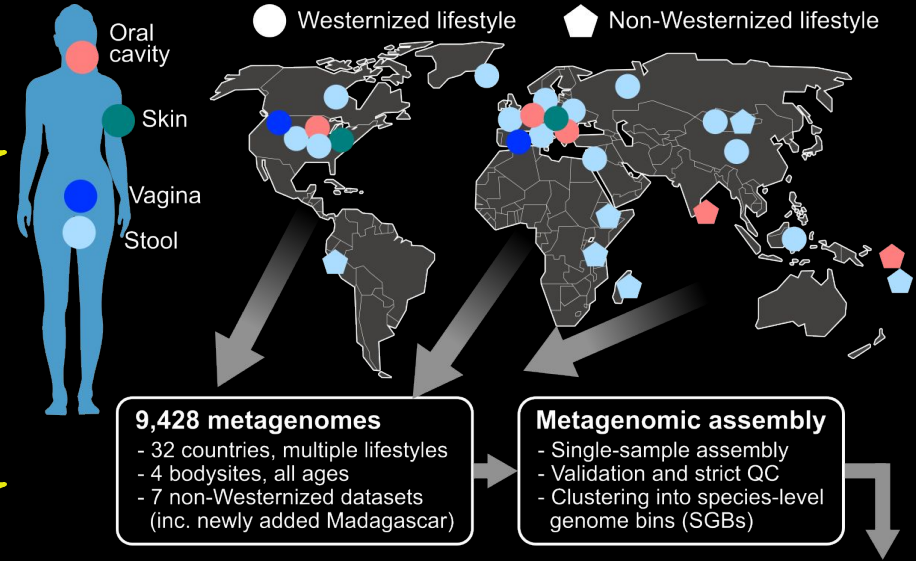
9,428 metagenomes

single-sample
assembly

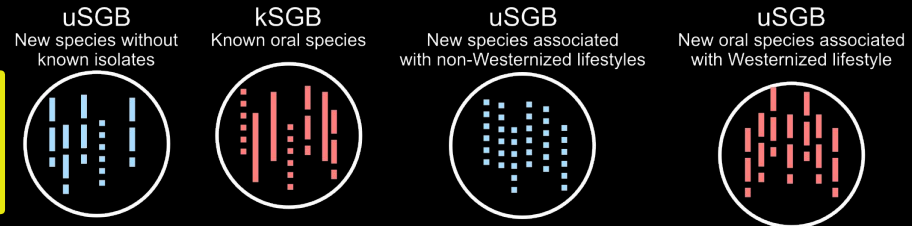
154,723 reconstructed genomes

Species-level Genome Bins
clustering

4,930 SGBs
(77% without a reference genome)



154,723 microbial genomes from metagenomes



Pasolli et al., *Cell* (2019)

— Genome from a Westernized population
- - - Genome from a non-Westernized population
— Available and annotated reference genome (usually from isolate sequencing)

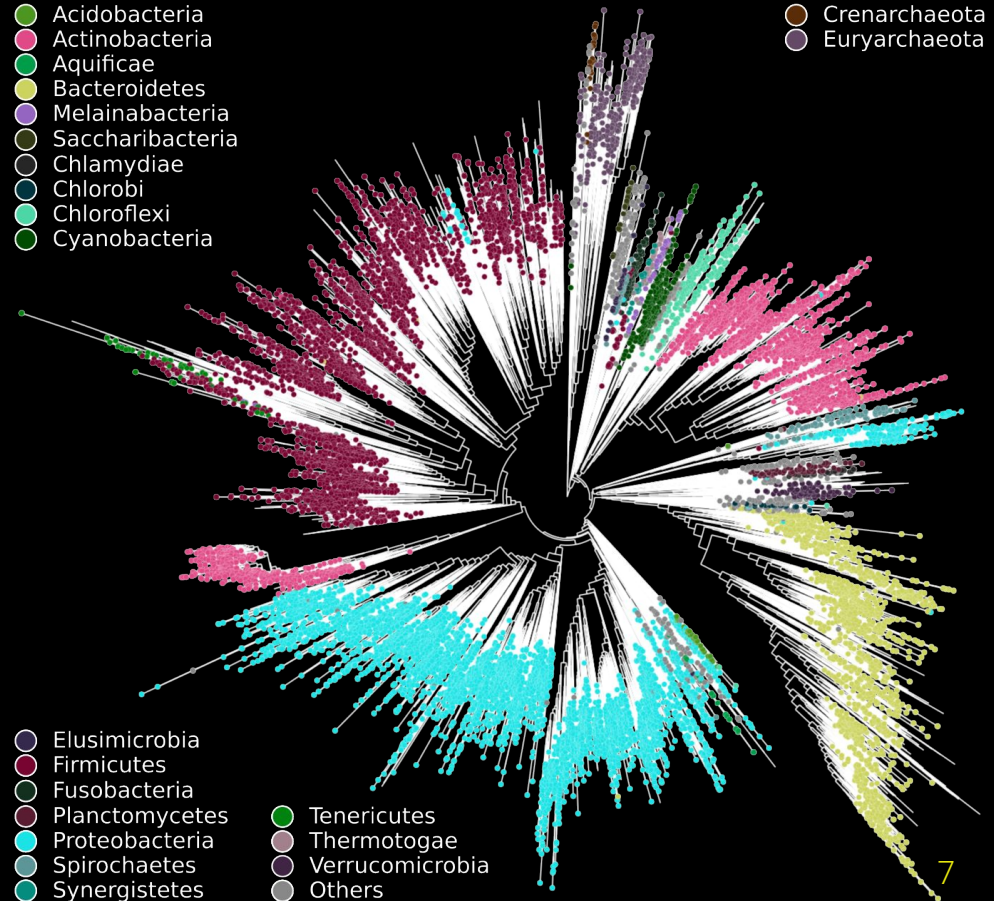
PhyloPhlAn for microbial tree-of-life phylogenies

Microbial tree of life of 17,672 SGBs

Including all isolate genomes in NCBI,
all MAGs from Parks (2017) and Pasolli
(2019)

Phylogeny based on the 400
PhyloPhlAn marker genes validated in
Qiyun (2019) with 4,522 AA positions

Building took 10 days and 15 hours
using 100 CPUs

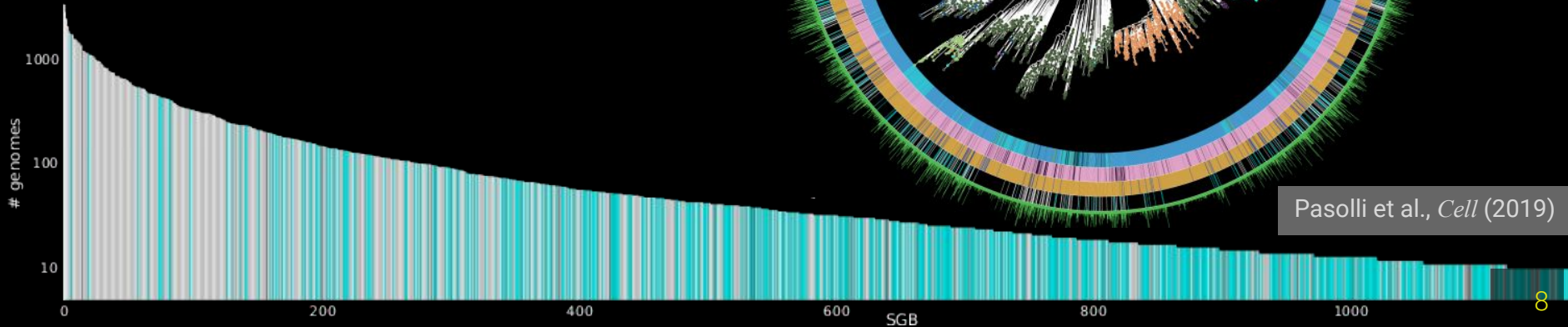


**Tree-of-life size
phylogenies**

Diversity and prevalence of human-associated SGBs

Phylogeny tree of the 4,930 human-associated SGBs

Distribution of the number of genomes for each SGB \Rightarrow several SGBs without a reference genome (uSGBs) are highly prevalent



Characterization of lower taxonomic clades

Ruminococcus spp.

A *Subdoligranulum variabile*
B *Gemmiger formicilis*

SGB ID: 15291
 # genomes: 137

SGB ID: 15287
 # genomes: 7

SGB ID: 15286
 # genomes: 1813

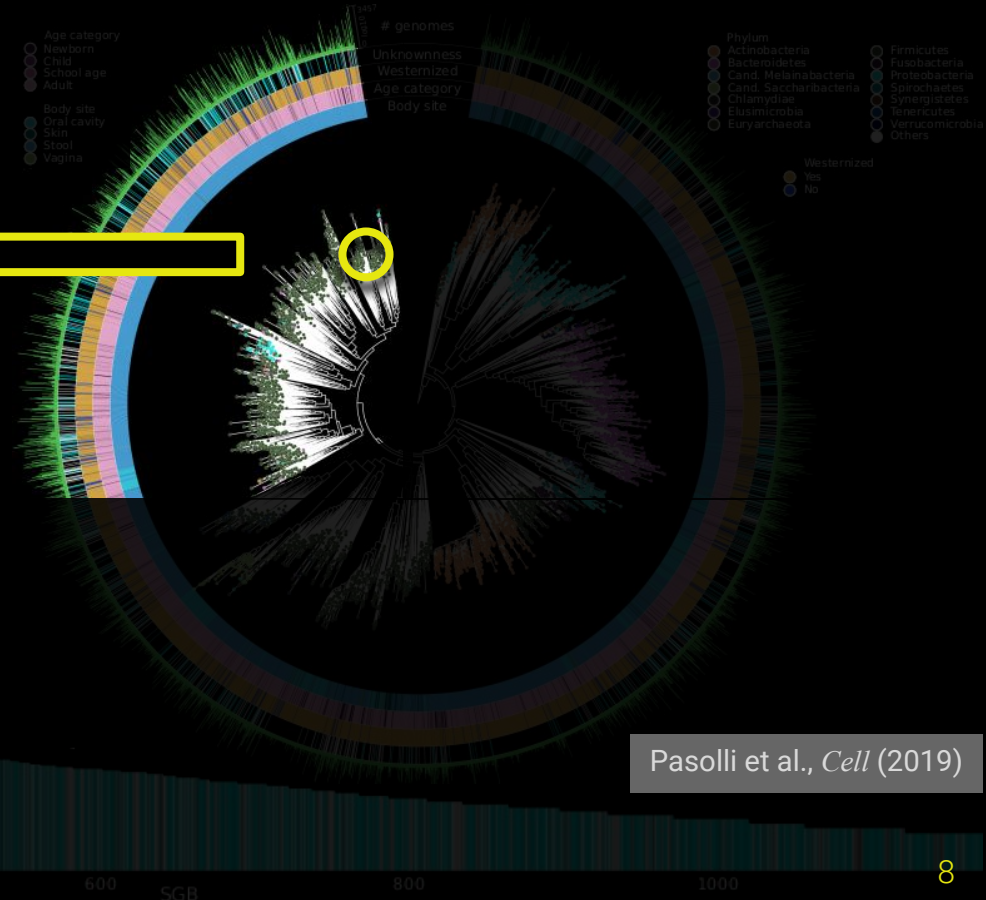
SGB ID: 15300
 # genomes: 1212

SGB ID: 15295
 # genomes: 99

SGB ID: 15292
 # genomes: 9

SGB ID: 15299
 # genomes: 99

Faecalibacterium prausnitzii



Pasolli et al., *Cell* (2019)

uSBG 15286: “*Candidatus Cibiobacter qucibialis*”

Pasolli et al., *Cell* (2019)

Ruminococcus spp.

A *Subdoligranulum variabile*
B *Gemmiger formicilis*

SGB ID: 15291
genomes: 137

SGB ID: 15287
genomes: 7

SGB ID: 15286
genomes: 1813

SGB ID: 15300
genomes: 1212

SGB ID: 15295
genomes: 99

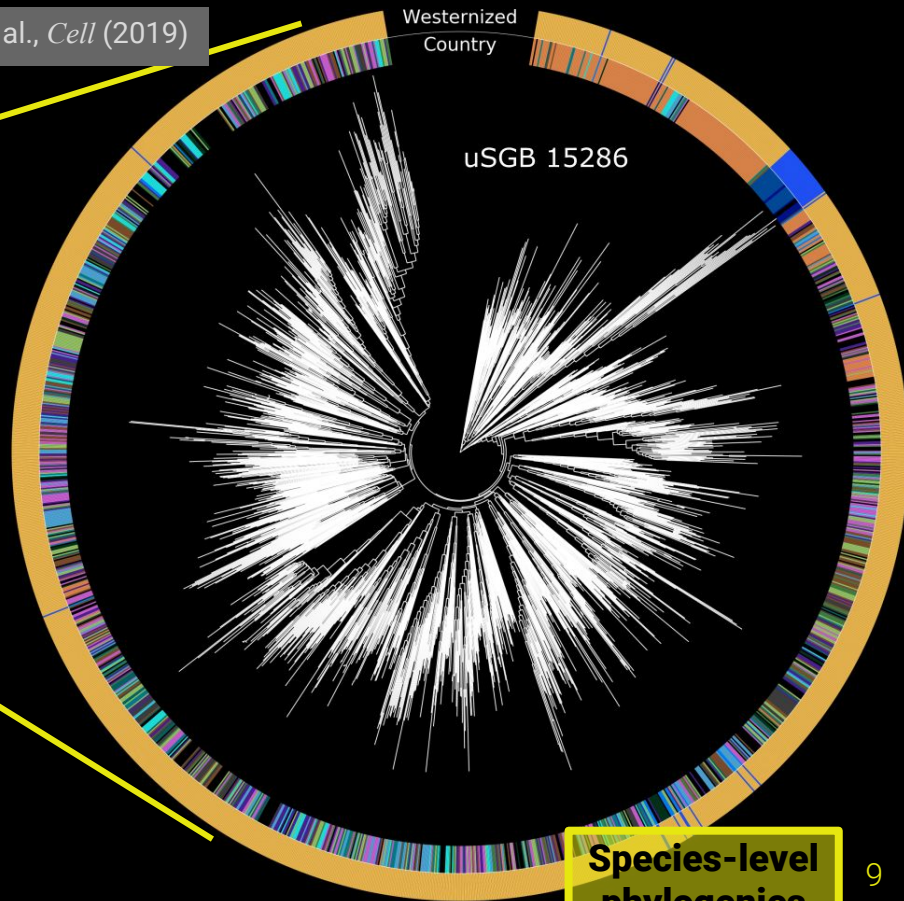
SGB ID: 15292
genomes: 9

SGB ID: 15299
genomes: 99

Faecalibacterium prausnitzii

Most prevalent
uSGB with 1,813
reconstructed
genomes

Subtrees
associated with
geography and
non-Westernized
populations



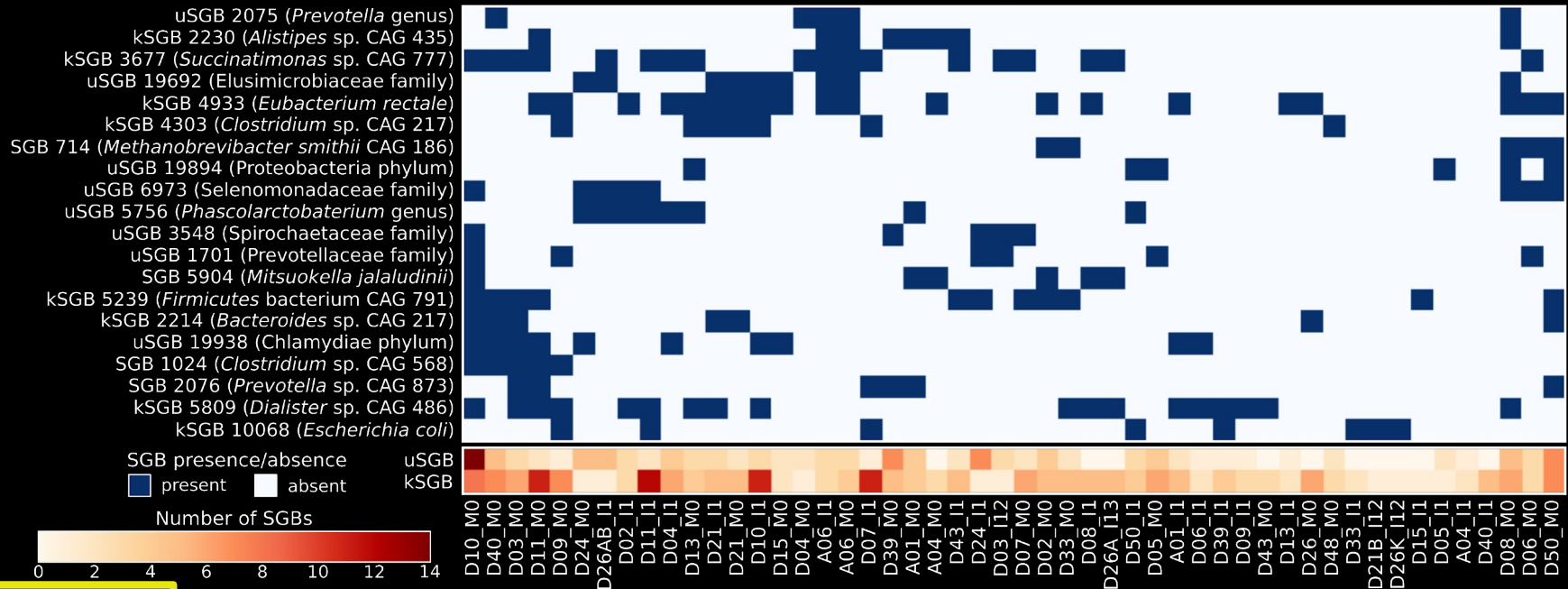
Species-level
phylogenies

Automatic SGBs assignment to new MAGs

New Ethiopian non-Westernized dataset of 50 metagenomes

in collaboration
with MC Collado

1. Which SGBs are present in this population?
2. Which new SGBs can be found in this population?



Metagenomic
application

Automatic SGBs assignment to new MAGs

New Ethiopian non-Westernized dataset of 50 metagenomes

1. Which SGBs are present in this population?
2. Which new SGBs can be found in this population?

in collaboration
with MC Collado

Focus on two SGBs:

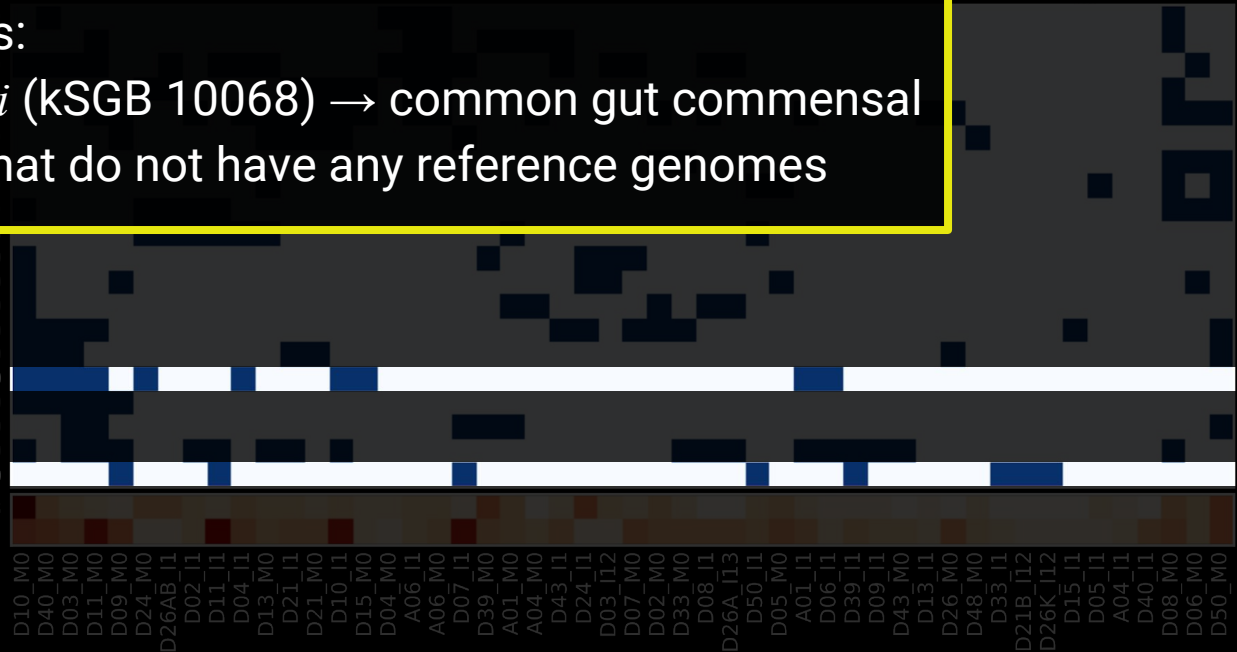
1. *Escherichia coli* (kSGB 10068) → common gut commensal
2. uSGB 19938 that do not have any reference genomes

SGB 7

uSGB 3136 (*Mitsucoccobacterium* genus)
uSGB 3548 (Spirochaetaceae family)
uSGB 1701 (Prevotellaceae family)
SGB 5904 (*Mitsuokella jalaludinii*)
kSGB 5239 (*Firmicutes* bacterium CAG 791)
kSGB 2214 (*Bacteroides* sp. CAG 217)
uSGB 19938 (Chlamydiae phylum)
SGB 1024 (*Clostridium* sp. CAG 568)
SGB 2076 (*Prevotella* sp. CAG 873)
kSGB 5809 (*Dialister* sp. CAG 486)
kSGB 10068 (*Escherichia coli*)

SGB presence/absence
■ present ■ absent

Number of SGBs



Metagenomic
application

kSGB 10068: *E. coli*

Automatic download of *E. coli* core set
of UniRef90 proteins

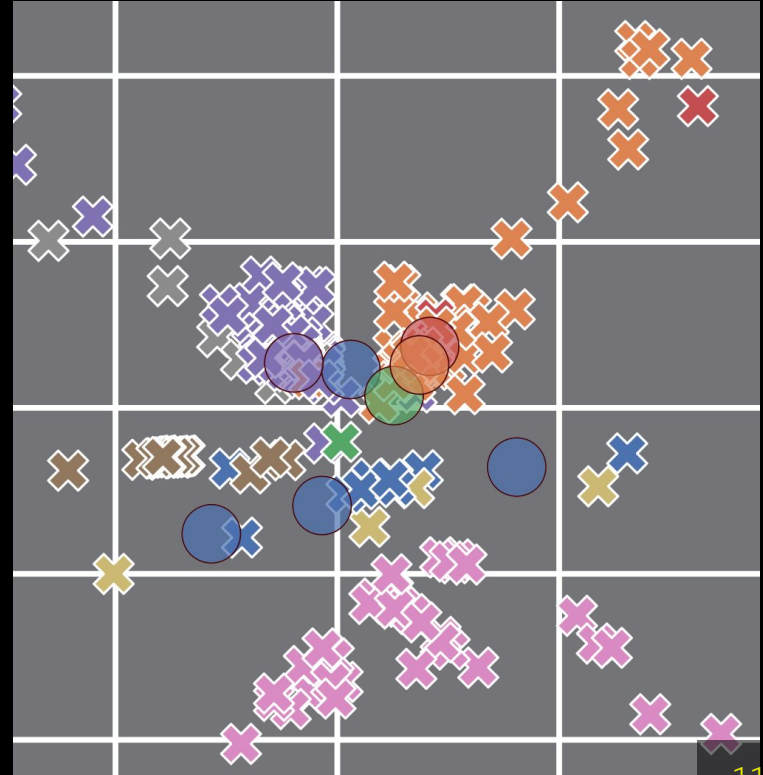
Automatic download of 1,000 *E. coli*
reference genomes

Ordination of the phylogenetic
distances of the *E. coli* phylogeny
comprising the **8 Ethiopian MAGs**

**Species-level
phylogenies**

Phylotype

- A
- B1
- B2
- C
- D
- E
- F
- U
- N/A
- Ethiopian
- ✕ Reference genomes

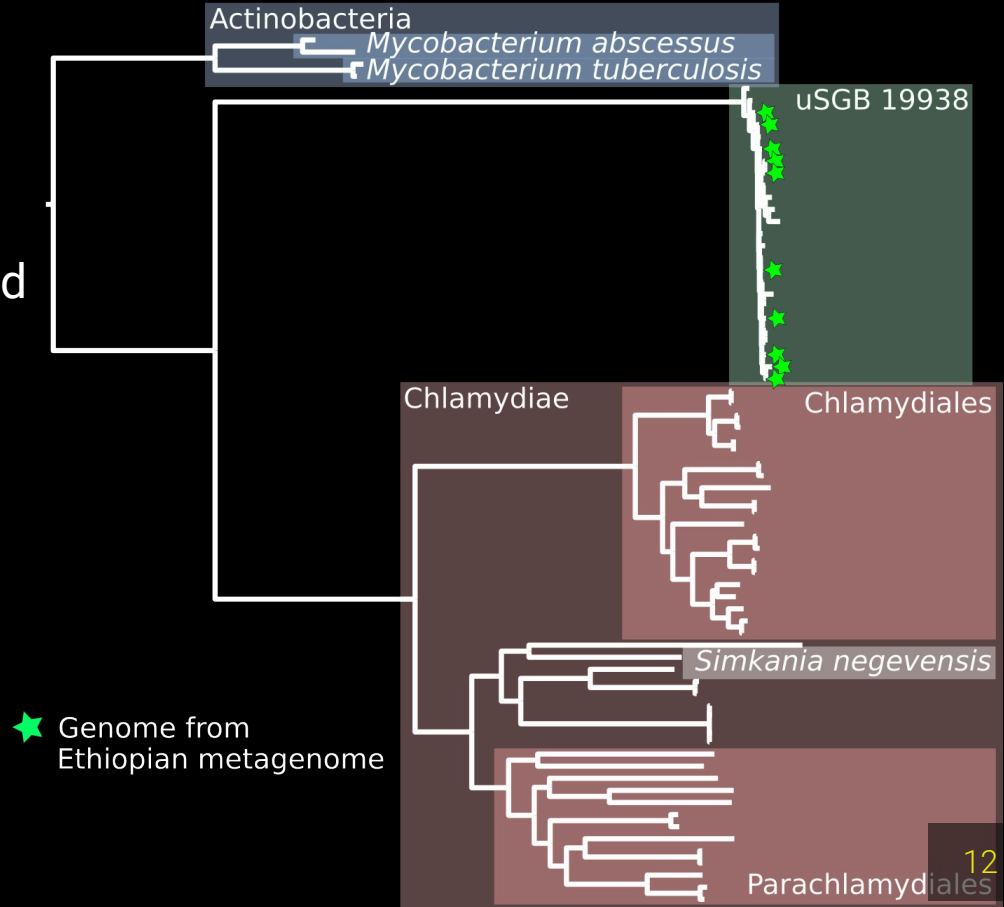


uSGB 19938: Chlamydiae (?)

The lower taxonomic assignment for uSGB 19938 is the Chlamydiae phylum

Automatic download of reference genomes for the Chlamydiae phylum and two *Mycobacterium* species (rooting)

Reconstruction of a phylum-level diversity phylogeny



Tree-of-life size
phylogenies

PhyloPhlAn (version 2)

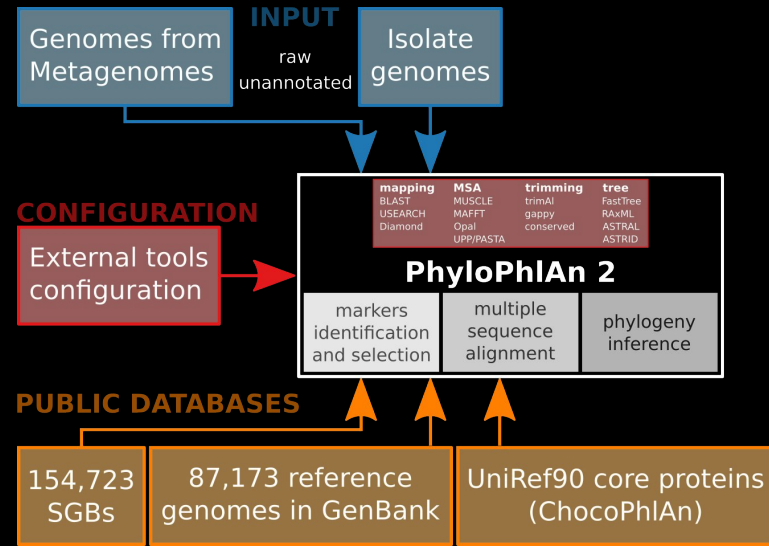
An integrated framework for phylogenetic analysis

- Reference genomes from isolates
- MAGs from metagenomes
- Clade-specific phylogenetic markers
- Retrieval of additional genomes & MAGs
- Taxonomic assignment of MAGs

Main features:

- Scalable, flexible, automatic, modular, customizable

- Available open-source in Bitbucket: <https://bitbucket.org/nsegata/phylophlan>
- Will be soon available in **BIOCONDA**
- Working on tutorials to describe the new functionalities by example
- The software can be used and any feedback is appreciated!



Integrated approach for metagenomic analysis

**Quantitative
taxonomic
profiling**

**Quantitative
functional
potential profiling**

Metagenomic Analysis

Integrated approach for metagenomic analysis

**Quantitative
taxonomic
profiling**

**Quantitative
functional
potential profiling**

**(novel) genome
reconstruction
+ PhyloPhlAn**

Metagenomic Analysis

Characterize new MAGs within the SGBs

Phylogenetic investigations including
available reference genomes

New data to update the SGB resource
to improve future analyses

Integrated approach for metagenomic analysis

Quantitative
taxonomic
profiling

Quantitative
functional
potential profiling

(novel) genome
reconstruction
+ PhyloPhlAn

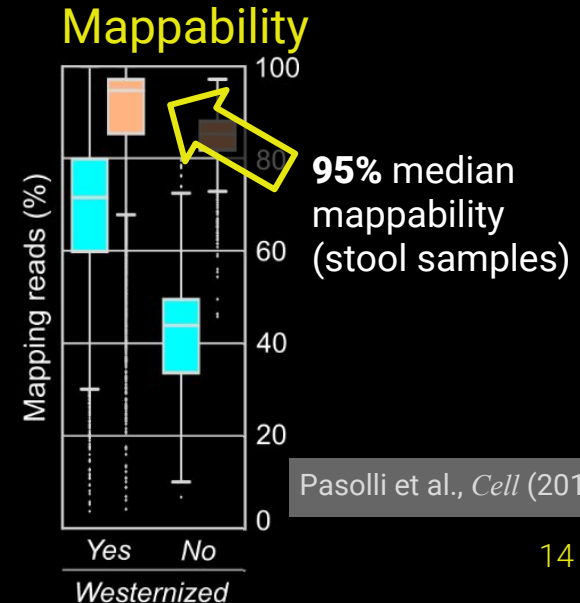
Feedback to improve
taxonomic and
functional analysis

Metagenomic Analysis

Characterize new MAGs within the SGBs

Phylogenetic investigations including
available reference genomes

New data to update the SGB resource
to improve future analyses



Thanks!

The Laboratory of Computational Metagenomics



Nicola Segata (PI)
Adrian Tett
Federica Pinto
Fabio Cumbo
Andrew Thomas
Giulia Masetti
Federica Armanini
Francesco Asnicar
Serena Manara
Paolo Ghensi
Moreno Zolfo
Francesco Beghini
Kun D. Huang
Nicolai Karcher
Paolo Manghi

<http://segatalab.cibio.unitn.it> - nicola.segata@unitn.it

Thanks to:

Curtis Huttenhower
Rob Knight
Siavash Mirarab
Qiyun Zhu
Maria C. Collado



ISTITUTO
G.B. MATTEI
RICERCA IN IDROLOGIA MEDICA
E MEDICINA TERMALE

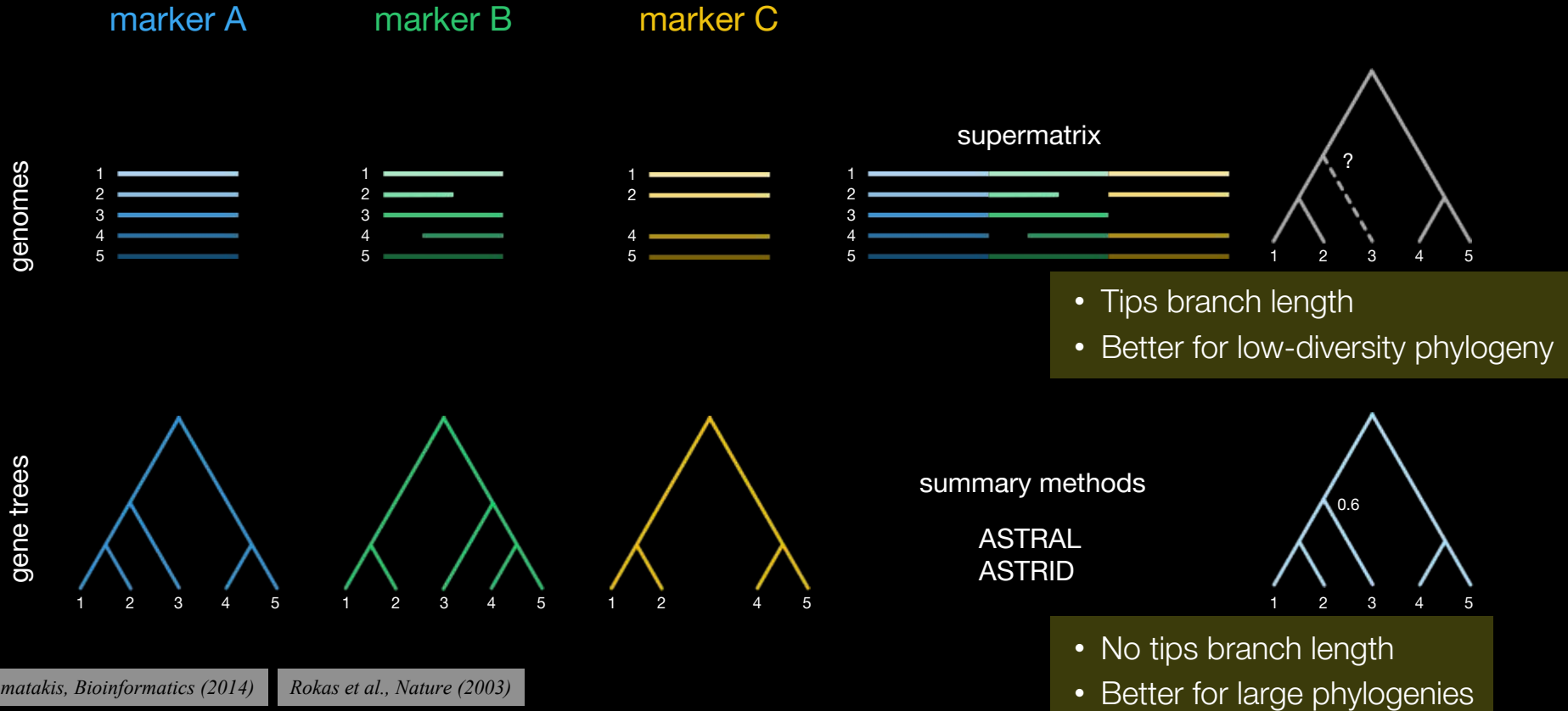


TERME DI COMANO

Interested?
We are recruiting!
nicola.segata@unitn.it



Supermatrix vs. Supertree



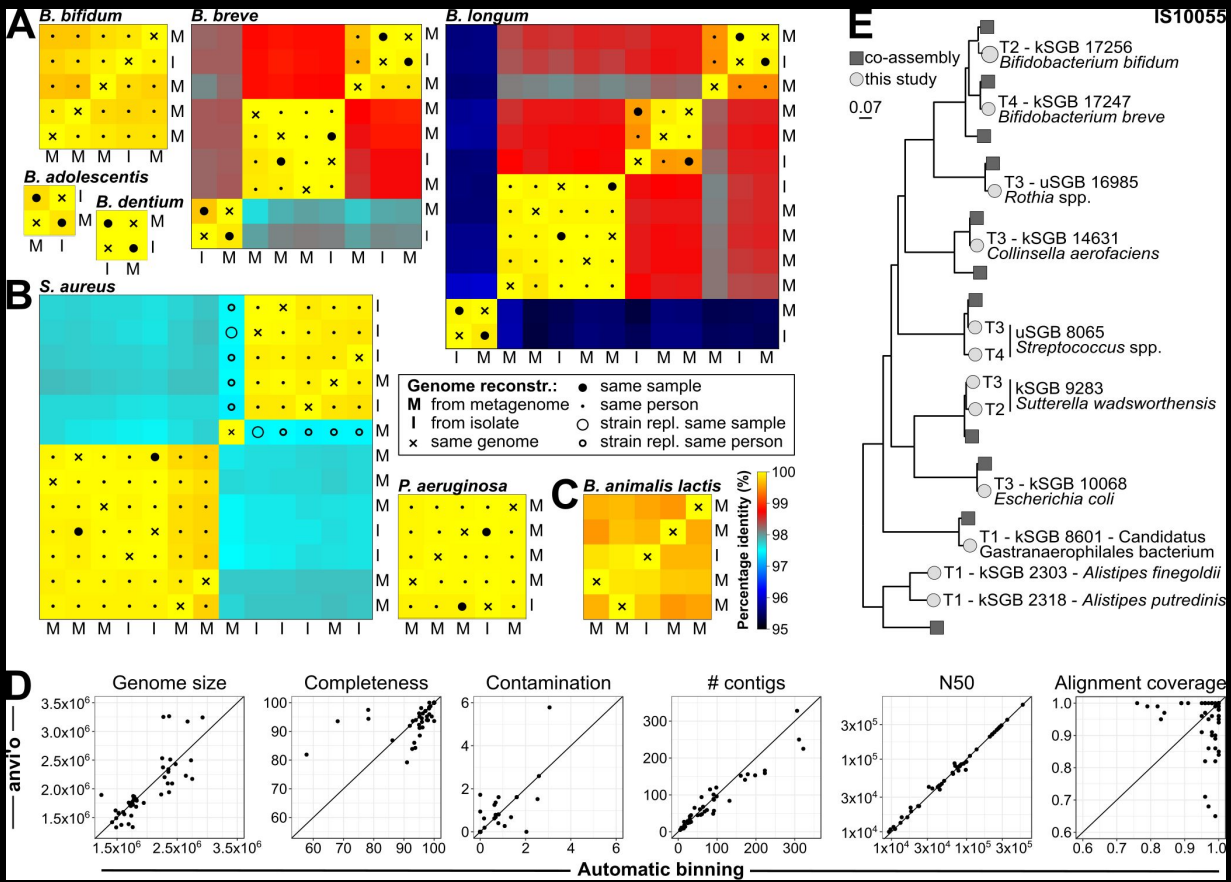
Stamatakis, Bioinformatics (2014)

Rokas et al., Nature (2003)

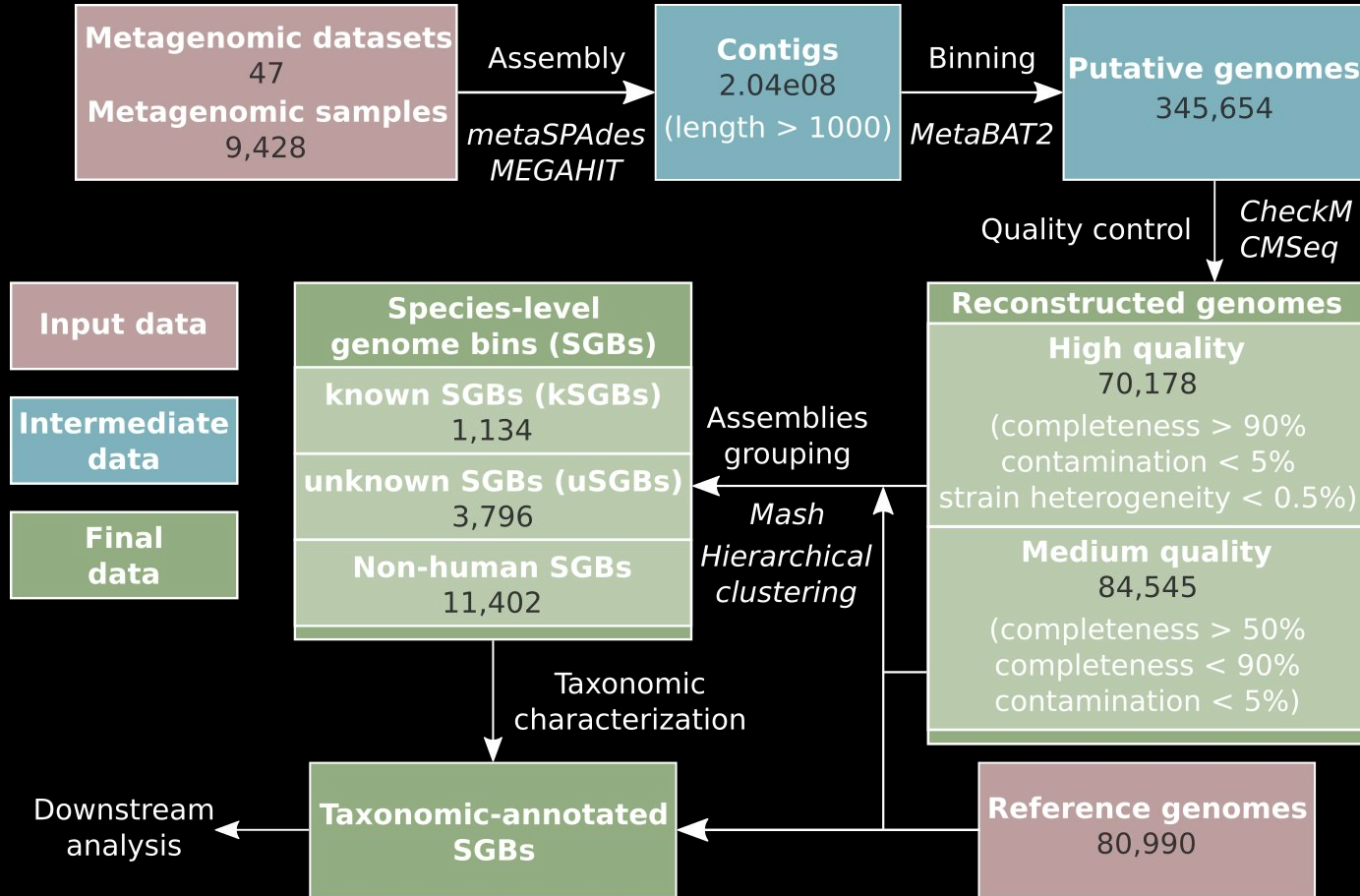
Siavash et al., Bioinformatics (2014)

Vachaspati and Warnow, BMC Genomics (2015)

Automatic and manual single-sample assembly and co-assembly evaluation



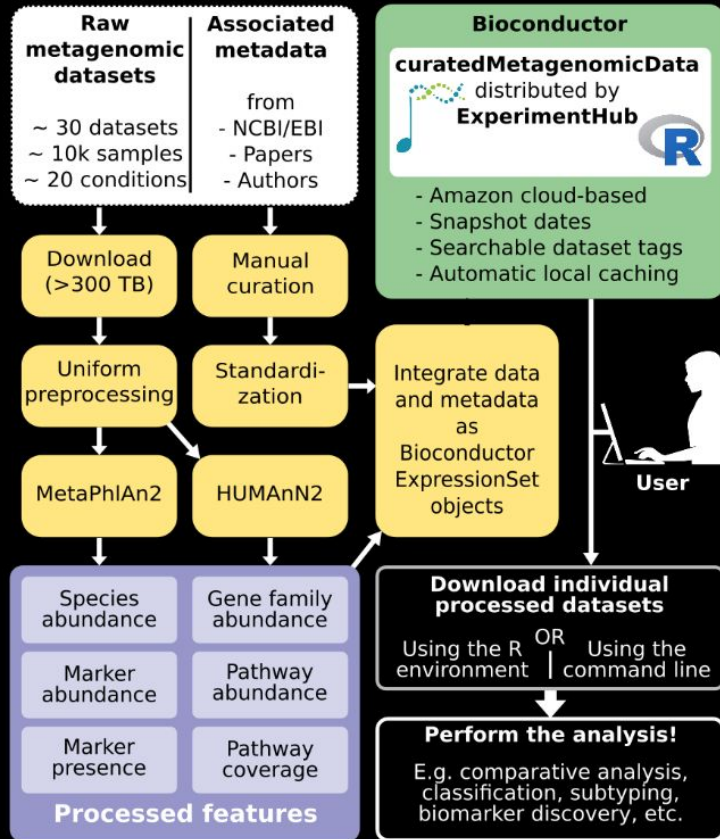
Large-scale single-sample assembly approach



CuratedMetagenomicData

Pasoli et al., *Nat Methods* (2017)

Offline high computational load pipeline (incrementally performed on new data)



Mandatory metadata fields:

```
sampleID DNA_extraction_kit
subjectID
sequencing_platform
body_site number_reads
country number_bases
antibiotics_current_use
minimum_read_length
study_condition median_read_length
disease NCBI_accession
age_category PMID
gender non_westernized
```

Optional metadata fields are... All the available ones!

• Currently >10,000 metagenome samples

• New available datasets are continuously included

<https://waldronlab.github.io/curatedMetagenomicData>