

Metagenomics 101

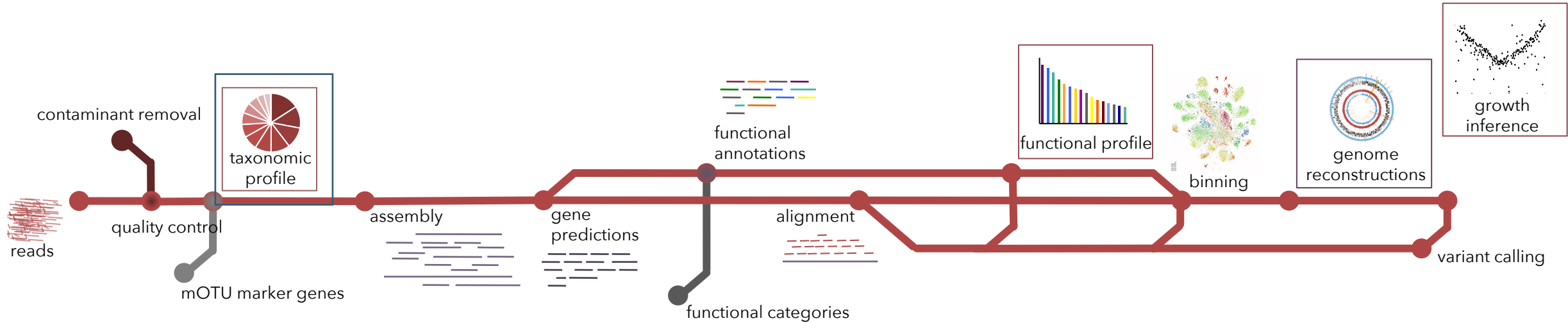
Session 4: Taxonomy I

Anna Heintz-Buschart

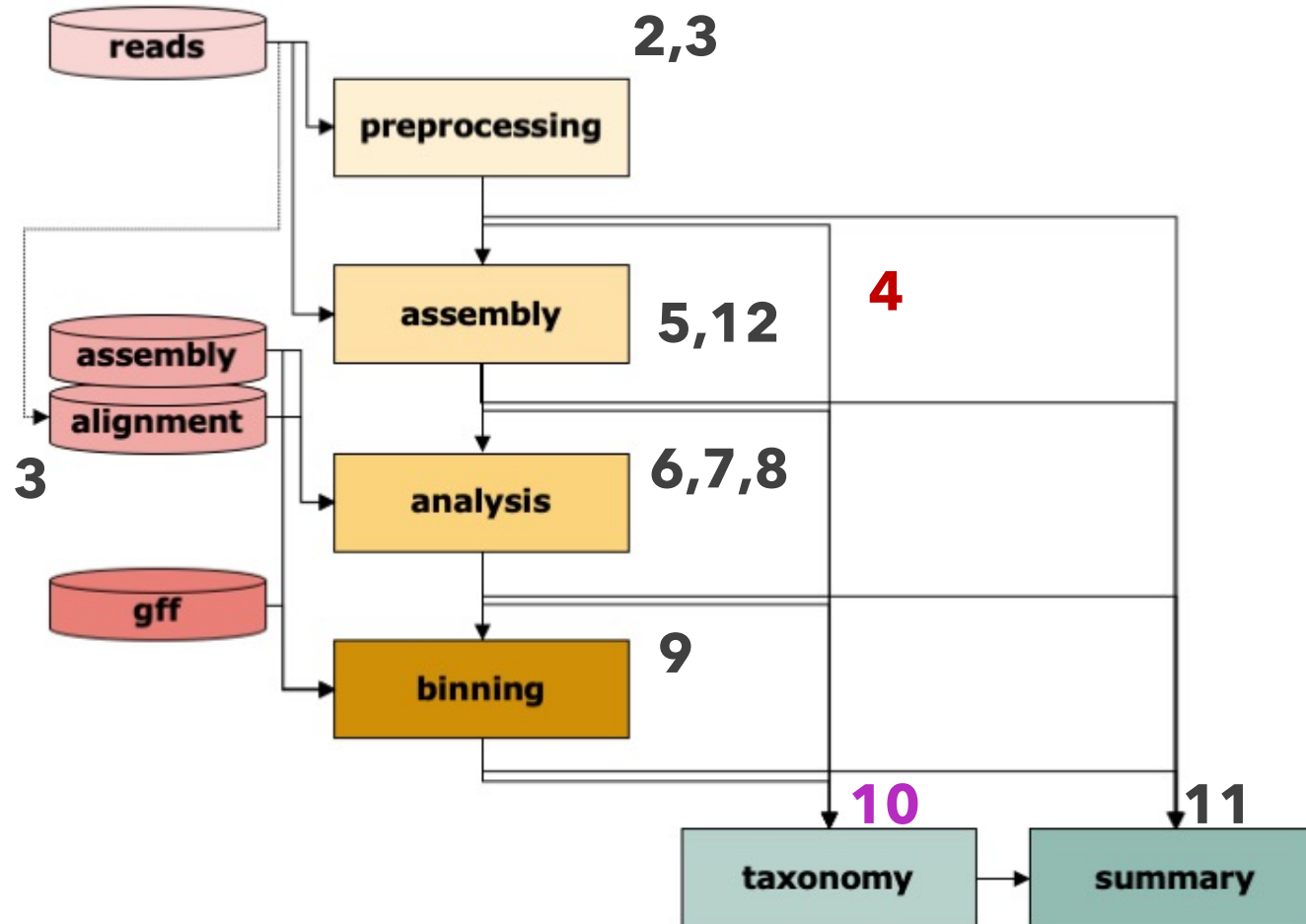
March 2022



Metagenomics (+ other omics) pipeline



Metagenomics (+ other omics) pipeline



Overview

- considerations: aims and requirements
- taxonomy
- methods:
 - reminder read mapping
 - non alignment based profilers
- benchmarking

Considerations

What do we want from taxonomy?

- accurate list of the taxa that are present
- accurate quantification of the present taxa
- robust semi-quantitative representation of present taxa

- accurate annotation of as many reads as possible
- accurate annotation of as many contigs as possible

What kind of taxa are we interested in?

- specific target taxa
 - known organisms
 - unknown organisms
-
- bacteria
 - archaea
 - eukaryotes – unicellular or multicellular
 - viruses

What resolution do we need?

- strains, sub-species
- species-like
- genera
- dynamic (can be a mix of species - phyla per sample)

What kind of microbiome do we have?

- (mostly) known organisms
- likely many unknowns
- unusual taxonomies
- many non-prokaryotic sequences

What kind of data do we have?

- reads
 - metagenomics
 - metatranscriptomics
- assembled contigs
- genome reconstructions

How much computational resources can we use?

- quick or slow?
 - scaling how with the number of reads?
- large memory or not?
- much space for our databases or less?

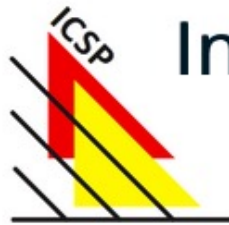
What kind of taxonomy do we want?

- anything, as long it's got the names
- a community standard
- a custom taxonomy
- we need to link it to other databases

What's a taxonomy, anyway?

taxonomy, in a broad sense the science of classification, but more strictly the classification of living and extinct organisms—i.e., biological classification. The term is derived from the Greek *taxis* ("arrangement") and *nomos* ("law"). ~~Taxonomy is, therefore, the methodology and principles of systematic botany and zoology and sets up arrangements of the kinds of plants and animals in hierarchies of superior and subordinate groups. Among biologists the Linnaean system of binomial nomenclature, created by Swedish naturalist Carolus Linnaeus in the 1750s, is internationally accepted."~~

Who makes the taxonomy?



International Committee on
Systematics of Prokaryotes

International Committee on Systematics of Prokaryotes (ICSP)

Who makes the taxonomy?

INTERNATIONAL
JOURNAL OF **SYSTEMATIC**
AND EVOLUTIONARY
MICROBIOLOGY

2019, volume 69, issue 1A, pages S1–S111



International Code of Nomenclature of Prokaryotes

Prokaryotic Code (2008 Revision)

Charles T. Parker¹, Brian J. Tindall² and George M. Garrity³ (Editors)

General Consideration 1

The progress of bacteriology can be furthered by a precise system of nomenclature accepted by the majority of bacteriologists of all nations.

General Consideration 4

Rules of nomenclature do not govern the delimitation of taxa nor determine their relations. The Rules are primarily for assessing the correctness of the names applied to defined taxa; they also prescribe the procedures for creating and proposing new names.

Why are there several taxonomies?

Published online 14 September 2021

Nucleic Acids Research, 2022, Vol. 50, Database issue D785–D794
<https://doi.org/10.1093/nar/gkab776>

GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy

Donovan H. Parks^{1*}, Maria Chuvpochina, Christian Rinke, Aaron J. Mussig²,
Pierre-Alain Chaumeil and Philip Hugenholtz¹

D590–D596 *Nucleic Acids Research*, 2013, Vol. 41, Database issue
[doi:10.1093/nar/gks1219](https://doi.org/10.1093/nar/gks1219)

Published online 28 November 2012

The SILVA ribosomal RNA gene database project: improved data processing and web-based tools

Christian Quast¹, Elmar Pruesse^{1,2}, Pelin Yilmaz¹, Jan Gerken^{1,2}, Timmy Schweer¹,
Pablo Yarza³, Jörg Peplies³ and Frank Oliver Glöckner^{1,2,*}



Database, 2020, 1–21
[doi: 10.1093/database/baaa062](https://doi.org/10.1093/database/baaa062)
Review



Review

NCBI Taxonomy: a comprehensive update on curation, resources and tools

Conrad L. Schoch^{*}, Stacy Ciufo, Mikhail Domrachev, Carol L. Hotton,
Sivakumar Kannan, Rogneda Khovanskaya, Detlef Leipe,
Richard Mcveigh, Kathleen O'Neill, Barbara Robbertse,
Shobha Sharma, Vladimir Sousoy, John P. Sullivan, Lu Sun,
Seán Turner and Ilene Karsch-Mizrachi



LPSN - List of Prokaryotic names with Standing in Nomenclature

Why are there several taxonomies?

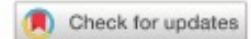
- need to include un-cultured taxa / taxa without specimen cultures
- drive to use more phylogenetic information
- phylogenies based on 16S rRNA or whole genomes
- split taxa based on phylogenetic information or based on literature weight?

Taxonomy, phylogeny, and un-cultured organisms

nature
biotechnology

RESOURCE

<https://doi.org/10.1038/s41587-020-0501-8>



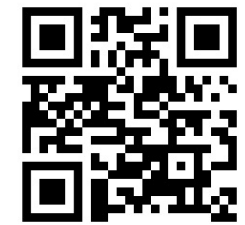
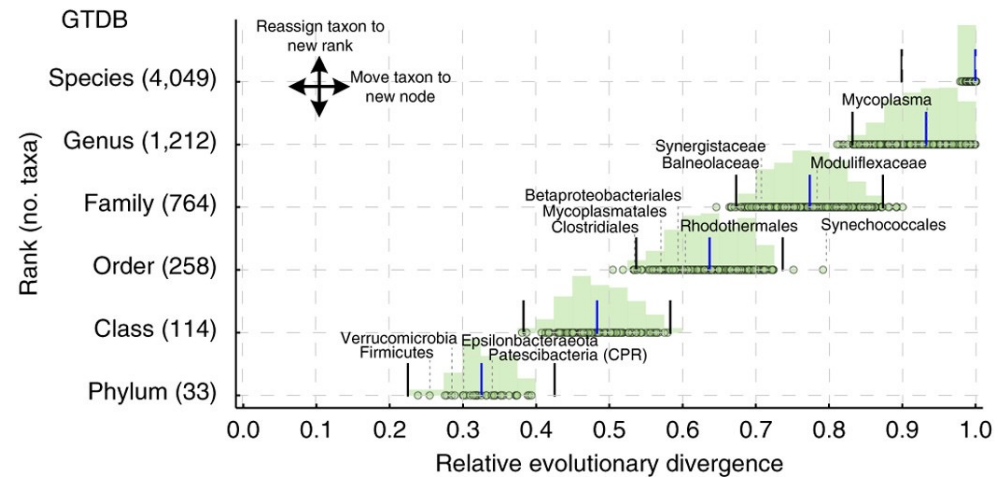
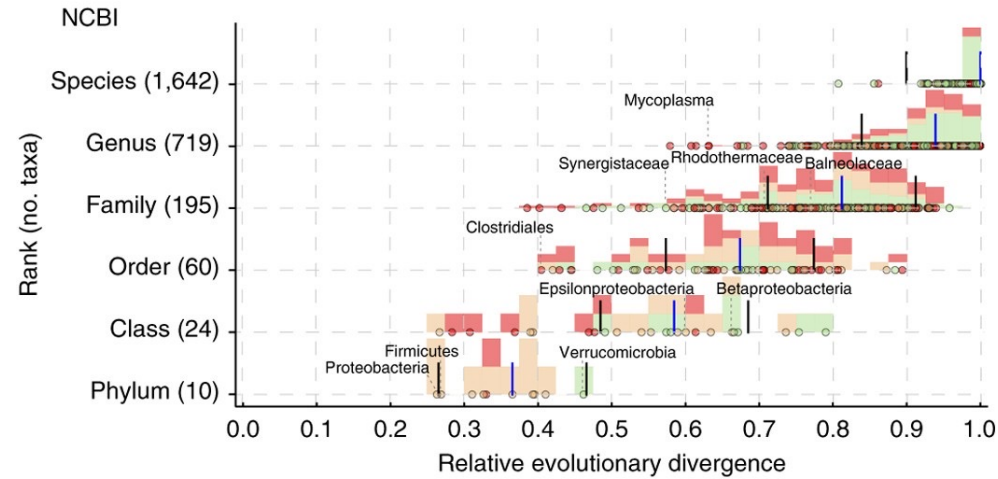
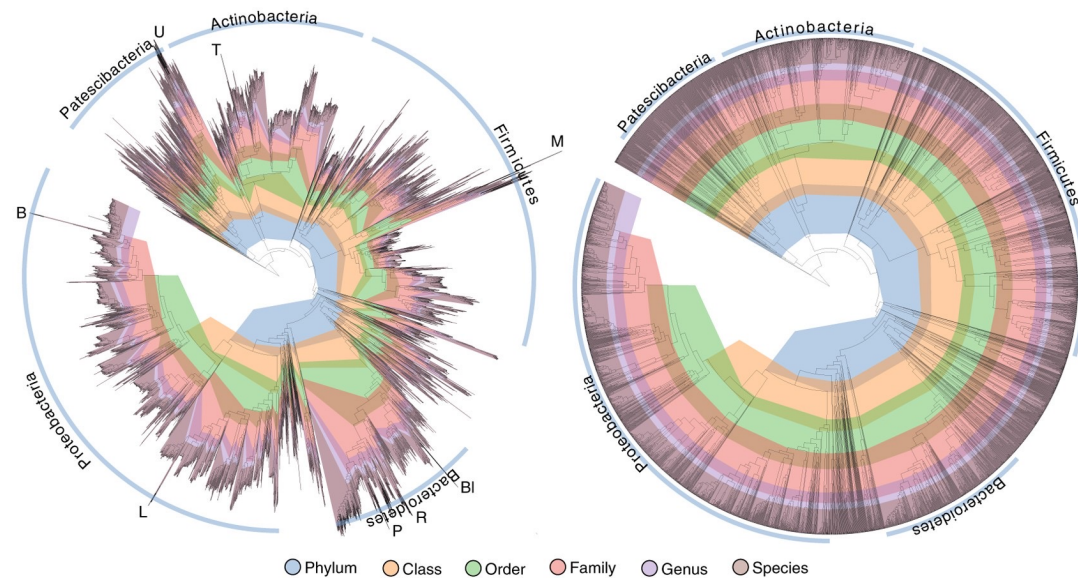
A complete domain-to-species taxonomy for Bacteria and Archaea

GTDB

Donovan H. Parks  , Maria Chuvochina, Pierre-Alain Chaumeil, Christian Rinke ,
Aaron J. Mussig  and Philip Hugenholtz 



Taxonomy, phylogeny, and un-cultured organisms



Taxonomy, phylogeny, and un-cultured organisms

GTDB
@ace_gtdb

We are proposing to reclassify *Shigella* species as synonyms of *E. coli* in the

Feedback on this reclassification is invited in the GTDB forums: [forum.gtdb.org/t/reclassification](#)

bioRxiv
THE PREPRINT SERVER FOR BIOLOGY

Reclassification of *Shigella* species as synonyms of *Escherichia coli* and are often considered to be atypical ...

5:57 PM · Sep 23, 2021 · Twitter Web App

Opinion

Microbial Taxonomy Run Amok

Trends in Microbiology

Robert A. Sanford,^{1,*} Karen G. Lloyd,² Konstantinos T. Konstantinidis,³ and Frank E. Löffler ^{2,4,5,6,7,*}

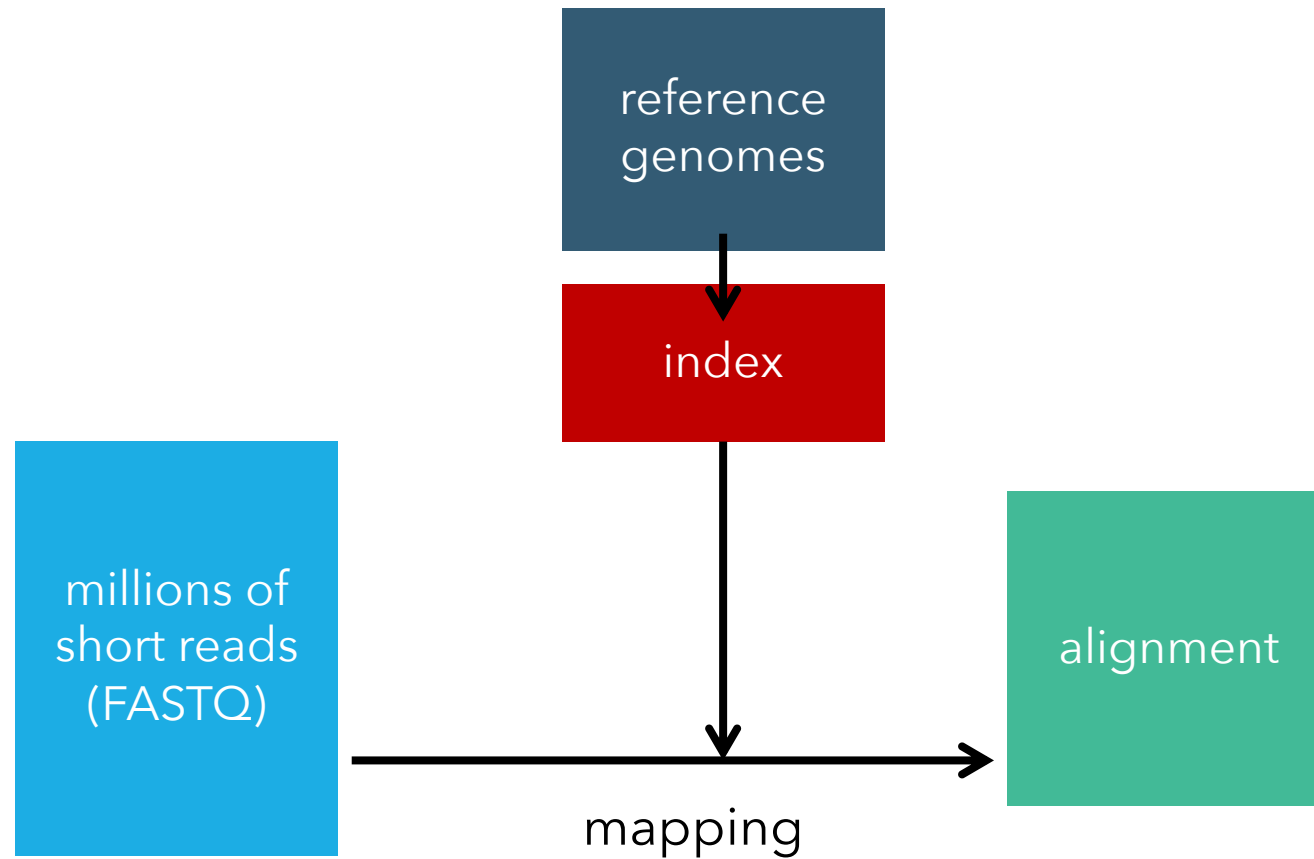


Methods

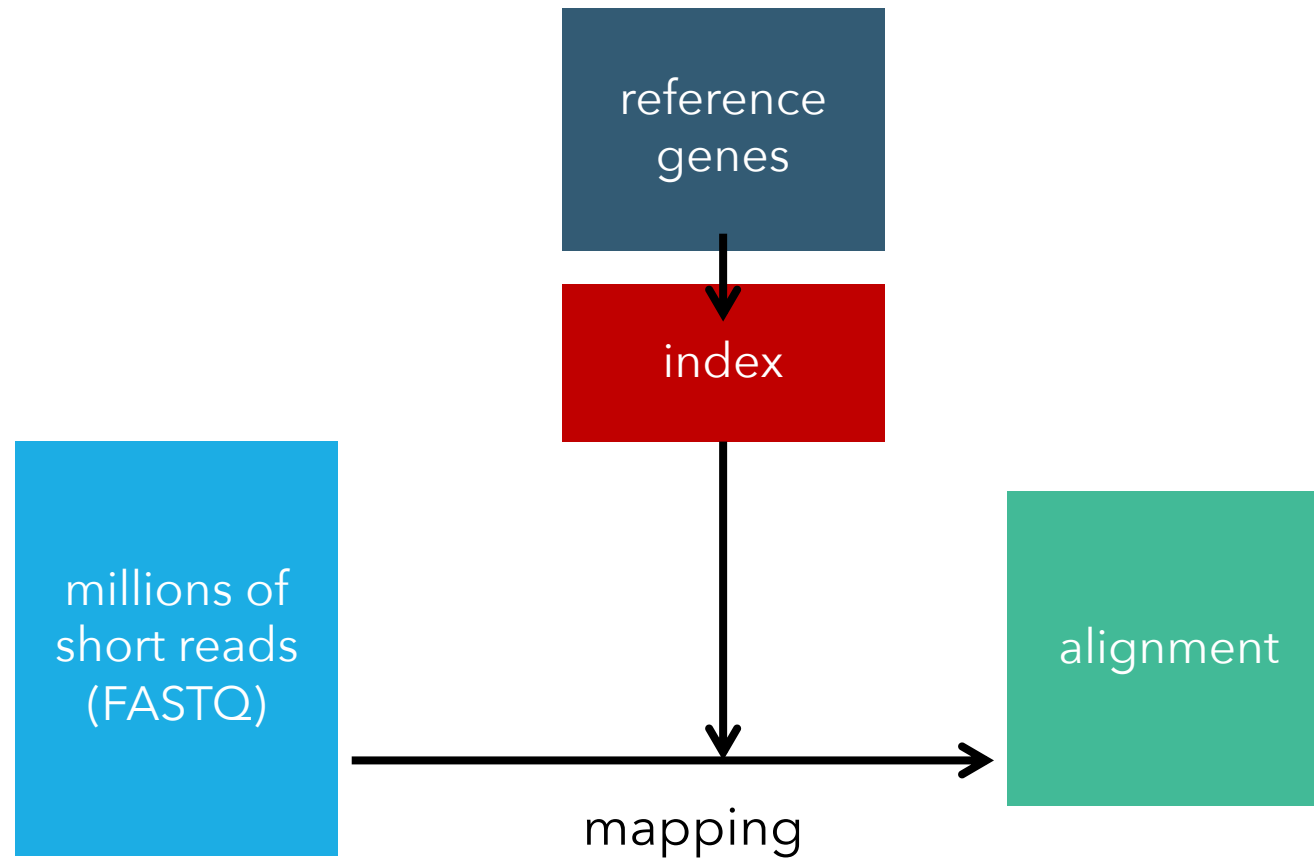
General profiling approaches

- align reads to taxonomically recognized genomes
- align reads to taxonomically annotated genes
- align reads to phylogenetic/specific marker regions
- match reads to taxonomically recognized genomes

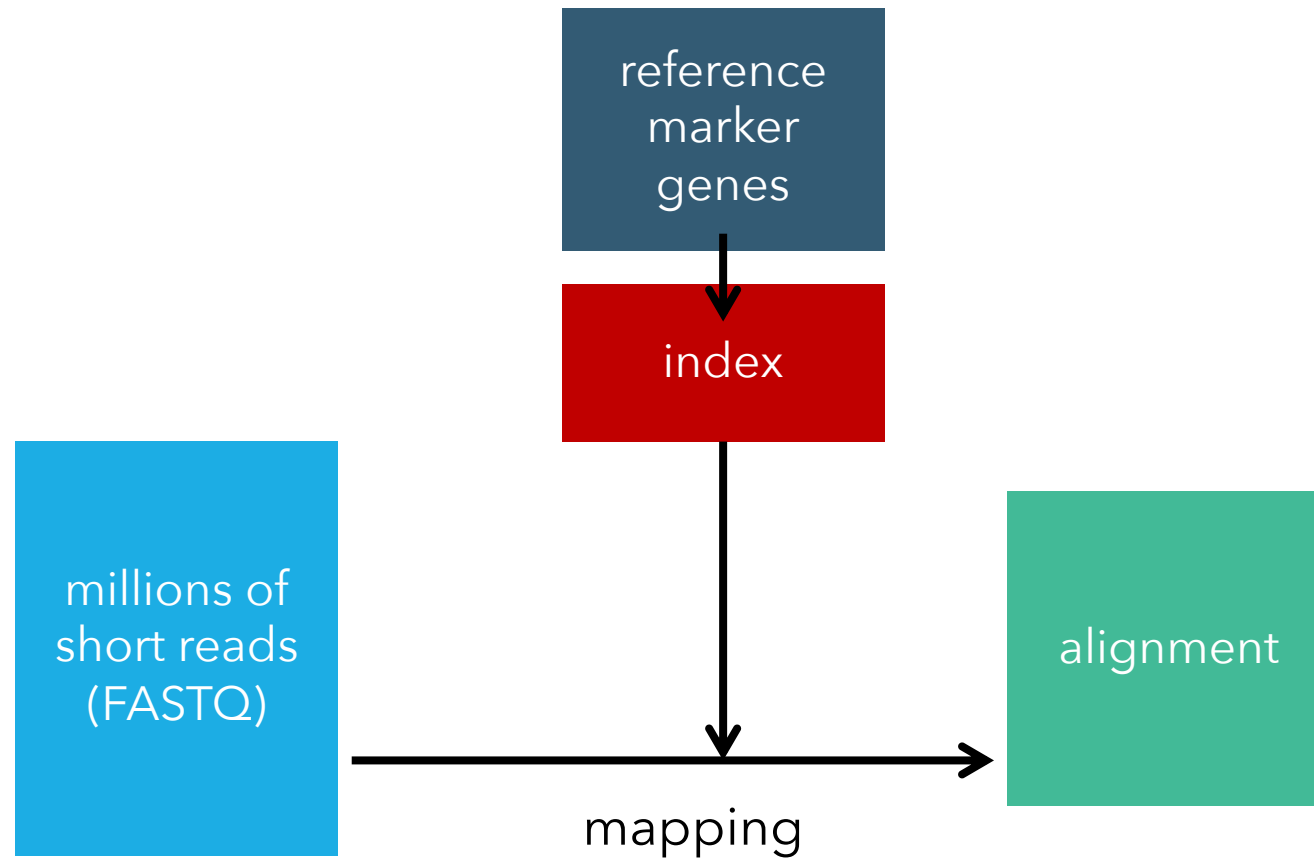
Taxonomic profilers - 1



Taxonomic profilers - 2

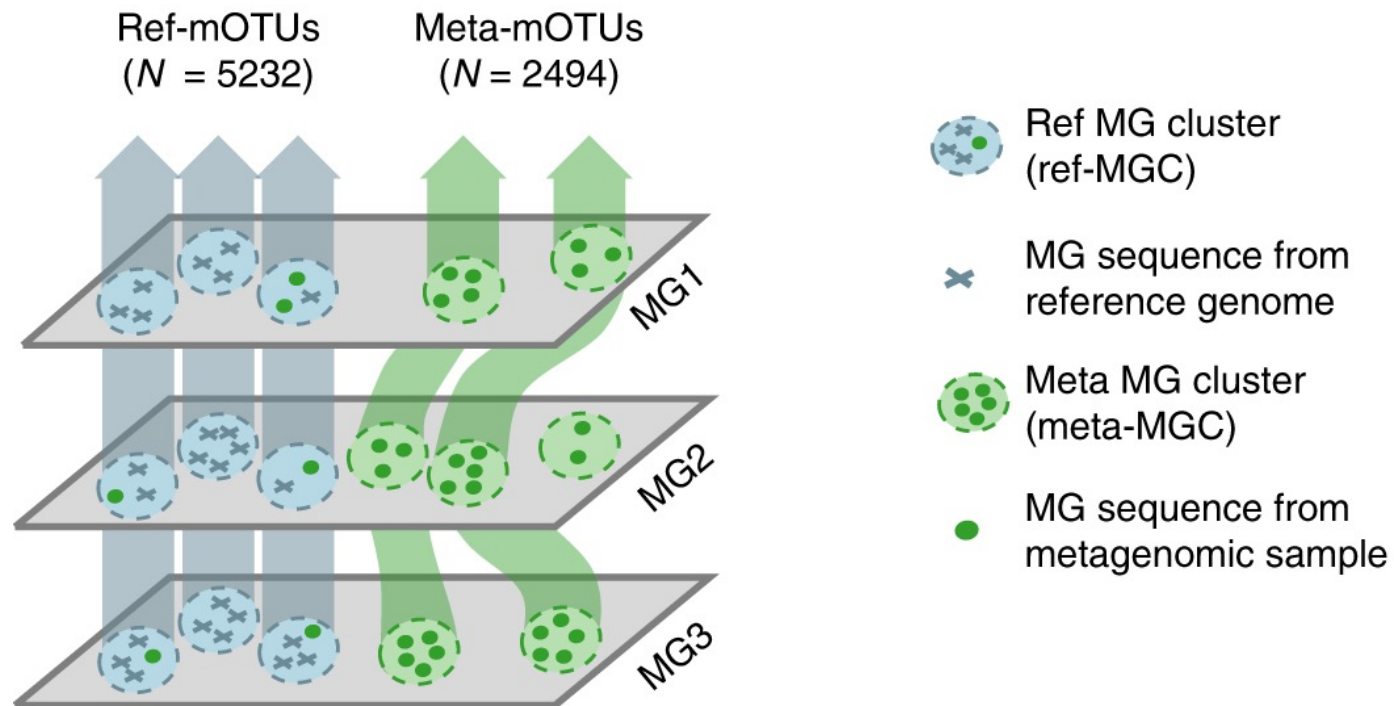


Taxonomic profilers - 3



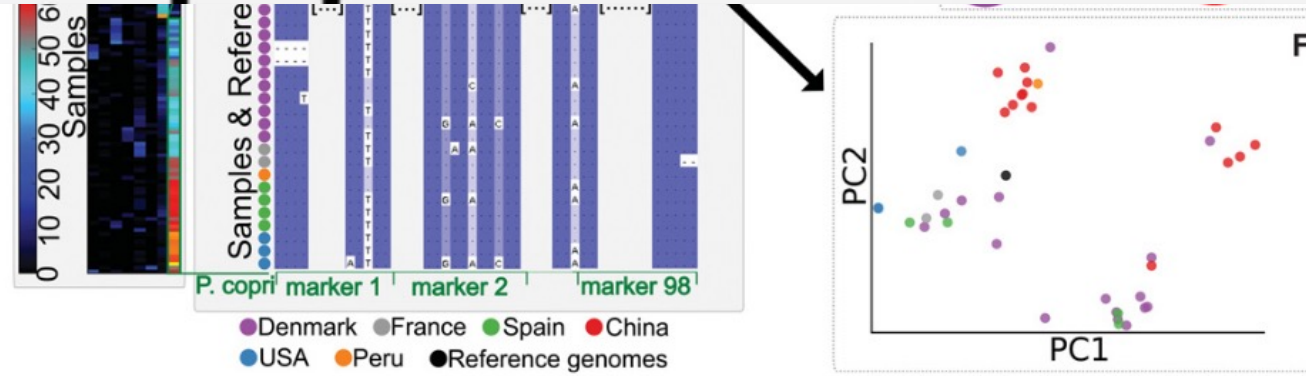
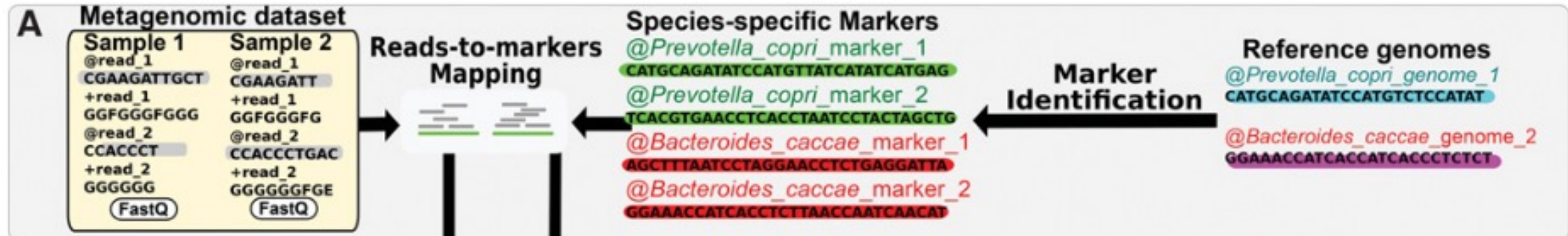
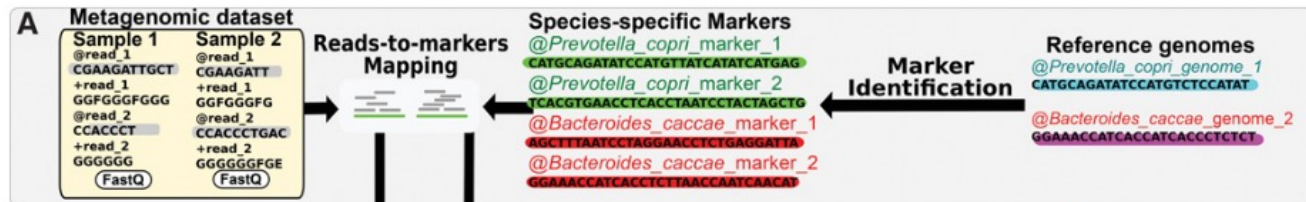
Taxonomic profilers - 3: marker gene/region approach

mOTUs:

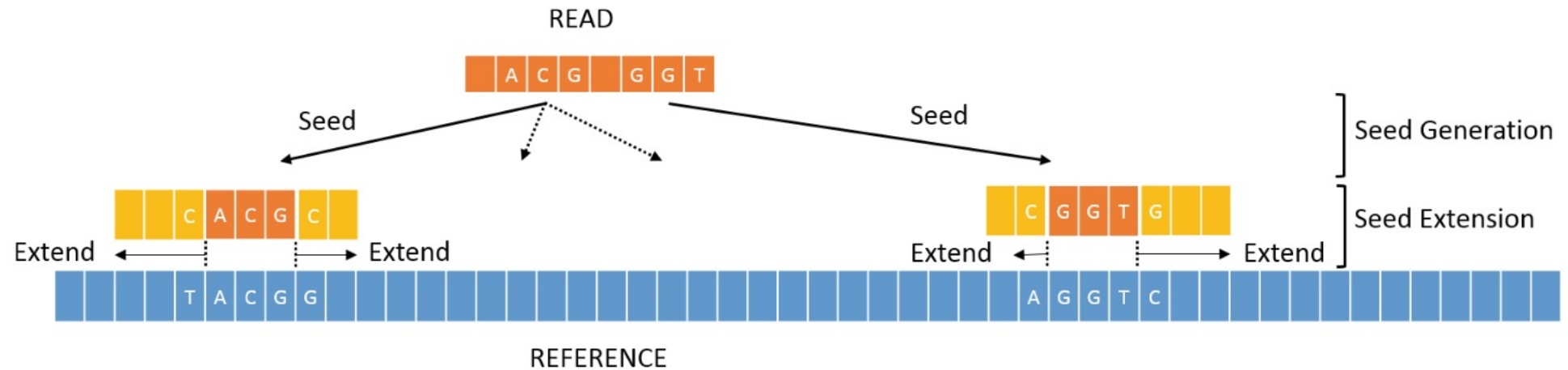


Taxonomic profilers - 3: marker gene/region approach

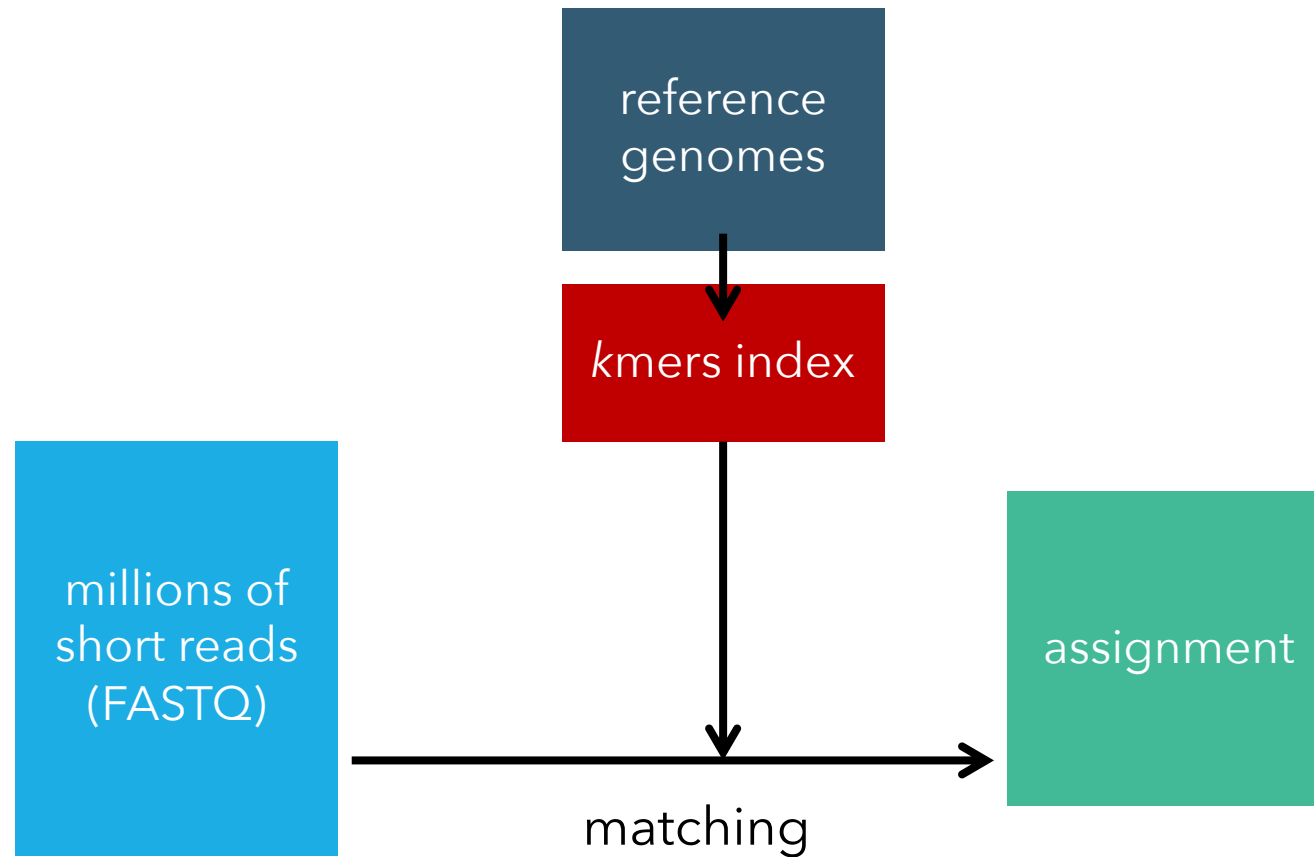
MetaPhlan:



Mapping reads = aligning reads

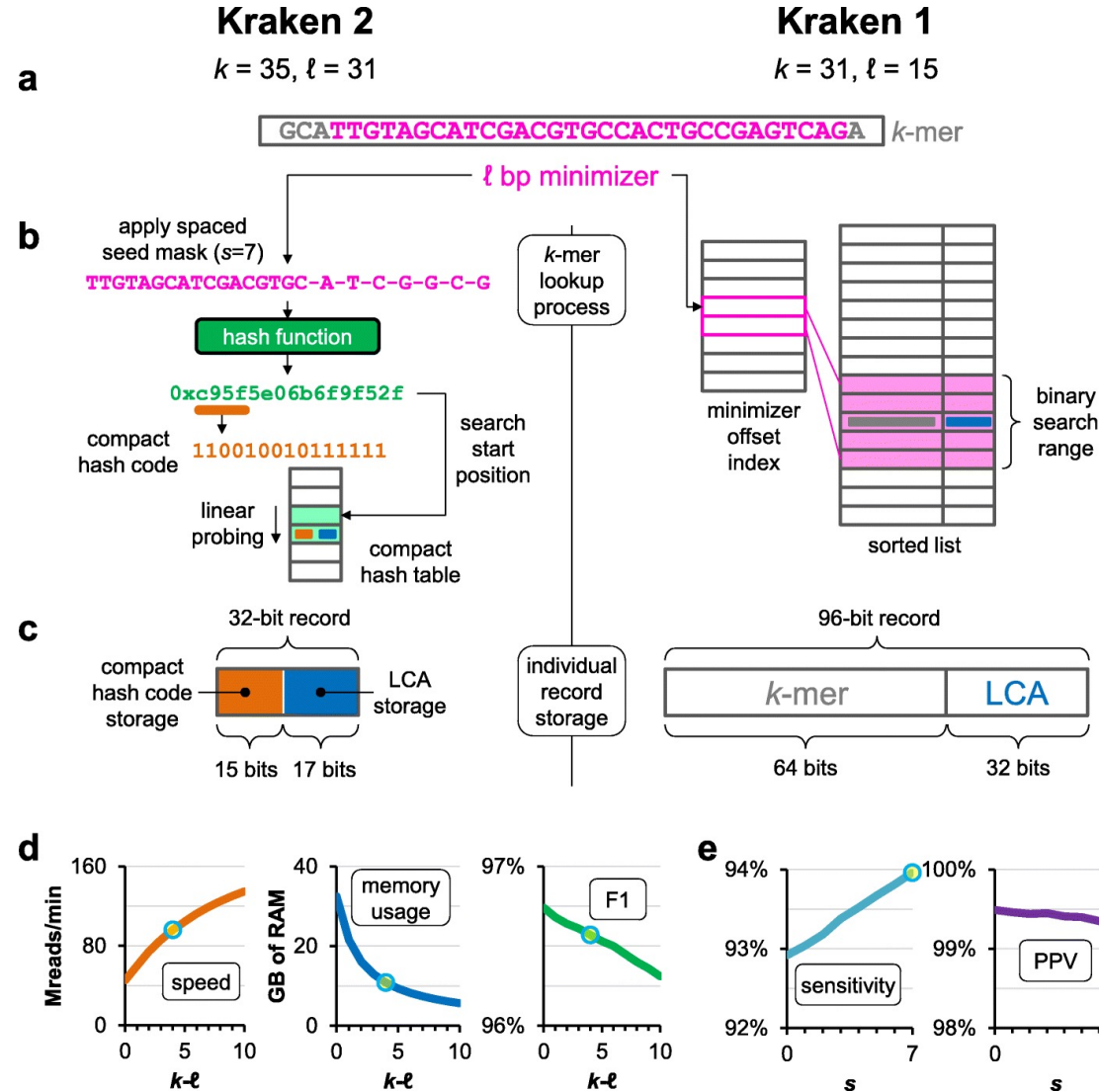


Taxonomic profilers - 4

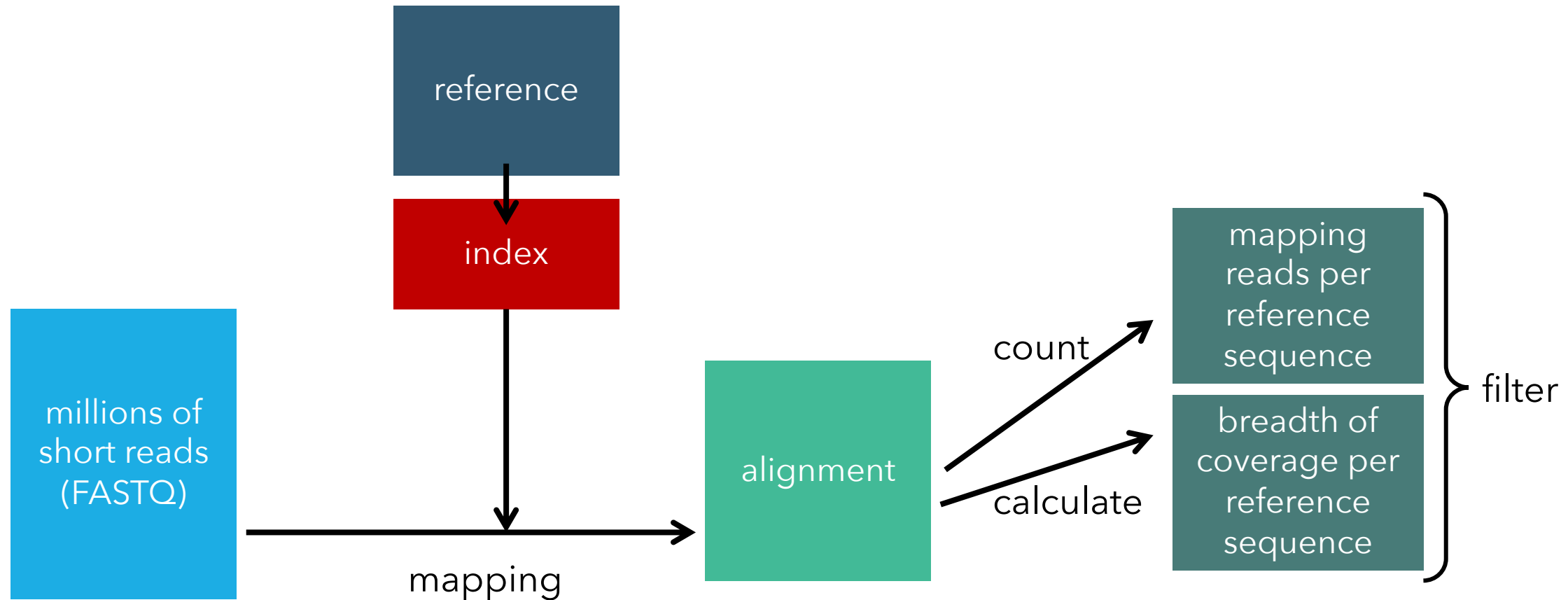


Taxonomic profilers - 4

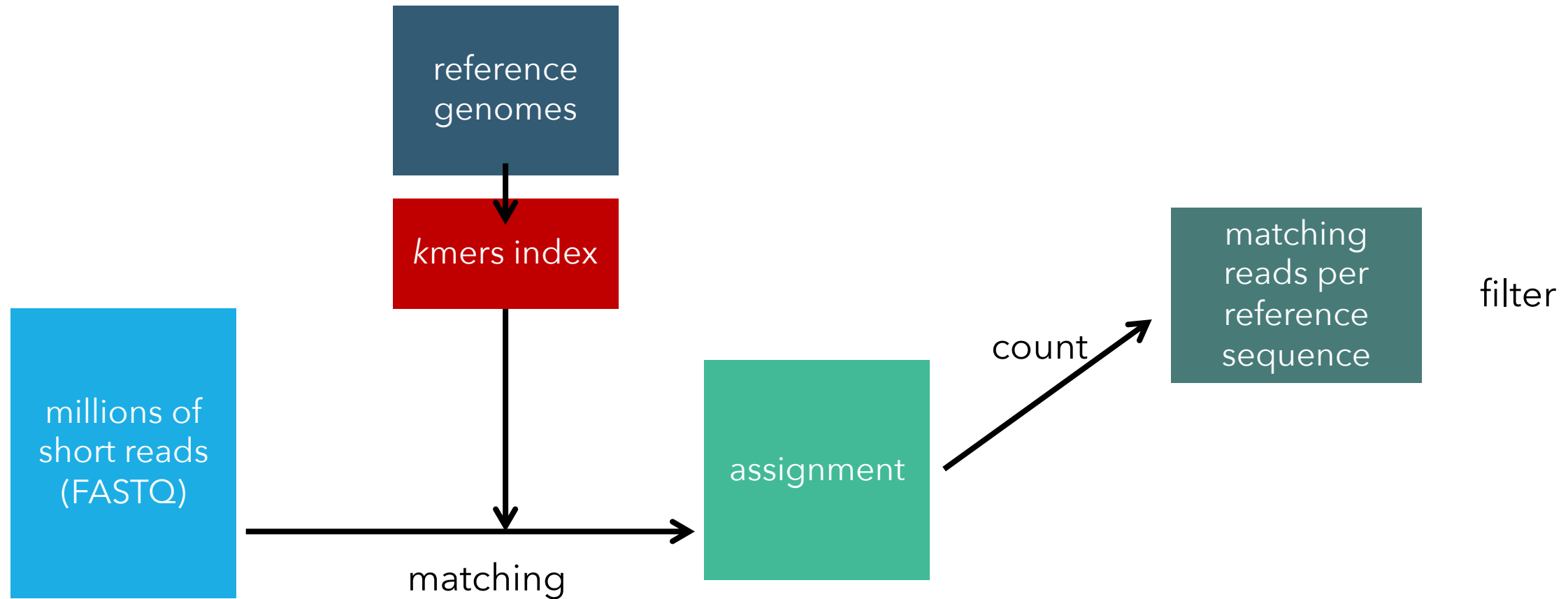
kraken2:



Mapping-based taxonomic profilers: quantification & specificity filtering



*k*mer-based taxonomic profilers : quantification & specificity filtering






How to choose a profiler? Benchmarks

ANALYSIS

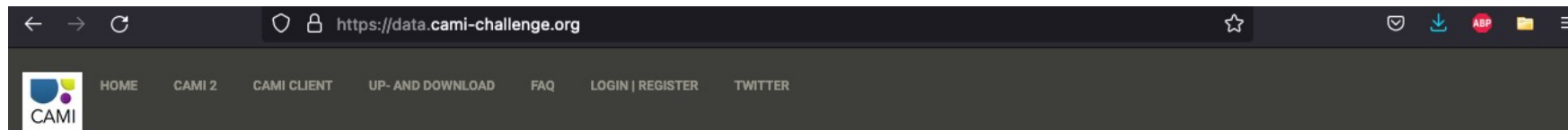
OPEN

Critical Assessment of Metagenome Interpretation—a benchmark of metagenomics software

Alexander Sczyrba^{1,2,48}, Peter Hofmann^{3-5,48}, Peter Belmann^{1,2,4,5,48}, David Koslicki⁶, Stefan Janssen^{4,7,8}, Johannes Dröge³⁻⁵, Ivan Gregor³⁻⁵, Stephan Majda^{3,47}, Jessika Fiedler^{3,4}, Eik Dahms³⁻⁵, Andreas Bremges^{1,2,4,5,9}, Adrian Fritz^{4,5}, Ruben Garrido-Oter^{3-5,10,11}, Tue Sparholt Jørgensen¹²⁻¹⁴, Nicole Shapiro¹⁵, Philip D Blood¹⁶, Alexey Gurevich¹⁷, Yang Bai^{10,47}, Dmitrij Turaev¹⁸, Matthew Z DeMaere¹⁹, Rayan Chikhi^{20,21}, Niranjana Nagarajan²², Christopher Quince²³, Fernando Meyer^{4,5}, Monika Balvočiūtė²⁴, Lars Hestbjerg Hansen¹², Søren J Sørensen¹³, Burton K H Chia²², Bertrand Denis²², Jeff L Froula¹⁵, Zhong Wang¹⁵, Robert Egan¹⁵, Dongwan Don Kang¹⁵, Jeffrey J Cook²⁵, Charles Deltel^{26,27}, Michael Beckstette²⁸, Claire Lemaitre^{26,27}, Pierre Peterlongo^{26,27}, Guillaume Rizk^{27,29}, Dominique Lavenier^{21,27}, Yu-Wei Wu^{30,31}, Steven W Singer^{30,32}, Chirag Jain³³, Marc Strous³⁴, Heiner Klingenberg³⁵, Peter Meinicke³⁵, Michael D Barton¹⁵, Thomas Lingner³⁶, Hsin-Hung Lin³⁷, Yu-Chieh Liao³⁷, Genivaldo Gueiros Z Silva³⁸, Daniel A Cuevas³⁸, Robert A Edwards³⁸, Surya Saha³⁹, Vitor C Piro^{40,41}, Bernhard Y Renard⁴⁰, Mihai Pop^{42,43}, Hans-Peter Klenk⁴⁴, Markus Göker⁴⁵, Nikos C Kyrpides¹⁵, Tanja Woyke¹⁵, Julia A Vorholt⁴⁶, Paul Schulze-Lefert^{10,11}, Edward M Rubin¹⁵, Aaron E Darling¹⁹ , Thomas Rattei¹⁸  & Alice C McHardy^{3-5,11} 



How to choose a profiler? Benchmarks

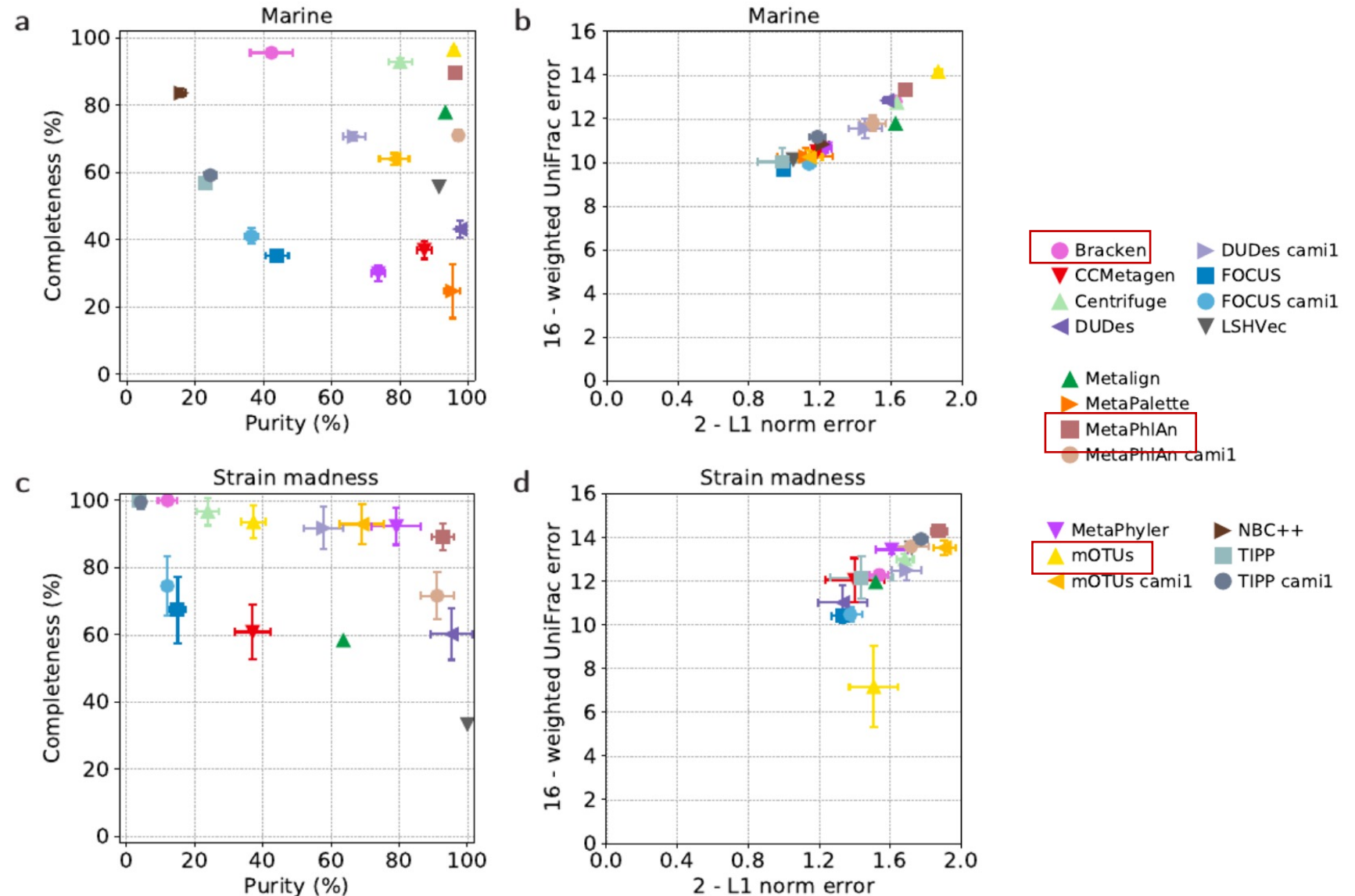


Critical Assessment of Metagenome Interpretation

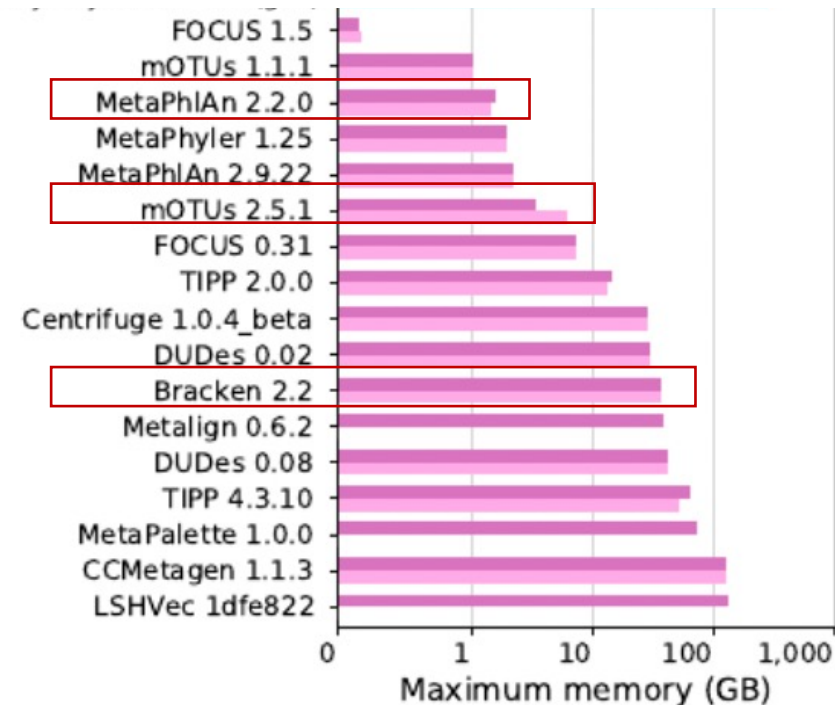
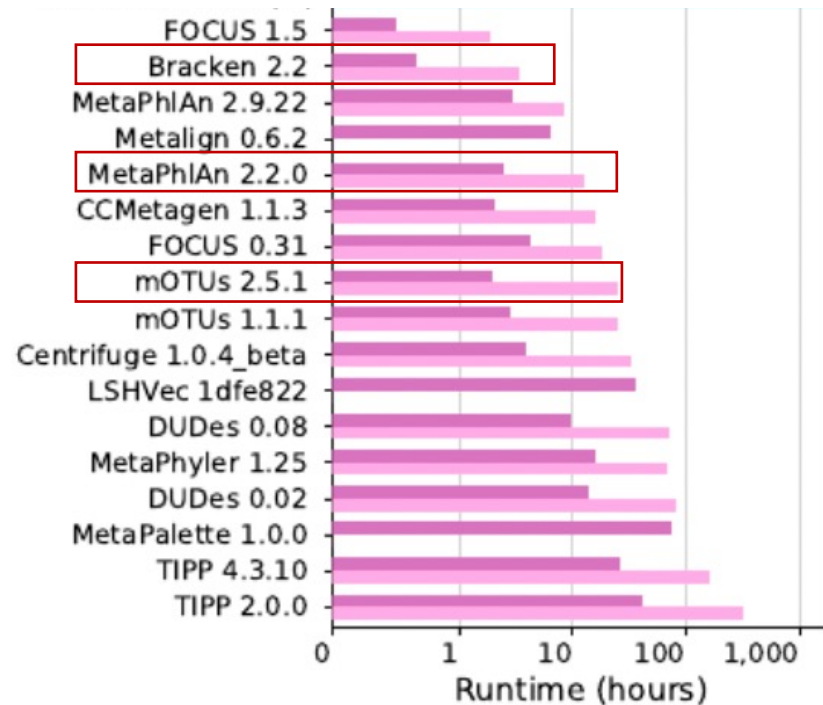
AGCTACG AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGTAACGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTACG AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGTAACGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG
AGCTAC AAAAGTACGATT TAACGTACCCCTACGTACGTACGT ACGTACGTAC ACGT CGTACGTACG



How to choose a profiler? Benchmarks



How to choose a profiler? Benchmarks



Taxonomic profiling
 ■ marine
 ■ strain madness



Thanks for your attention!



a.u.s.heintzbuschart@uva.nl

SP C2.205



github.com/a-h-b



twitter.com/_a_h_b_

