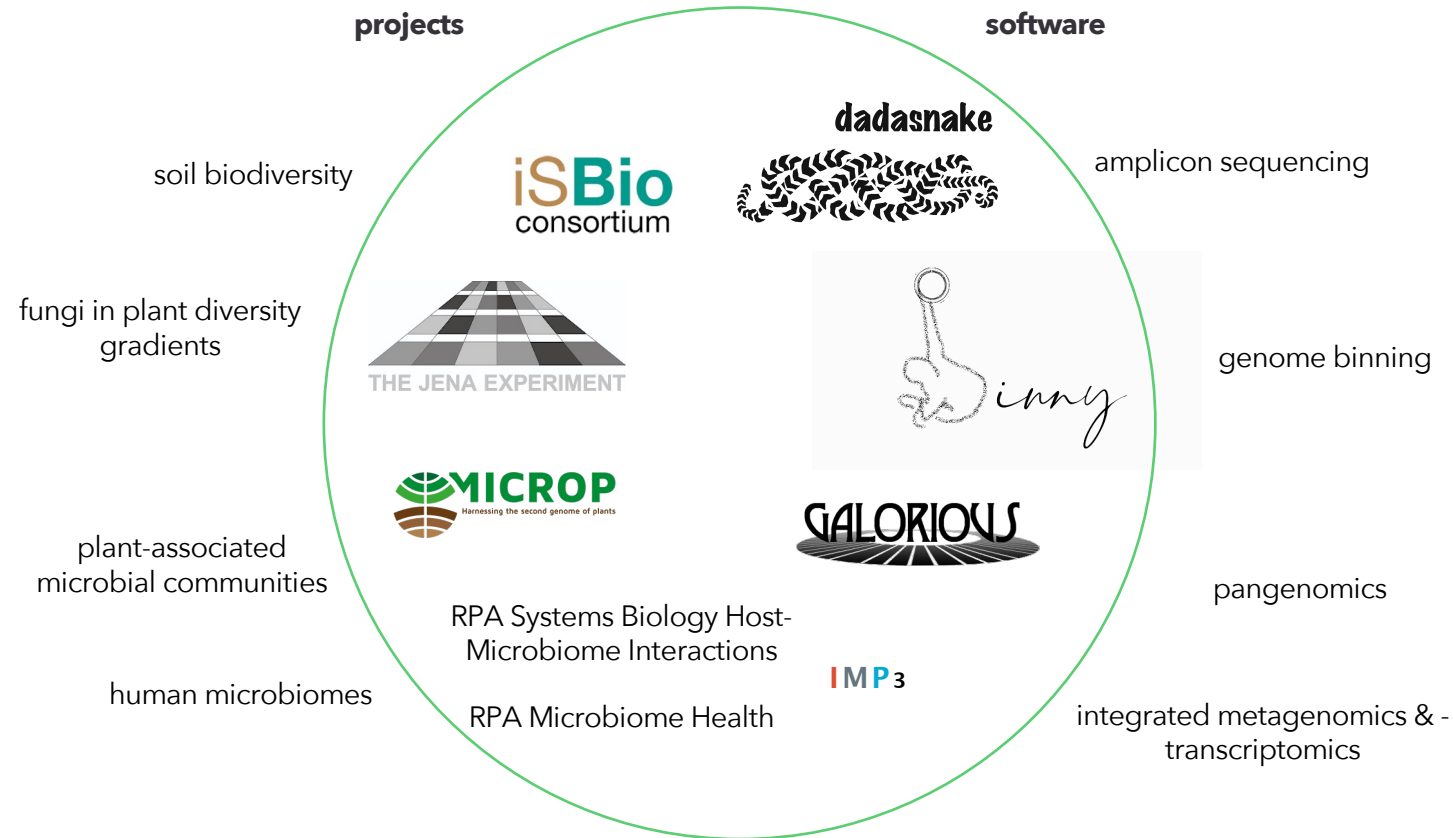


# Metagenomics 101

Anna Heintz-Buschart  
February 2022

# Who is here?

- A.H.-B.:



2008

MSc Biology (Microbiology, Botany, Molecular & Cell Biology)



PhD: Fungal human pathogen

- compound screening, mode-of action
- gene expression analysis



2011

Postdoc: Gene regulatory network modelling

2012

Postdoc: Integrated meta-omics

- human microbiome, wastewater treatment
- metagenomics, metatranscriptomics, metaproteomics
- lab automation
- bioinformatics pipelines



2017

Metagenomics support:

- biodiversity
- soil, plants, animal microbiomes
- bioinformatics pipelines
- data integration



2021

Assistant Prof Microbial Metagenomics

- meta-omics integration
- human and plant microbiomes



# Who is here?

# "Metagenomics"

- "**directly** accessing the **genomes** of [...] organisms that cannot be, or have not been, cultured **by isolating their DNA**" (Handelsman *et al.* Chem Biol. 1998)
- "accessing" (nowadays): by sequencing
- uncultured organisms: usually more than one



← Tweet



**Irene Newton**  
@chicaScientific

16S rRNA gene sequencing is NOT metagenomics,  
people

4:50 PM · Dec 29, 2020 · Twitter W

← Tweet



**Willem van Schaik**  
@WvSchaik

Replying to @WvSchaik

And 16S is not metagenomics

6:31 PM · Nov 9, 2015 · Twitter for Android

← Tweet

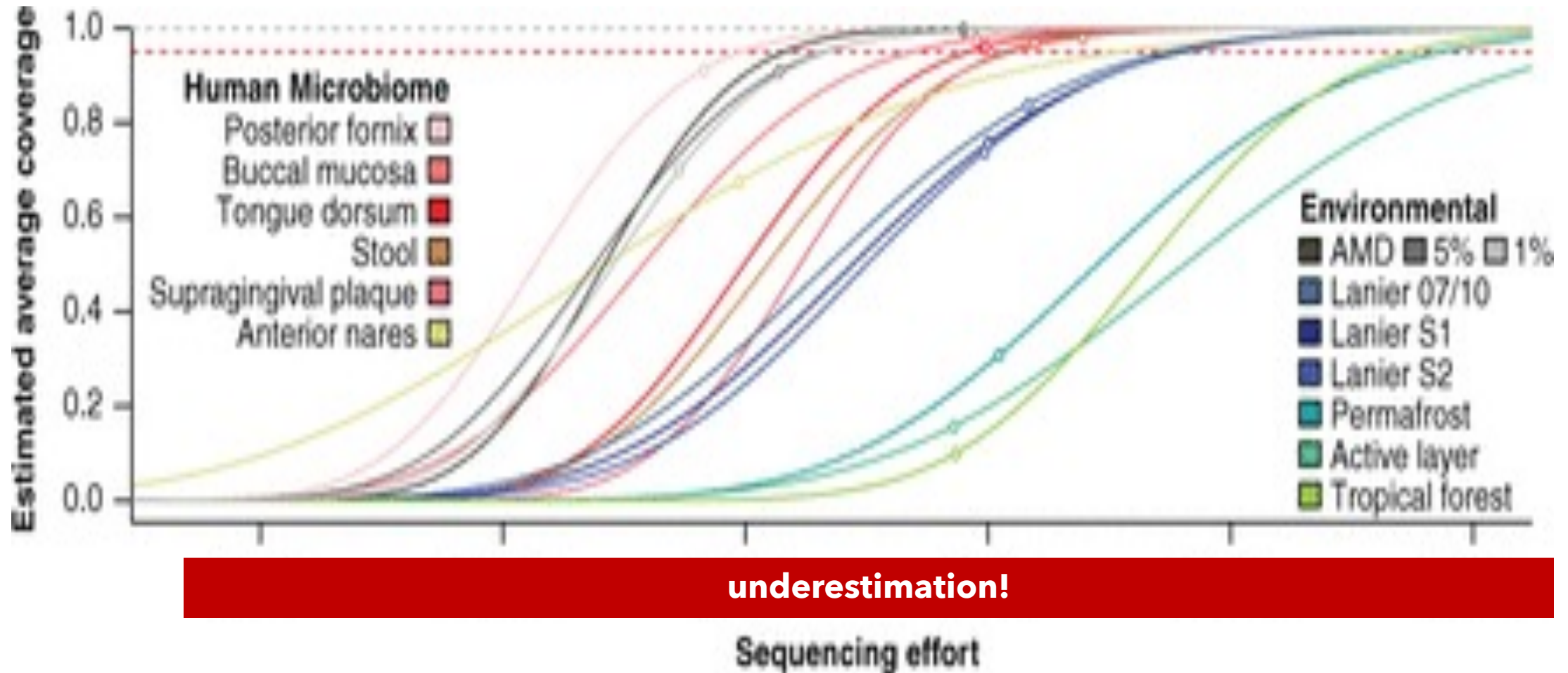


**Ken McGrath**  
@DrKenMcGrath

Hey @illumina - 16S profiling is not metagenomics.  
Clear distinction here should avoid confusion for  
people starting out in the field.

This course is about shotgun  
metagenomics

# WHOLE metagenome shotgun sequencing?



# Sequencing depth and metagenomic coverage



$$P(B = k) = \binom{R}{k} \sum_{\beta=k}^{\eta} \binom{R-k}{\beta-k} (-1)^{\beta-k} \alpha^{\beta} (1 - \beta\varphi)^{\beta-1} (1 - \beta\varphi\alpha)^{R-\beta}$$

$P(B=k)$ : probability of  $k$  gaps

$\alpha$ : relative abundance

$R$ : number of reads

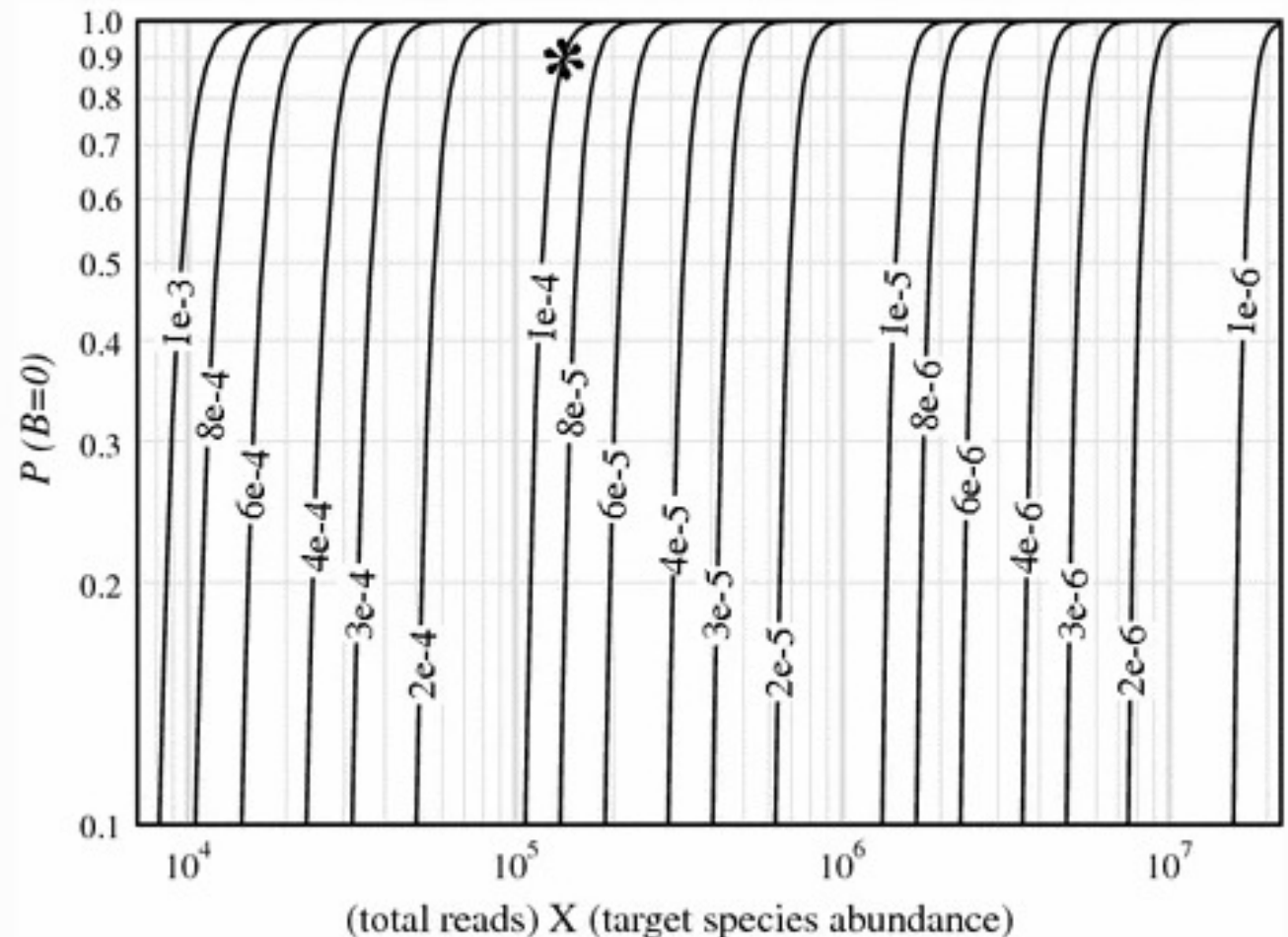
$\varphi$ : probability of a position being covered =  $L/\gamma$

$L$ : read length

$\gamma$ : genome size

$\eta$ : the smaller of  $R$  or the maximum

number of non-overlapping reads on the genome



# Sequencing depth and metagenomic coverage

- **beware**: these calculations are still often underestimations, because of:
  - sequencing biases
  - uneven copy numbers during replication
  - regions with high inter-species similarity and horizontal gene transfer
  - repeats
  - micro-diversity / strains



# Short & biased history of metagenomics

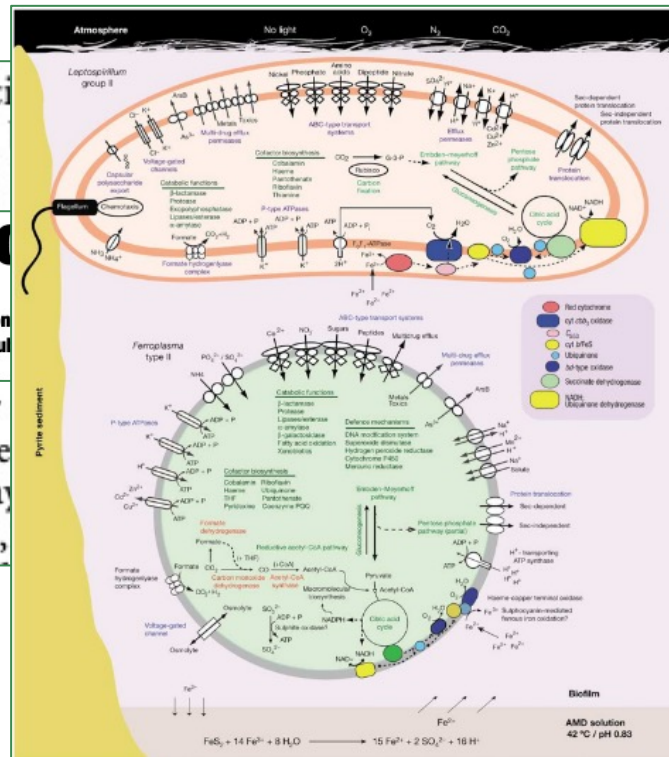
## • 2004 beginnings

sequenced  
million  
103,462

genome

Gene W. Tyson  
Edward M. Rubin

constraints. Over  
scaffolds longer  
region that may  
length of the 1,



## articles

A total of 76.2  
generated from  
r read). Analysis

ment

Richardson<sup>4</sup>, Victor V. Soloviyev<sup>4</sup>,

Physics, University of California,

ed into  
genomic  
mbined  
assembly



R  
collec  
of nor

TIGR role category	Total genes
Amino acid biosynthesis	37,118
Biosynthesis of cofactors, prosthetic groups, and carriers	25,905
Cell envelope	27,883
Cellular processes	17,260
Central intermediary metabolism	13,639
DNA metabolism	25,346
Energy metabolism	69,718
Fatty acid and phospholipid metabolism	18,558
Mobile and extrachromosomal element functions	1,061
Protein fate	28,768
Protein synthesis	48,012
Purines, pyrimidines, nucleosides, and nucleotides	19,912
Regulatory functions	8,392
Signal transduction	4,817
Transcription	12,756
Transport and binding proteins	49,185
Unknown function	38,067
Miscellaneous	1,864
Conserved hypothetical	794,061
Total number of roles assigned	1,242,230
Total number of genes	1,214,207

ing for roughly 410,000 reads, or 25% of the

LE

A total of 1.045 billion base pairs  
otated, and analyzed to elucidate



semblies. Our  
l-sampled ge-  
s with at least  
333 scaffolds  
panning 30.9  
S3), account-

# Short & biased history of metagenomics

- **2007-2012 'large-scale' projects**

OPEN ACCESS Freely available online

PLOS BIOLOGY

Community Page

## The Human Microbiome Project: A Community Resource for the Healthy Human Microbiome

**Dirk Gevers<sup>1</sup>, Rob Knight<sup>2,3</sup>, Joseph F. Petrosino<sup>4,5,6</sup>, Katherine Huang<sup>1</sup>, Amy L. McGuire<sup>7</sup>, Bruce W. Birren<sup>1</sup>, Karen E. Nelson<sup>8</sup>, Owen White<sup>9</sup>, Barbara A. Methé<sup>8\*</sup>, Curtis Huttenhower<sup>1,10\*</sup>**

<sup>1</sup> The Broad Institute of MIT and Harvard, Cambridge, Massachusetts, United States of America, <sup>2</sup> Department of Chemistry and Biochemistry, University of Colorado, Boulder, Colorado, United States of America, <sup>3</sup> Howard Hughes Medical Institute, Boulder, Colorado, United States of America, <sup>4</sup> Human Genome Sequencing Center, Baylor College of Medicine, Houston, Texas, United States of America, <sup>5</sup> Molecular Virology and Microbiology, Baylor College of Medicine, Houston, Texas, United States of America, <sup>6</sup> Alkek Center for Metagenomics and Microbiome Research, Baylor College of Medicine, Houston, Texas, United States of America, <sup>7</sup> Center for Medical Ethics and Health Policy, Baylor College of Medicine, Houston, Texas, United States of America, <sup>8</sup> J. Craig Venter Institute, Rockville, Maryland, United States of America, <sup>9</sup> Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, Maryland, United States of America, <sup>10</sup> Biostatistics, Harvard School of Public Health, Boston, Massachusetts, United States of America

# Short & biased history of metagenomics

- around 2014 tailored tools

*Bioinformatics*, 31(10), 2015, 1674–1676  
doi: 10.1093/bioinformatics/btv033  
Advance Access Publication Date: 20 January 2015  
Applications Note

OXFORD

## BRIEF COMMUNICATIONS

### Binning metagenomic contigs by coverage and composition

Johannes Alneberg<sup>1,8</sup>, Brynjar Smári Bjarnason<sup>1,8</sup>,  
Ino de Bruijn<sup>1,2</sup>, Melanie Schirmer<sup>3</sup>, Joshua Quick<sup>4,5</sup>,  
Umer Z Ijaz<sup>3</sup>, Leo Lahti<sup>6,7</sup>, Nicholas J Loman<sup>4</sup>,  
Anders F Andersson<sup>1,9</sup> & Christopher Quince<sup>3,9</sup>

PeerJ

### MetaBAT, an efficient tool for accurately reconstructing single genomes from complex microbial communities

Dongwan D. Kang<sup>1,2</sup>, Jeff Froula<sup>1,2</sup>, Rob Egan<sup>1,2</sup> and Zhong Wang<sup>1,2,3</sup>

<sup>1</sup> Department of Energy Joint Genome Institute, Walnut Creek, CA, USA

<sup>2</sup> Genomics Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA

<sup>3</sup> School of Natural Sciences, University of California at Merced, Merced, CA, USA

Sequence analysis

### MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct *de Bruijn* graph

Dinghua Li<sup>1,†</sup>, Chi-Man Liu<sup>2,†</sup>, Ruibang Luo<sup>2,†</sup>, Kunihiro Sadakane<sup>3</sup> and Tak-Wah Lam<sup>1,2,\*</sup>



# Short & biased history of metagenomics

- since 2015 pipelines & platforms

Bioinformatics, 32(16), 2016, 2520–2523  
doi: 10.1093/bioinformatics/btw183  
Advance Access Publication Date: 8 April 2016  
Applications Note



TECHNOLOGY REPORT  
published: 24 January 2019  
doi: 10.3389/fmicb.2018.03049



Uritskiy et al. *Microbiome* (2018) 6:158  
<https://doi.org/10.1186/s40168-018-0541-1>

Genome analysis

## MOCAT2: a metagenomic assembly, annotation and profiling framework

Jens Roat Kultima<sup>1</sup>, Luis Pedro Coelho<sup>1</sup>, Kristoffer Forslund<sup>1</sup>,  
Jaime Huerta-Cepas<sup>1</sup>, Simone S. Li<sup>1,2</sup>, Marja Driessen<sup>1</sup>,  
Anita Yvonne Voigt<sup>1,3</sup>, Georg Zeller<sup>1</sup>, Shinichi Sunagawa<sup>1</sup> and  
Peer Bork<sup>1,3,4,5,\*</sup>

## SqueezeMeta, A Highly Portable, Fully Automatic Metagenomic Analysis Pipeline

Javier Tamames\* and Fernando Puente-Sánchez



TOOLS AND RESOURCES



## Integrating taxonomic, functional, and strain-level profiling of diverse microbial communities with bioBakery 3

Francesco Beghini<sup>1†</sup>, Lauren J McIver<sup>2†</sup>, Aitor Blanco-Míguez<sup>1</sup>, Leonard Dubois<sup>1</sup>,  
Francesco Asnicar<sup>1</sup>, Sagun Maharjan<sup>2,3</sup>, Ana Mailyan<sup>2,3</sup>, Paolo Manghi<sup>1</sup>,  
Matthias Scholz<sup>4</sup>, Andrew Maltez Thomas<sup>1</sup>, Mireia Valles-Colomer<sup>1</sup>,  
George Weingart<sup>2,3</sup>, Yancong Zhang<sup>2,3</sup>, Moreno Zolfo<sup>1</sup>, Curtis Huttenhower<sup>2,3\*</sup>,  
Eric A Franzosa<sup>2,3\*</sup>, Nicola Segata<sup>1,5\*</sup>



## Anvi'o: an advanced analysis and visualization platform for 'omics data

A. Murat Eren<sup>1,2</sup>, Özcan C. Esen<sup>1</sup>, Christopher Quince<sup>3</sup>,  
Joseph H. Vineis<sup>1</sup>, Hilary G. Morrison<sup>1</sup>, Mitchell L. Sogin<sup>1</sup> and  
Tom O. Delmont<sup>1</sup>

<sup>1</sup> Josephine Bay Paul Center, Marine Biological Laboratory, Woods Hole, MA, United States  
<sup>2</sup> Department of Medicine, The University of Chicago, Chicago, IL, United States  
<sup>3</sup> Warwick Medical School, University of Warwick, Coventry, United Kingdom

Microbiome

SOFTWARE

Open Access



## MetaWRAP—a flexible pipeline for genome-resolved metagenomic data analysis

Gherman V. Uritskiy, Jocelyne DiRuggiero\* and James Taylor\*

Narayanasamy et al. *Genome Biology* (2016) 17:260  
DOI 10.1186/s13059-016-1116-8

Genome Biology

SOFTWARE

Open Access



## IMP: a pipeline for reproducible reference-independent integrated metagenomic and metatranscriptomic analyses

Shaman Narayanasamy<sup>1†</sup>, Yohan Jarosz<sup>1†</sup>, Emilie E. L. Muller<sup>1,2</sup>, Anna Heintz-Buschart<sup>1</sup>, Malte Herold<sup>1</sup>,  
Anne Kaysen<sup>1</sup>, Cédric C. Laczný<sup>1,3</sup>, Nicolás Pinel<sup>4,5</sup>, Patrick May<sup>1</sup> and Paul Wilmes<sup>1\*</sup>

# Short & biased history of metagenomics

- since 2017 diversification of tools

Published online 13 January 2018

*Nucleic Acids Research*, 2018, Vol. 46, No. 6 e35  
doi: 10.1093/nar/gkx1321

## PlasFlow: predicting plasmid sequences in metagenomic data using genome signatures

Pawel S. Krawczyk<sup>1,2,\*</sup>, Leszek Lipinski<sup>1</sup> and Andrzej Dziembowski<sup>1,2</sup>

ANALYSIS

nature  
biotechnology

Measurement of bacterial replication rates in microbial communities

Christopher T Brown<sup>1</sup>, Matthew R Olm<sup>1</sup>, Brian C Thomas<sup>2</sup> & Jillian F Banfield<sup>2-4</sup>



## MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies

Dongwan D. Kang<sup>1</sup>, Feng Li<sup>2</sup>, Edward Kirton<sup>1</sup>, Ashleigh Thomas<sup>1</sup>, Rob Egan<sup>1</sup>, Hong An<sup>2</sup> and Zhong Wang<sup>1,3,4</sup>

Lind and Pollard *Microbiome* (2021) 9:58  
<https://doi.org/10.1186/s40168-021-01015-y>

Microbiome

METHODOLOGY

Open Access

## Accurate and sensitive detection of microbial eukaryotes from whole metagenome shotgun sequencing

Abigail L. Lind<sup>1</sup> and Katherine S. Pollard<sup>1,2,3,4,5\*</sup>



# Short & biased history of metagenomics

## • discoveries!



Cell Host & Microbe  
Article

### The *Prevotella copri* Complex Comprises Four Distinct Clades Underrepresented in Westernized Populations

Adrian Tett,<sup>1,\*</sup> Kun D. Huang,<sup>1,2</sup> Francesco Asnicar,<sup>1</sup> Hannah Fehlner-Peach,<sup>3</sup> Edoardo Pasolli,<sup>1,19</sup> Nicolai Karcher,<sup>1</sup> Federica Armanini,<sup>1</sup> Paolo Manghi,<sup>1</sup> Kevin Bonham,<sup>4,6</sup> Moreno Zolfo,<sup>1</sup> Francesca De Filippis,<sup>5,7</sup> Cara Magnabosco,<sup>8</sup> Richard Bonneau,<sup>8,9</sup> John Lusingu,<sup>10</sup> John Amuasi,<sup>11</sup> Karl Reinhard,<sup>12</sup> Thomas Rattei,<sup>13</sup> Fredrik Boulund,<sup>14</sup> Lars Engstrand,<sup>15</sup> Albert Zink,<sup>16</sup> Maria Carmen Collado,<sup>16</sup> Dan R. Littman,<sup>3</sup> Daniel Eibach,<sup>17,18</sup> Danilo Ercolini,<sup>5,7</sup> Omar Rota-Stabelli,<sup>2</sup> Curtis Huttenhower,<sup>4,6</sup> Frank Maixner,<sup>16</sup> and Nicola Segata<sup>1,20,\*</sup>

### ARTICLES

<https://doi.org/10.1038/s41564-020-00840-5>

nature  
microbiology

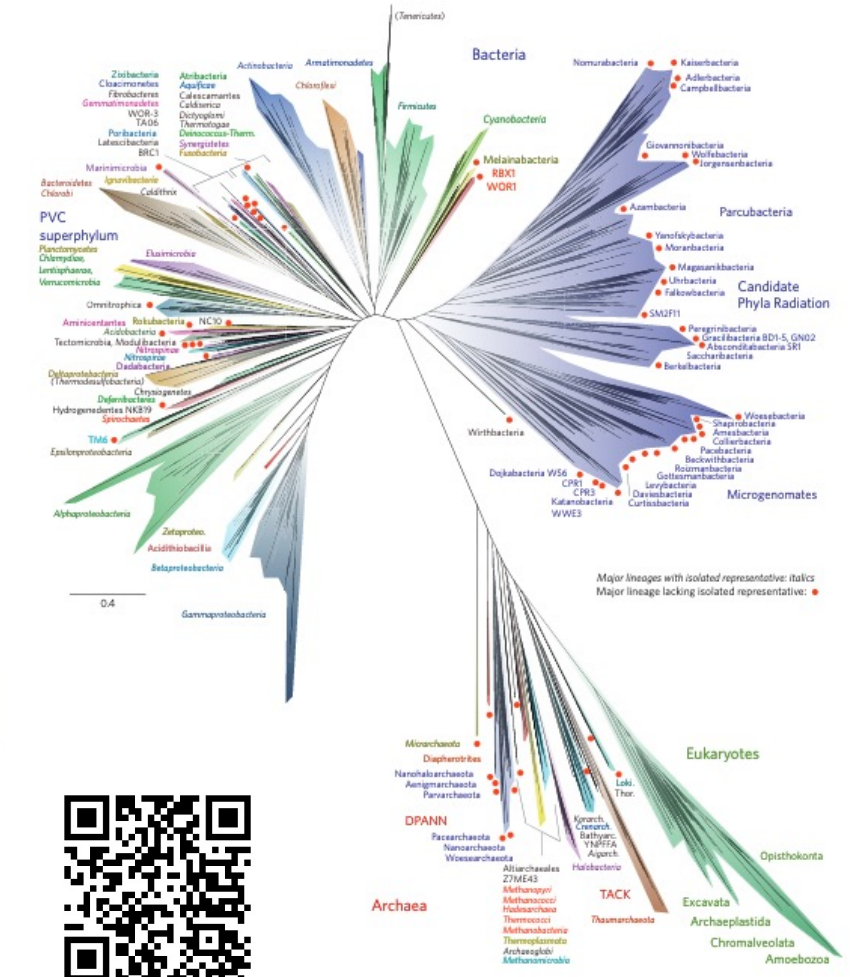


### OPEN

Genome-resolved metagenomics reveals site-specific diversity of episymbiotic CPR bacteria and DPANN archaea in groundwater ecosystems

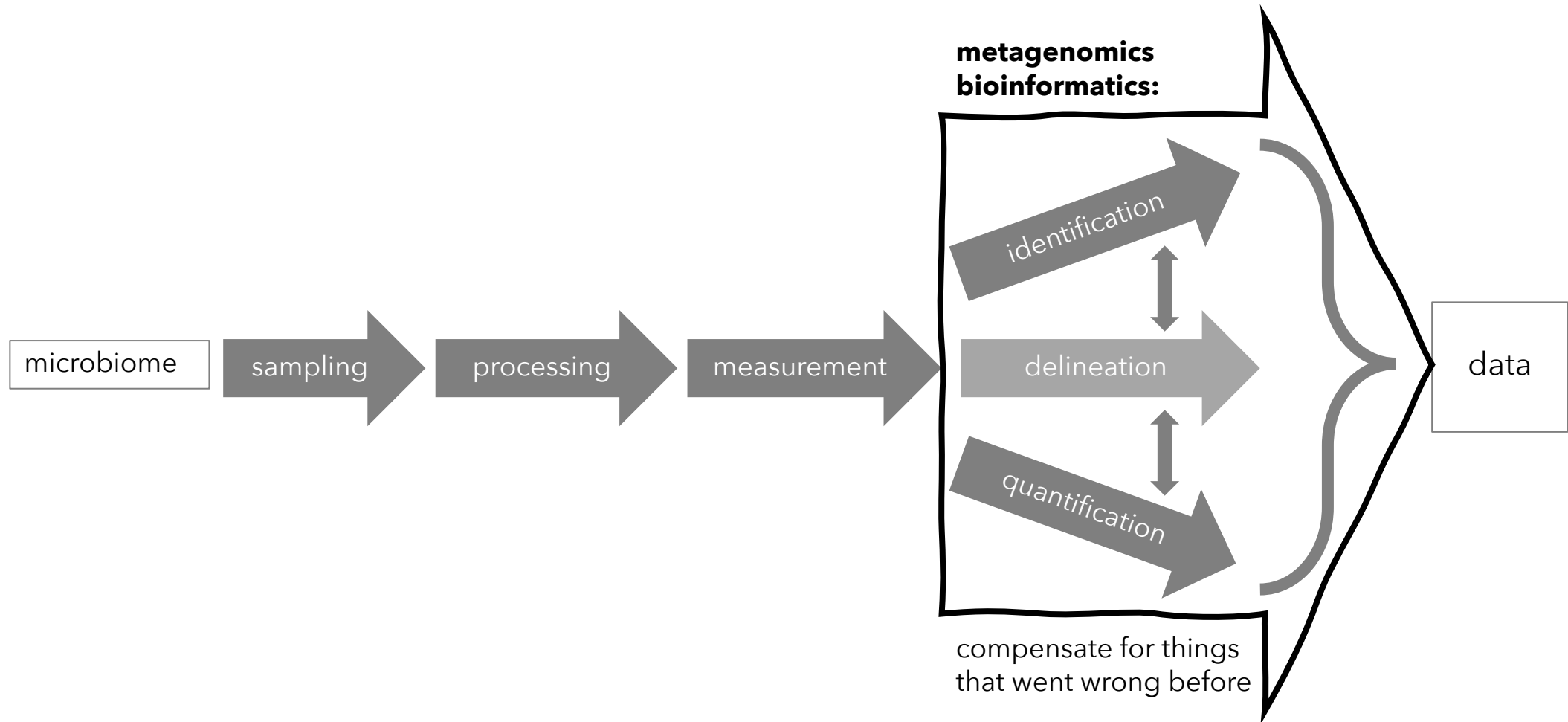
Christine He<sup>1</sup>, Ray Keren<sup>2</sup>, Michael L. Whittaker<sup>3,4</sup>, Ibrahim F. Farag<sup>5</sup>, Jennifer A. Doudna<sup>1,5,6,7,8</sup>, Jamie H. D. Cate<sup>1,5,6,7</sup> and Jillian F. Banfield<sup>1,4,9</sup>✉

LETTERS NATURE MICROBIOLOGY DOI: 10.1038/NMICROBIOL.2016.48



# Why do you study microbiomes?

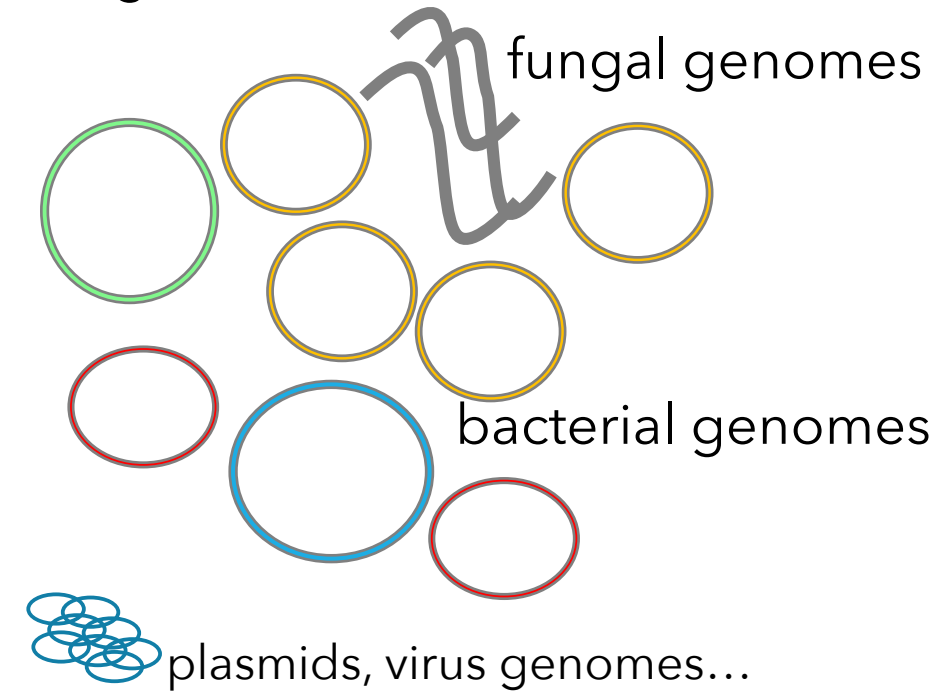
# Omics paradigm





# Measuring microbiomes: DNA based methods

the 'metagenome':



Line 1: Name

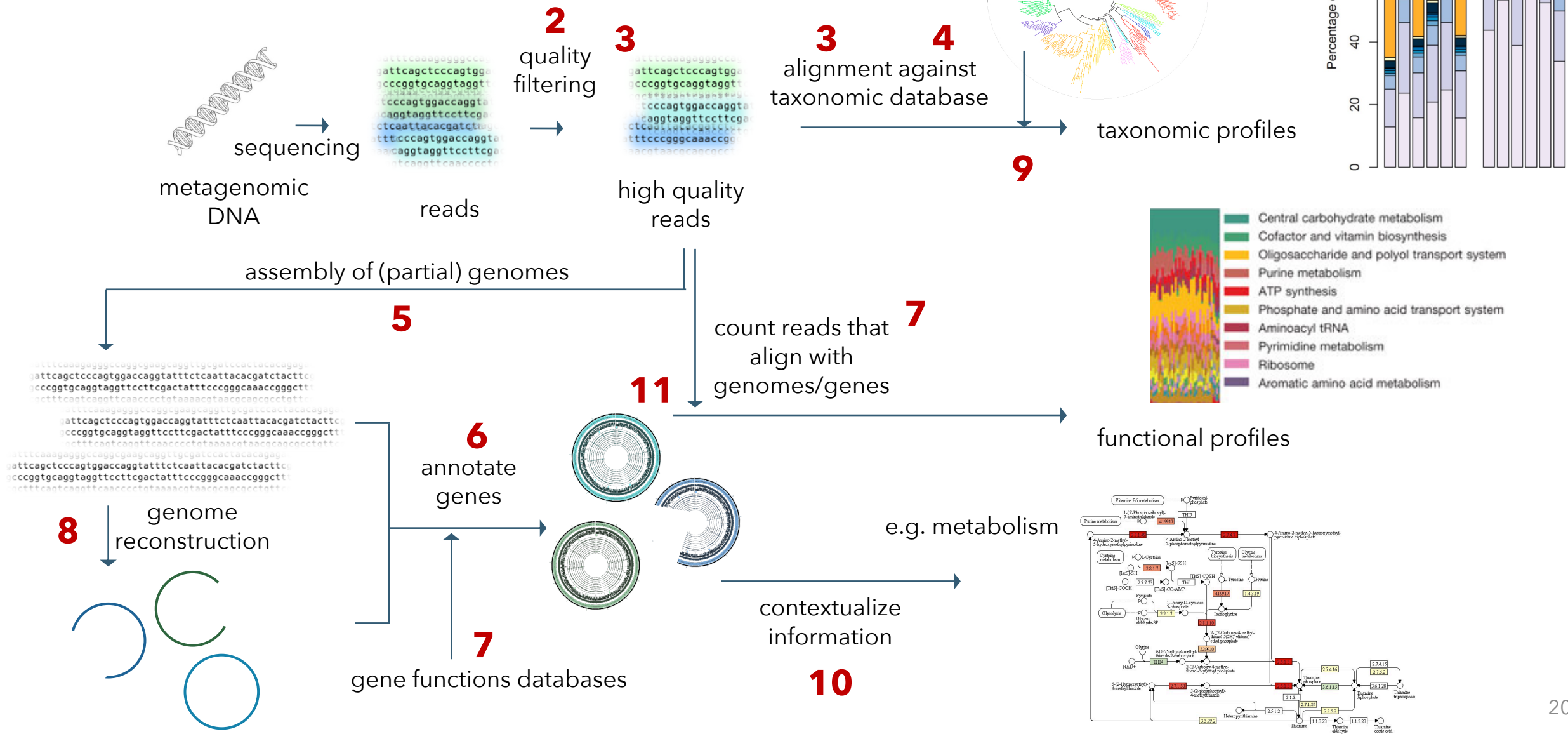
Line 2: Sequence

Line 3: anything

Line 4: Quality at each position

as many as we have reads (forward- & reverse files)

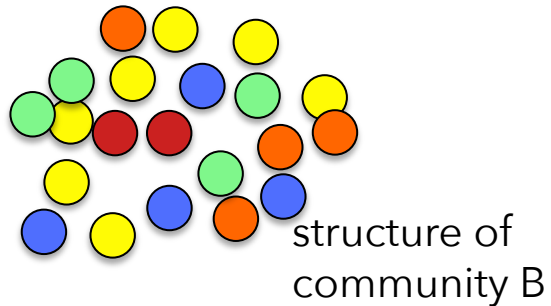
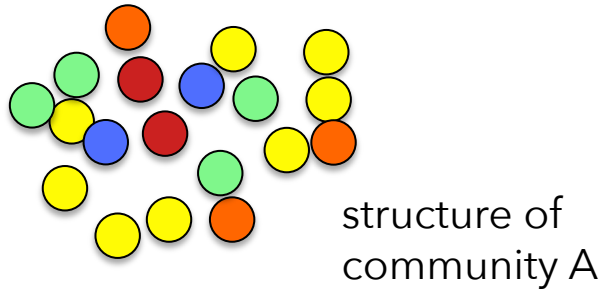
# Metagenomics overview





# Metagenomics (+ other omics) pipeline

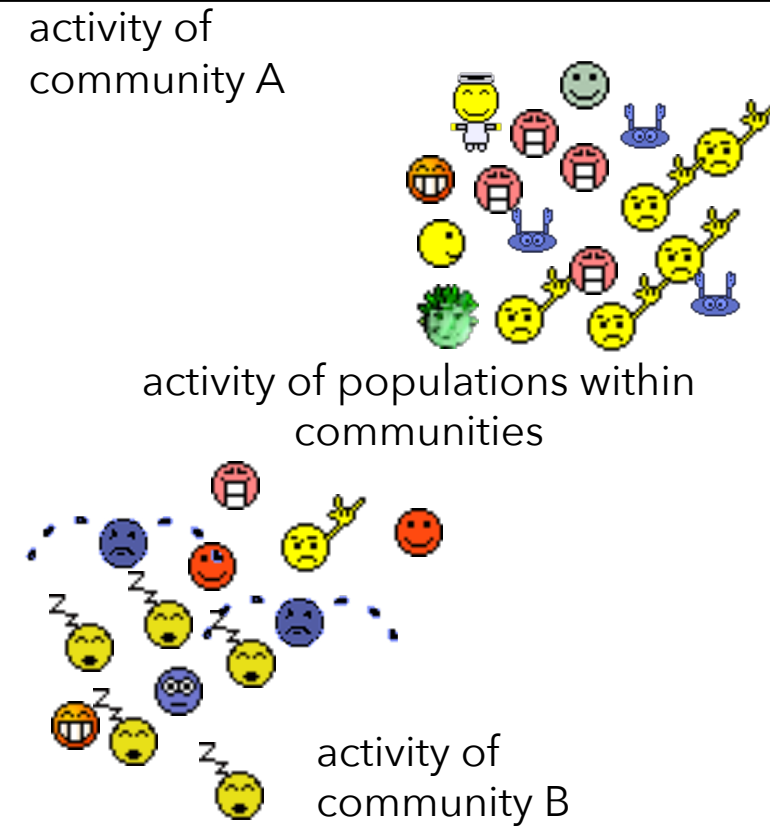
## metabarcoding:



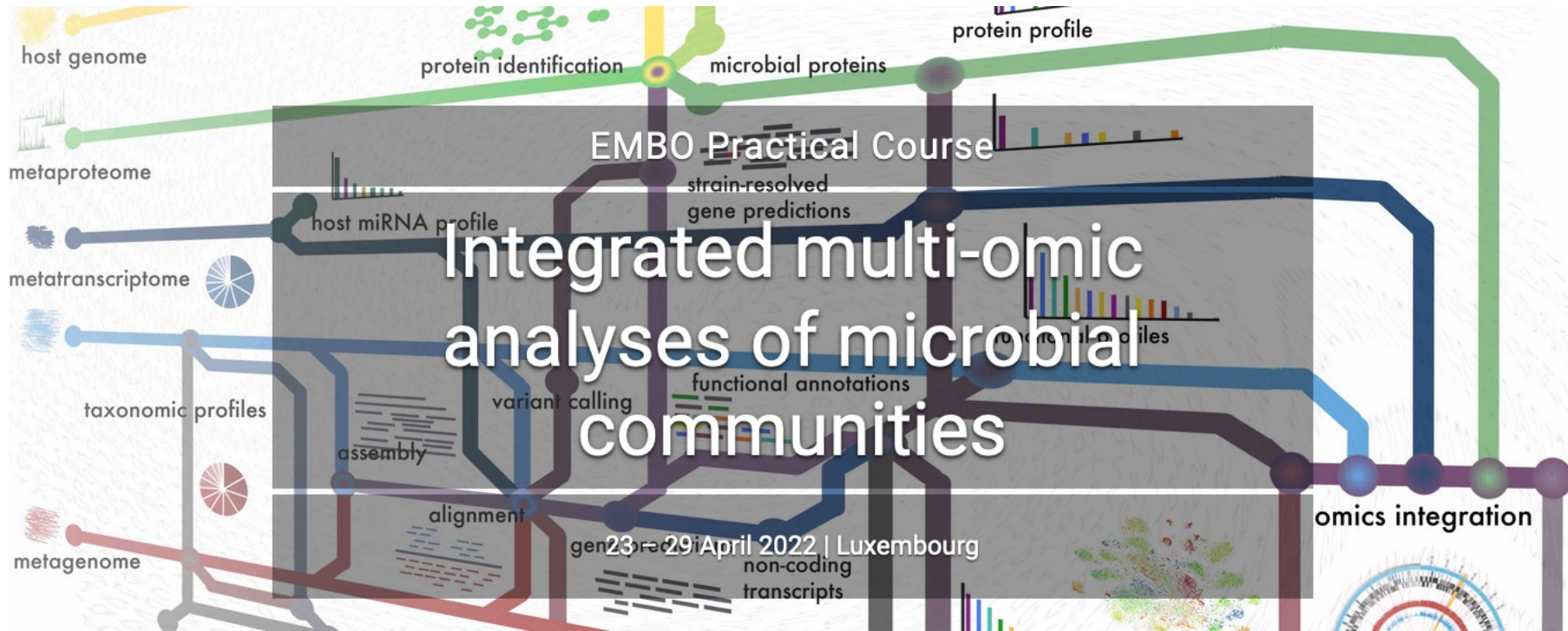
## metagenomics:



## functional omics:



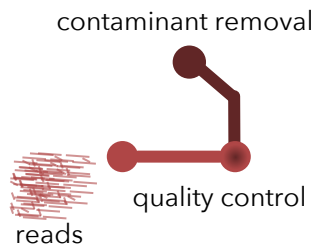
# Advertisement Break



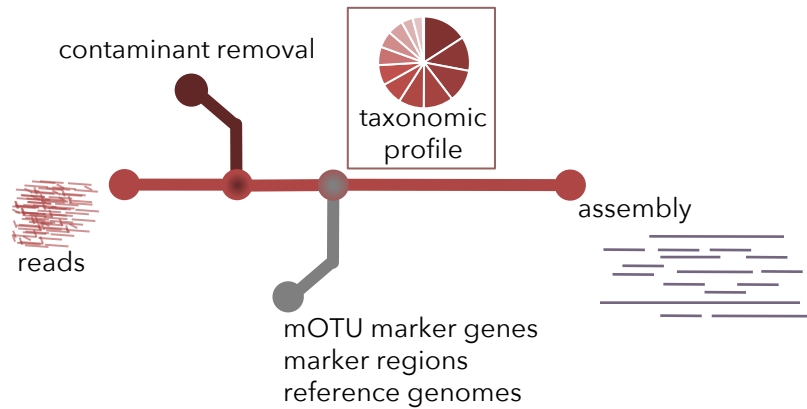
- 23rd-29th April, Luxembourg
- application deadline: 24th February



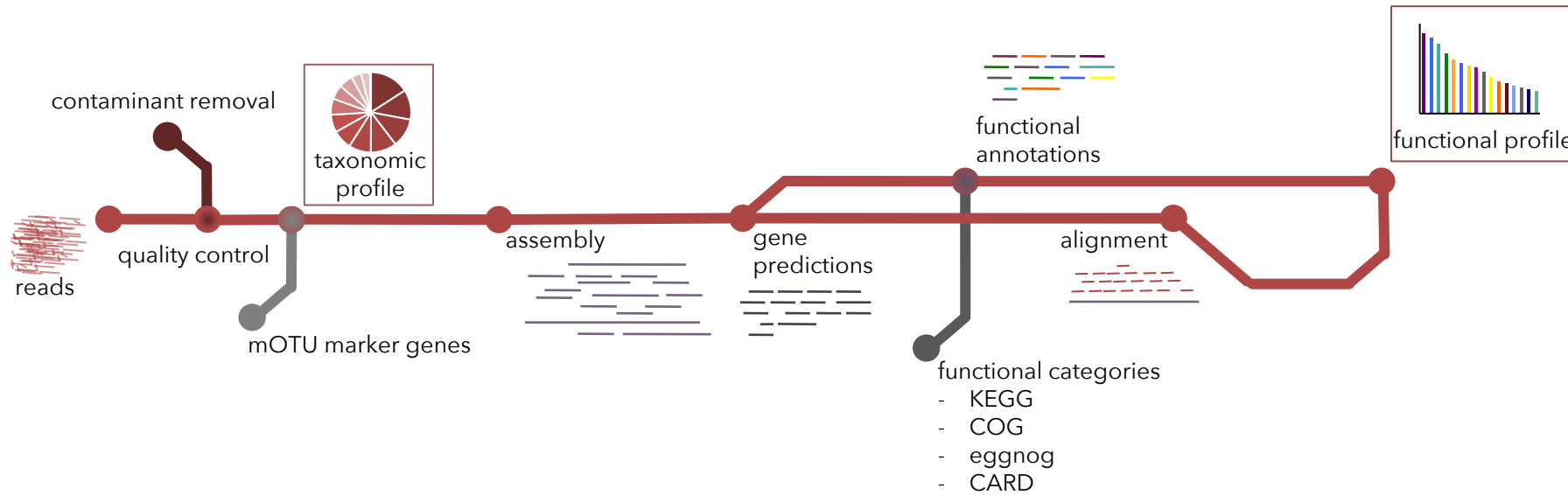
# Metagenomics (+ other omics) pipeline



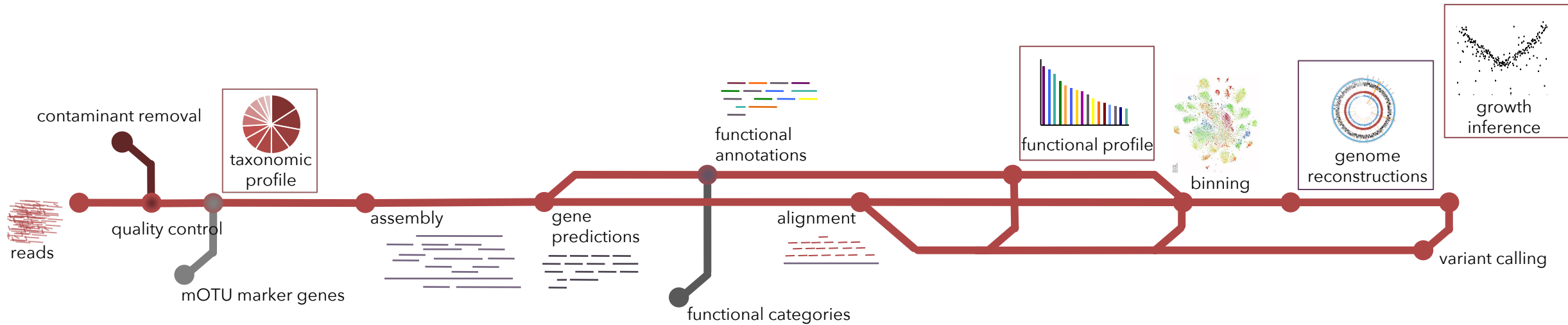
# Metagenomics (+ other omics) pipeline



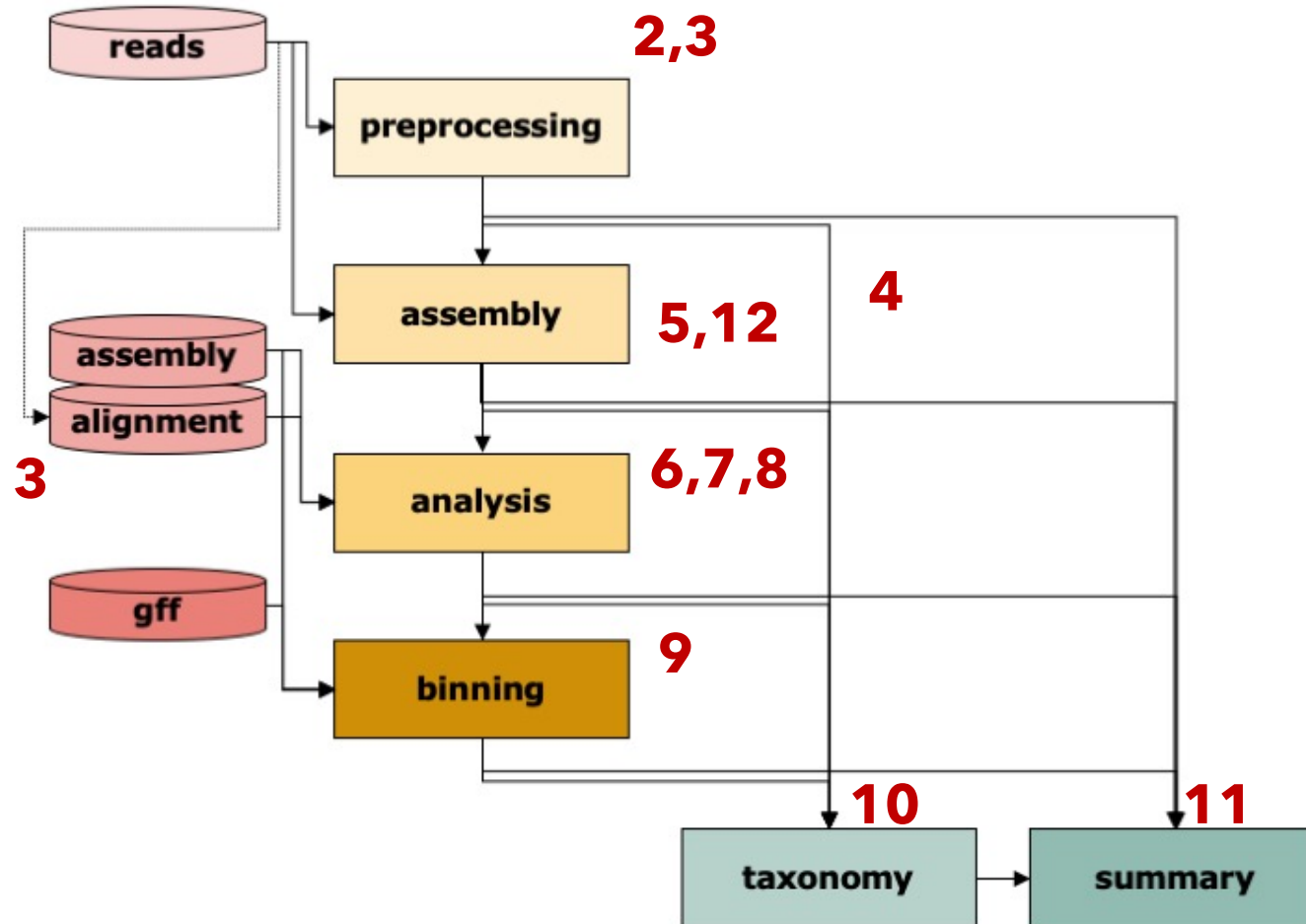
# Metagenomics (+ other omics) pipeline



# Metagenomics (+ other omics) pipeline



# Metagenomics (+ other omics) pipeline



# Thanks for your attention!



[a.u.s.heintzbuschart@uva.nl](mailto:a.u.s.heintzbuschart@uva.nl)

SP C2.205



[github.com/a-h-b](https://github.com/a-h-b)



[twitter.com/\\_a\\_h\\_b\\_](https://twitter.com/_a_h_b_)

