



Universidade de São Paulo
Instituto de Química



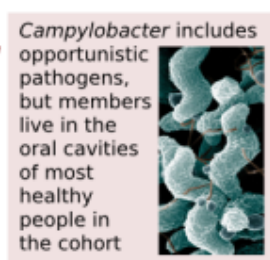
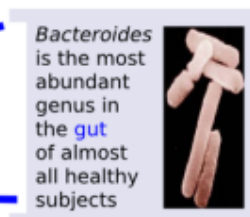
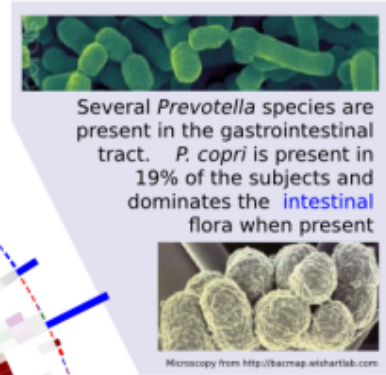
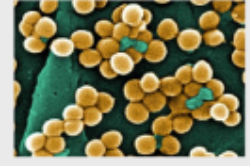
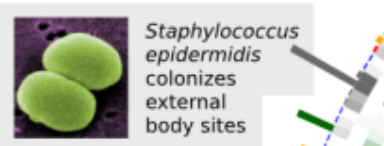
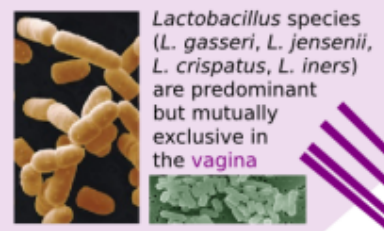
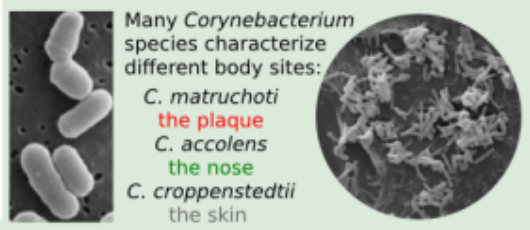
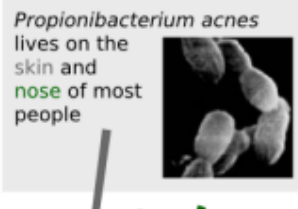
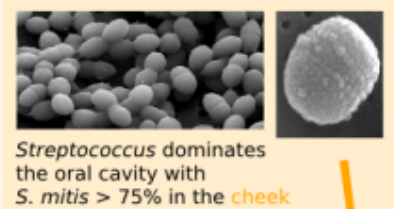
Análise de Microbiomas

João Carlos Setubal

Os microorganismos estão por toda parte

- São responsáveis por muitos processos **fundamentais para a vida do planeta** em geral e para **a vida dos seres humanos** em particular

A map of diversity in the human microbiome



○ Commensal microbes
 ☆ Potential pathogens

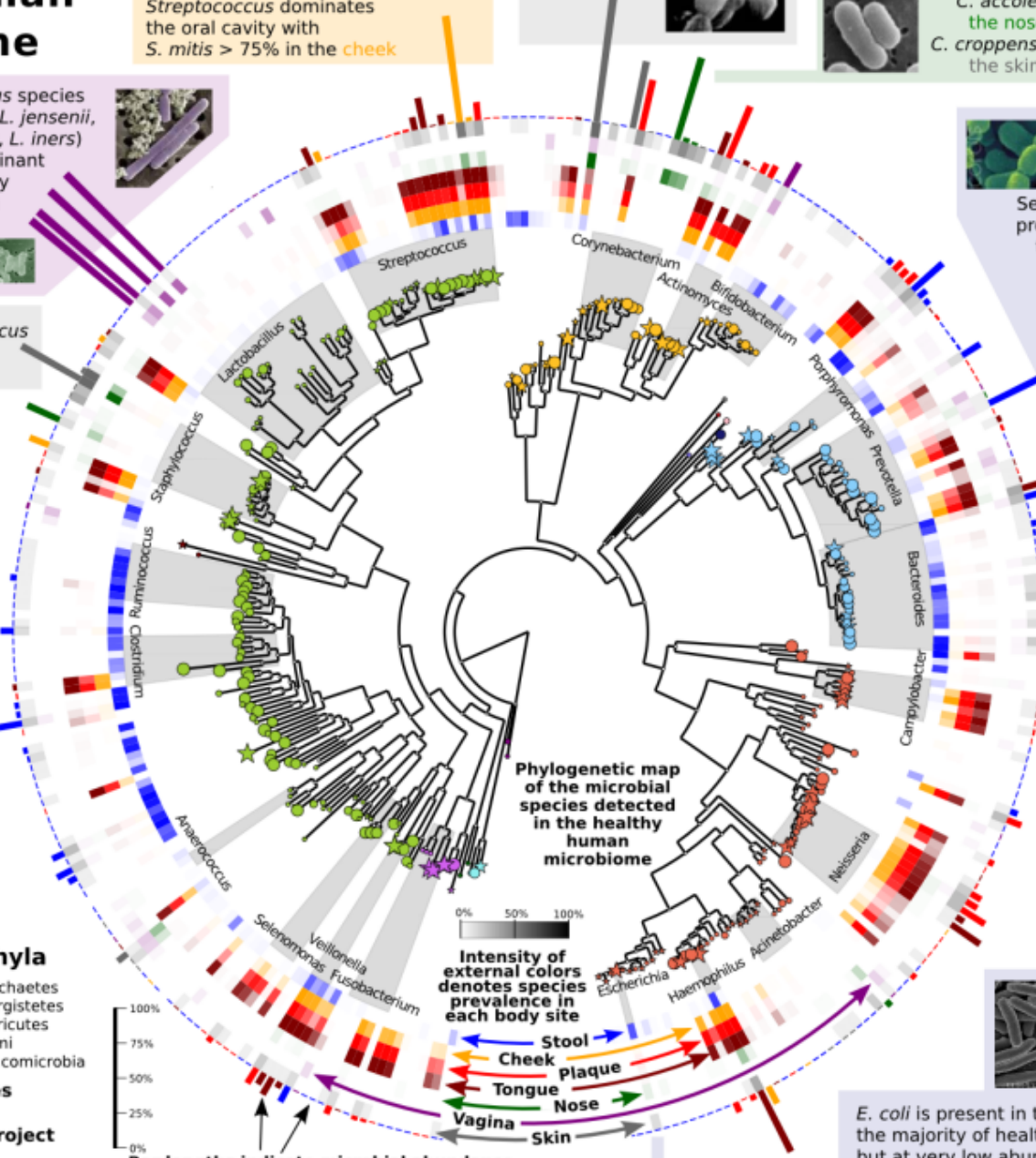
The four most abundant phyla

- Actinobacteria
- Bacteroidetes
- Firmicutes
- Proteobacteria

Low abundance phyla

- Chloroflexi
- Cyanobacteria
- Euryarchaeota
- Fusobacteria
- Lentisphaerae
- Spirochaetes
- Synergistetes
- Tenericutes
- Thermi
- Verrucomicrobia

National Institutes of Health Human Microbiome Project



E. coli is present in the **gut** of the majority of healthy subjects but at very low abundance

Projeto Microbioma Humano



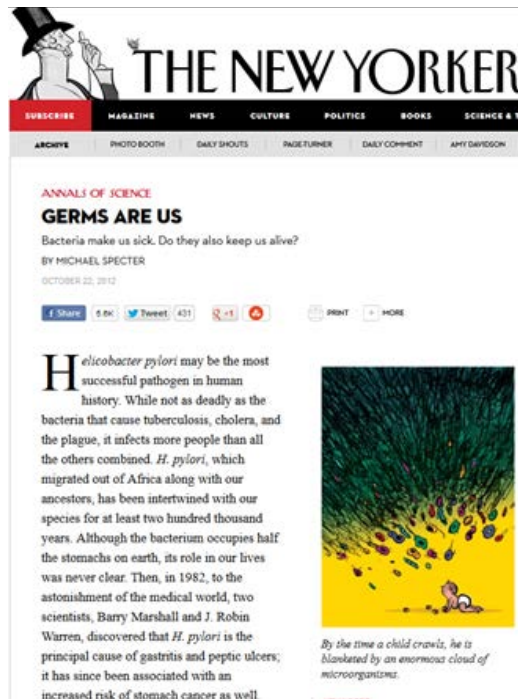
junho 2012

outubro 2012



My Microbiome and Me

Science 8 June 2012:



June 2012 Issue



maio 2013

www.earthmicrobiome.org



Home Defining the Tasks Getting Involved EMP Protocols and Standards Affiliations Publications Meetings EMP Logo

No categories



The Earth Microbiome Project is a systematic attempt to characterize the global microbial taxonomic and functional diversity for the benefit of the planet and mankind

Constructing the Microbial Biomap for Planet Earth

The Earth Microbiome Project is a proposed massively multidisciplinary effort to analyze microbial communities across the globe. The general premise is to examine microbial communities from their own perspective. Hence we propose to characterize the Earth by environmental parameter space into different biomes and then explore these using samples currently available from researchers across the globe. We will analyze 200,000 samples from these communities using metagenomics, metatranscriptomics and amplicon sequencing to produce a global Gene Atlas describing protein space, environmental metabolic models for each biome,

Meetings

There are currently no EMP centric meetings planned, however we will update this space as soon as the next meeting is organized.

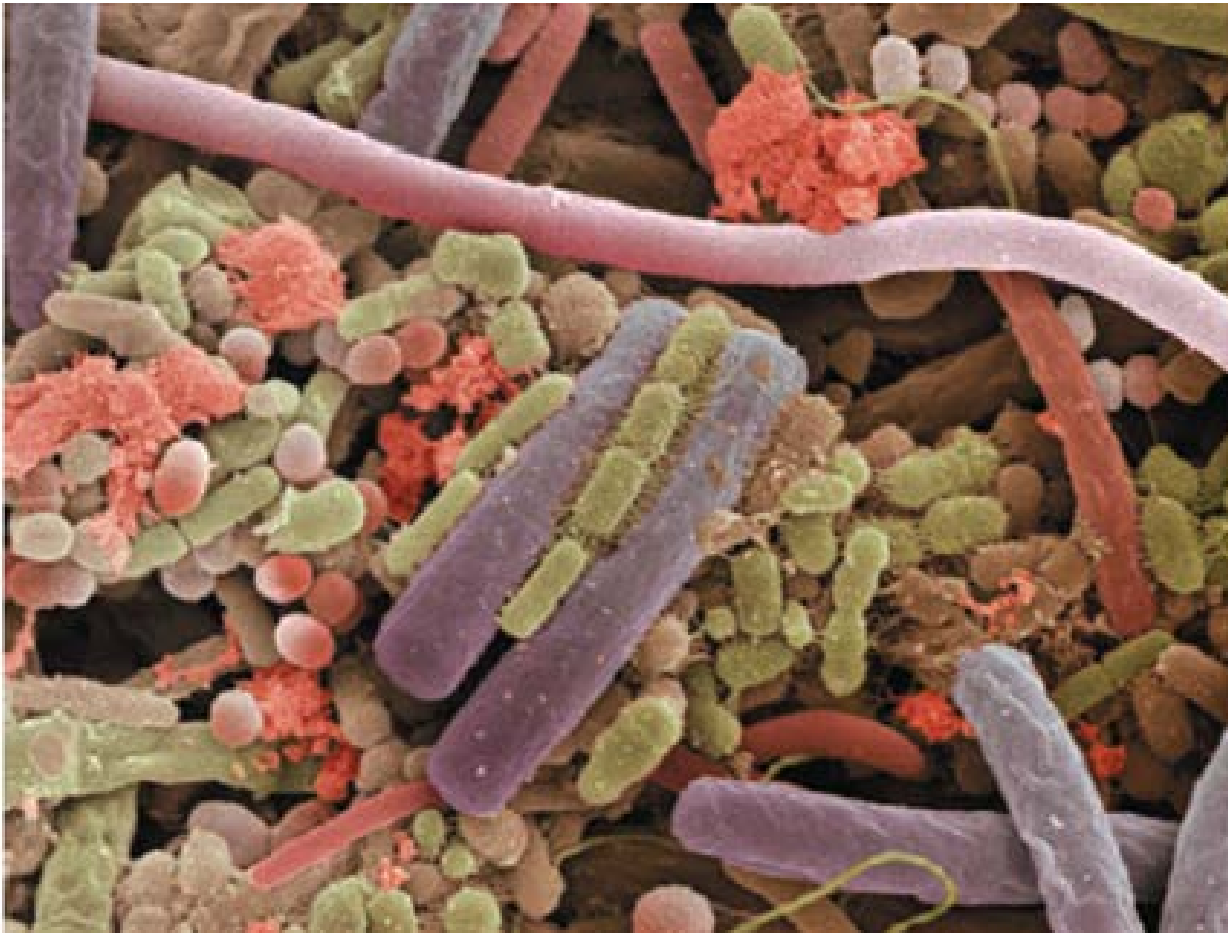
News

Earth Microbiome Project:
Rick Stevens at
TEDxMaperville

Há uma certa confusão

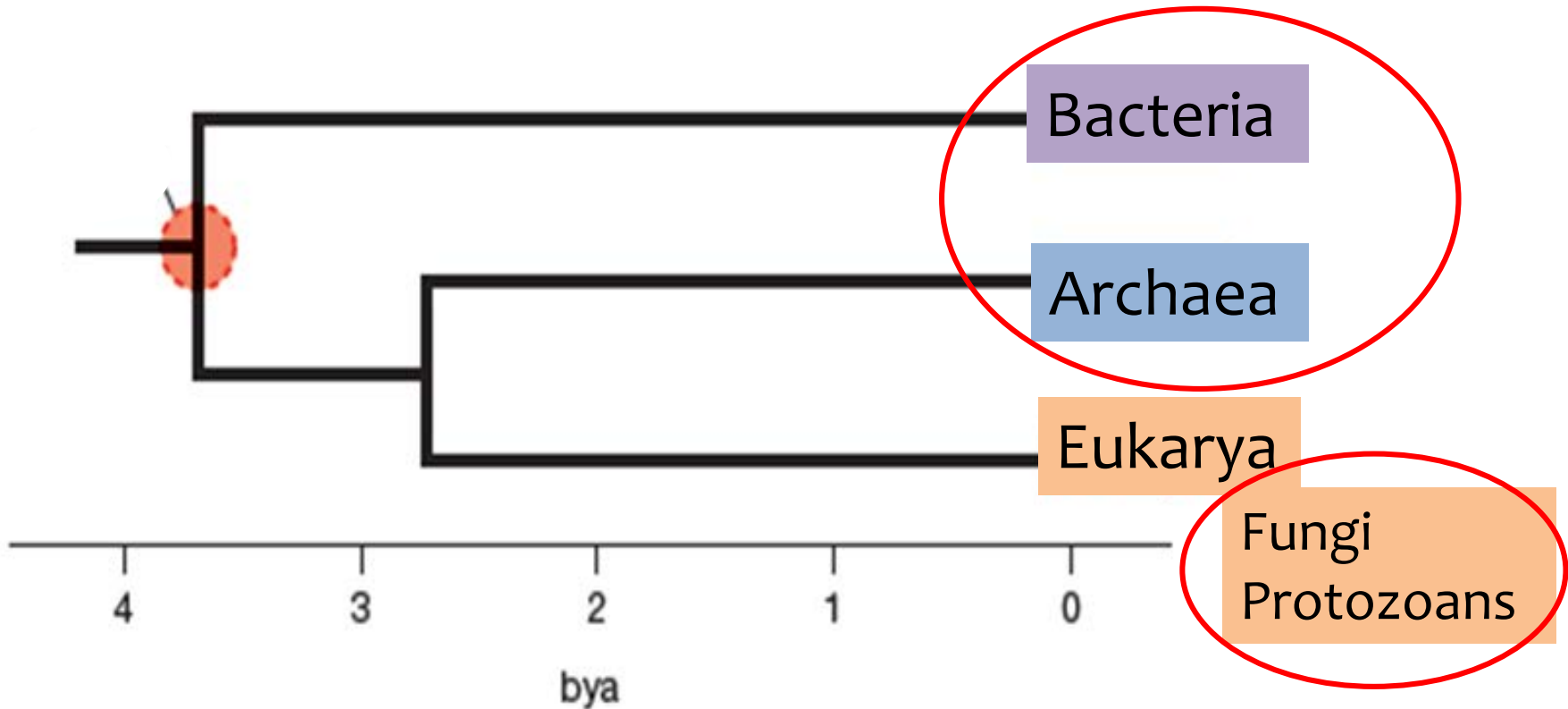
- **Earth Biogenome Project (EBP)**
- Projeto lançado em 2017 que pretende sequenciar “**all life on Earth**”
 - voltado para eucariotos

Comunidades microbianas –**Microbiotas**– são típicas de cada ambiente



ecossistema
microbiano

Microbiotas contêm variedade de microrganismos



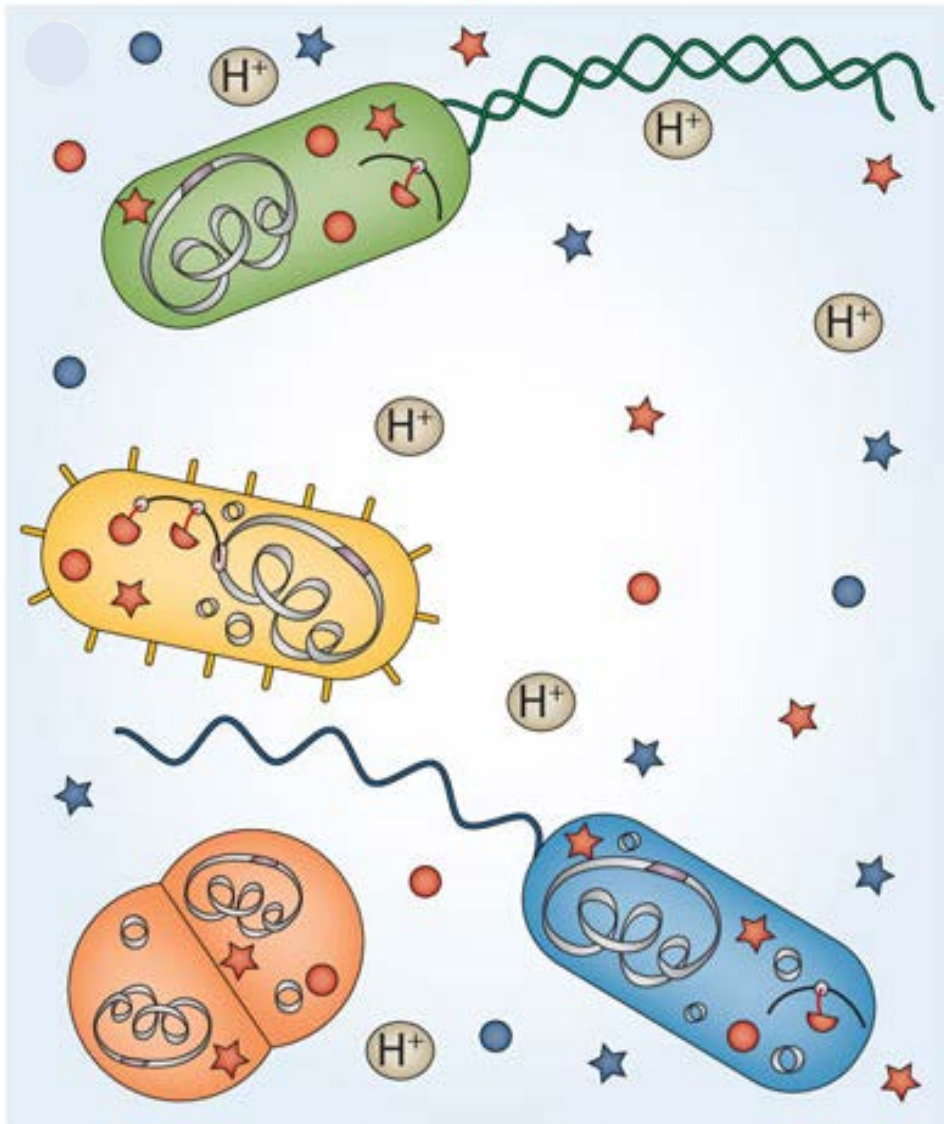
+ Vírus e Bacteriófagos

Microbioma

Genes, Genomas,
Proteínas e Metabólitos da
Microbiota

+

Proteínas e Metabólitos da resposta
do Hospedeiro à interação com a
microbiota



- ★ Metabólito da microbiota
- ★ Metabólitos do hospedeiro
- Proteína da microbiota
- Proteínas do hospedeiro

Como acessar essa extraordinária riqueza microbiológica?



Abordagens dependentes de cultivo

Cultivo de bactérias em meio sólido

Porém...

A **fração cultivável** da vasta riqueza microbiana da biosfera é muito pequena (estimada em 1%)

Como acessar a extraordinária **maioria invisível**?

→ Abordagens **independentes do cultivo**

MetaGenômica

revela as **espécies**, os **genes** e **genomas** de comunidades microbianas

MetaTranscritômica

revela os **genes expressos** (microbiota ativa)

MetaProteômica

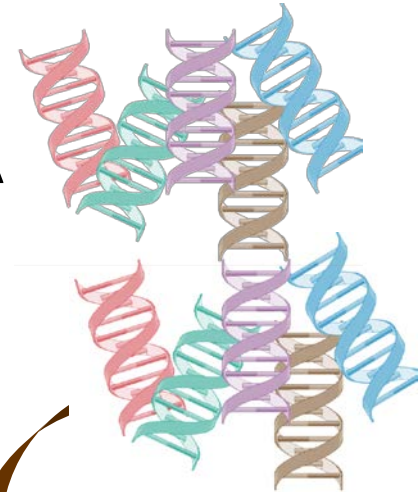
revela as **proteínas expressas** (microbiota ativa)

MetaGenômica e MetaTranscritômica

Amostra ambiental



Extrair o DNA
(ou RNA)



Sequenciar



Analisar as sequências de
DNA: metagenômica
cDNA: metatranscritômica



Sequenciamento de DNA
alto-desempenho

Tecnologias de sequenciamento

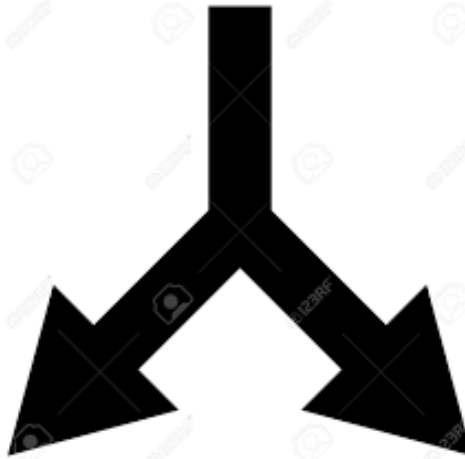
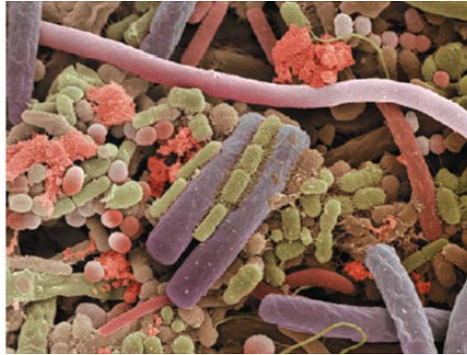
- NGS – next generation sequencing
 - Illumina
 - 90% do mercado
 - Em metagenômica talvez seja perto de 100%
 - PacBio
 - Long reads
 - Nanopore
 - Long reads



Big Data

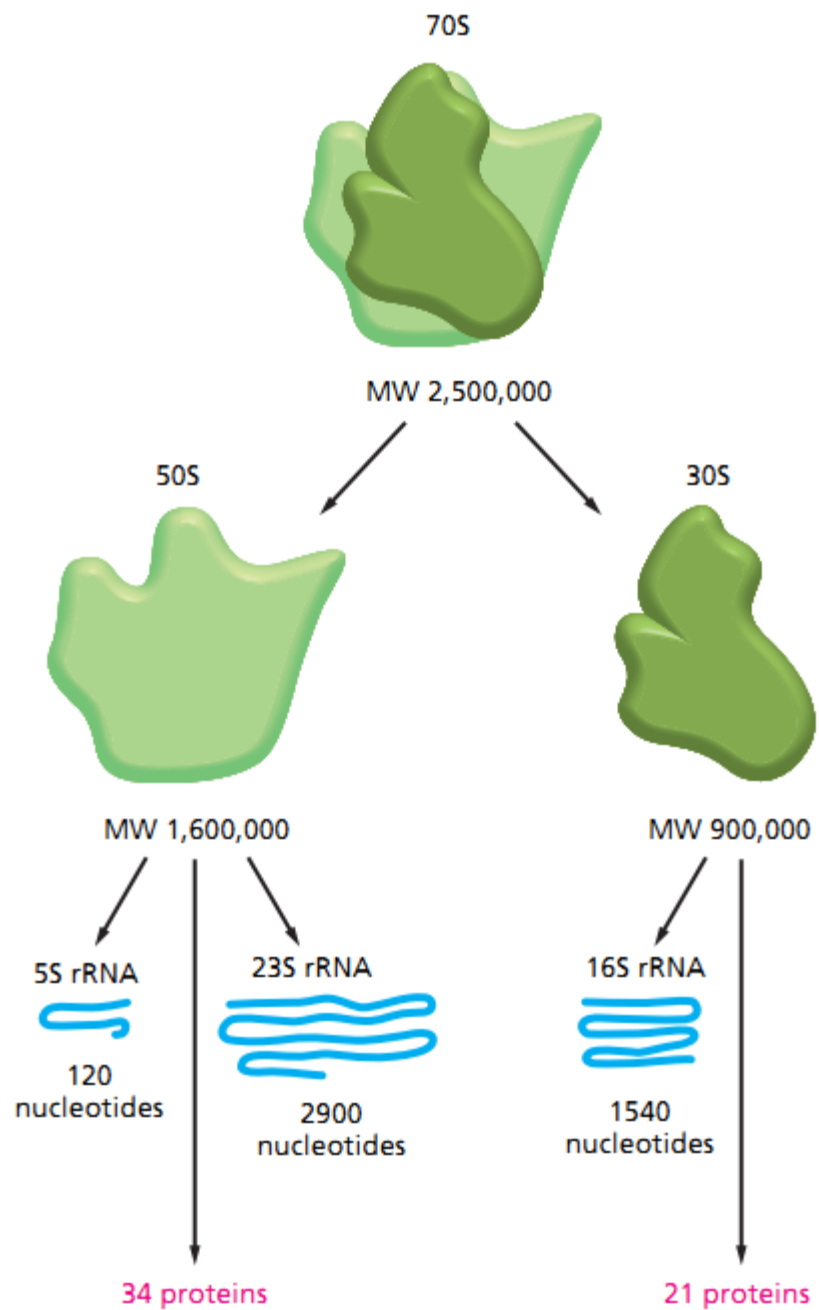
- Milhões de reads
- Que significa isto?
- Supondo
 - cada read com 300 bp
 - 10 milhões de reads para **uma amostra**
 - $10 \times 10^6 \times 300 = 3 \times 10^9$ bp
 - Um genoma bacteriano: 5×10^6 bp
 - Equivalente a **600 genomas bacterianos**
- **A bioinformática é essencial**

Metagenômica: tipos de Dados

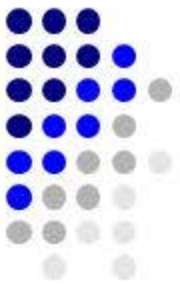


16S / 18S

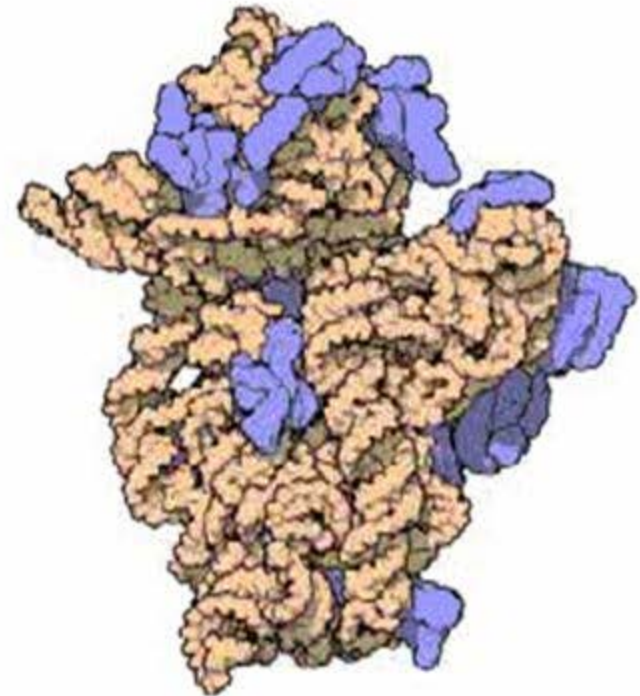
shotgun

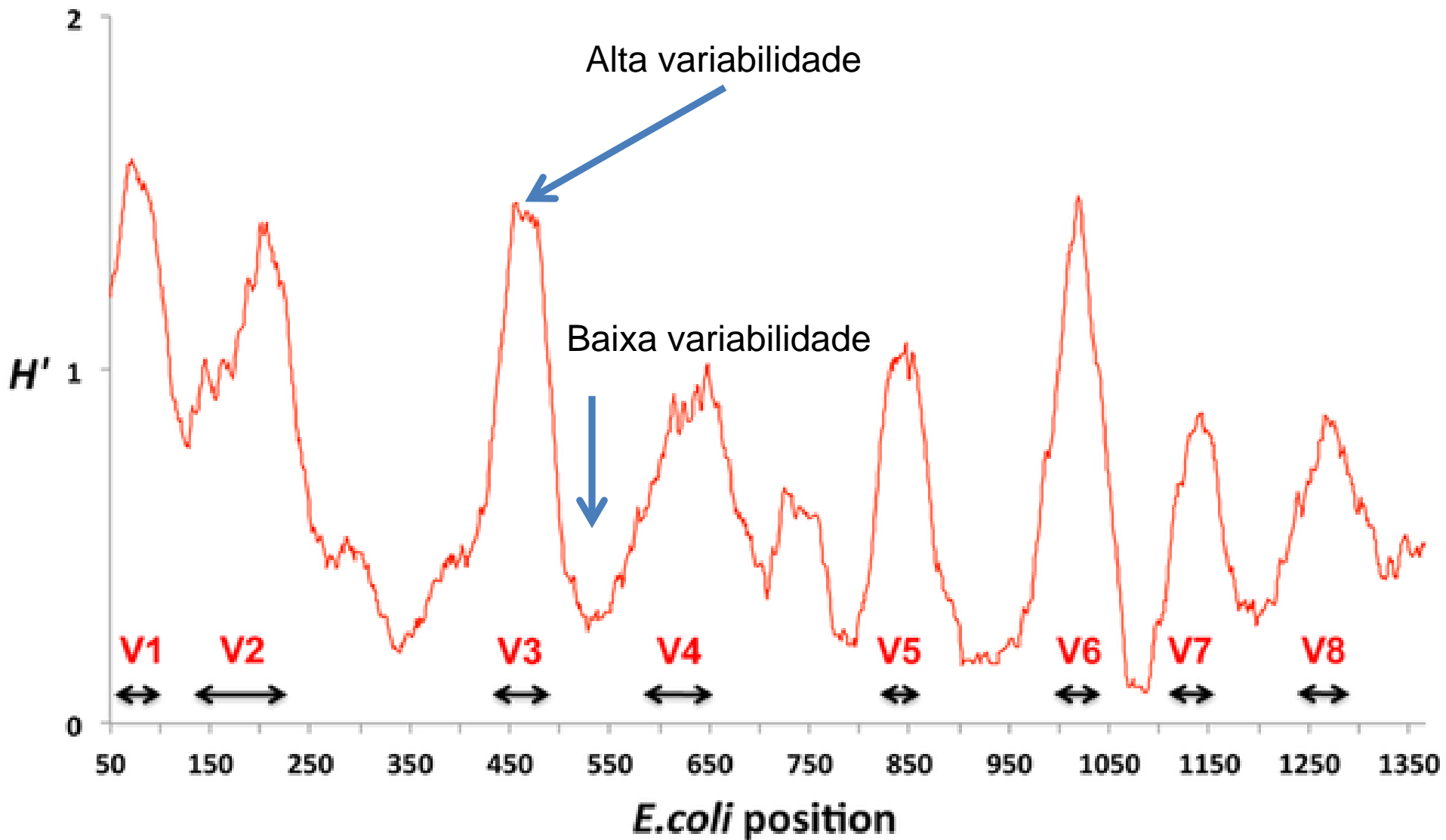


16S rRNA



- 16S
 - Ribosomal RNA
 - Large RNA component of the small subunit of the ribosome
 - Phylogenetic Markers
 - Species Identification
 - 1542 bp





⇒ Primers “universais”

DNA shotgun

- Sequenciar o **DNA total** da amostra
- Resultado
 - Milhões de fragmentos
 - Mistura dos DNAs dos diversos organismos presentes

16S vs. shotgun: objetivos

- 16S
 - Composição e estrutura da microbiota
 - “perfil taxonômico”
- Shotgun
 - Resultados mais detalhados
 - Perfil taxonômico
 - Funções
 - genomas

16S e shotgun: positivos e negativos

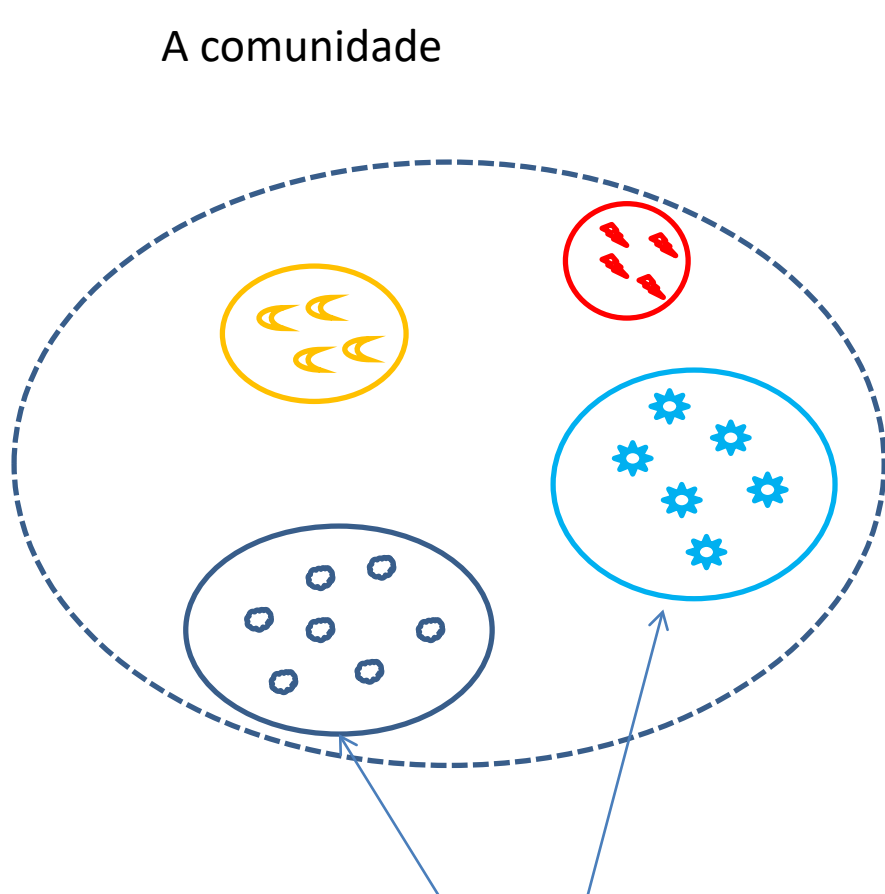
| | 16S | shotgun |
|--------------------------|--|--|
| custo | Mais baixo | Mais alto |
| Vieses (biases) | Menor chance de ser representativo | Maior chance de “pegar tudo” |
| Bancos de dados | Maior cobertura | Menor cobertura |
| Identificação taxonômica | Menos precisa (em geral, não mais do que gênero) | Mais precisa, podendo chegar a espécie, e talvez cepas |

Que perguntas queremos fazer?

Quem está na amostra?

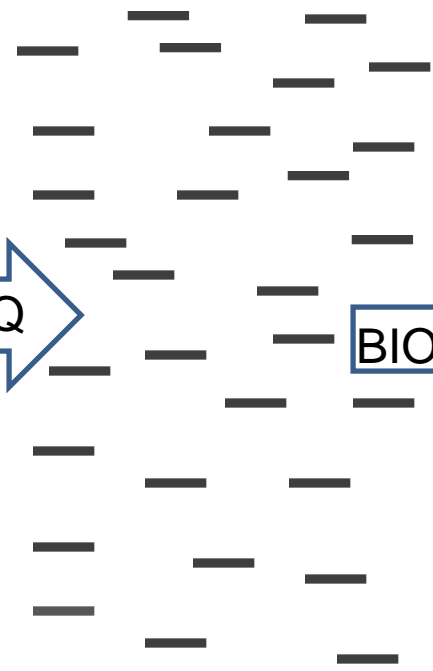
- Identificação taxonômica (16S, shotgun)
- Recuperação de genomas (shotgun)

A comunidade

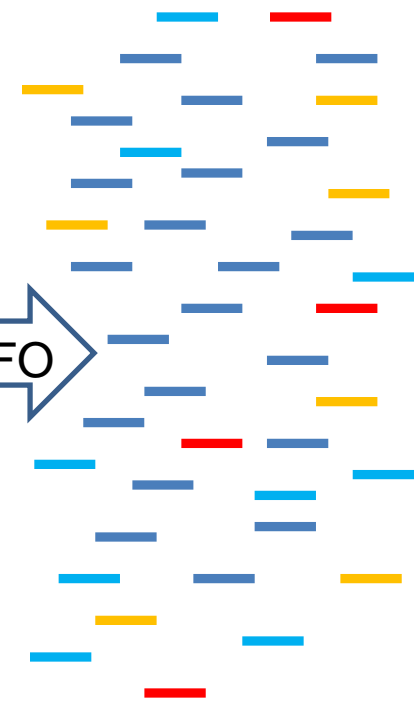


populações

SEQ

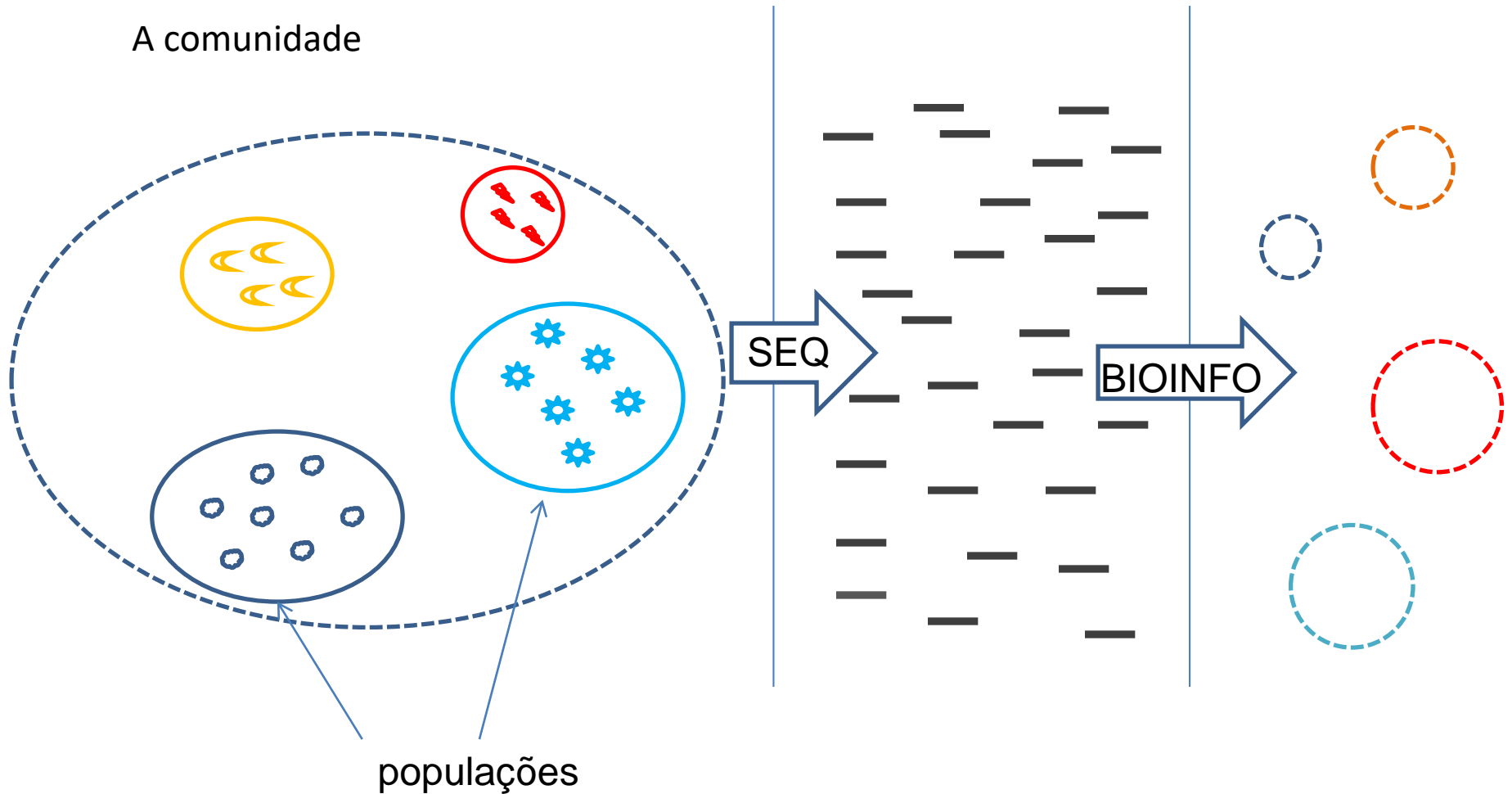


BIOINFO



Recuperação de genomas

A comunidade



Identificação taxonômica depende
de bancos de dados

Bancos de dados de 16S



GREENGENES
The 16S rRNA Gene Database and Tools

The Greengenes Database

While we are setting up our site, please visit the [download](#) area to obtain files.



The Greengenes Database by The Greengenes Database Consortium is licensed under a [Creative Commons Attribution-ShareAlike 3.0 Unported License](#).

SILVA

Welcome to the SILVA rRNA database project

A comprehensive on-line resource for quality checked and aligned ribosomal RNA sequence data.

SILVA provides comprehensive, quality checked and regularly updated datasets of aligned small (16S/18S, SSU) and large subunit (23S/28S, LSU) ribosomal RNA (rRNA) sequences for all three domains of life (*Bacteria*, *Archaea* and *Eukarya*).

SILVA are the official databases of the software package ARB.

For more background information → [Click here](#)

SILVAngs



Check out our new service for Next Generation Amplicon data

SILVA Tree Viewer

The SILVA Tree Viewer is a web application to browse and query the SILVA guide trees.

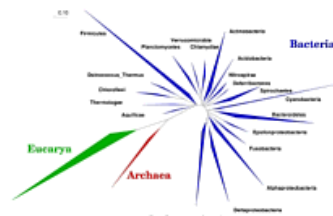
A technical preview is available at www.arb-silva.de/treeviewer



ARB

The software package ARB represents a graphically-oriented, fully-integrated package of cooperating software tools for handling and analysis of sequence information.

The ARB project has been started more than 15 years ago by Wolfgang Ludwig at the Technical University in Munich, Germany, see www.arb-home.de.



News

23.11.2017

The 10th de.NBI Quaterly Newsletter published



Good news for de.NBI, the German Network for Bioinformatics Infrastructure: In September, de.NBI has passed successfully the midterm evaluation in Berlin. The international evaluation panel stated that de.NBI is working successfully from the beginning and that it should be continued.

09.11.2017

Call for Action - We need your Help!



The UniEuk project needs your help to launch EukBank 1.0.

28.10.2017

de.NBI Handbook ready for Download



The Handbook is the first comprehensive document that lists the work and effort of all de.NBI partners. Content: How de.NBI is structured, Presentations of all Partners, Index of Persons/Contact Details.

05.10.2017

SILVA TreeViewer published



SILVA TreeViewer: interactive web browsing of the SILVA phylogenetic guide trees now published in BMC Bioinformatics.

[go to Archive ->](#)

User satisfaction survey

SILVA is now part of the German Network for Bioinformatics Infrastructure de.NBI.



To evaluate and improve our quality of service we need your feedback. Please help us by participating in this short [survey](#).

SILVA SSU / LSU 128 - full release

SSU Parc SSU Ref SSU Ref NR 99 LSU Parc LSU Ref



ANNOUNCEMENTS

RDP News

11/10/2017 [myRDP login problem fixed!](#)

11/09/2017 [Apologize for the problem with myRDP login.](#)
Our team is working to fix it as soon as possible.

05/16/2017 [Apology for slow/NO connection to RDP tools today](#)
Thanks to Alex/Brian, etc. for working things out in the server room

05/10/2017 [RDP Director at GSC 19, May 14-17](#)
Genomic Standards Consortium Meeting, Brisbane, Queensland, Australia

05/10/2017 [Possible Friday, May 12, morning interruptions](#)
Emergency Generator Testing from 9-10 A.M.

12/13/2016 [Most Highly Cited Researchers](#)
Congratulations to RDP Director James Cole

09/30/2016 [RDP Release 11.5 available](#)
Updated 16S rRNA training set to training set No. 16.

08/16/2016 [Possible Friday morning interruptions](#)
Building electrical testing/maintenance

06/30/2016 [RDP Classifier Updates](#)
The Classifier 16S training set and Fungal ITS Warcup set have been updated

06/03/2016 [RDP staff on the road!](#)
Teaching in China, Genomic Standards Consortium meeting in Crete, special ASM Microbe events in Boston



RDP Release 11, Update 5 :: September 30, 2016

3,356,809 16S rRNAs :: 125,525 Fungal 28S rRNAs
Find out what's new in RDP Release 11.5 [here](#).

Cite RDP's latest tool articles.

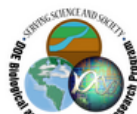
RDP provides quality-controlled, aligned and annotated Bacterial and Archaeal 16S rRNA sequences, and Fungal 28S rRNA sequences, and a suite of analysis tools to the scientific community. New to RDP release 11:

- RDP tools have been updated to work with the new fungal 28S rRNA sequence collection.
- A new Fungal 28S Aligner and updated Bacterial and Archaeal 16S Aligner. We optimized the parameters for these secondary-structure based Infernal aligners to provide improved handling for partial sequences.
- Updated RDPipeline offers extended processing and analysis tools to process high-throughput sequencing data, including single-strand and paired-end reads.
- Most of the RDP tools are now available as open source packages for users to incorporate in their local workflow.



RDP's mission and funding:

Part of RDP's mission is to provide support to our users. Email and phone contacts are available on the [contacts page](#).

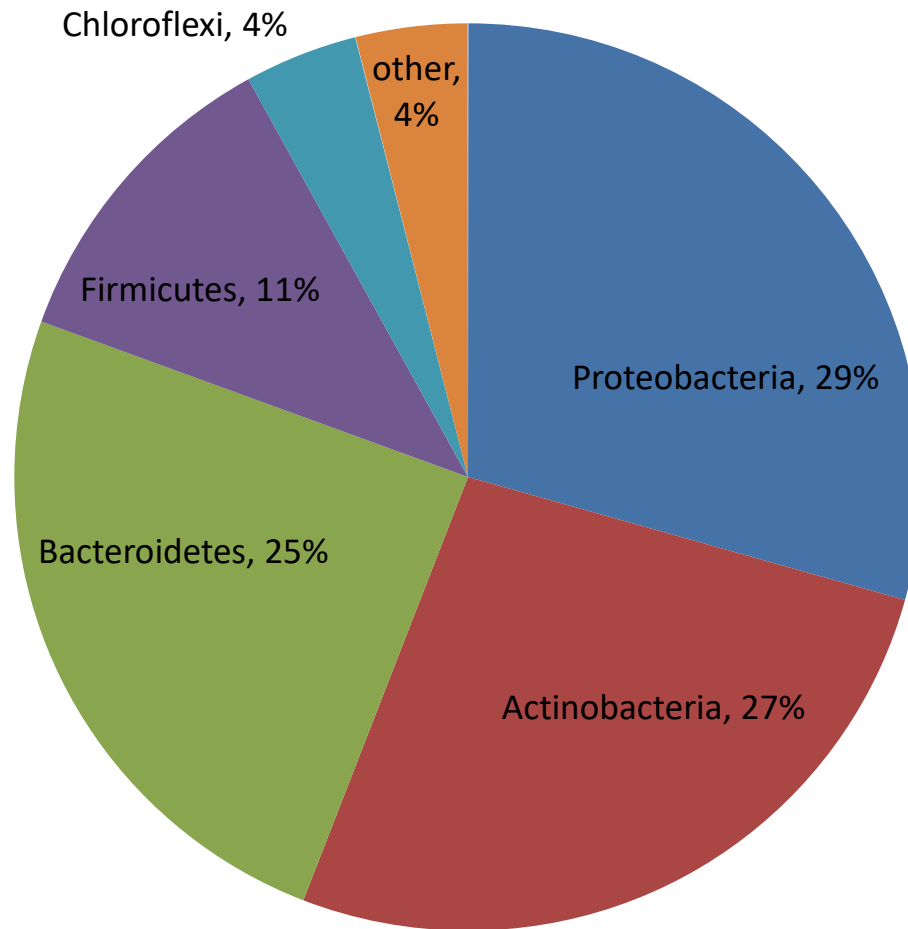


Bancos de datos para DNA total

- GenBank
 - nt
 - nr
 - env_nr
 - refSeq
 - WGS

Quais são as abundâncias relativas?

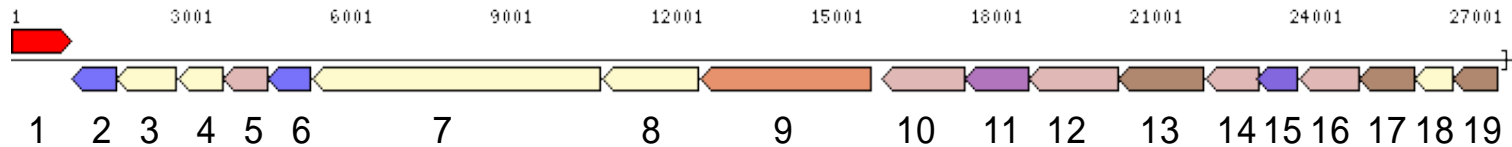
16S e shotgun



Quais funções estão presentes?

- Em genes (shotgun)
- Em genes expressos (metaTranscritômica)

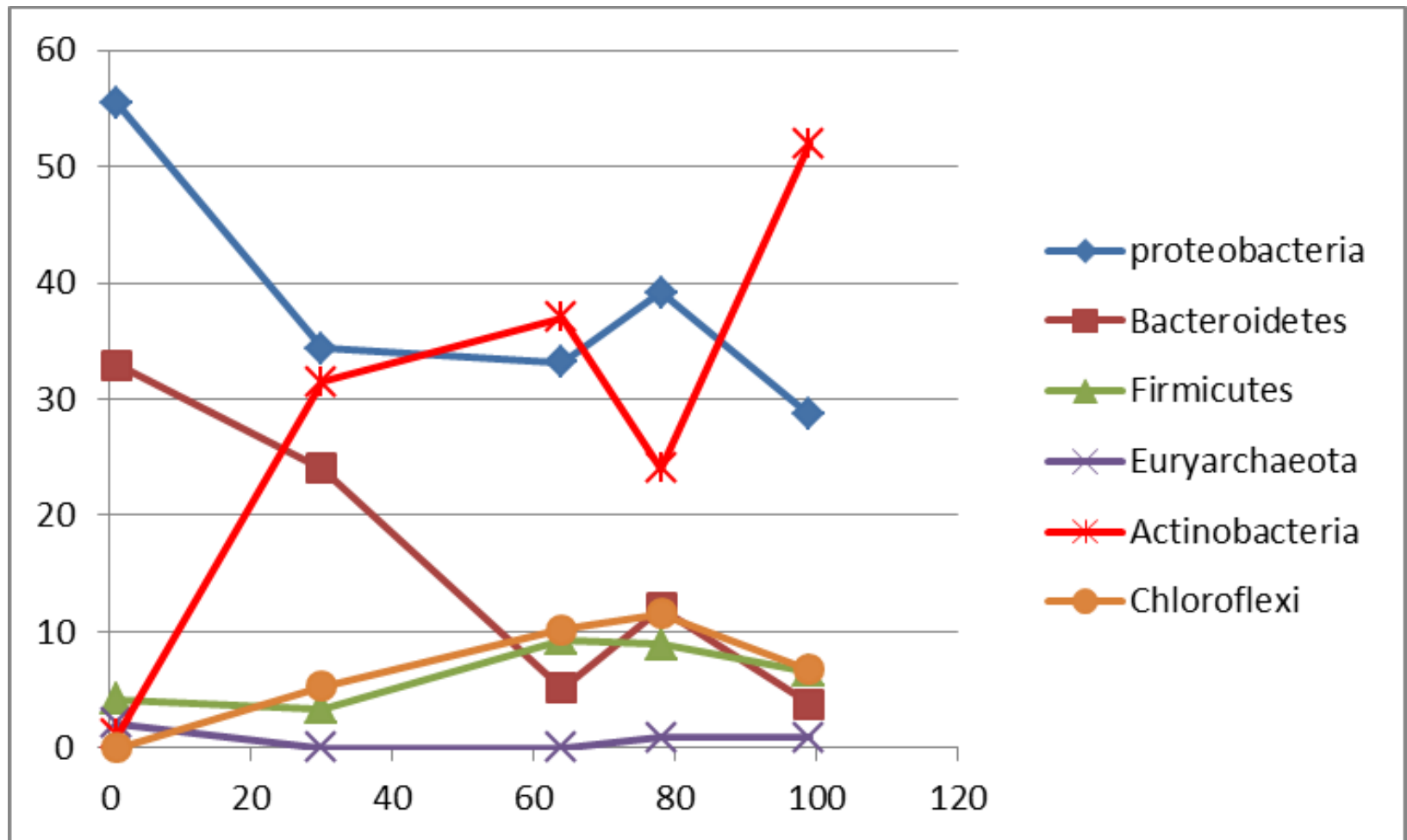
ZC1 contig00009.9 (27,919 bp)



1. Beta-xylosidase (376aa, COG3507)
2. Dehydrogenases (280aa, COG1028)
3. hypothetical protein (379aa);
4. hypothetical protein (283aa)
5. 5-keto 4-deoxyuronate isomerase (280aa, COG3717)
6. Dehydrogenases (267aa, COG1028)
7. hypothetical protein (1799aa)
8. SusD family protein (606aa, pfam07980)
9. TonB-linked outer membrane protein (1068aa, COG4771);
- 10. Pectate lyase (518aa, COG3866)**
11. Predicted unsaturated glucuronyl hydrolase
- 12. Pectin methylesterase (568aa, COG4677)**
- 13. Endopolygalacturonase (523aa, COG5434)**
14. Nucleoside-diphosphate-sugar epimerase (326aa, COG0451)
15. Nucleoside-diphosphate-sugar pyrophosphorylase (249aa, pfam00483)
16. Galactokinase (377aa, COG0153)
17. Soluble lytic murein transglycosylase (347aa, COG0741)
18. hypothetical protein (235aa)
19. Predicted UDP-glucose 6-dehydrogenase (283aa, COG1004).

Metagenômica comparativa

Mesmo local, variação no tempo

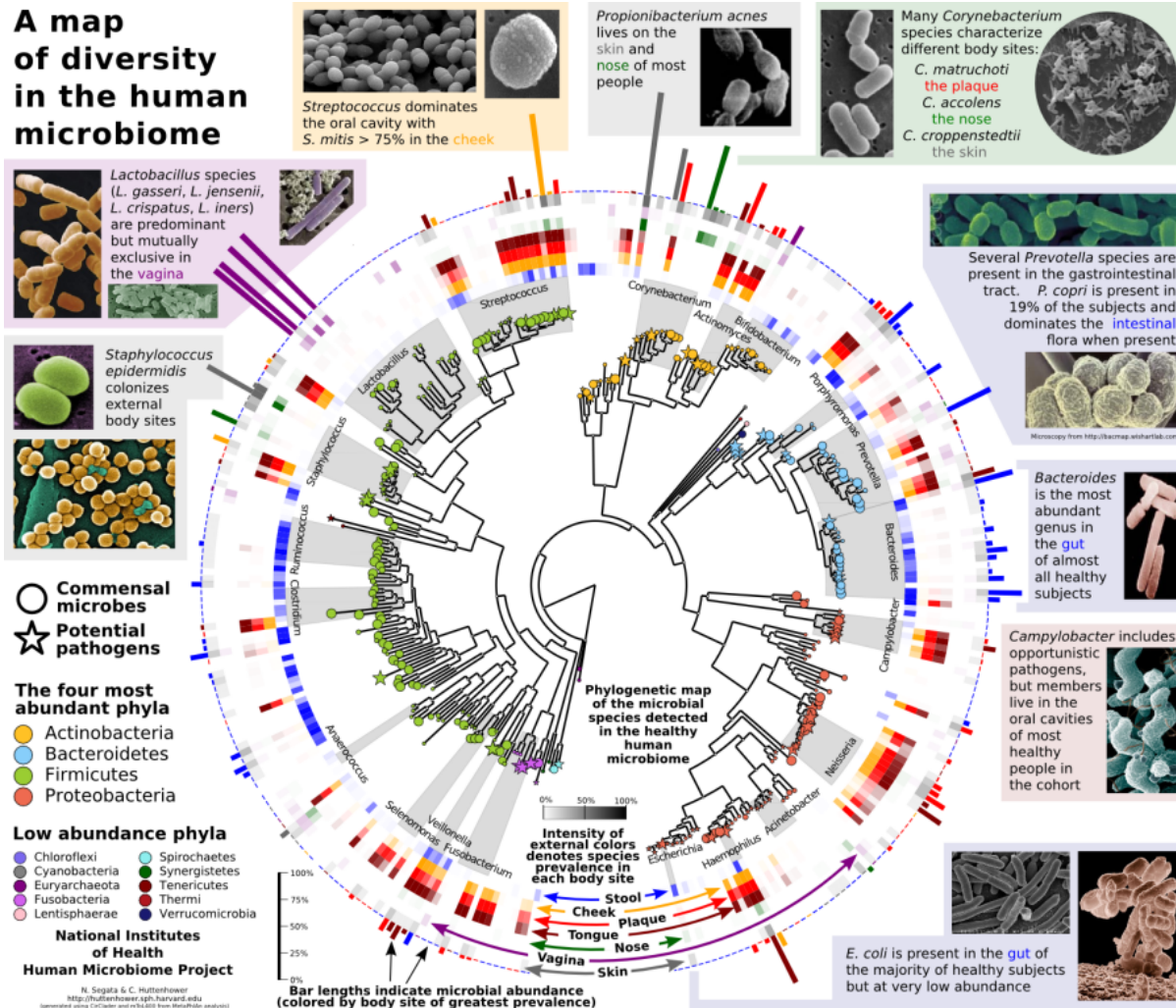


Mesmo local, variação de indivíduos

- Amostras da boca
 - Indivíduos que fumam
 - Indivíduos que não fumam

Diferentes locais, mesmo indivíduo

A map of diversity in the human microbiome



Taxonomia

- *Xanthomonas citri*
- Filo: **proteobacteria**
 - Classe: **proteobacteria gama**
 - Ordem: **xanthomonadales**
 - Família: **xanthomonadacea**
 - » Gênero: **xanthomonas**
 - Espécie: **citri**

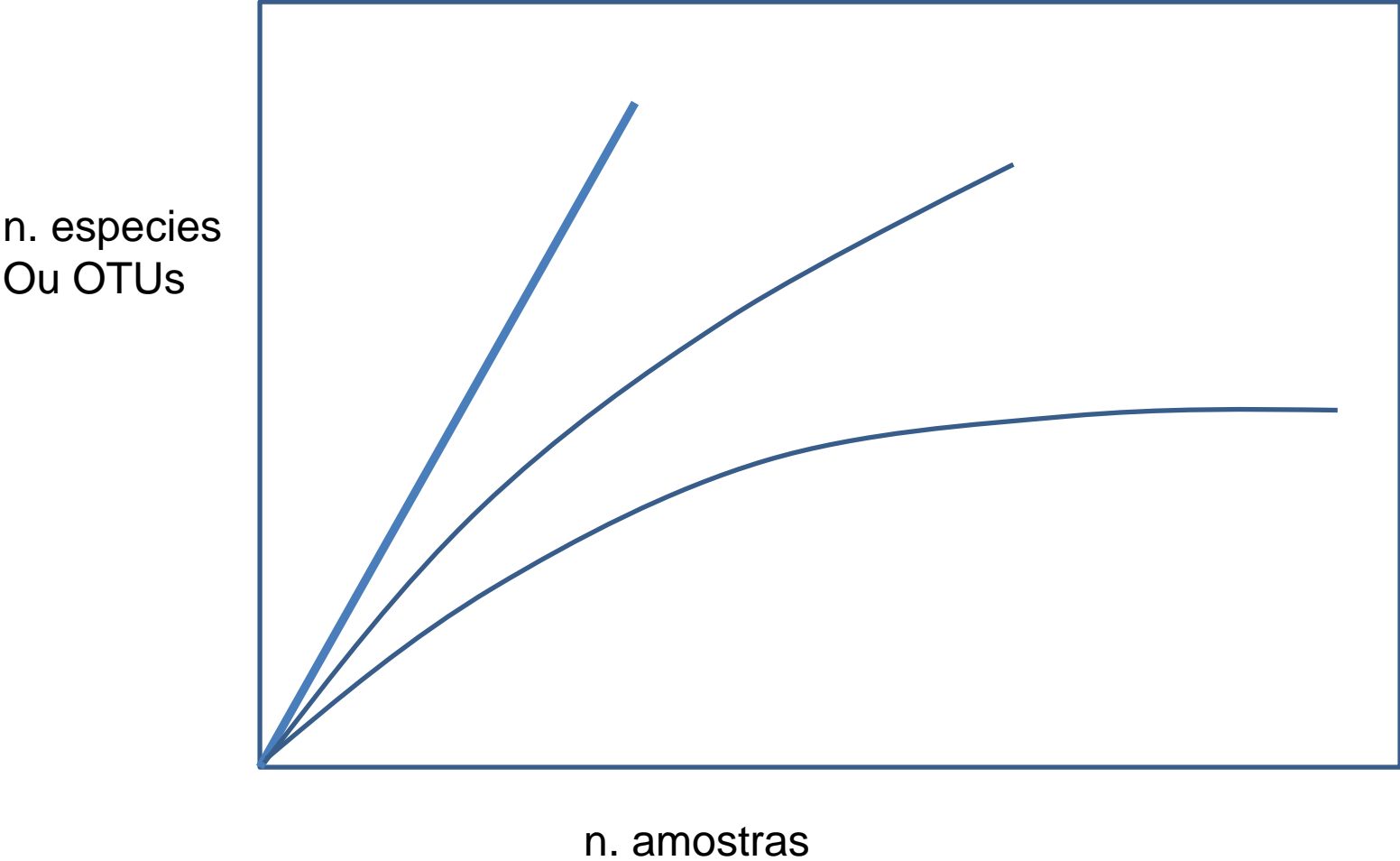
OTU

- Unidade taxonômica operacional
- Se for **conhecida**, leva um rótulo padronizado
 - *Xanthomonas citri*
- Mas pode ser **desconhecida**
 - Nesse caso, recebe um número, que varia de análise para análise

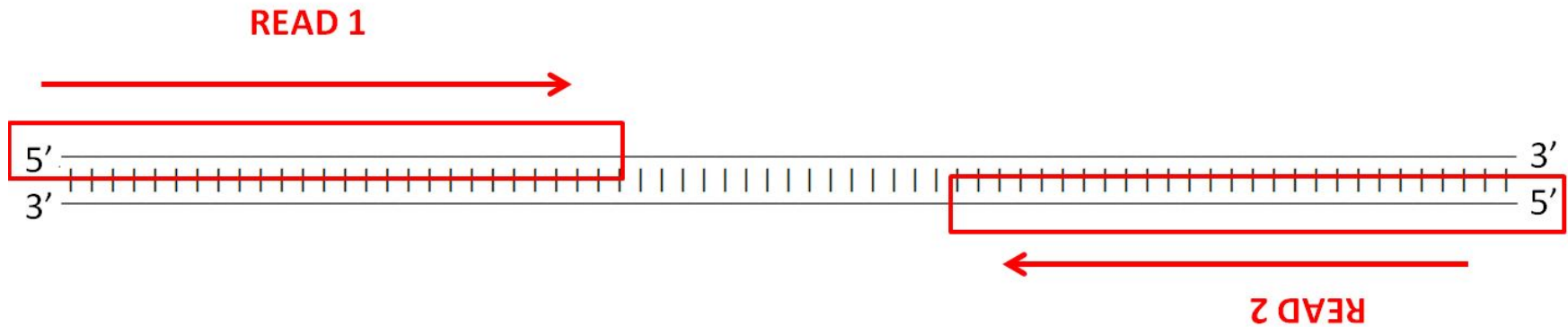
A amostra é representativa?

- Curvas de rarefação

Curvas de rarefação (ou saturamento)



Single-end and Paired-end reads



Muitas fontes de erro

- Amostragem
- Preparação da biblioteca
- Sequenciamento
- Tamanho da sequência (pode ser curta demais)
- Programas
- Viéses dos bancos de dados

Classificação de reads de DNA total

- **Similaridade** com sequências de origem conhecida
 - BLAST
- Propriedades intrínsecas de cada sequência
 - **Assinaturas genômicas**
 - Adequado para binning

Classificação com base na frequência de palavras de k bases

$k = 4$: AAAA, AAAC, AAAG, AAAT, CAAA, etc...

Dada uma janela de x kb, podemos contar as ocorrências de cada uma dessas palavras dentro da janela

Exemplo:

AG**ATTA**GCGACT**ATT**ATAGCCTAGATCGATC**ATTA**CC

AGAT ocorre 2 vezes

ATTA ocorre 3 vezes

etc

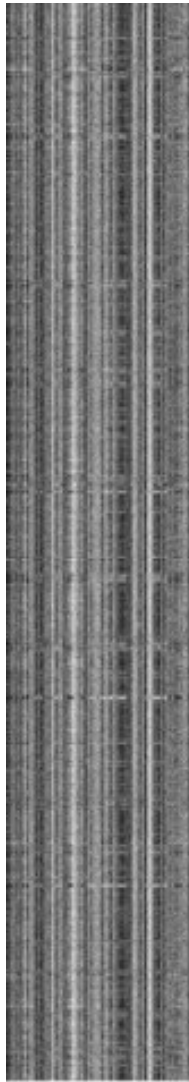
Palavras de k bases: k -mers (kâmeros)

Matriz de frequências

| janela | AAAA | AAAC | AAAG | AAAT | ACAA | ACAC | ACAG | ACAT |
|--------|------|------|------|------|------|------|------|------|
| 1 | 15 | 2 | | | | | | |
| 2 | 16 | 3 | | | | | | |
| 3 | 14 | 0 | | | | | | |
| 4 | 13 | 2 | | | | | | |
| 5 | 15 | 4 | | | | | | |
| 6 | 12 | 0 | | | | | | |
| 7 | 18 | 1 | | | | | | |
| 8 | 17 | 3 | | | | | | |
| 9 | 16 | 1 | | | | | | |

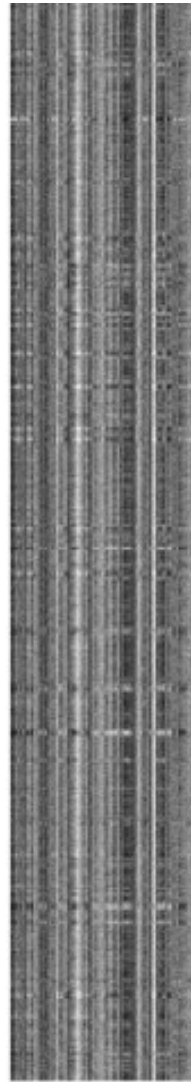
Genome "barcodes"

Burkholderia pseudomallei



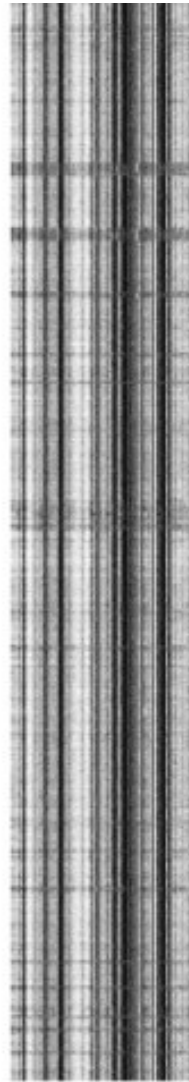
(a)

E. coli K12



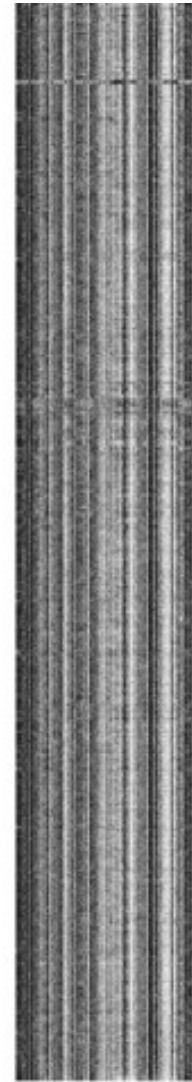
(b)

E. coli O157



(c)

Pyrococcus furiosus



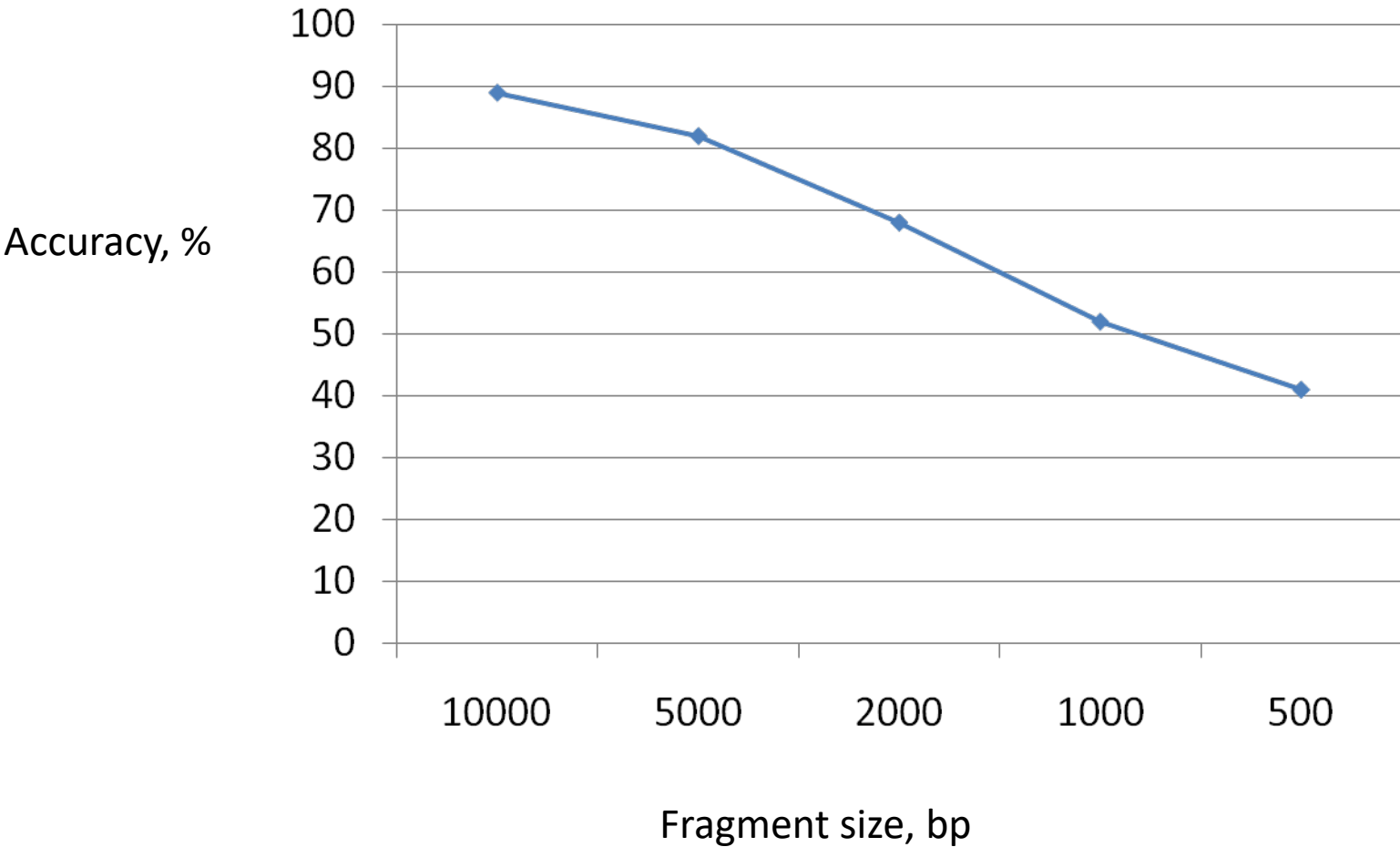
(d)



(e)

random

Não funciona bem com fragmentos curtos

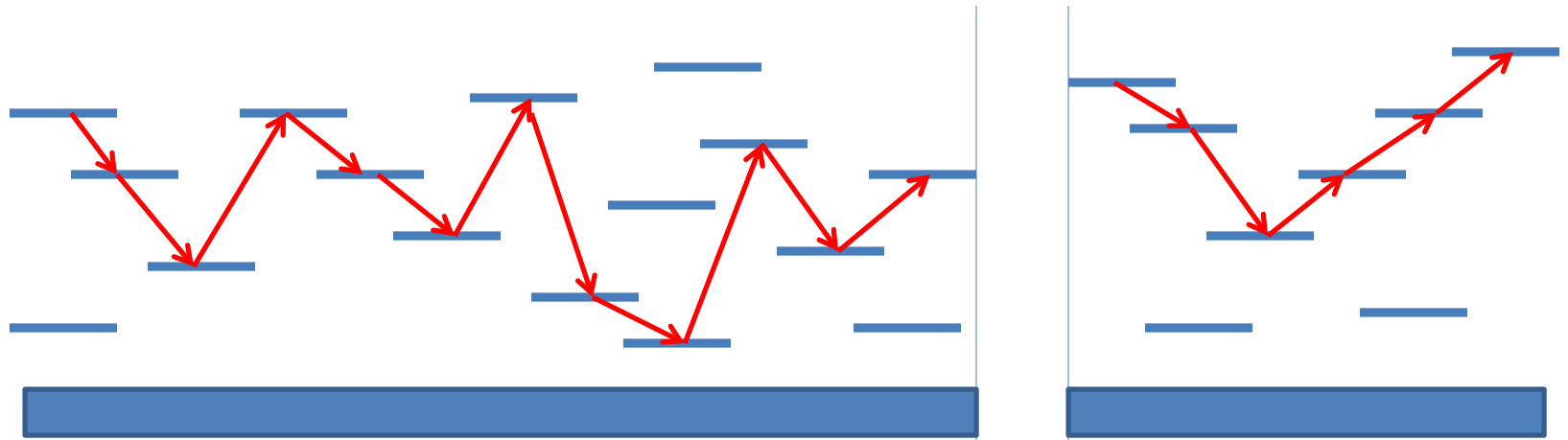


Zhou et al, 2009 simulated data

Exercício

- $S_1 = \text{TTCTACTACT}$
- $S_2 = \text{TTGTACTAGG}$
- $S_3 = \text{ACTTCTACTA}$
- Contar palavras de tamanho 2

Montagem de genomas



contig

```
...ACCGTAAATGGGCTGATCATGCTTAAA  
TGATCATGCTTAAACCCTGTGCATCCTACTG...
```

buraco

Montagem

- Em genomas bacterianos isolados, é um processo **razoavelmente bem compreendido**
- Em metagenomas há velhas e novas dificuldades
 - Mistura de organismos
 - Quimeras
 - Transferência lateral
 - Repetições
 - Tamanho dos conjuntos de dados
 - Chegando a **bilhões** de reads

Exemplo de quimerismo

genes

contig

g1

g2

g3

g4

g5



chlorobium

firmicutes

euryarch.

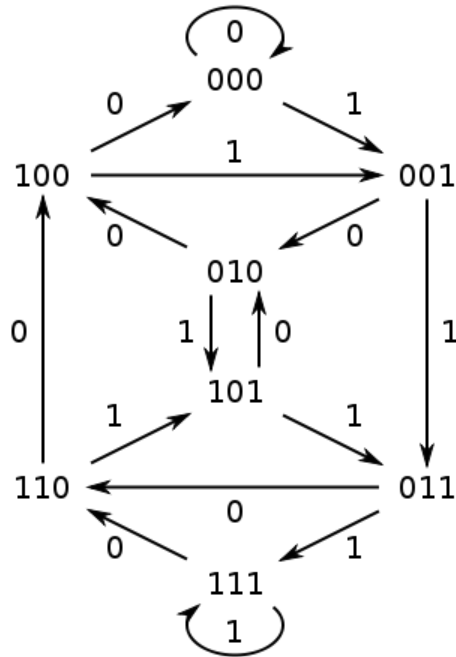
γ proteob.

crenarch.

Paradigmas de montagem

- OLC
 - overlap, layout, consensus
 - mais rigoroso, mas mais lento
- k-meros + grafos de de Bruijn
 - menos rigoroso, mas muito mais rápido
 - mais apropriado para metagenômica

grafos de de Bruijn



Sobreposição de k -mers

$k = 1$

Grafo de de Bruijn em montagem

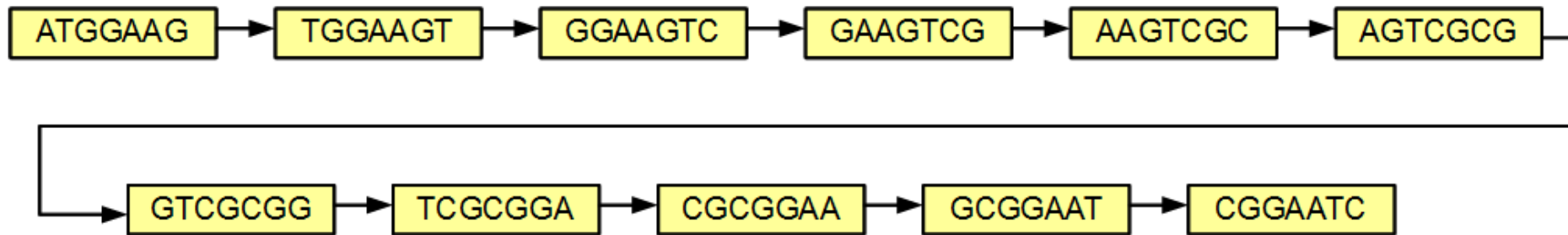
sequence

ATGGAAGTCGCGGAATC

7mers

ATGGAAG
TGGAAGT
GGAAGTC
GAAGTCG
AAGTCGC
AGTCGCG
GTCGCGG
TCGCGGA
CGCGGAA
GCGGAAT
CGGAATC

de Bruijn graph



Anotação funcional

- Pipeline para genomas completos pode ser usado
 - Exemplo: IMG/M
- Revejam aula sobre anotação de genomas

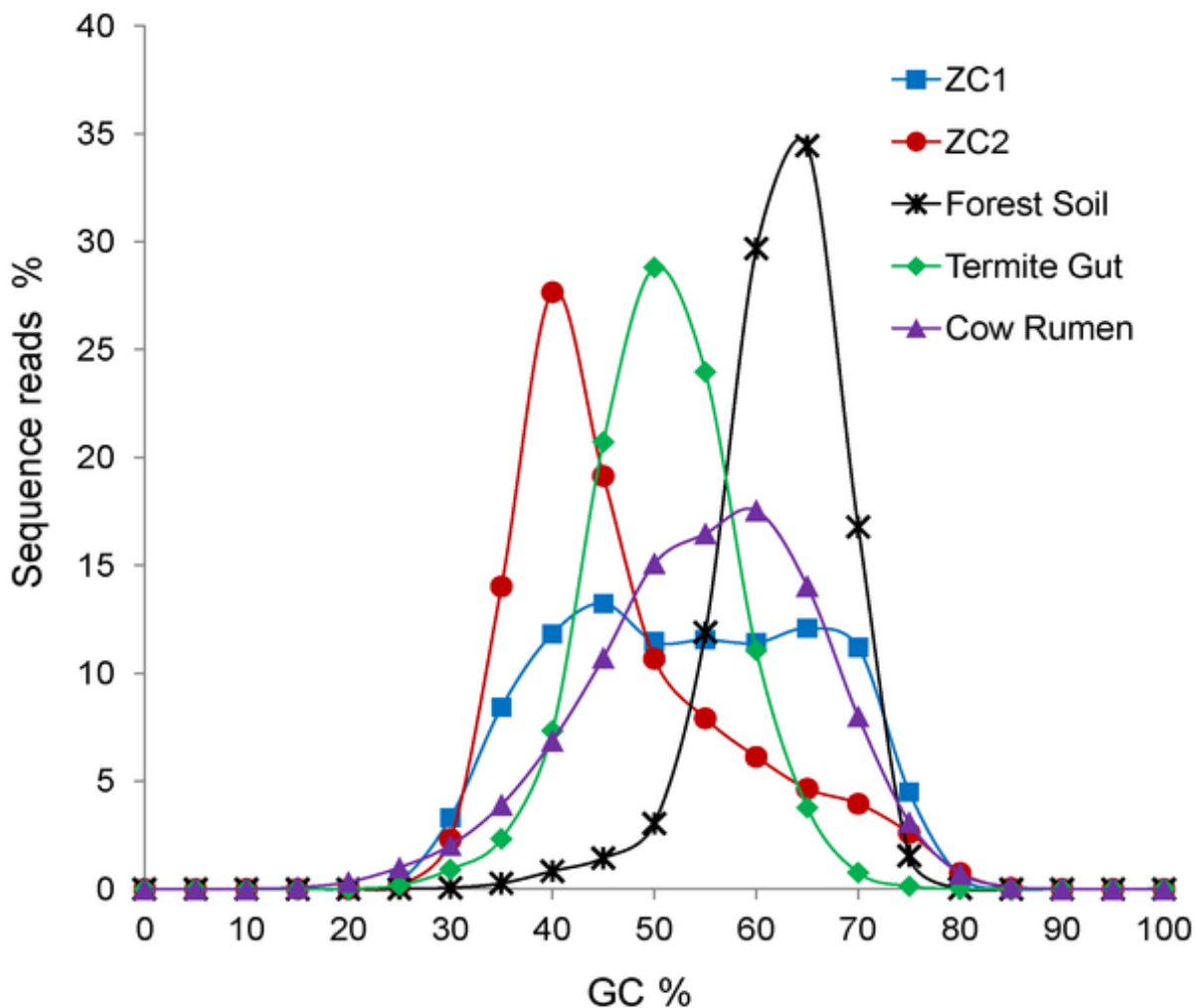
Cobertura

- Quanto cada genoma é coberto pelos reads obtidos
- Ambientes de grande riqueza: **cobertura baixa**
- Cobertura baixa cria **contigs pequenos**
 - maioria das ORFs são parciais
 - Dificulta atribuição de função
 - Potencial gerador de erros

Comparação de metagenomas

- Genomicamente
- Taxonomicamente
- Funcionalmente
- Recursos oferecidos pelo IMG/M

Figure 1. Distribution of the GC content percentage for ZC1 and ZC2 compared with selected metagenomes.



Martins LF, Antunes LP, Pascon RC, de Oliveira JCF, Digiampietri LA, et al. (2013) Metagenomic Analysis of a Tropical Composting Operation at the São Paulo Zoo Park Reveals Diversity of Biomass Degradation Functions and Organisms. PLoS ONE 8(4): e61928. doi:10.1371/journal.pone.0061928

<http://127.0.0.1:8081/plosone/article?id=info:doi/10.1371/journal.pone.0061928>

Genome clustering (IMG/M)

Clustering Type:

By Function:

- COG
- Pfam
- KO

By Taxonomy:

- Class
- Family
- Genus

By Function Category:

- COG Categories
- COG Pathways
- KEGG Pathway Categories (KO)
- KEGG Pathway Categories (EC)
- KEGG Pathways (KO)
- KEGG Pathways (EC)
- Pfam Categories

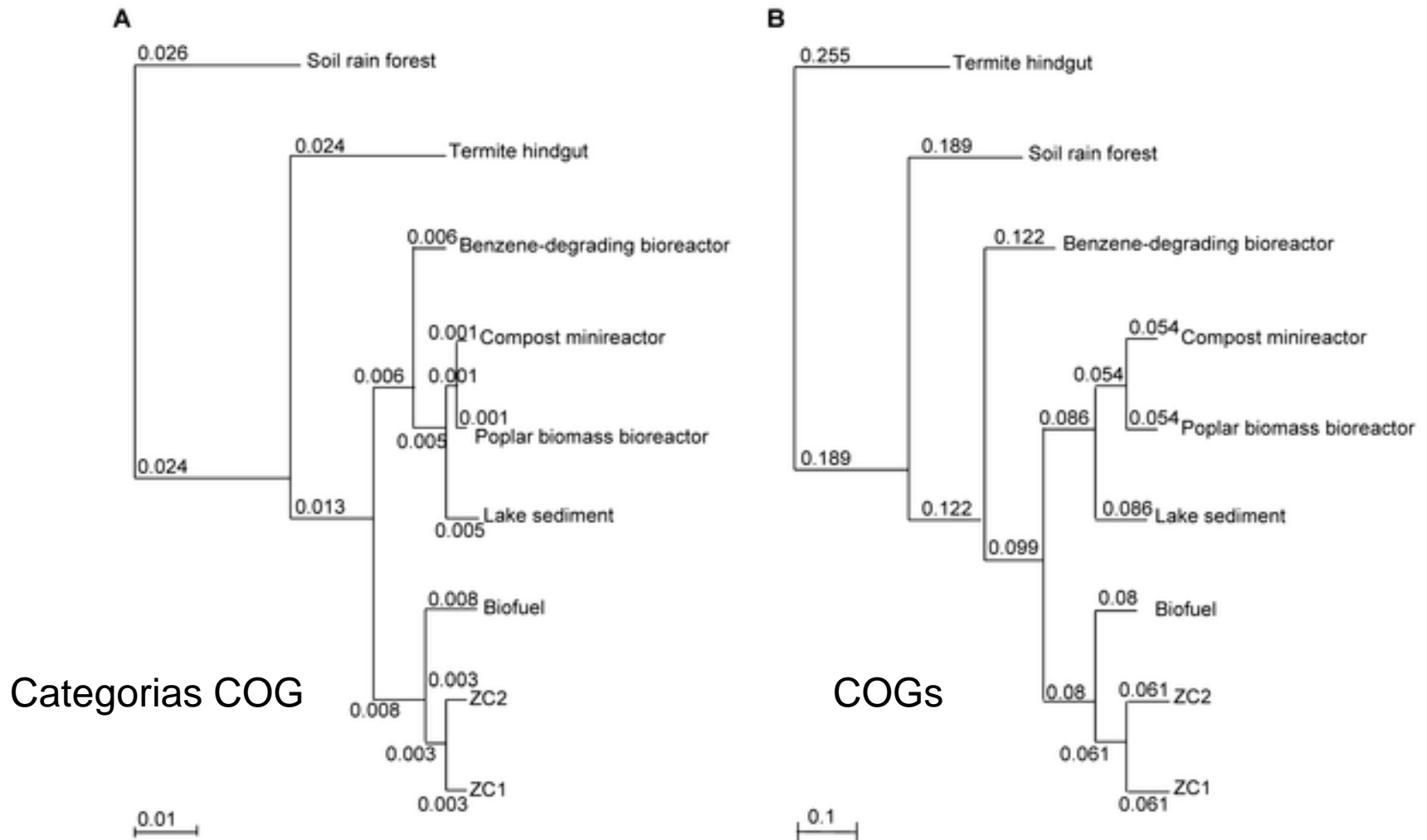
Clustering Method:

- Hierarchical Clustering
- Principal Components Analysis (PCA)
- Principal Coordinates Analysis (PCoA)
- Non-metric MultiDimensional Scaling (NMDS)
- Correlation Matrix

Go

Reset

Figure 8. Hierarchical clustering of functional gene groups of ZC1 and ZC2 and seven public metagenomes.



Martins LF, Antunes LP, Pascon RC, de Oliveira JCF, Digiampietri LA, et al. (2013) Metagenomic Analysis of a Tropical Composting Operation at the São Paulo Zoo Park Reveals Diversity of Biomass Degradation Functions and Organisms. PLoS ONE 8(4): e61928. doi:10.1371/journal.pone.0061928

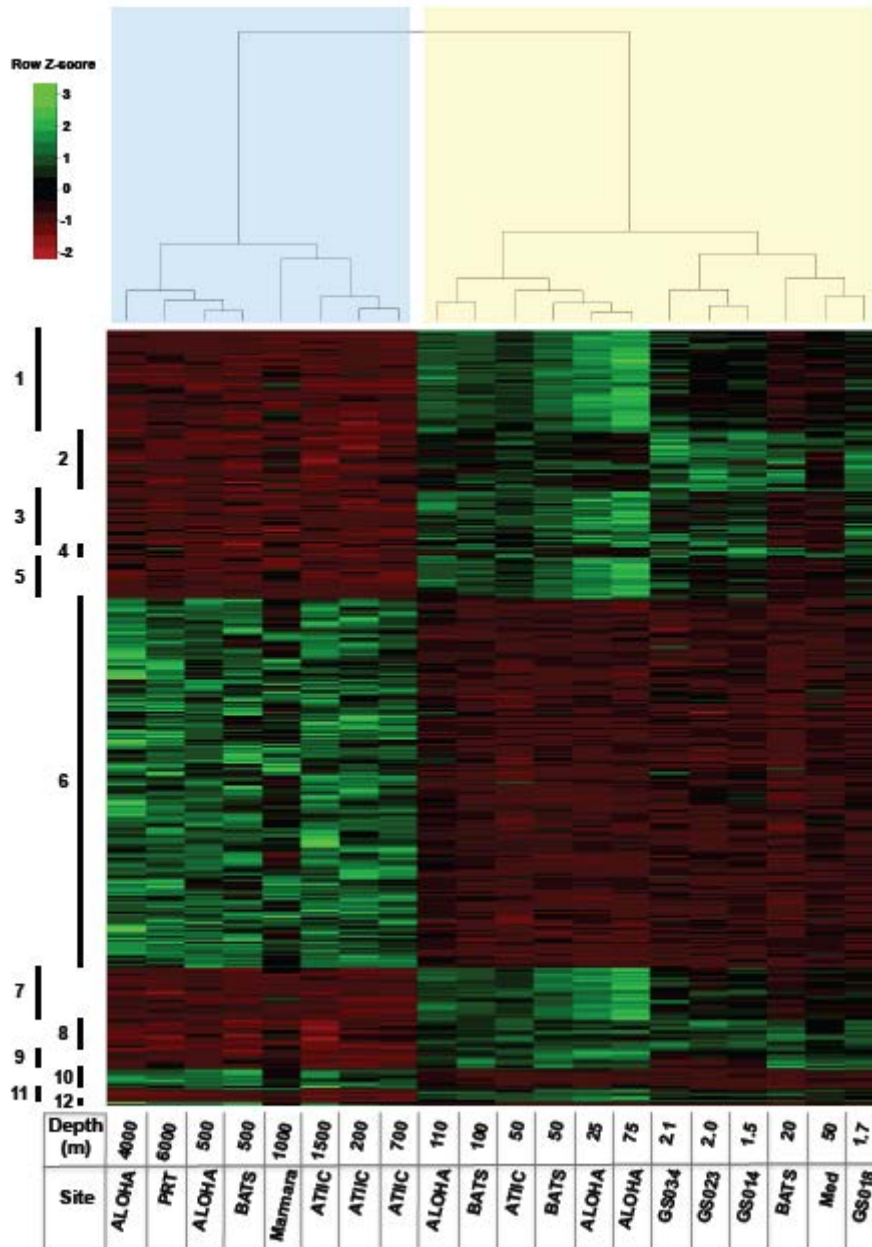
<http://127.0.0.1:8081/plosone/article?id=info:doi/10.1371/journal.pone.0061928>

Abundância de funções

- mapeamento de reads em ORFs anotadas

Abundância relativa espacial

- **COGs diferencialmente representados**
- Semelhante a genes diferencialmente expressos
- Heat maps, clusterização hierárquica



Based on 386 COGs shared by ATIIC, Aloha, BATS with differential representation

← COGs

Iquique not included

Plataformas web de processamento

- Laboratórios governamentais
- Serviços padronizados de processamento

MG-RAST

metagenomics analysis server

LOGIN



[Browse Metagenomes](#)

search for metagenomes



[Register](#)



[Contact](#)



[Help](#)



[Upload](#)*



[News](#)

About

MG-RAST (the Metagenomics RAST) server is an automated analysis platform for metagenomes providing quantitative insights into microbial populations based on sequence data.

| | |
|-------------------------|----------------|
| # of metagenomes | 77,307 |
| # base pairs | 25.81 Tbp |
| # of sequences | 236.94 billion |
| # of public metagenomes | 12,527 |

The server primarily provides upload, quality control, automated annotation and analysis for prokaryotic metagenomic shotgun samples. MG-RAST was launched in 2007 and has over 8000 registered users and 77,307 data sets. The current server version is 3.3.3.3. We suggest users take a look at [MG-RAST for the impatient](#).

[Updates](#)

[MG-RAST 3.2.4 release notes \[October 2012\]](#)

* login required

This project has been funded in part with Federal funds from the National Institute of Allergy and Infectious Diseases, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN272200900040C.

This work was supported in part by the Office of Advanced Scientific Computing Research, Office of Science, U.S. Department of Energy, under Contract DE-AC02-08CH11357.

[cite MG-RAST](#)

Microbiome Details (Assembled Data)

Add to Genome Cart

 Browse Genome

 ATC BLAST Genome

About Genome

- [Overview](#)
- [Statistics](#)
- [Genes](#)

Overview

| | |
|----------------------------|---------------------------------------|
| Proposal Name | Sao Paulo Zoo Compost |
| Sample Name | Sample C4 |
| Taxon Object ID | 2156126000 |
| IMG Submission ID | 2671 |
| GOLD ID in IMG Database | Project Id: Gm0002180 |
| External Links | |
| Genome type | metagenome |
| Sequencing Status | Draft |
| IMG Release | |
| Comment | |
| Sample Information | |
| Sample Site | Sao Paulo Zoo composting operation |
| Sample Collection Date | January 26, 2011 |
| Isolation Country | Brazil |
| Sampling Strategy | 8 days after composting started |
| Sample Isolation | done 8 days after composting started |
| Temperature Range | Thermophile |
| Sample Assembly Method | newbler |
| Sample Geographic Location | Sao Pulo Zoo |
| Longitude | -46.62 |
| Latitude | -23.65 |



Easy submission



Manually supported submission process, with help available for meta-data provision. Accepted data formats include SFF (454) and FASTQ (Illumina and IonTorrent).

[Find out more](#)

Powerful analysis



Functional analysis of metagenomic sequences using InterPro - a powerful and sophisticated alternative to BLAST-based analyses. Taxonomy diversity analysis is performed using Qiime.

[Find out more](#)

Data archiving



Data automatically archived at the Sequence Read Archive (SRA), ensuring accession numbers are supplied - a prerequisite for publication in many journals.

[Find out more](#)


Projects

Latest public projects (Total: 37)

Metatranscriptomics of the marine sponge *Geodia barretti*: Tackling phylogeny and function of its microbial community.

Geodia barretti is a marine cold-water sponge harbouring high numbers of microorganisms. ...

[View more](#) - 1 sample

A metagenomics transect into the deepest point of the Baltic Sea reveals clear stratification of microbial functional capacities

The Baltic Sea is characterized by hyposaline surface waters, hypoxic and anoxic deep waters and ...

[View more](#) - 6 samples

Gut metagenome in European women with normal, impaired and diabetic glucose control

Type 2 diabetes (T2D) is a result of complex gene-environment interactions, and several risk ...

[View more](#) - 147 samples

Samples

Latest public samples (Total: 1053)

Fecal sample from Crohn's patient 1

Fecal sample from Crohn's patient 1 ...

[View more](#) - Taxonomy | Function results |

Fecal sample from Crohn's patient 10

Fecal sample from Crohn's patient 10 ...

[View more](#) - Taxonomy | Function results |

Fecal sample from Crohn's patient 2

Fecal sample from Crohn's patient 2 ...

[View more](#) - Taxonomy | Function results |

Fecal sample from Crohn's patient 3

Fecal sample from Crohn's patient 3 ...

[View more](#) - Taxonomy | Function results |

Fecal sample from Crohn's patient 4

Fecal sample from Crohn's patient 4 ...

[View more](#) - Taxonomy | Function results |

Data content

1053 public samples (37 public projects)

191 private samples (13 private projects)

News & events

Tweets

[Follow @EBImetagenomics](#)



EBI Metagenomics @EBImetagenomics 30 Sep

Check out our new analysis page, using improved data visualisation (Google & Krona charts), and with taxonomic info:

ebi.ac.uk/metagenomics/

Expand



EBI Metagenomics @EBImetagenomics 8 Aug

The poster we presented at #SMBECCB is now available at F1000 posters and describes the EBI metagenomics pipeline: [f1000.com/poster/f1000100](https://doi.org/10.1093/f1000/poster/f1000100)

Sugestão de leitura

Analytical Tools and Databases for Metagenomics in the Next-Generation Sequencing Era

Mincheol Kim¹, Ki-Hyun Lee¹, Seok-Whan Yoon¹, Bong-Soo Kim², Jongsik Chun^{1,2}, Hana Yi^{3,4,5*}

¹School of Biological Sciences & Institute of Bioinformatics (BIOMAX), Seoul National University, Seoul 151-742, Korea,

²Chunlab Inc., Seoul National University, Seoul 151-742, Korea, ³Department of Environmental Health, Korea University, Seoul 136-703, Korea, ⁴Department of Public Health Sciences, Graduate School, Korea University, Seoul 136-703, Korea,

⁵Korea University Guro Hospital, Korea University College of Medicine, Seoul 136-703, Korea

Metagenomics has become one of the indispensable tools in microbial ecology for the last few decades, and a new revolution in metagenomic studies is now about to begin, with the help of recent advances of sequencing techniques. The massive data production and substantial cost reduction in next-generation sequencing have led to the rapid growth of metagenomic research both quantitatively and qualitatively. It is evident that metagenomics will be a standard tool for studying the diversity and function of microbes in the near future, as fingerprinting methods did previously. As the speed of data accumulation is accelerating, bioinformatic tools and associated databases for handling those datasets have become more urgent and necessary. To facilitate the bioinformatics analysis of metagenomic data, we review some recent tools and databases that are used widely in this field and give insights into the current challenges and future of metagenomics from a bioinformatics perspective.

Keywords: computational biology, high-throughput nucleotide sequencing, metagenomics