# **Microbial diversity profiling**

# using

# **Next Generation Sequencing**

prof. Matteo Ramazzotti, PhD



The prokaryotes consist of

#### many millions

of genetically distinct organisms.

They are phylogenetically distinct from eukaryotes and comprise Bacteria and Archea.

Less than 1% are known because most of them cannot be isolated and cultured.



A combination of techniques is needed to give them a (partial) classification:

- Morphology
- Motility
- Structural features
- Oxygen requirements
- Metabolism

- Gram stain
- Temperature resistance
- GC%
- Physiology
- Genetics

Woese C PNAS 1990 www.textbookofbacteriology.net



The definition of "species" has evolved with scientific knowledge and is far to have a stable definition. According to John Wilkins:

- ✓ There is one species concept (and it refers to real species).
- There are two explanations of why real species are species: ecological adaptation and reproductive reach.
- ✓ There are seven distinct definitions of species plus 27 variations.
- ✓ And there are n+1 definitions of "species" in a room of n biologists.

**Taxonomy** is the practice and science of classification. It organizes *taxa* (singular *taxon*) in hierarchical, progressively inclusive classes termed **ranks**.

**Taxonomic profiling** is the characterization the taxonomic diversity of a population, e.g. a microbial community.





# The diversity of prokaryotes: taxonomy

Since January 2000, names of prokaryotes change at a rate approaching 750 validly published names every year.

"List of Prokaryotic names with Standing in Nomenclature" [...] provides accurate information about the current status of a name, synonyms, and other useful information."



# Genus Abiotrophia Family: "Aerococcaceae" - (List of genera included in the family) Suborder: Order: "Lactobacillales" - (List of genera included in the order - List of families included in the order) Subclass: Class: "Bacilli" - (List of orders and suborders included in the class) Division or phylum: "Firmicutes" - (see Phylum "Firmicutes" in the file Classification of domains and phyla) Domain or empire: "Bacteria" - (see Domain "Bacteria" in the file Classification of domains and phyla) Species Abiotrophia defectiva (Bouvet et al. 1989) Kawamura et al. 1995, comb. nov. (Type species of the genus).

Species Abiotrophia defectiva (Bouvet et al. 1989) Kawamura et al. 1995, comb. nov. (Type species of the genus).
 Type strain (see also StrainInfo.net): strain SC10 = ATCC 49176 = CIP 103242 = CCUG 27639 = CCUG 27804 = DSM 9849 = LMG 14740.
 GenBank/EMBL/DDBJ accession number for the 16S rRNA gene sequence of the type strain: D50541.
 Basonym: © Streptococcus defectivus Bouvet et al. 1989.
 Etymology: L. fem. adj. defectiva, deficient.
 Reference: KAWAMURA (Y.), HOU (X.G.), SULTANA (F.), LIU (S.), YAMAMOTO (H.) and EZAKI (T.): Transfer of Streptococcus adjace Abiotrophia gen. nov. as Abiotrophia adiacens comb. nov. and Abiotrophia defectiva comb. nov., respectively. Int. J. Syst. Bacteriol., 1995, 45 Original article in USEM Online



INIVERSIT

LPSN is accessible at www.bacterio.net

# The diversity of prokaryotes in databanks (10 years later)

2007 2017 domain Bacteria (0/703222/0) domain Bacteria (0/1502575/0) **RDP11** RDP8 phylum "Acidobacteria" (0/7334/0) phylum "Actinobacteria" (0/226279/0) phylum "Firmicutes" (0/256748/0) phylum "Aquificae" (0/1013/0) phylum "Bacteroidetes" (0/180267/0) phylum "Bacteroidetes" (0/69546/0) phylum "Caldiserica" (0/230/0) phylum "Actinobacteria" (0/119442/0) phylum "Chlamydiae" (0/800/0) phylum "Aquificae" (0/851/0) phylum "Chlorobi" (0/409/0) phylum "Caldiserica" (0/44/0) phylum "Chloroflexi" (0/22363/0) phylum "Chlamydiae" (0/307/0) phylum "Chrysiogenetes" (0/12/0) phylum "Chlorobi" (0/319/0) phylum "Deferribacteres" (0/420/0) phylum "Chloroflexi" (0/3764/0) phylum "Deinococcus-Thermus" (0/2186/0) phylum "Chrysiogenetes" (0/9/0) phylum "Dictyoglomi" (0/20/0) phylum "Deferribacteres" (0/276/0) phylum "Elusimicrobia" (0/403/0) phylum "Deinococcus-Thermus" (0/851/0) phylum "Fibrobacteres" (0/850/0) phylum "Dictyoglomi" (0/21/0) phylum "Fusobacteria" (0/11979/0) phylum "Fibrobacteres" (0/235/0) phylum "Gemmatimonadetes" (0/1638/0) phylum "Fusobacteria" (0/2213/0) phylum "Lentisphaerae" (0/1798/0) phylum "Gemmatimonadetes" (0/738/0) phylum Nitrospirae (0/1802/0) phylum "Lentisphaerae" (0/176/0) phylum "Planctomycetes" (0/16438/0) phylum "Nitrospira" (0/818/0) phylum "Proteobacteria" (0/437996/0) phylum "Spirochaetes" (0/11723/0) phylum "Planctomycetes" (0/4311/0) phylum "Synergistetes" (0/1333/0) phylum "Proteobacteria" (0/202798/0) phylum "Tenericutes" (0/5184/0) phylum "Spirochaetes" (0/3354/0) phylum "Thermodesulfobacteria" (0/119/0) phylum "Synergistetes" (0/778/0) phylum "Thermotogae" (0/728/0) phylum "Tenericutes" (0/2233/0) phylum BRC1 (0/406/0) phylum "Thermodesulfobacteria" (0/94/0) phylum Parcubacteria (0/174/0) phylum "Thermotogae" (0/478/0) phylum Microgenomates (0/127/0) phylum "Verrucomicrobia" (0/4194/0) phylum SR1 (0/343/0) phylum Bacteria incertae sedis (0/271/0) phylum Candidatus Saccharibacteria (0/2670/0) phylum OP10 (0/202/0) phylum Latescibacteria (0/559/0) phylum OP11 (0/65/0) phylum "Armatimonadetes" (0/1178/0) phylum OD1 (0/103/0) phylum "Verrucomicrobia" (0/10731/0) phylum BRC1 (0/56/0) phylum "Acidobacteria" (0/16286/0) phylum Cyanobacteria (0/9562/0) phylum Firmicutes (0/482511/0) phylum Cyanobacteria/Chloroplast (0/26471/0) phylum SR1 (0/2/0) phylum Marinimicrobia (0/1065/0) phylum WS3 (0/119/0) phylum Aminicenantes (0/1543/0) phylum TM7 (0/805/0) phylum Omnitrophica (0/21/0) phylum Acetothermia (0/40/0) phylum Poribacteria (0/105/0) phylum Atribacteria (0/69/0) phylum Cloacimonetes (0/220/0) phylum Candidatus Calescamantes (0/3/0) phylum candidate division WPS-1 (0/311/0) phylum candidate division WPS-2 (0/183/0) phylum Hydrogenedentes (0/468/0) phylum candidate division ZB3 (0/73/0)

http://rdp.cme.msu.edu/hierarchy/hierarchy\_browser.jsp?



phylum Ignavibacteriae (0/870/0)

phylum Nitrospinae (0/619/0)





1 rulu

# The diversity of prokaryotes: guidelines

Macroscopic morphology	Microscopic morphology	Cell component	Growth characteristics	Metabolism	Molecular genetics	Biochemistry
Type of growth	Shape of the cell	Cell wall	Athmospheric requests	Carbon sources	GC%	Exposed antigens
Shape of the colony	Dimension of the cell della cell	Gram staining	PH tolerance	Nitrogen sources	DNA sequences	Enzymes
Pigments	Internal structures	Capsule	Growth temperature	Sulfur sources	rDNA sequences	Toxins
	Accessory structures		Lifestyle	Type of fermentation	Molecular probes	
			Sensitivity to antibiotics	Type of respiration	PCR	
				End products		



1\_ rulu

# The current definition of a prokaryotic species

A bacterial species is defined as a group of strains having

- 1. > 70% DDH (DNA/DNA Hybridization) similarity
- 2. < 5°C ∆tm
- 3. < 5% mol G + C difference of total genomic DNA
- 4. > 98% 16S rRNA identity (full length gene).

Modern genomics have proposed two additional measures

5. AAI : average amino acid identity.

6. ANI: average nucleotide identity (>94% ANI ~ > 70% DDH)

7. digital DDH (based on complete genomes) (> 70% dDDH)

In prokaryotes, species definition is "operational", used to rapidly establish a taxonomic framework based on microbial evolution.

#### No biological criteria are involved in the definition of a prokaryote species

Thomposon CC et al. BMC Genomics 2013, 14:913



1\_ rule

#### FOR OUR PURPOSES:

= 1 sequence and one species

#### **REAL BACTERIA**

One species can have multiple, different copies of a gene

Two similar species may share an identical sequence











FIRENZE

# Chain Reaction: an amplification cascade



# DNA sequencing with the Sanger method



Capillary electrophoresis

DNA flowing down



Frederick Sanger Nobel prize 1958/1980





nain

detector

Template Sequence

3'GAGCAAATTCCGATACATTATTGT... 5' Primer

5'CTCGTTTAAG... 3'

CTCGTTTAAGGGTATGTA CTCGTTTAAGGGTA CTCGTTTAAGGGTA CTCGTTTAAGGGTATG CTCGTTTAAGGGTATG CTCGTTTAAGGGTATGT CTCGTTTAAGGGTATGTA CTCGTTTAAGGGTATGTA





Sampling ANNUM CONTRACTOR Isolation methods of TTT Culturing DNA extraction PCR evade NGS Sanger sequencing sequencing

#### **Metagenomics**

only needs DNA (whole extract or amplified via PCR) sequencing and occurs in parallel: for need no isolation.

Traditional microbial ecology require organisms collected from an environment to be cultured in the laboratory. From 90% to 99.9% of the cells defy cultivation, and therefore identification.







Platform	Instrument	Unit	Reads/unit	Read Length (bp)	Read Type	Error Type
PacBio	PacBio RS II	SMRT Cell	47000	14000	SR	indel
PacBio	PacBio RS	SMRT Cell	22000	4600	SR	indel
Roche 454	GS FLX+ / FLX	1 PTP	700000	700	SR	indel
Roche 454	GS FLX+ / FLX	1/2 PTP	350000	700	SR	indel
Roche 454	GS FLX+ / FLX	1/4 PTP	125000	700	SR	indel
Roche 454	GS FLX+ / FLX	1/8 PTP	50000	700	SR	indel
Roche 454	GS FLX+ / FLX	1/16 PTP	20000	700	SR	indel
Illumina	MiSeq v3	Lane	25000000	600	SR & PE	substitution
Illumina	MiSeq v2 Nano	Lane	1000000	500	SR & PE	substitution
Ion	PGM 318	Chip	4000000	400	SR	indel
Ion	PGM 316	Chip	2000000	400	SR	indel
Ion	PGM 314	Chip	400000	400	SR	indel
Roche 454	GS FLX+ / FLX	1 PTP	70000	400	SR	indel
Illumina	HiSeq X	Lane	375000000	300	PE	substitution
Illumina	HiSeq 3000/4000	Lane	312500000	300	SR & PE	substitution
Illumina	NextSeq 500 High-Output	Run	40000000	300	SR & PE	substitution
Illumina	NextSeq 500 Mid-Output	Run	13000000	300	PE	substitution
Illumina	HiSeq Rapid Run	Lane	150696000	300	SR & PE	substitution
Illumina	GAIIx	Lane	42075000	300	SR & PE	substitution
Illumina	MiSeq v2 Micro	Lane	4000000	300	SR & PE	substitution
Illumina	HiSeq High-Output v4	Lane	250000000	250	SR & PE	substitution
Illumina	HiSeq High-Output v3	Lane	186048000	250	SR & PE	substitution
Illumina	MiSeq v2	Lane	16000000	250	SR & PE	substitution
Illumina	MiSeq	Lane	5000000	250	SR & PE	substitution
Illumina	HiScanSQ	Lane	93024000	200	SR & PE	substitution
Ion	Proton I	Chip	6000000	200	SR	indel
SOLiD	5500×I W	Lane	266666667	100	SR & PE	A/T Bias
SOLiD	5500 W	Lane	266666667	100	SR & PE	A/T Bias
SOLiD	5500	Lane	81500000	100	SR & PE	A/T Bias
SOLiD	5500xl	Lane	81500000	100	SR & PE	A/T Bias

Taken from http://genohub.com/ngs-instrument-guide/



FASTA format

#### >SEQ\_ID GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCAT

FASTQ format





www.bioinformatics.babraham.ac.uk/projects/fastqc



1\_ rulu

Introduced by Phil Green's group at the University of Washington in the 1990s, the phred quality estimator **calls each sequenced base with an integer related to the probability for a base to be incorrectly identified**. Born **for ABI chromatograms**, it has been reused in NGS base-callers, with technology specific modifications.

$$p = 10^{\frac{q}{-10}}$$
  $q = -10\log_{10}(p)$ 

p = error	probability	for the	base
-----------	-------------	---------	------

if p = 0.01 (1% chance of error), then q = 20

e.g. q = 32 p =  $10^{(32/-10)} = 0.00063$ 

Phred quality values are rounded to the nearest integer

Quality Value	Error Probability	Probability Called Base is Correct	Description
10	0.1	0.9	error rate of 1 in 10
20	0.01	0.99	error rate of 1 in 100
30	0.001	0.999	error rate of 1 in 1000
40	0.0001	0.9999	error rate of 1 in 10000



### A good molecular classifier is a regoin in the DNA that:

- $\checkmark$  must be present in all organisms of interest.
- ✓ cannot be subject to transfer between species (lateral transfer).
- must display an appropriate level of sequence conservation for the divergences of interest.
- must be sufficiently large to contain a record of the historical information.
- ✓ (optional) singolarity, i.e. being present in only one copy per genome ...

We need a genetic element that is involved in a central dogma of biology and that may accommodate a significant number of changes though maintaining its fundamental function...



1\_ rile

Ribosomes are RNA/protein complexes that act as the universal protein factory in all living cells. They have been conserved trough billions of years. So some regions have sequences identical in all the three kingdoms of life.

human	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAATTC</mark> CAGCTCCAAT <mark>AG</mark> CGTATATTAAAG <mark>TT</mark> GC <mark>T</mark> GCAGT <mark>TAAA</mark> AAG
yeast	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAATTC</mark> CAGCTCCAAT <mark>AG</mark> CGTATATTAAAG <mark>TT</mark> GT <mark>T</mark> GCAGT <mark>TAAA</mark> AAG
corn	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAATTC</mark> CAGCTCCAAT <mark>AG</mark> CGTATATTTAAG <mark>TT</mark> GT <mark>T</mark> GCAGT <mark>TAAA</mark> AAG
Escherichia coli	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAAT</mark> ACGGAGGGTGCA <mark>AG</mark> CGTTAATCGGAA <mark>TT</mark> AC <mark>T</mark> GGGCG <mark>TAAA</mark> GCG
Anacystis nidulans	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAAT</mark> ACGGGAGAGGCA <mark>AG</mark> CGTTATCCGGAA <mark>TT</mark> AT <mark>T</mark> GGGCG <mark>TAAA</mark> GCG
Thermotoga maratima	<mark>GTG</mark> CCAGCA <mark>GCCGCGGTAAT</mark> ACGTAGGGGGCA <mark>AG</mark> CGTTACCCGGAT <mark>TT</mark> AC <mark>T</mark> GGGCG <mark>TAAA</mark> GGG
Methanococcus vanniel	ii <mark>GTG</mark> CCAGCA <mark>GCCGCGGTAAT</mark> ACCGACGGCCCG <mark>AG</mark> TGGTAGCCACTC <mark>TT</mark> AT <mark>T</mark> GGGCC <mark>TAAA</mark> GCG
Thermococcus celer	<mark>GTG</mark> GCAGCC <mark>GCCGCGGTAAT</mark> ACCGGCGGCCCG <mark>AG</mark> TGGTGGCCGCTA <mark>TT</mark> AT <mark>T</mark> GGGCC <mark>TAAA</mark> GCG
Sulfolobus sulfotaricus	<mark>GTG</mark> TCAGCC <mark>GCCGCGGTAAT</mark> ACCAGCTCCGCG <mark>AG</mark> TGGTCGGGGTGA <mark>TT</mark> AC <mark>T</mark> GGGCC <mark>TAAA</mark> GCG

#### Ribosomal RNAs in Prokaryotes

Name 5S	<b>Size (bp)</b> 120	<b>Location</b> Large subunit of ribosome	Carl Woese in 1990 realized
<b>16S</b>	1500	Small subunit of ribosome	→ that IOS IDINA could be
23S	2900	Large subunit of ribosome	phylogenesis.



16S rDNA is a perfect candidate because it is:

- present in all prokaryotic organisms
- sufficiently conserved in some regions
- highly variable in other regions

A Simplified map of the 16S rRNA molecule





CONSERVED REGIONS: unspecific applications VARIABLE REGIONS: group or species-specific applications



INIVERSIT







The general definition of distance for biological sequences is

#### d = fraction of mismatches at aligned positions

So sequences need to be aligned, compared and differences counted and expressed as a fraction of the sequence length. Alternative definitions exists, based e.g. on the number of k-mer shared (no alignment is necessary in this case...)

Substitution models can be introduced now, but no generalized consensus exists.

Problems arise when gaps have to be taken into account



0 0.143 0.077 0.111 0.05 0.038 0.091 0.076 0.083 0.077 0.088 0 0.232 0.233 0.177 0.177 0.191 0.209 0.215 0.213 0 129 0.136 0.127 0.143 0.157 0.129 0.127 0 0.143 0.149 0.176 0.136 0.143 0.143 0 0.052 0.054 0.149 0.072 0.215 0.188 0.183 0.141 0.157 0.149 0.052 0.155 0.073 0.15 0.054 0.055 0.213 0.207 0.177 0.129 0.154 0 0.147 0.062 0 126 0 133 0 145 0 149 0 155 0 147 0 146 0 189 0.16 0.174 0.072 0.073 0.053 0.053 0.162 0.05 0.162 0.147 0.066 0.139 0.127 0.137 0.149 0.139 0.156 0.145 0.063 0.059 0.062 0.18 0.232 0.205 0.178 0.183 0.186 0.08 0.082 0.083 0.187 0.092 0.284 0.256 0.239 0.202 0.231 0.225 0.133 0.134 0.139 0.224 0.161 0.212 0.176 0.164 0.154 0.094 0.098 0.102 0.164 0.111 0.093 0.127 0.113 0.116 0.131 0.117 0.14 0.143 0.162 0.183 0.164 0.11 0.099 0.127 0.137 0.148 0.047 0.049 0.043 0.198 0.165 0.176 0.162 0.106 0.098 0.096





Each OTU will have a variable number of sequences (i.e. individuals) assigned.





1 rulu









ΟΤυ	ASV
Can be subject to reference bias	Reference is not used until taxonomy assignment
OTU tables cannot be combined between studies	ASV tables can be compared across studies
Represented by a consensus sequence	Represented by an exact sequence
Can represent multiple species with different sequences	If it represents multiple species, it is because they share the sequence
Subject to chimeric sequences	Subject to chimeric sequences
Chimera detection can be complex and may require reference bias	Chimera detection is simple and reference-free



\_\_\_\_



**Chao1** richness index is a estimation of species abundance taking into account rare sequences.

$$S_{chao1} = S_{obs} + \frac{n_1 (n_1 - 1)}{2 (n_2 + 1)}$$

**Shannon** diversity index is related to the effective species number given an OTU definition

$$\hat{H}_{shannon} = \sum_{i=1}^{S_t} \frac{\hat{C}\pi_i \ln\left(\hat{C}\pi_i\right)}{1 - \left(1 - \hat{C}\pi_i\right)^N}$$

#### $\beta$ -diversity: between samples

**Jaccard** dissimilarity index evaluate the fraction of OTUs shared by two communities

$$D_{Jaccard} = \frac{S_{AB}}{S_A + S_B - S_{AB}}$$

Yue & Clayton theta similarity coefficient (dissimilarity index) evaluate the fraction of OTUs shared by two communities

$$D_{\Theta_{YC}} = 1 - \frac{\sum_{i=1}^{S_T} a_i b_i}{\sum_{i=1}^{S_T} (a_i - b_i)^2 + \sum_{i=1}^{S_T} a_i b_i}$$



# OTU-based analysis: the rarefaction curves

Rarefaction is a technique to compare population estimators (e.g. richness). Rarefaction curves are plots of value of the estimator as a function of the number of individuals sampled.

To build a rarefaction curve the following procedure is used:

- given N sequences (the sample) and
- an initial richness value of R=0:
- 1. Take a random sample with n% individuals
- 2. Calculate richness R'
- 3. Accumulate richness with R = R+R'

4. go to 1.

On the left the steep slope indicates that a large fraction of the species diversity remains to be discovered.

On the right, a flat curve indicate richness saturation and a growing curve indicate an insufficient sampling.







Number of reads sampled

Insufficient !!!



l n



RDP Release 11, Update 4 :: May 26, 2015

3,224,600 16S rRNAs :: 108,901 Fungal 28S rRNAs



RDP offers a number of online tools for the analysis of 16S reads, from classification to tree building and probe design.

1\_ rule

Each sequence is framed into its taxonomic assignment using infernal SS aligner. The taxonomy is kept consistent by periodical revision.



BLAST allows a fast and efficient euristic search into 16S rDNA databases. The best-hit method is usually applied.



with identical scores, so the degree and confidence of the assignment depend on the the taxonomic spread of the matches.



## Taxonomic assignment with k-mer search

One of the simplest and fastest way or searching in a database is to pre-index it using words of defined length. Length should be optimized for sensitivity and speed.

Given the size length of the word (k) the search is space is  $4^{k}$ .

							5	1024
Each sequence	ce of	E r	o rank Root (0/20/7136	97) (selected/	match/total	RDP sequences)	6	4096
, databaaa	of	+	domain Bacteria (0.	/20/703222)	indeen/ cocat	tor sequences,	7	16384
a ualabase	01	+	phylum "Actinoba	cteria" (0/20/	119442)	-	8	65536
16 <u>5</u> 500000			class Actinobac	teria (0/20/1	(0/20/1134	46)	9	262144
TOS Sequent	~ <b>C</b> 2	÷	order Bifi	dobacteriales	(0/20/924)	40)	10	1048576
1		+	family	Bifidobacteria	ceae (0/20	/924)	10	1040370
		+	gen	us Bifidobacter	rium (0/20/	842)		
				S000008715	1.000 1367	Bifidobacterium breve; JCM 701	6;	
	ΑΤΛΟΟΟΛΤΟΟΛΟΤ			5000015908	1.000 1367	Bifidobacterium breve; JCM 701	7;	
	ATACUCATCUACT			S000136536	1.000 1377	Bifidobacterium longum subsp. lo	ongum biovar Lon	igum; Y10; .
	ATACCCATCCACT			5000381784	1.000 1429	Bifidobacterium longum (T): ATCC	. 15700;	
	ATACUCATCUACT			5000414117	1.000 1412	Bifidobacterium longum subsp. //	ongum by Infanti	is: BG5:
$\sim$	ΛΤ <b>Λ<u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u><u></u></b>			S000443702	1.000 1429	Bifidobacterium breve: BR2:	ongum ov. miana	3, 003,
	ATACOCATCOACT			5000468460	1.000 1431	Bifidobacterium breve; BGM6;		
	ATACGCATCGACT			S000536096	1.000 1393	uncultured bacterium; rRNA019;		
				S000536186	1.000 1417	uncultured bacterium; rRNA109;		
	ATAC <b>GCATCGA</b> CT			S000536224	1.000 1413	uncultured bacterium; rRNA147;		
				S000536312	1.000 1388	uncultured bacterium; rRNA235;		
	ATACG <b>CATCGAC</b> T			5000536425	1.000 1420	Bifidebactorium longum subsp. 4		
/			i i i i i i i i i i i i i i i i i i i	5000806300	1.000 1319	Bifidobacterium longum subsp. lo	ongum by Infanti	is: THT-010201:
/	ATACGC <b>ATCGACT</b>			S001052051	1.000 1418	uncultured bacterium: rRNA187:	ongum ov. mund	3, 1111-010201,
				S001239302	1.000 1416	Bifidobacterium longum subsp. in	nfantis ATCC 156	97;
	Kmer decomposition			5001239305	1.000 1416	Bifidobacterium longum subsp. in	nfantis ATCC 156	97;
				5001239307	1.000 1416	Bifidobacterium longum subsp. i	nfantis ATCC 156	97;
I				S001239309	1.000 1416	Bifidobacterium longum subsp. in	nfantis ATCC 1569	97;
Linoccianod								
Unassigned				Coorc	databa	a aguanaa hira	umbor of	
soquonco				Score	ualana	se sequences by n		
Sequence				shared	kmers	or a fraction over t	he length	
				21100.00				

INIVERSIT

Length

1 2

3

4

Space

4

16

64

256

# Taxonomic assignment with naïve bayesian classifiers

Developed at RDP, is it based on the conditional probability of assigning a sequence to a rank given a set of rank specific sequences in a training set. The probability is based on the frequencies of  $w_{mers}$  (optimal length 8).

Basically, we want to calculate P(G|S) i.e. the probability of having e.g. the sequence S in the genus G.

We can use the Bayes Theorem P(G|S) = P(S|G)

Given a dataset of N sequences, count w-mer occurrences (n) and derive the **prior** probability of a w-mer to occur in N.

Given a dataset of M sequences assigned to a genus G, count w-mer occurrences (m) and derive the genus conditioned probability.

Given a sequence S, count w-mer occurrences (v) and derive the genus conditioned probability of that sequence.

matteo.ramazzotti@unfi.it

(U) 
$$(U_{j,1}E_{j}) \leq P(E_{j})$$
  
(U)  $(U_{j,1}E_{j}) \leq P(E_{j})$   
(of  
Siven the database, they are  
constant and can be removed

B, (A)= P(B;)P(AB)

$$P(w-mer) = \frac{n(w-mer) + 0.5}{N + 1}$$
$$P(w-mer|G) = \frac{m(w-mer) + P(w-mer)}{M + 1}$$

 $P(S|G) = \prod P(w-mer|G)$ 

#### Do this for each genus, sort by P and retain the most probable genus



## Counts at different taxonomic ranks

The main application of taxonomic assignment is to evaluate the samples in terms of rank counts. When done for all the ranks, one can appreciate the actual biological variability, if any, in the sample of interest.





- ✓ A wide variability in copy number has been observed even within closely related organisms
- ✓ Different amplified regions have different classification performances.
- ✓ Some primers couples covers more "unknown" organisms than other.
- ✓ It has been well documented that evolutionarily distant SSU rRNA genes that are similar in nucleotide composition have been consistently but nevertheless incorrectly placed close together in phylogenetic trees.

And anyway,

✓ Inferring the phylogeny of organisms from any single gene carries some risks and should be corroborated by the use of other phylogenetic markers.











# Finding genes in untargeted reads

Several programs have been developed to address the issue of finding actual genes in the context metagenomics. Such programs incorporate HMM, GC content models and k-mer composition to predict:

- 1. Ribosome binding sites
- 2. Start and stop codons
- 3. Composition of genes (and non-coding regions)
- 4. Sequencing errors

FragGeneScan have been specifically developed for predicting genes in short reads (< 100 bp) so it can be used directly on raw (preprocessed) data !!!

#### Complete HMM of FragGeneScan



Noguchi et al. DNA Res. 2008, 15, 387-396 Rho M et al. Nucl. Acids Res 2010, gkq747 Zhu W et al. Nucl. Acids Res 2010, 38(12):e132 Kelley DR et al. Nucl. Acids Res 2012, 40(1): e9



# Many probable hits...should we take the top?

The lower the E-va The better the targe	 :?	The lo the	onger th better t	ne alig he ta	gnme Irget?	ent ?	-	The hig the b	gher the ide etter the tar	ntity get?	%
Q.name	Q.length	E-val	A.Length	Q.start	Q.end	gaps	strand	A.id%	T.name	T.st	T.en
HWI-ST170:227:5:1:2755-2167	77nt	4.9e-24	76	77	2	0	-	85.53%	YP_003142869.1	833	908
HWI-ST170:227:5:1:2755-2167	77nt	2.3e-15	43	77	35	0	-	88.37%	YP_606609.1	776	818
HWI-ST170:227:5:1:2755-2167	77nt	8.5e-15	33	75	43	0	-	93.94%	YP_004033514.1	823	855
HWI-ST170:227:5:1:2755-2167	77nt	8.5e-15	33	75	43	0	-	93.94%	YP_618644.1	823	855
HWI-ST170:227:5:1:2755-2167	77nt	8.5e-15	33	75	43	0	-	93.94%	YP_812565.1	823	855
HWI-ST170:227:5:1:2755-2167	77nt	1.2e-13	43	77	35	0	-	86.05%	YP_001667032.1	776	818
HWI-ST170:227:5:1:2755-2167	77nt	1.2e-13	40	75	36	0	-	87.50%	YP_001093961.1	763	802
HWI-ST170:227:5:1:2755-2167	77nt	6.6e-12	31	77	47	0	-	90.32%	YP_002354340.1	788	818
HWI-ST170:227:5:1:2755-2167	77nt	3.6e-10	31	77	47	0	-	87.10%	YP_002435333.1	803	833
HWI-ST170:227:5:1:2755-2167	77nt	1.7e-12	41	75	35	0	-	85.37%	YP_003912539.1	808	848
HWI-ST170:227:5:1:2755-2167	77nt	1.2e-13	31	77	47	0	-	93.55%	YP_004195006.1	806	836
HWI-ST170:227:5:1:2755-2167	77nt	3.6e-10	31	77	47	0	-	87.10%	YP_004368623.1	728	758
HWI-ST170:227:5:1:2755-2167	77nt	3.6e-10	31	77	47	0	-	87.10%	NP_903294.1	755	785

Targets are different, and probably in different organisms...but what if they have all the same function?



1 rile



This algorithm is of paramount importance in metagenomics since it allows to attribute functions and taxonomies to ambiguous assignments, that are frequent due to the shortness of the reads and orthology of genes.



MEGAN (now in its 6<sup>th</sup> version) analyzes BLAST-like output for each read to establish a taxonomy of the read, and then show the result in the taxonomy tree. Multiple samples can also be visualized.



LCA can be applied every time we need to quantify something that can be placed in a predefined, hierarchical stricture such as taxonomy.

Huson DH et al Genome Research 2011, 21:1552-1560 Mitra S et al Bioinformatics 2009, 25:1849–1855



# Exponential need of calculation power







# **Other alternatives?**



# Taxonomic markers from multiple alignments

We all know that proteins can be grouped in many different ways according to function, structure, orthology. In general, sequences that share common features can be represented as multiple sequence alignments.



A large number of databases exist with different scopes that collect MSAs.

			10	20	30	40	50	60	
				*		* <mark> </mark>	*	*	
Feature	1		# #	# #		# #	#	#	
1TOT_A		7	YTCNECKh HVe ·	TRWHCTVCe	DYDL	_CINCYNTk -	SH	THKMVKW	47
gi 70230	994	95	ISCDGCDeIAp-	-wHRYRCLQCs	DMDL	_CKTCFLGgv	/kpe <b>GH</b> gd	DHEMVNM	142
gi 50750	9334	201	VRCRVCKt fpITg-	· LRYRCLKCl	NFDL	_CQVCFFTgr	hsk <b>PH</b> ks	SHPVVEH	249
gi 50257	7626	582	SECTICLtalFS-	NRFKCVSCp	KFDL	CRSCYQKvd	leIHp-	AHAFLSL	626
gi 47222	2763	98	IICDSCKkhgIMg-	MRWKCKVCf	DYDL	_CTQCYMNn -	KHdl	SHAFERY	142
gi 40743	3717	1022	RVCNNCLk - eFDe ·	-gKMVSCADCd	DFDL	_CITCILGhk	:hg <b>HH</b> p-	SHTFVLL	1068
gi 51261	627	15	PPCKGCSs-yLMe-	PYIKCAECgp -	pEFLL	_CLQCFSGfe	:yk <b>KH</b> q≤	DHSYEIM	63
gi 16944	4480	367	RTCNCCIq-dLPe-	-aEFVHCQTCd	DFDL	_CKVCFAKnr	·h <mark>GH</mark> hp	KHAFSPI	413
gi 42546	5497	336	RTCNCCVq-eHPe-	-aEFLHCRMCe	DFDL	_CQSCFARds	; h <b>GH</b> hp	KHSFAPA	382
gi 40745	5179	373	IICDGCNaegLA-	VQYHCADCe	DYDL	_CQSCYKAgt	rc-gyk <b>GH</b> t-	YHLEFNA	421
gi 49648	3418	505	FVCDYCLe-pISe-	ARFHCQSCv	DFD\	/CSSCYPSra	ikSHaq	KHGEVQL	550
gi 56512	2242	87	YLCDECQy-aIMpl	-sYVFECNTCe	NFTL	_CKKCFKKg-	KH	EHPLKKM	130
gi 42555	5336	1607	SFCDACLm - nNYg-	LRFHCDVCw	DFDL	_CFKCYRSqr	ıi <b>IH</b> p-	KHSFTNP	1651
gi 47207	7916	134	VTCDGCEg - pVVg	TRFKCSVCp	NYDL	_CSACQAKg -	THt -	EHPLLPI	176
gi 42551	1863	1612	MFCDCCLl - dIYg-	FYYTCSTCf	ECDL	CDKCYLSvs	;kI <b>H</b> p-	AHSSFMK	1656
gi 40744	4696	1669	YSCDICLs-tIWg-	· PVYECETCl	DFT/	ACKKCHGRin	1 <b>lYH</b>	GHLRLEN	1712
gi 30725	5524	289	IRCDGCGvlpITg-	PRFKSKVKe	DYDL	_CTICYSVm-	GN	EGDYTRM	331
gi 75127	750	167	FKCDNCGiepIQg-	VRWHCQDCppe	msldFCDS	SCSDCLHEt -	dIHke	DHQLEPI	218
gi 70230	994	144	FTCDHCQg-lIIg-	RRMNCNVCd	DFDL	_CYGCYAAkk	: y s y GHl p	THSITAH	191
gi 47223	3744	1889	YACDHCQg-lIVg-	SRINCNVCe	DFDL	_CFGCYNAkk	:ypd <b>SH</b> lp	THRITVY	1936
gi 60467	7539	3292	FSCDLCNinpITg-	KRWNCSNCg	DFDL	_CNQCYQNpe	e k DHpk	DHIFKEF	3338
gi 30678	3519	2616	YCCDGCSt vpILr+	RRWHCTVCp	DFDL	_CEACYEV1 d	ladrlpp <b>PH</b> tr	DHPMTAI	2667
gi 50758	3066	1840	YACDHCQg-vIIg-	RRMNCNVCd	DFDL	_CYGCYSAkk	:ysd <b>SH</b> lp	THSITVY	1887
gi 47585	520	2706	VTCDGCQmfpINg-	SRFKCRNCd	DFDF	FCETCFKTk -	KHnt	RHTFGRI	2750
gi 60465	5335	40	YSCNGCGs-eIWpp	kqERYACNECs	NFDL	CSECYRKer	ıil <b>IN</b> gt	QEEKDKL	89
gi 17530	949	281	NSCAGCRkehIVg-	IRFRCQVCr	DISL	_CLPCFAVgf	agg <b>RH</b> ep	GHRMCEV	329
gi 11760	914	258	YICHTCGn-eSIn-	VRYHNLRAr	DTNL	_CSRCFQEgh	1fg <mark>AN</mark>	FQSSDFI	302
gi 46852	2178	7	VSCDACLkgnFRg-	RRYKCLICy	DYDL	CASCYESga	ttt <b>RH</b> tt	DHPMQCI	55

Multiple Sequence Alignments allow to catch sequence conservation in the aligned sequences and identify columns that contribute to define the reason for which they are grouped together.





Model positions trough a series of interconnected states with pre-computed emission probabilities



http://pfam.sanger.ac.uk/



# Clade specific marker genes: improved taxonomy

MetaPhIAn2 (2012) relies on ~1M unique clade-specific marker genes pre-identified from ~17,000 reference genomes (~13,500 bacterial and archaeal, ~3,500 viral, and ~110 eukaryotic), allowing to compute abundance at the "species"-level for bacteria, archaea, eukaryotes and viruses.





Segata Net al. Nat Methods. 2012, 9: 811-814

# Using 16S in non-targeted metagenomics

RiboFrame solves the issue of bias-byamplification when using 16 rDNA.

Reads overlapping the 16S rDNA genes are identified using Hidden Markov Models and a taxonomic assignment is obtained by naïve Bayesian classification.

All reads identified as ribosomal are coherently positioned in the 16S rDNA gene, allowing the use of the topology of the gene (i.e., the secondary structure and the location of variable regions) to guide the abundance analysis.





INIVERSIT

FIRENZE

# Estimating community functions from genes



matteo.ramazzotti@unfi.it

università degli studi FIRENZE

XIPARTIMENTO DI SCIENZE BIOMEDICHE

# Estimating community functions from taxonomy



Langille MG et al. Nat Biotechnol. 2013, 31(9):814-21

UNIVERSITÀ Degli stud

FIRENZE

# MGnify: the EBI metagenomics resource



		1763	studies	112090
fν	7	10988	34 sample	s 7026
- J		14186	analyse	es 2035
				17961
				1739
				P
		63 <sup>#</sup>	~~~	Th
n	Digestive	Aquatic	Marino	Plants
4)	system	(19005)	(1597/1)	(10/09)
+)	(31053)	(13003)	(15574)	(10403)
	(51055)			
10				
	. Su	•••••••••••••••••••••••••••••••••••••••	•	
	Digestive	Skin	Wastewater	Food
9)	system	(4337)	(2473)	production
	(6782)			(1245)

1 700

110000 amplicon assemblies metabarcoding metagenomes metatranscriptom







Soil (10389



