

Genomes, metagenomes and environments:
a perspective

Alessandra Carbone
Laboratoire de Génomique des Microorganismes
UMR7238 CNRS-Université Pierre et Marie Curie

The sargasso sea **Santa Cruz whale carcass bone**

Different profiles of gene enrichment in environment specific-functions




The sargasso sea

Nutrient poor environment, and its genes for ABC-type transporters dedicated to amino-acids transport and metabolism are translationally optimized.



Santa Cruz whale carcass bone

Microbes live in an abundant food source. Translational optimization is shown in energy production and conversion genes.



This difference can reflect functional adaptation of microbes to different environmental conditions

The sargasso sea

Nutrient poor environment, and its genes for ABC-type transporters dedicated to amino-acids transport and metabolism are translationally optimized.



Santa Cruz whale carcass bone

Microbes live in an abundant food source. Translational optimization is shown in energy production and conversion genes.

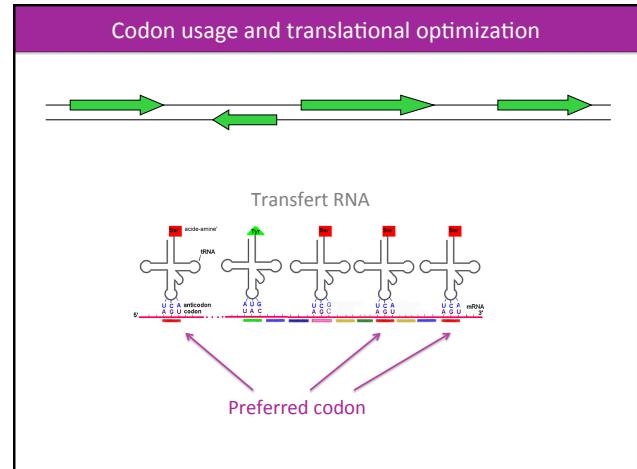


This difference can reflect functional adaptation of microbes to different environmental conditions

Three main results are shown for microbes:

1. Genome coding (codon bias) can be automatically studied and a set of optimized (most biased) genes can be identified for each genome
2. There exists a genome organization based on codon bias that reflects environmental living conditions
3. Codon biases reflect metabolic processes important for an organism

Are these statements true for microbial communities?

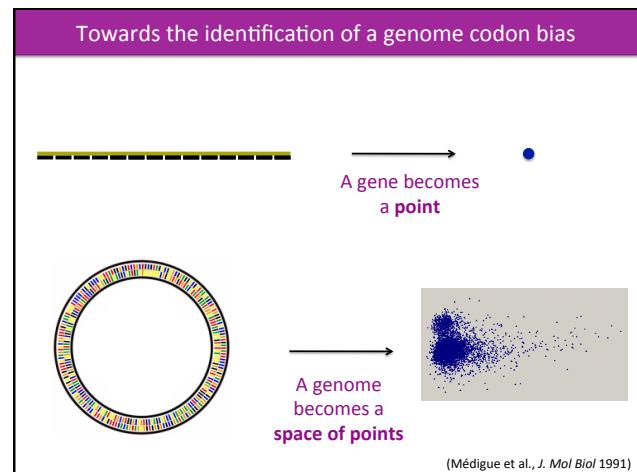


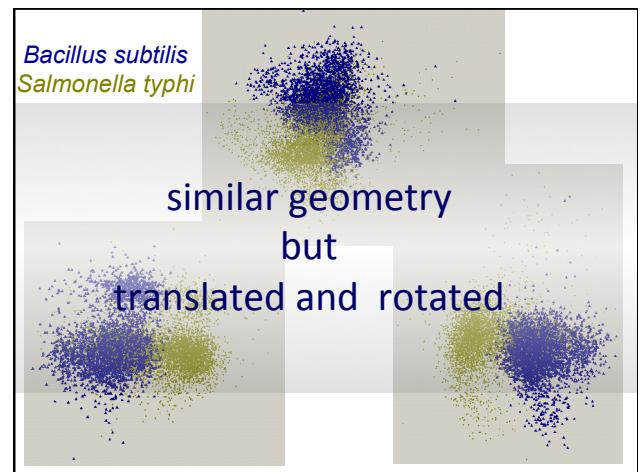
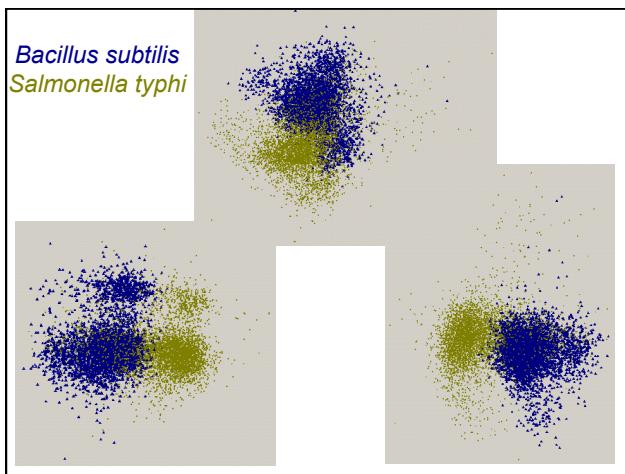
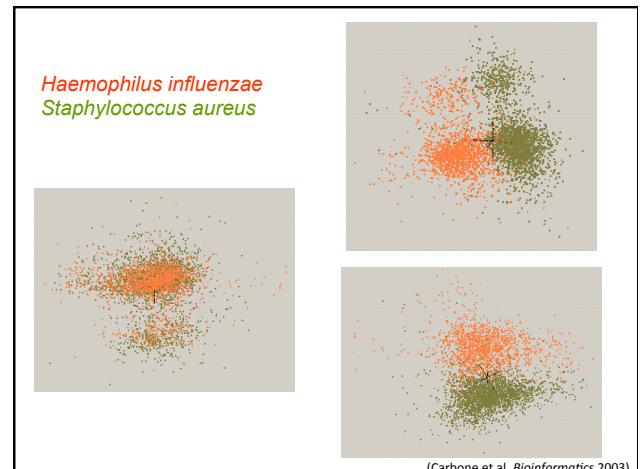
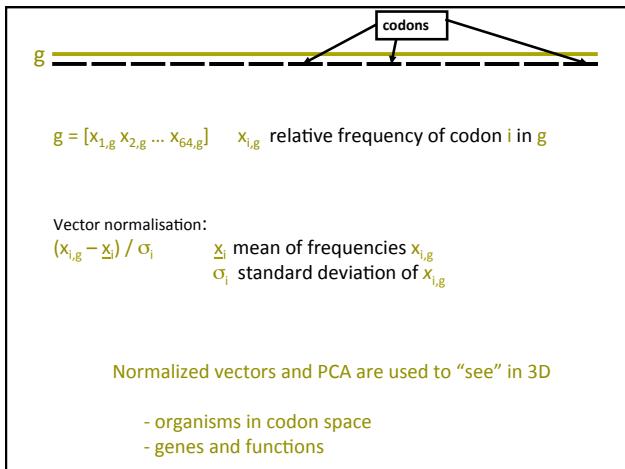
Codon usage and translational optimization

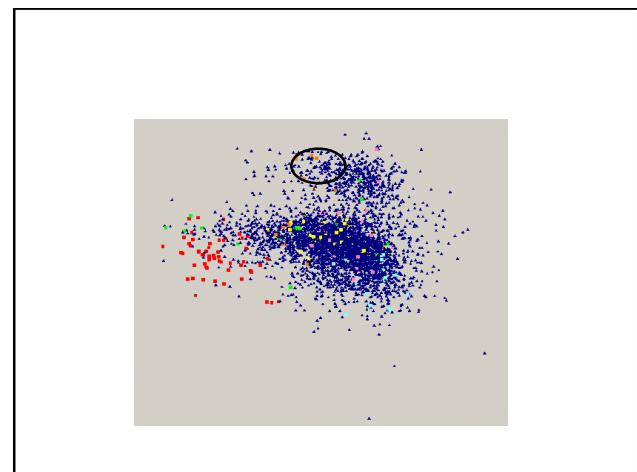
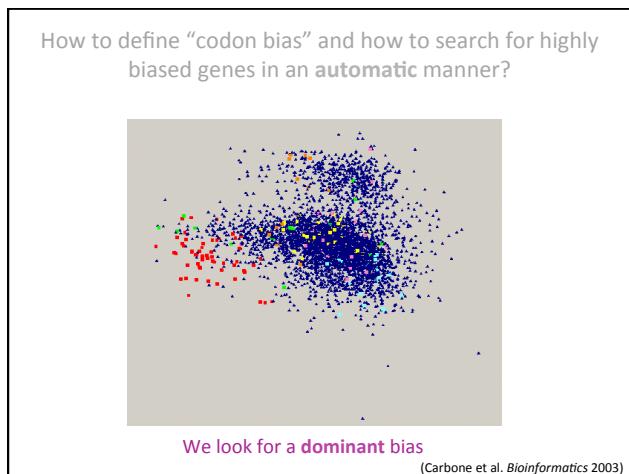
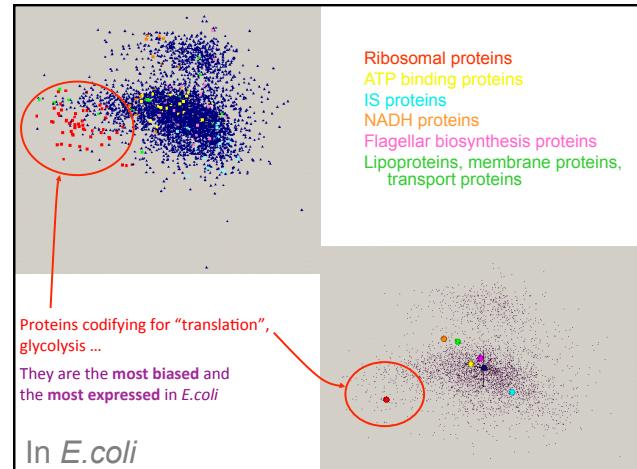
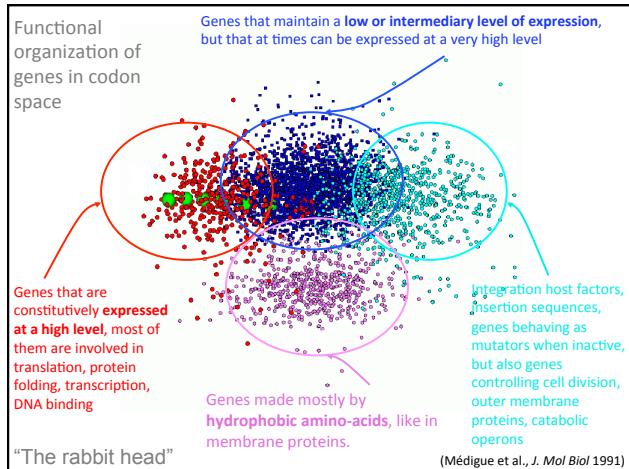
In *E.coli* and other organisms that reproduce rapidly

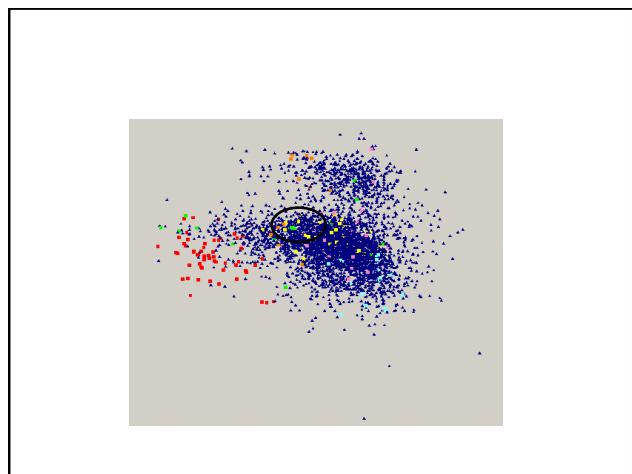
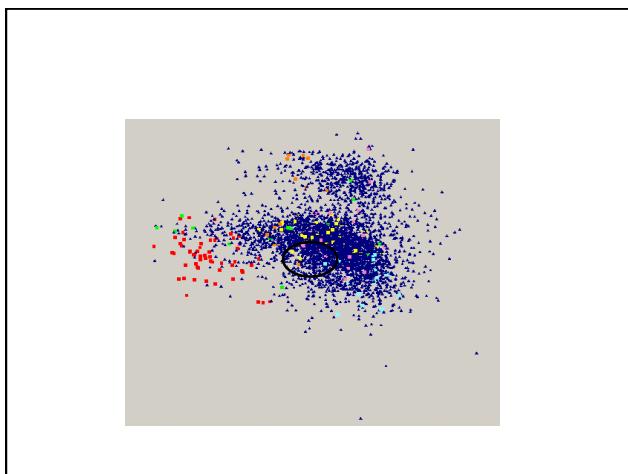
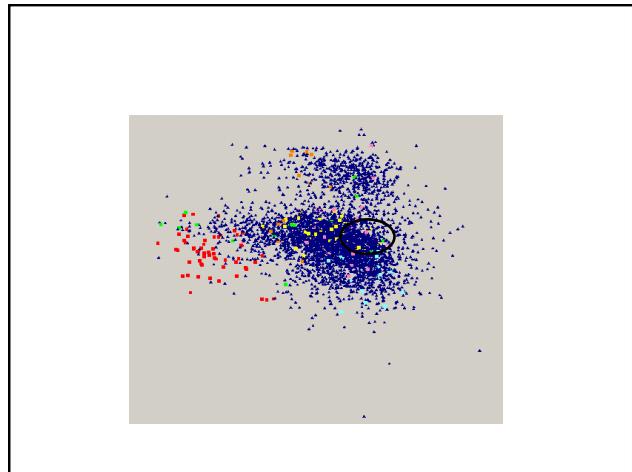
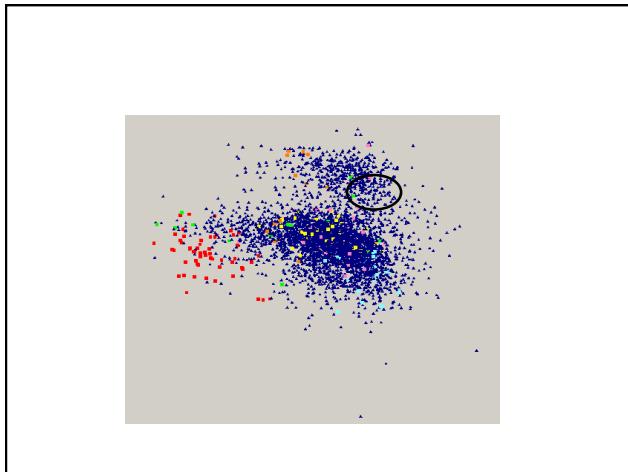
high tRNA number	correlated to	codon preference
high expression	(experimentally)	

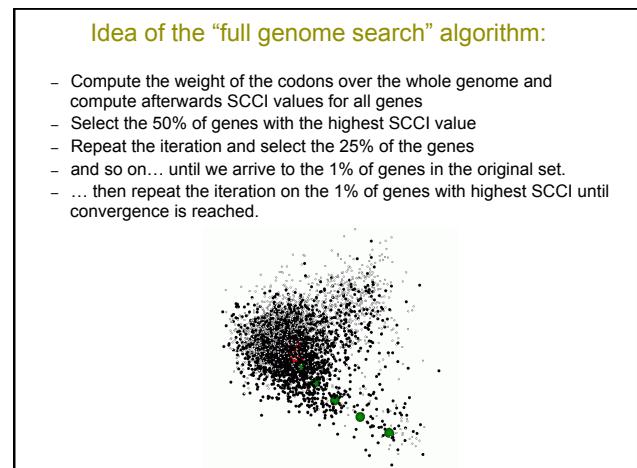
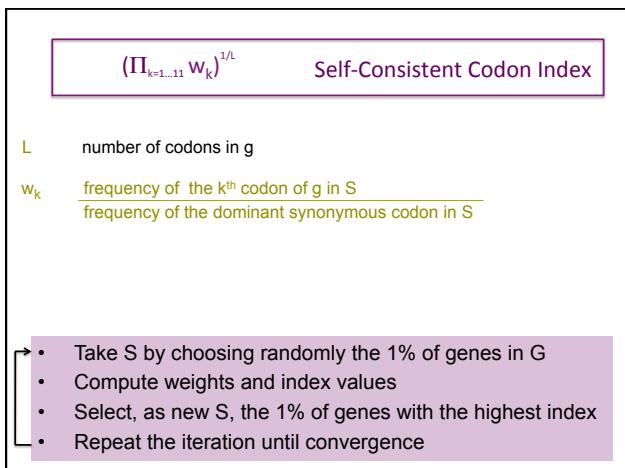
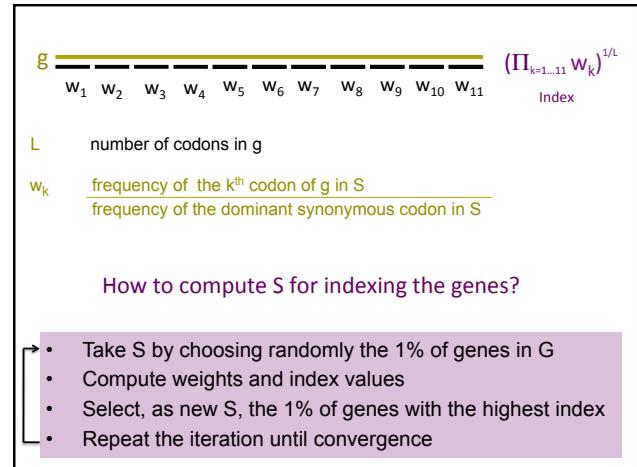
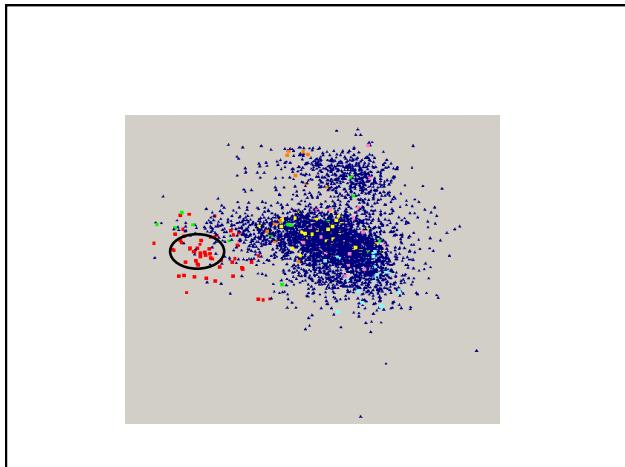
Codon preference and tRNA : Ikemura, 1985; Bennetzen and Hall, 1982; Bulmer, 1987; Gouy and Gautier, 1982.
tRNA and elongation rate : Varenne et al., 1984.
High expression and codon preference : Grantham et al., 1980; Wada et al., 1990; Sharp and Li, 1987;
Sharp et al., 1986; Medigue et al., 1991; Shields and Sharp, 1987; Sharp et al., 1988; Stenico et al., 1994.





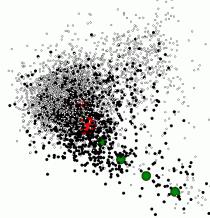






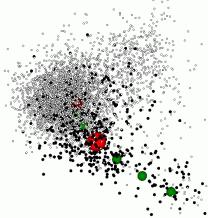
Idea of the algorithm:

- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



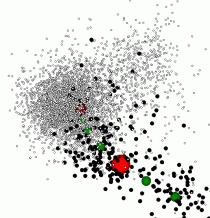
Idea of the algorithm:

- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



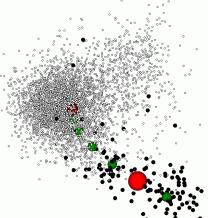
Idea of the algorithm:

- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



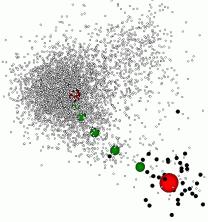
Idea of the algorithm:

- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



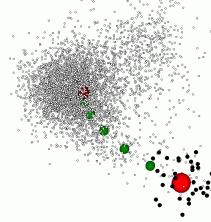
Idea of the algorithm:

- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



Idea of the algorithm:

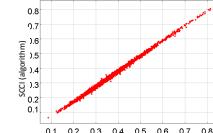
- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat the iteration on the 1% of genes with highest SCCI until convergence is reached.



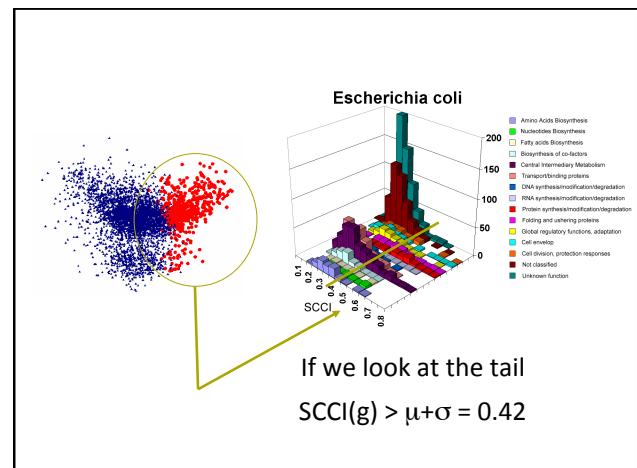
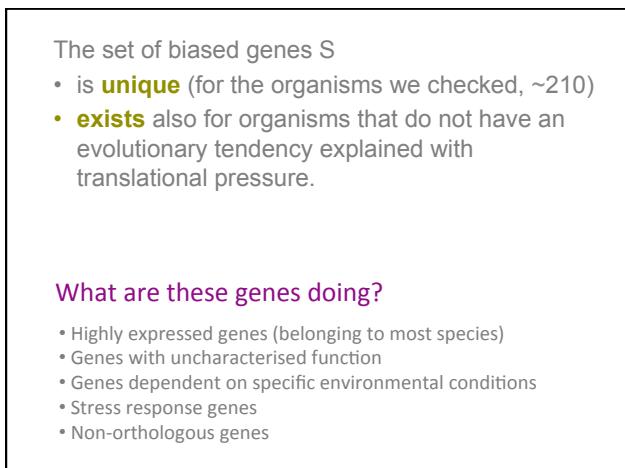
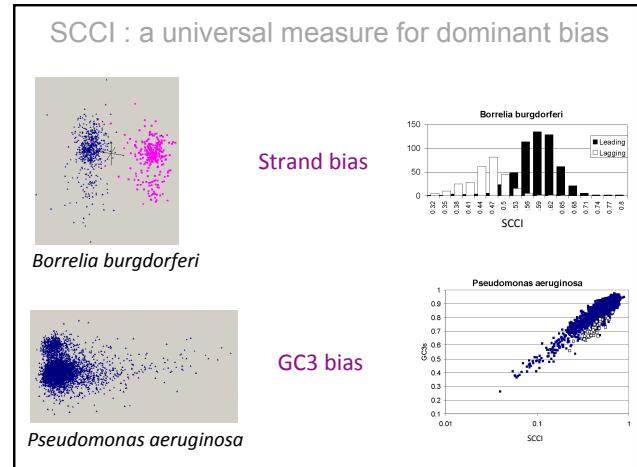
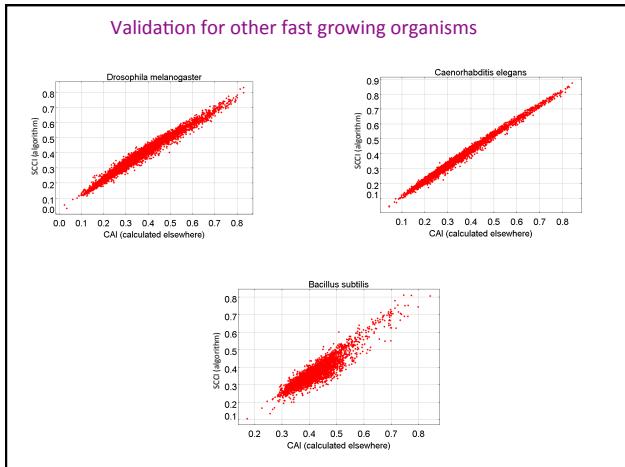
The algorithm associates to each genome a **vector of weights**, one for each codon, representing the occurrence of the codon within the most biased set of genes of the genome.

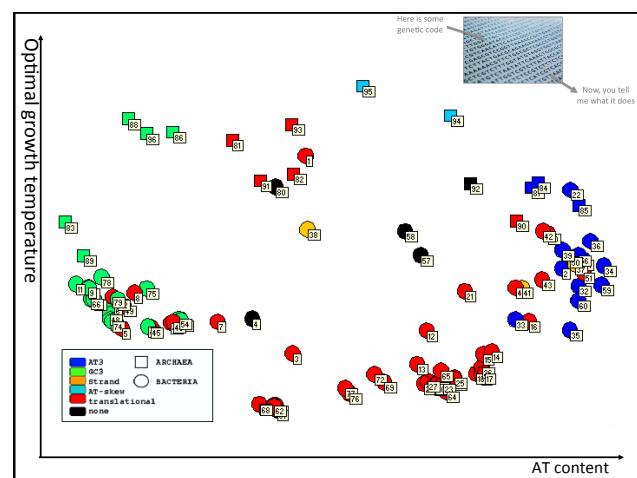
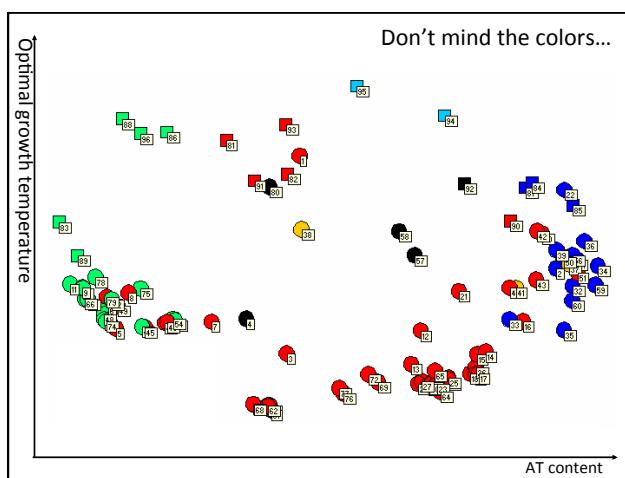
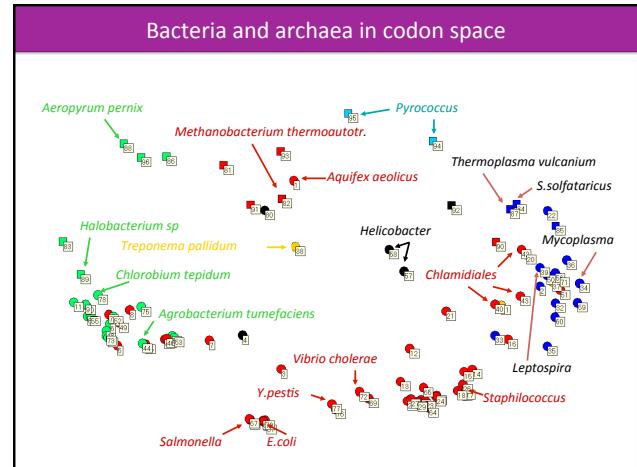
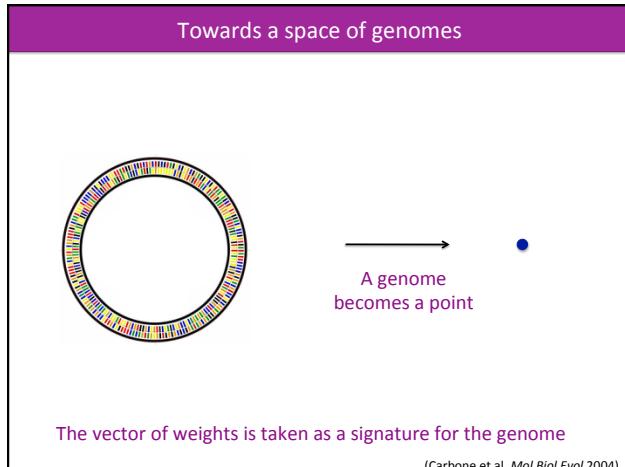
Most biased genes in *E.coli*

(*E.coli* reproduce rapidly)

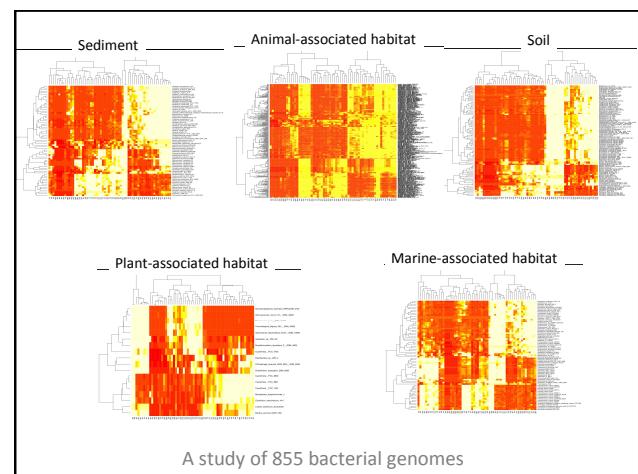
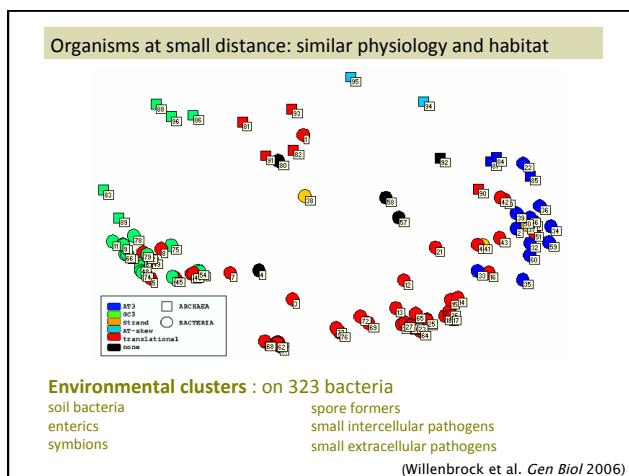
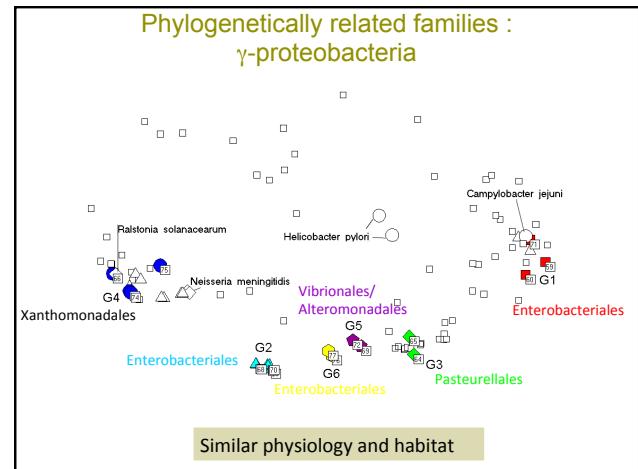


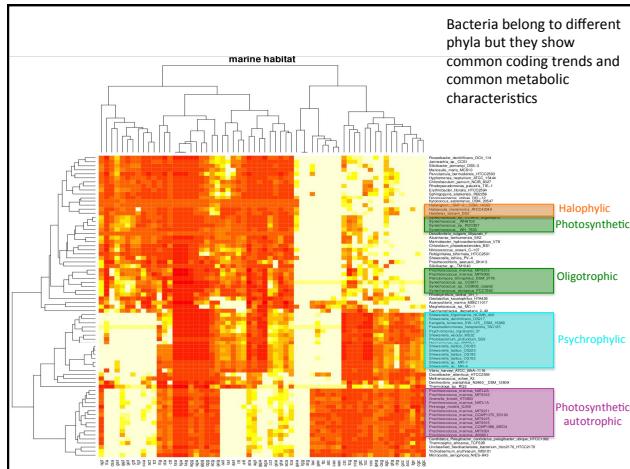
Gene	Annotation
rifB	protein chain elongation factor EF-Tu
tufB	protein chain elongation factor EF-Tu
tsf	protein chain elongation factor EF-Ts
fusA	GTP-binding protein chain elongation factor EF-G
ompA	chaperone GroEL
dsbA	heat shock protein DnaK
espA	cold shock protein 7.4
tig	trigger factor
ompA	outer membrane protein
ompB	outer membrane protein
ompC	outer membrane protein
lpp	mannin lipoprotein
pal	peptidoglycan-associated protein
yadD	putative flagellin associated protein
eneD	putative formate acetyltransferase
envA	dihydroxyacetonephosphate
tpuA	trisphosphate isomerase
tpuB	phosphotriose-isomerase
gapA	glyceraldehyde-3-phosphate dehydrogenase A
fba	fructose-biphosphate aldolase class II
pkfF	pyruvate kinase I
palP	formate:acetyl-phosphate acetyltransferase C22 subunit
sodA	alkyl hydroperoxide reductase C22 subunit
tskA	RNA transcriptase 1/2 isozyme
rncC	RNA polymerase beta prime subunit
rncD	30S ribosomal subunit protein S9
rncA	30S ribosomal subunit protein S1
rncB	30S ribosomal subunit protein S2
rncC	30S ribosomal subunit protein S3
rncD	30S ribosomal subunit protein S9
rncA	50S ribosomal subunit protein L1
rncY	50S ribosomal subunit protein L25
rncL	50S ribosomal subunit protein L9
rncL	50S ribosomal subunit protein L7/L12
rncC	50S ribosomal subunit protein L3
rncE	50S ribosomal subunit protein L31
rncB	50S ribosomal subunit protein L2
rncK	50S ribosomal subunit protein L11
rncM	50S ribosomal subunit protein L1
rncA	50S ribosomal subunit protein L27
rncD	50S ribosomal subunit protein L4, regulates expression of S10 operon





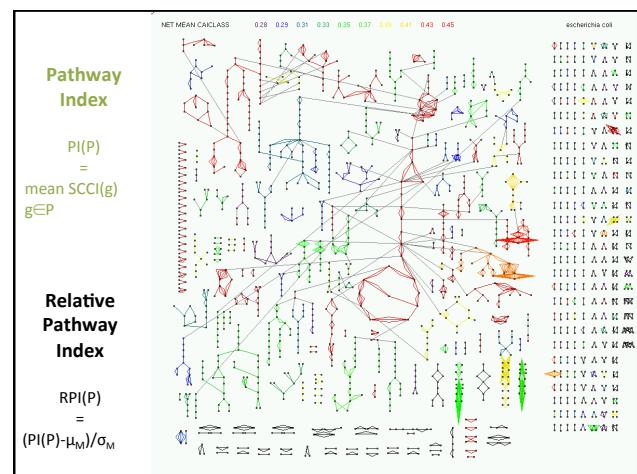
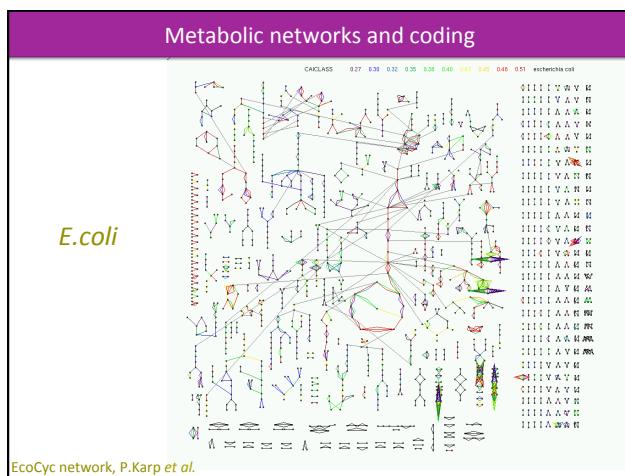
Can we exploit the geometry of the space to derive functional characteristics of groups of organisms?

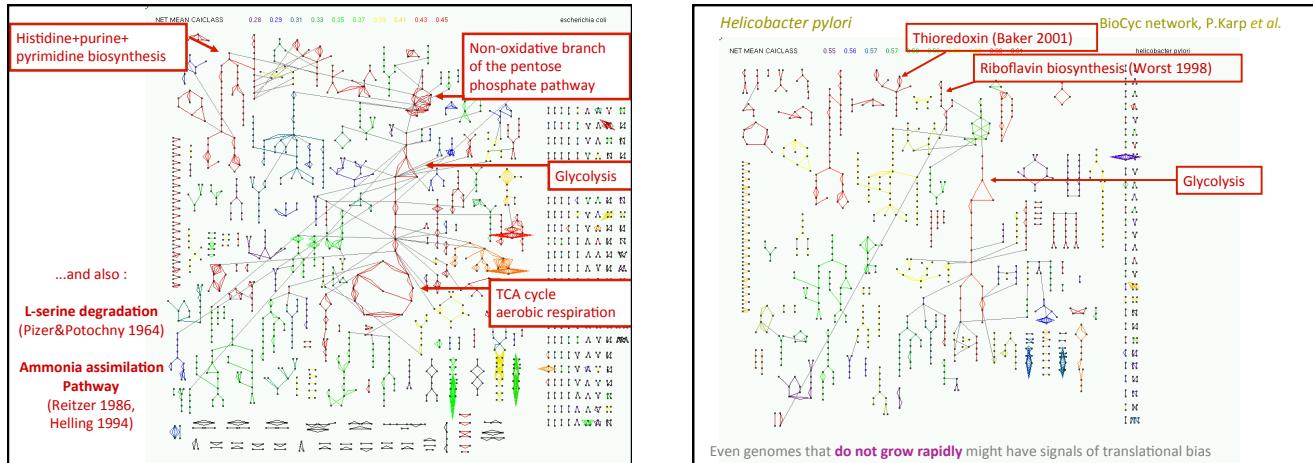




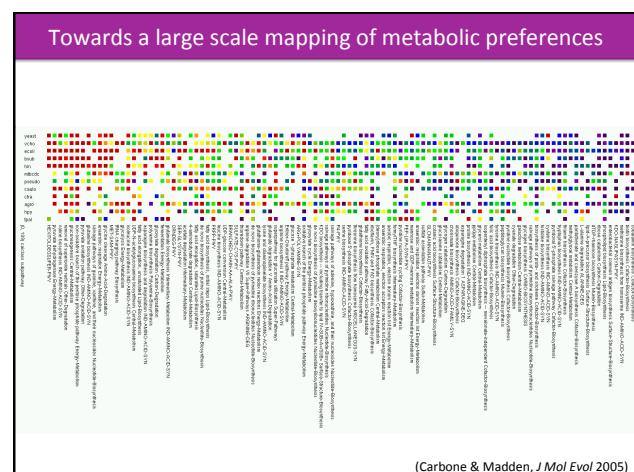
Can we use this signal to identify the most important **metabolic networks** in an organism?

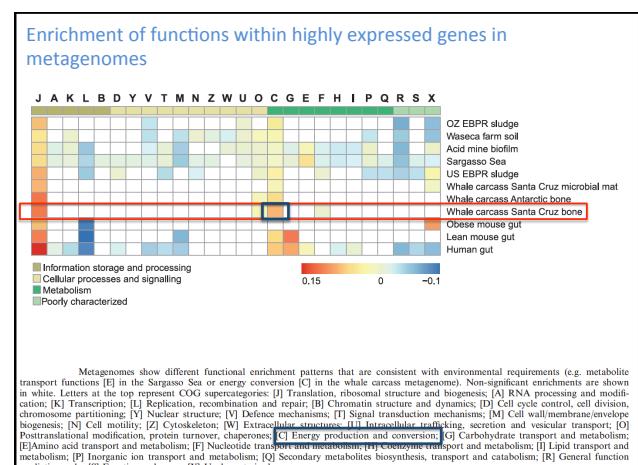
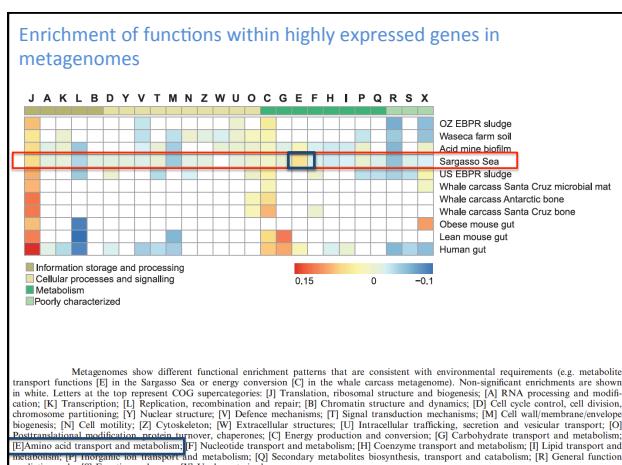
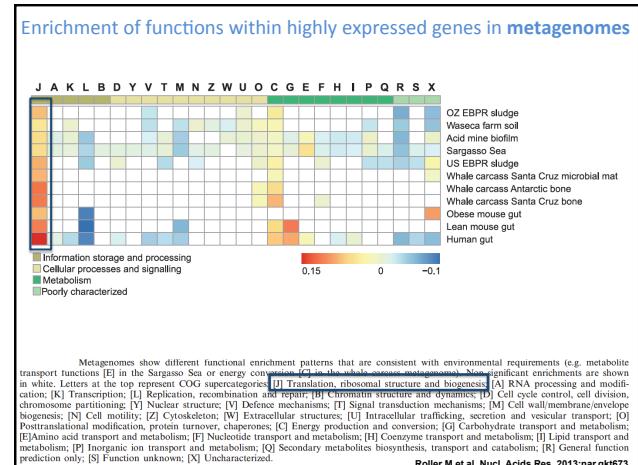
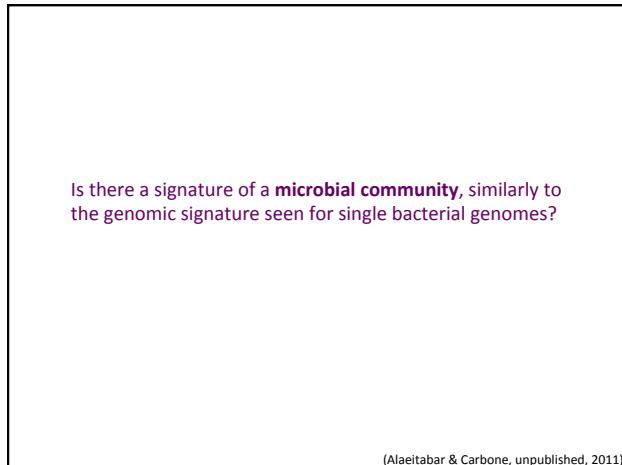
(Carbone & Madden, *J Mol Evol* 2005)





Metabolic pathways essential to <i>Mycobacterium tuberculosis</i>	
Essential to <i>M.tuberculosis</i> but not to other bacteria	
Biotin synthesis	(Norman et al. 1994)
Chorismate biosynthesis	(Parish and Stoker 2002)
Asparagine degradation	(Sassetti et al. 2003)
Pyridoxal 5'phosphate biosynthesis	(Sassetti et al. 2003)
Valine degradation	(Sassetti et al. 2003)
Leucine biosynthesis	(Sassetti et al. 2003)
ppGpp	(Primm et al. 2000)





Conclusions

- These findings suggest that microbial communities are representable by **genomic signatures**, specific to different communities, as single genomes are.
- This might be true for bacterial communities but also for eukaryotic ones. For these latter, **assembly** is much harder and our approach does not ask for large contigs for the analysis.
- The community-wide “optimization effect” is an important metagenomic feature with **predictive power**: genes with **unknown function** that are potentially important for the community can be identified

Conclusions

- Likely, we will be able to rank metabolic functions and orthologous groups of genes at the **system level**. Such effort is likely to be important to understand the **adaptation of the entire metagenome** to its particular environment.

Can we draw a metabolic map for communities?

Conclusions

Analysis at the systemic level should go parallel to an improvement of **gene annotation**.

We should work at the **annotation level**:

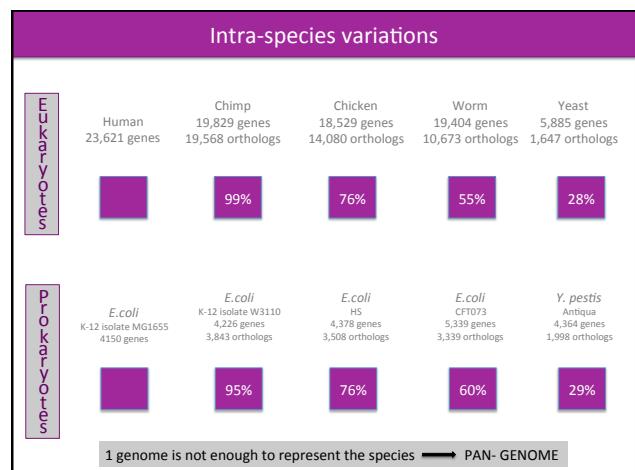
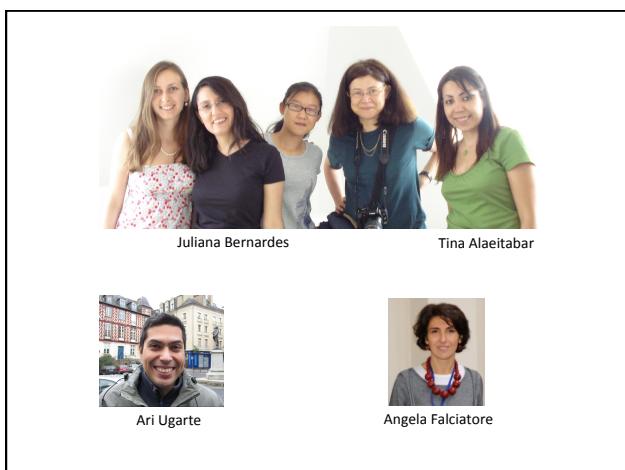
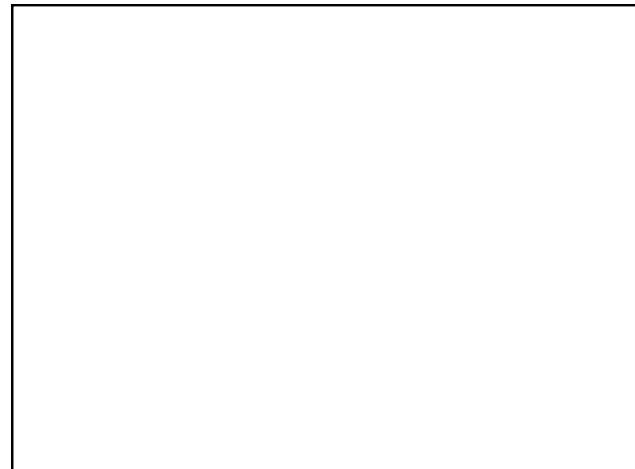
On a test realized on 51 bacterial genomes containing 159 930 CDS and 189 726 annotated domains (by Pfam) we found 28 107 new domains.

(Bernardes et al. submitted, 2013)
(Ugarte et al, in preparation, 2013)

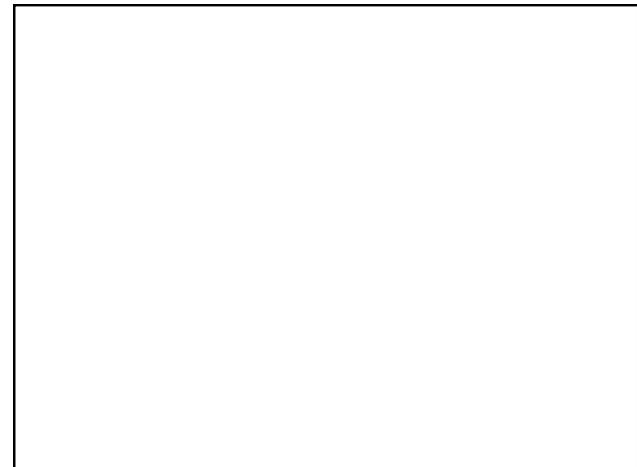
Conclusions

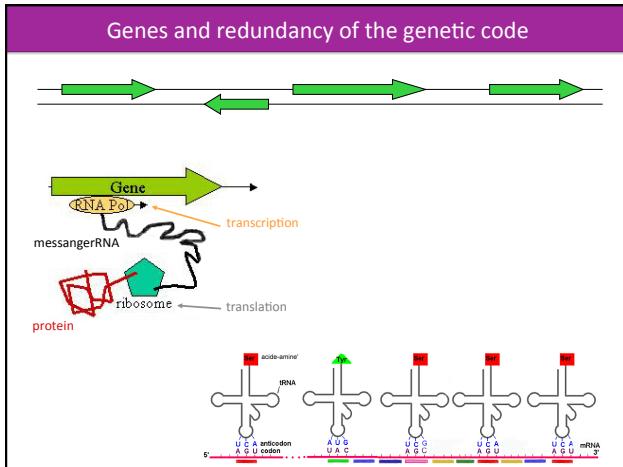
We should work at the **annotation level**:

Kingdom	Clade	Species	number of proteins	% of unknown proteins
Eukaryota	Metazoa	<i>Bryozoa matayi</i>	11472	20
		<i>Ctenophoridellus elegans</i>	26047	24
		<i>Drosophila melanogaster</i>	2012	15
		<i>Anopheles Gambiae</i>	14576	22
		<i>Ciona intestinalis</i>	4122	25
	Fungi	<i>Saccharomyces cerevisiae</i>	6607	23
		<i>Entamoeba histolytica</i>	8201	40
		<i>Paramecium tetraurelia</i>	1241	30
		<i>Oryctostilium discoloratum</i>	12646	25
		<i>Diplomonadida</i>	5012	49
Bacteria	Proteobacteria	<i>Giardia intestinalis</i>	24822	51
		<i>Giardia lamblia</i>	21000	47
		<i>Cryptophyta</i>	1000	20
		<i>Stromomonoglypha</i>	59681	43
		<i>Parabasalia</i>	4720	33
	Archaea	<i>Elasmicronema viride</i>	1593	32
		<i>Elasmicronema minutum</i>	1548	35
		<i>Elasmicronema curvum</i>	2450	65
		<i>Firmicutes</i>	2278	27
		<i>Streptococcus pneumoniae</i>	1990	23
Euryarchaeota	<i>Mycobacterium tuberculosis</i> T46	4134	42	
	<i>Cyanobacteria</i>	5356	50	
	<i>Leptothrix sp. str. methanilis</i>	2710	22	
	<i>Methanococcus maripaludis</i>	1807	31	
	<i>Halobacterium salinarum</i>	2749	47	
	<i>Korarchaeota</i>	1612	25	
	<i>Thaumarchaeota</i>	2017	48	
	<i>Gemmatarchaeota</i>	2869	38	



Percentage of proteins with unknown function in different genomes				
Kingdom	Clade	Species	number of proteins	% of unknown proteins
Eukaryota	Metazoa	<i>Brugia malayi</i>	11472	20
		<i>Caenorhabditis elegans</i>	26047	24
		<i>Drosophila melanogaster</i>	27752	19
		<i>Anopheles Gambiae</i>	14576	22
	Fungi	<i>Ciona intestinalis</i>	4122	25
		<i>Saccharomyces cerevisiae</i>	6607	23
	Amoebozoa	<i>Entamoeba histolytica</i>	8201	44
		<i>Plasmodium falciparum</i>	5491	30
	Diplomonadida	<i>Dictyostelium discoideum</i>	12646	25
	Cryptophyta	<i>Giardia intestinalis</i>	5012	49
Bacteria	Cryptophyta	<i>Giardia theta</i>	24822	51
	Stramenopiles	<i>Bgelovella natans</i>	21000	47
	Proteobacteria	<i>Phaeodactylum tricornutum</i>	10408	20
	Fusimicrobia	<i>Trichomonas vaginalis</i>	59681	43
		<i>Helicobacter pylori</i>	4720	33
	Firmicutes	<i>Elusimicrobium minutum</i>	1593	32
	Actinobacteria	<i>Staphylococcus aureus</i>	1548	35
		<i>Enterococcus faecalis</i>	2620	62
	Cyanobacteria	<i>Streptococcus pneumoniae</i>	3278	27
	Euryarchaeota	<i>Mycobacterium tuberculosis Td6</i>	1990	23
Archaea	Euryarchaeota	<i>Mycobacterium tuberculosis Td6</i>	4134	42
		<i>Microcystis aeruginosa</i>	5356	50
		<i>Methanobrevibacter smithii</i>	1710	25
	Korarchaeota	<i>Methanococcus maripaludis</i>	1807	31
		<i>Halobacterium salinarum</i>	2749	47
		<i>Candidatus Korarchaeum cryptofilum</i>	1612	25
Thaumarchaeota	Crenarchaeota	<i>Cenarchaeum symbiosum A</i>	2017	48
		<i>Pyrobaculum oguniense TE7</i>	2869	38





$$g = w_1 \ w_2 \ w_3 \ w_4 \ w_5 \ w_6 \ w_7 \ w_8 \ w_9 \ w_{10} \ w_{11} \quad (\prod_{k=1 \dots 11} w_k)^{1/11}$$

$$\text{CAI}(g) = (\prod_{k=1 \dots L} w_k)^{1/L} \quad (\text{Sharp \& Li, 1987})$$

Codon Adaptation Index

L number of codons in g

w_k $\frac{\text{frequency of the } k^{\text{th}} \text{ codon of } g \text{ in } S}{\text{frequency of the dominant synonymous codon in } S}$

proteines codifying for “translation”, glycolysis ...

Let S be a set of genes and g be a gene

$$\text{CAI}(g) = (\prod_{k=1 \dots L} w_k)^{1/L} \quad (\text{Sharp \& Li, 1987})$$

Codon Adaptation Index

L number of codons in g

w_k $\frac{\text{frequency of the } k^{\text{th}} \text{ codon of } g \text{ in } S}{\text{frequency of the dominant synonymous codon in } S}$

proteines codifying for “translation”, glycolysis ...

X Let S be a set of genes and g be a gene

$$\text{CAI}(g) = (\prod_{k=1 \dots L} w_k)^{1/L} \quad (\text{Sharp \& Li, 1987})$$

Codon Adaptation Index

L number of codons in g

w_k $\frac{\text{frequency of the } k^{\text{th}} \text{ codon of } g \text{ in } S}{\text{frequency of the dominant synonymous codon in } S}$

we compute S

Let S be a set of genes and g be a gene

$$\text{SCCI}(g) = (\prod_{k=1 \dots L} w_k)^{1/L}$$

Self Consistent Codon Index

L number of codons in g

w_k frequency of the k^{th} codon of g in S
frequency of the dominant synonymous codon in S

we compute S

Let S be a set of genes and g be a gene

$$\text{SCCI}(g) = (\prod_{k=1 \dots L} w_k)^{1/L}$$

Self Consistent Codon Index

Self consistency condition

SCCI values on genes in S are maximal :
 $\text{SCCI}(G/S) \leq \text{SCCI}(S)$, G is the set of all genes

Can we classify microbial organisms starting from codon optimization?

Here is some genetic code



Now, you tell me what it does

