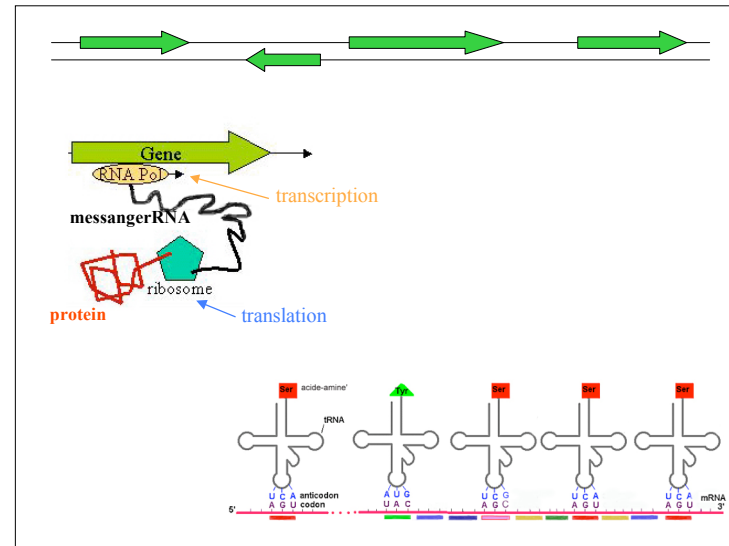


## M2 - BIM

## EVOL

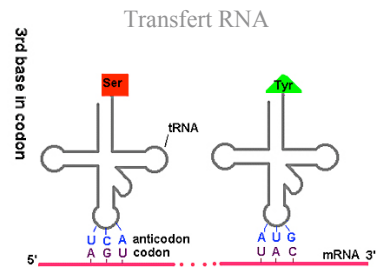
Algorithme de recherche de biais des codons dans les génomes, espaces de genes et espaces d'organismes

Alessandra Carbone  
Université Pierre et Marie Curie



## Redundancy of the genetic code

		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu	Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys Trp STOP	U C A G
	C	Leu Leu Leu	Pro Pro Pro	His His Gln	Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr	Asn Asn Lys	Ser Ser Arg Arg	U C A G
G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G	



## Bias on codon usage & preferred codons

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; Médigue *et al.*, 1991; Shields and Sharp, 1987; Sharp *et al.*, 1988; Stenico *et al.*, 1994.

**Visualisation** of genes in a genome  
and identification of their codon bias



$$g = [x_{1,g} \ x_{2,g} \ \dots \ x_{64,g}] \quad x_{i,g} \text{ relative frequency of codon } i \text{ in } g$$

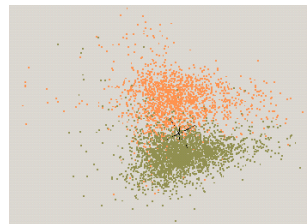
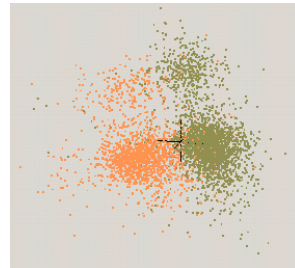
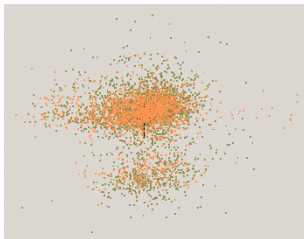
Vector normalisation:

$$(x_{i,g} - \bar{x}_i) / \sigma_i \quad \begin{array}{l} \bar{x}_i \text{ mean of frequencies } x_{i,g} \\ \sigma_i \text{ standard deviation of } x_{i,g} \end{array}$$

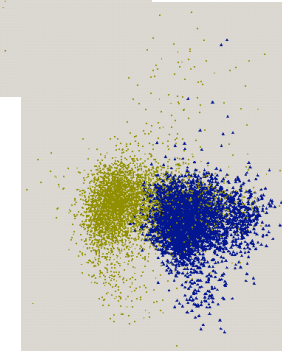
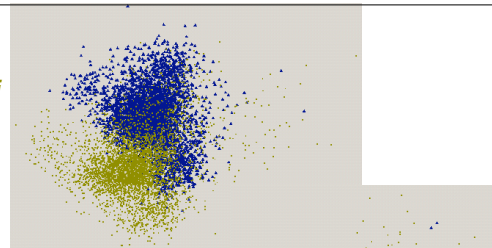
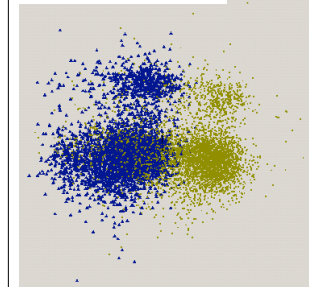
Normalized vectors and PCA are used to “see”

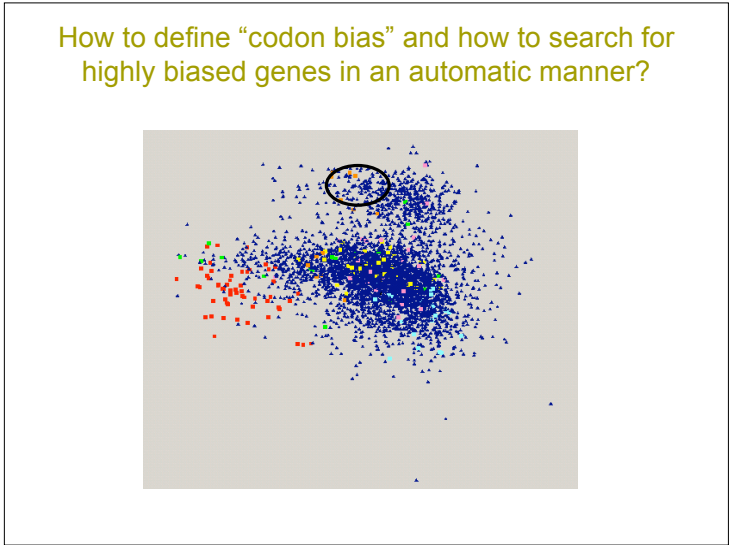
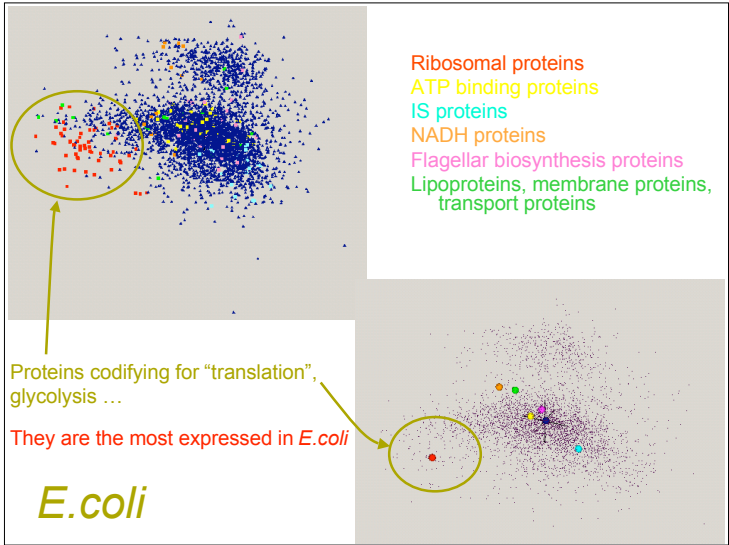
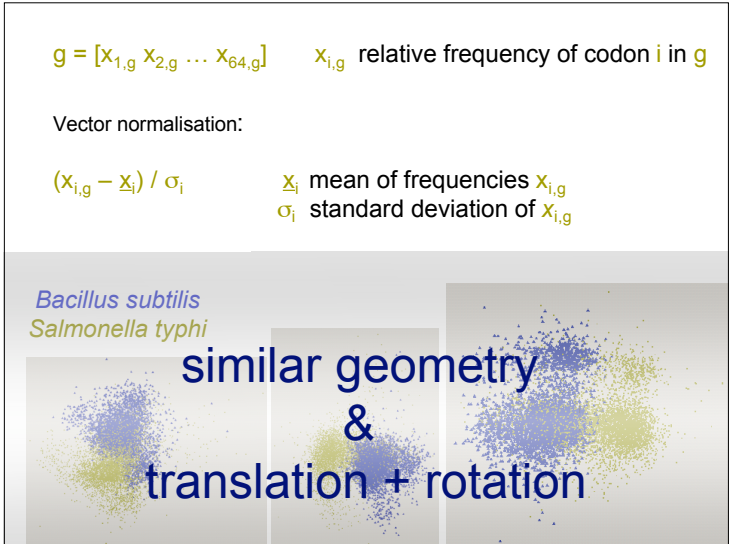
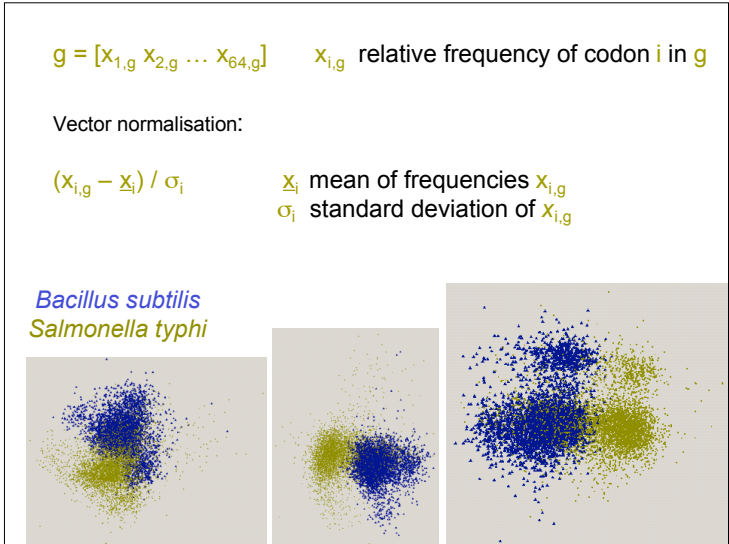
- organisms in codon space
- genes and functions

*Haemophilus influenzae*  
*Staphylococcus aureus*

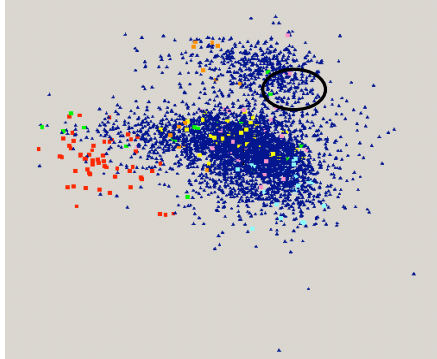


*Bacillus subtilis*  
*Salmonella typhi*

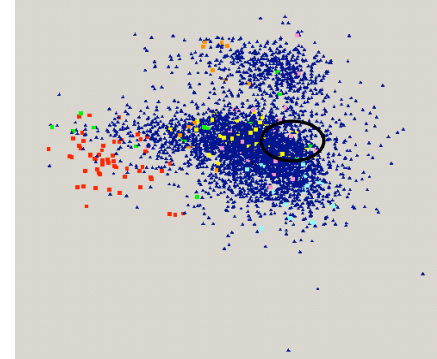




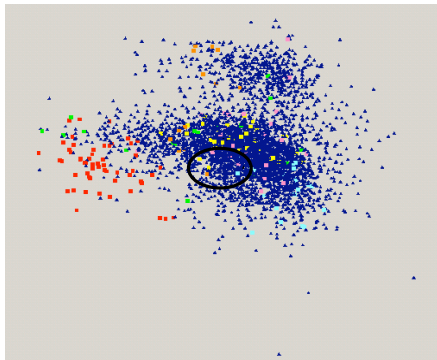
How to define "codon bias" and how to search for highly biased genes in an automatic manner?



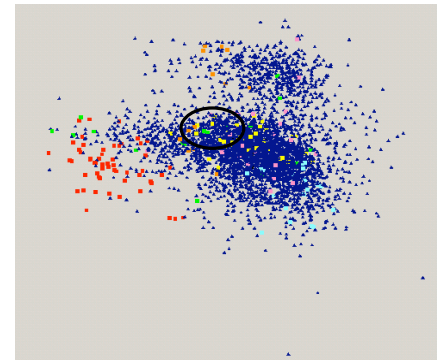
How to define "codon bias" and how to search for highly biased genes in an automatic manner?



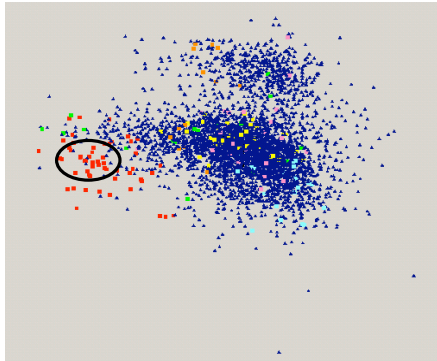
How to define "codon bias" and how to search for highly biased genes in an automatic manner?



How to define "codon bias" and how to search for highly biased genes in an automatic manner?



How to define “codon bias” and how to search for highly biased genes in an automatic manner?



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

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

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

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

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

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

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

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

Self Consistent Codon 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

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

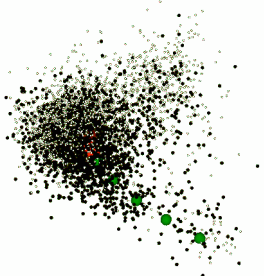
Self Consistent Codon Index

Self consistency condition

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

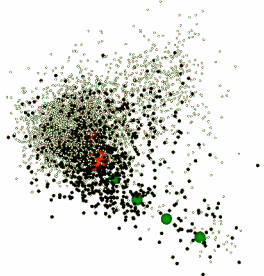
Idea of the algorithm:

- 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.



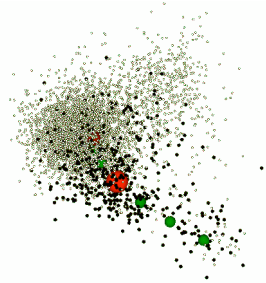
Idea of the algorithm:

- 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.



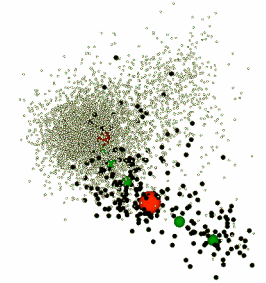
### Idea of the algorithm:

- 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.



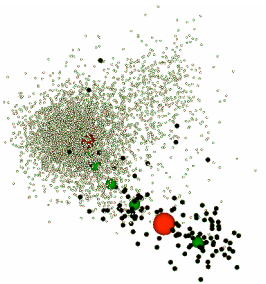
### Idea of the algorithm:

- 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.



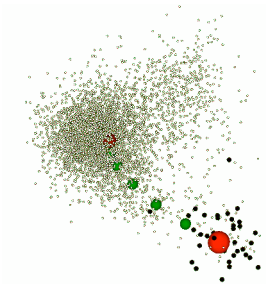
### Idea of the algorithm:

- 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:

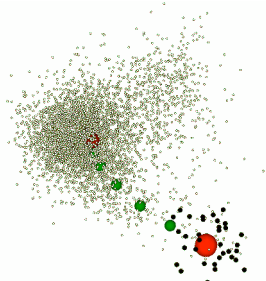
- Compute the weight of the codons over the whole genome and compute afterwards SCCI values for all genes
- Select the 50% of genes with the highest SCCI value
- Repeat the iteration and select the 25% of the genes
- and so on... until we arrive to the 1% of genes in the original set.
- ... then repeat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





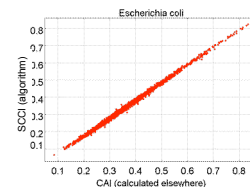
### Idea of the algorithm:

- 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.



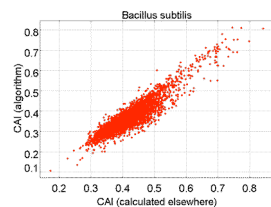
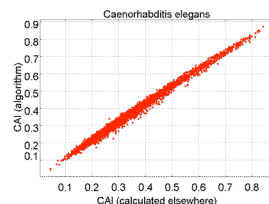
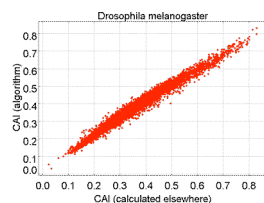
### S found by the algorithm: *E.coli*

(*E.coli* reproduce rapidly)

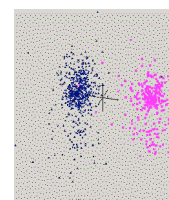


Gene	Annotation
tufA	protein chain elongation factor EF-Tu
tufB	protein chain elongation factor EF-Tu
tsf	protein chain elongation factor EF-Ts
fusA	GTP-binding protein chain elongation factor EF-G
mopA	chaperonin GroEL
dnaK	heat shock protein DnaK
ospA	cold shock protein 7.4
tig	trigger factor
ompA	outer membrane protein
ompX	outer membrane protein
ompC	outer membrane protein
lpp	murein lipoprotein
pal	peptidoglycan-associated lipoprotein
yafU	putative flagellin structural protein
yfjD	putative formate acetyltransferase
eno	diadenosine tetraphosphatase
tpiA	triosephosphate isomerase
pgk	phosphoglycerate kinase
ospA	glyceroldehydes-3-phosphate dehydrogenase A
fla	fructose-bisphosphate aldolase class II
pykF	pyruvate kinase I
pfbB	formate acetyltransferase 1
ahpC	alkyl hydroperoxide reductase C22 subunit
scsIA	superoxide dismutase ScsIA
tktA	transketolase 1/2 isozyme
rpoC	RNA polymerase beta prime subunit
rpsI	30S ribosomal subunit protein S9
rpsA	30S ribosomal subunit protein S1
rpsD	30S ribosomal subunit protein S2
rpsC	30S ribosomal subunit protein S3
rpsU	30S ribosomal subunit protein S21
rplA	50S ribosomal subunit protein L1
rplY	50S ribosomal subunit protein L25
rplI	50S ribosomal subunit protein L9
rplL	50S ribosomal subunit protein L7/L12
rplC	50S ribosomal subunit protein L3
rpmE	50S ribosomal subunit protein L31
rplB	50S ribosomal subunit protein L2
rplK	50S ribosomal subunit protein L11
rplM	50S ribosomal subunit protein A
rpmA	50S ribosomal subunit protein L27
rplD	50S ribosomal subunit protein L4, regulates expression of S10 operon

### Validation for other fast growing organisms

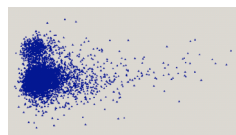
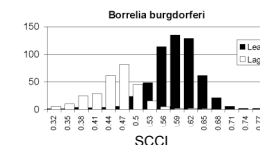


### SCCI : a universal measure



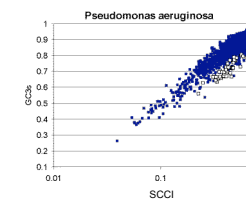
*Borrelia burgdorferi*

Strand bias



*Pseudomonas aeruginosa*

GC3 bias





The set of biased genes

- is **unique** (for the organisms we checked, ~210)
- **exists** also for organisms that do not have an evolutionary tendency explained with translational pressure.

For **any** bacteria we can compute:

- + dominant bias: strand bias, GC3, AT, ...
- + numerical criteria to determine the strength of translational bias

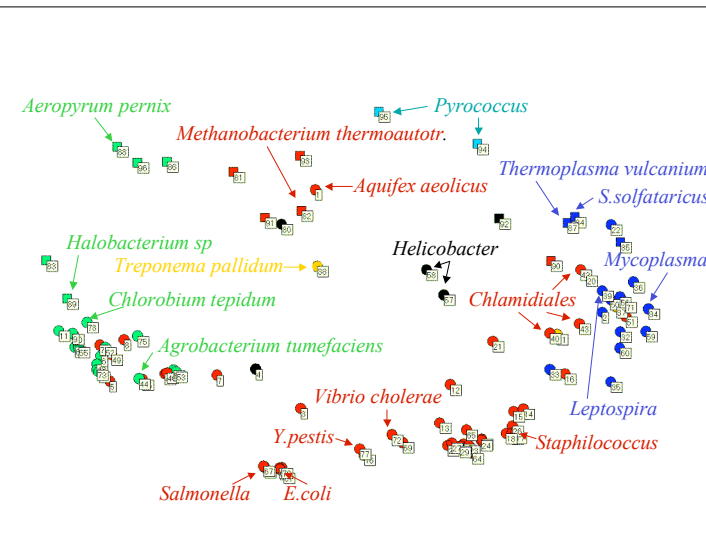
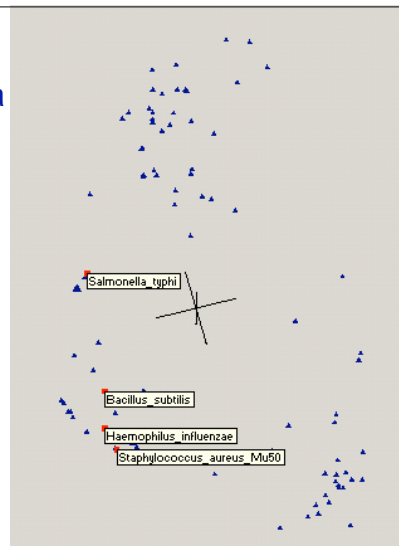
## Random version of the algorithm

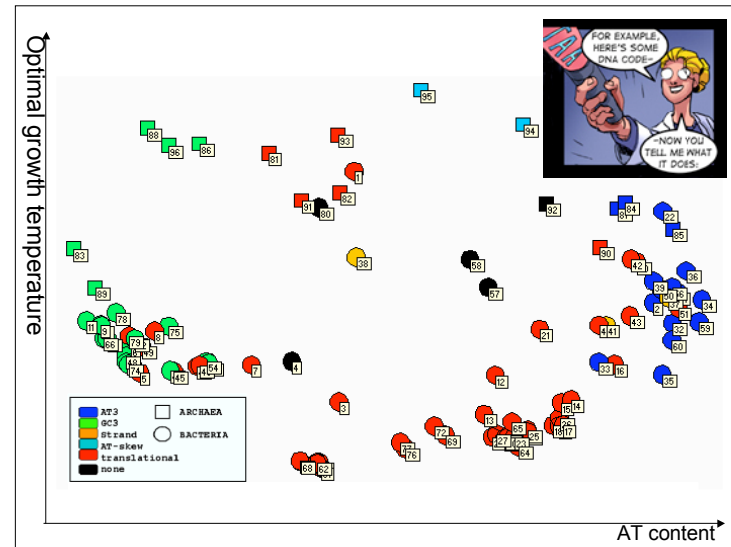
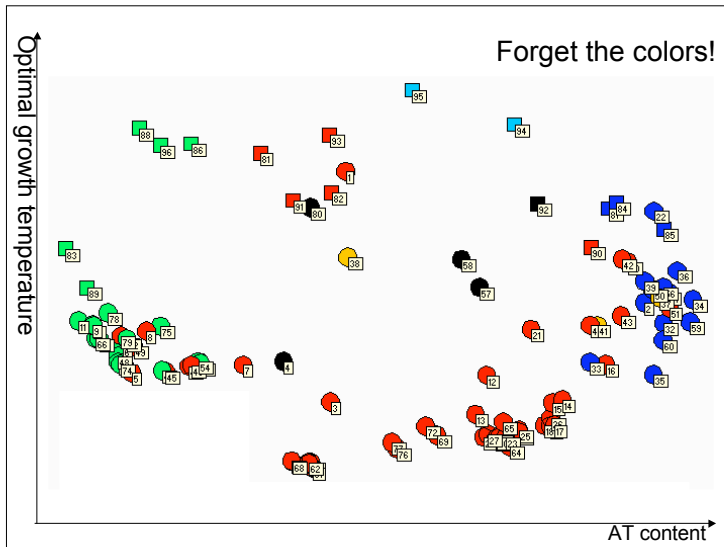
- Choose randomly the 1% of genes in G
- Compute weights and CAI values
- Select the 1% of genes with the highest CAI
- Repeat the iteration until convergence

## Bacteria and Archaea in SCCI codon space

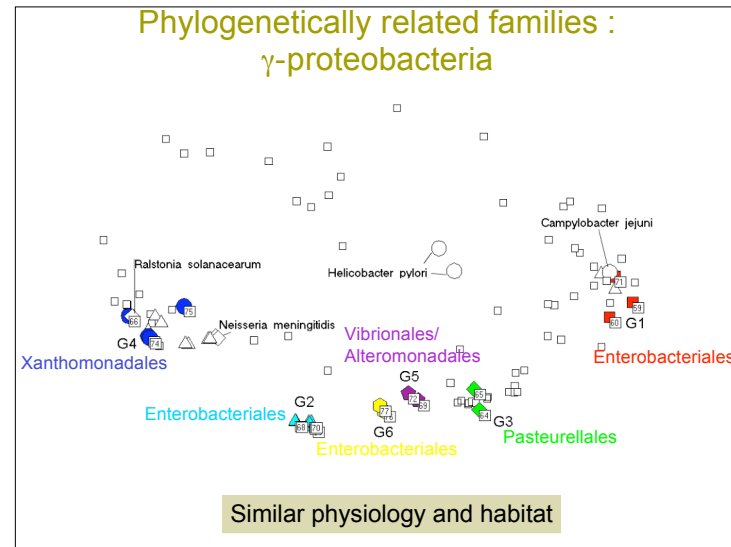
An organism is a 64-dim vector where

coordinate  
=  
SCCI codon weight

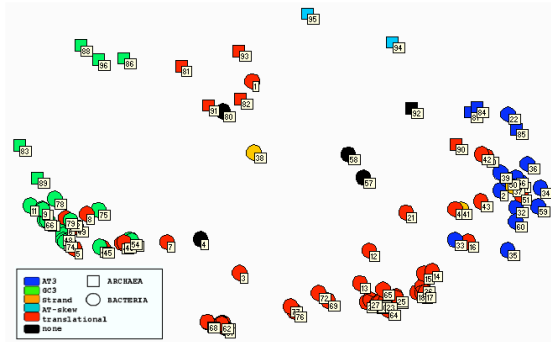




Can we exploit the geometry of the space to derive functional characteristics of groups of organisms?



Organisms at small distance: similar physiology and habitat



**Environmental clusters :**

soil bacteria  
enterics  
symbions

spore formers  
small intercellular pathogens  
small extracellular pathogens

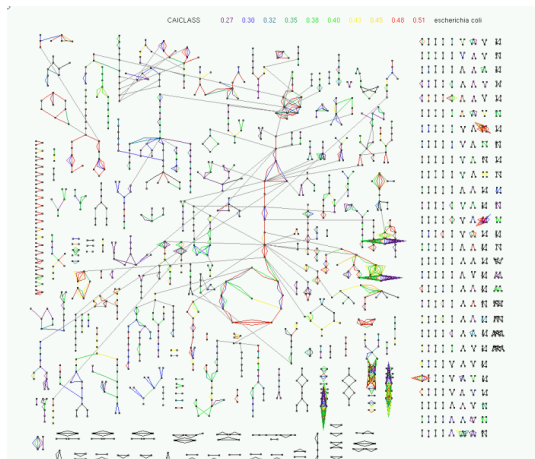
Coherence in the organisms space  
based on SCCI

Can we use this signal to deduce some  
more biological information ?

Can we determine the most important metabolic networks  
in a (translationally biased) organism ?

## Metabolic networks

*E. coli*



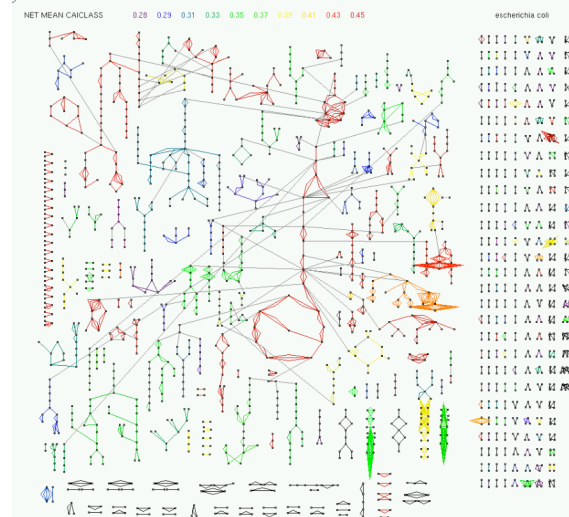
EcoCyc network, P.Karp *et al.*

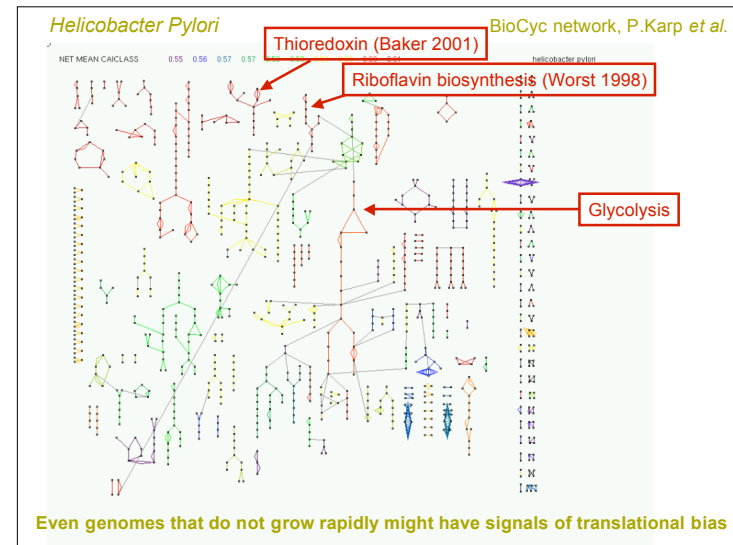
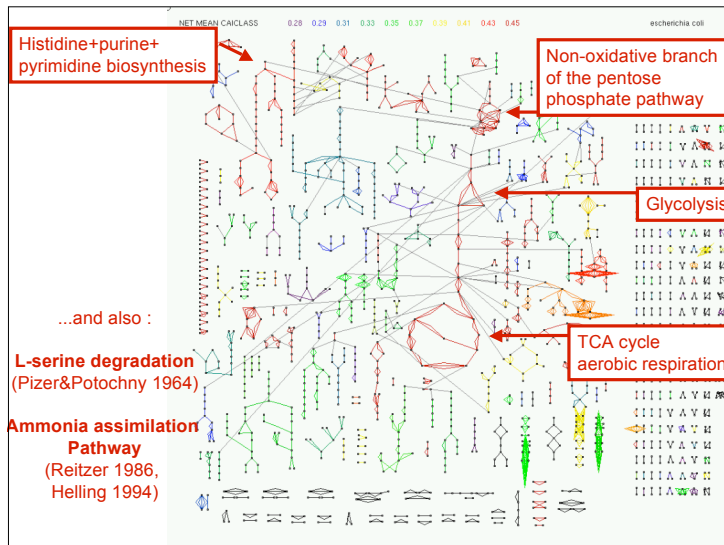
**Pathway  
Index**

$$PI(P) = \text{mean SCCI}(g)_{g \in P}$$

**Relative  
Pathway  
Index**

$$RPI(P) = \frac{PI(P) - \mu_M}{\sigma_M}$$





### Metabolic pathways essential to *Mycobacterium tuberculosis*

Essential to *M.tuberculosis* but not to other bacteria

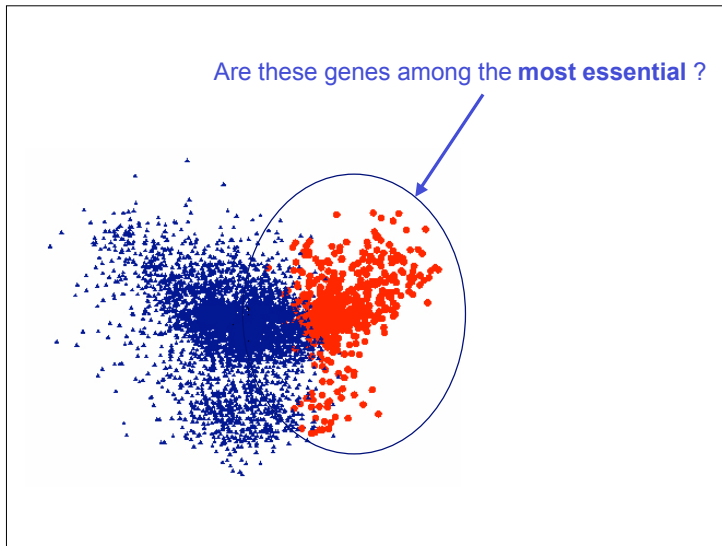
Biotin synthesis	(Norman et al. 1994)
Chorismate biosynthesis	(Parish and Stoker 2002)
Asparagine degradation	(Sasseti et al. 2003)
Pyridoxal 5'phosphate biosynthesis	(Sasseti et al. 2003)
Valine degradation	(Sasseti et al. 2003)
Leucine biosynthesis	(Sasseti et al. 2003)
ppGpp	(Primm et al. 2000)

### Coherence in the organisms space based on SCCI

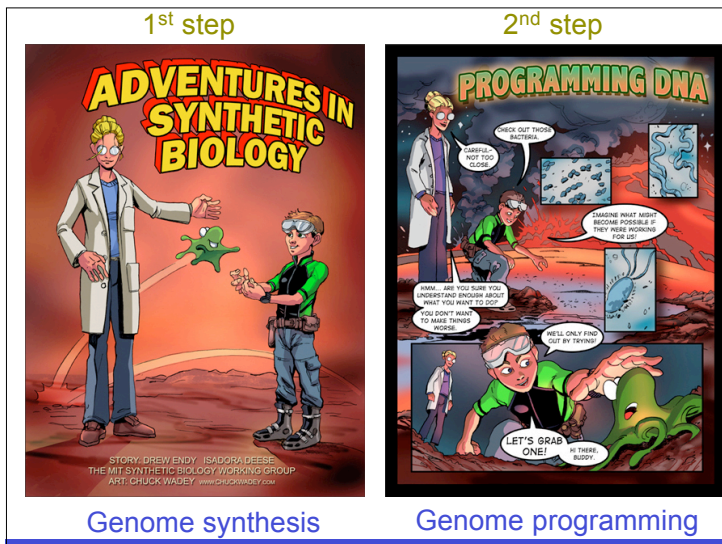
Can we use this signal to deduce some more biological information ?

We determined the most important metabolic networks in a (translationally biased) organism

Can we determine genes belonging to **minimal gene sets** ?



A parenthesis on synthetic biology



**Craig Venter, November 2002**  
**Synthesis of a bacterial genome**  
 the chromosome will be inserted in a living cell (whose genetic material has been removed) to verify if it can direct normal functional activities of the organism.

**Clyde Hutchison, 1999** (*Science* 286, 2165-2169):  
 Gene knock out (517) of *Mycoplasma genitalium* (580kb), and estimation of how many genes are necessary to life over 517: about 300 to survive.

**Eckard Wimmer, 2002** (*Science* 297, 1016-1018):  
 Synthesis of a poliovirus that infects cells! (~7500b)

**Venter, Hutchison, Smith, 2008** (*Science* 319, 1215-1220)  
**Synthesis of *Mycoplasma genitalium***

## Search for a minimal genome

### Why to do this :

Add genes to transform *Mycoplasma* in a “useful” bacteria

Remedy against environmental pollution, new industrial chemical substances production, insuline production...

## To search for a minimal set is not easy...

### Experiments : transposomal mutagenesis & RNA silencing

<i>B.subtilis</i> 300 genes/~4000 (Itaya, 1995) 248 genes/~4100 (Kobayashi, 2003)	<i>M.genitalium</i> 265 genes / 482 (Hutchison et al., 1999) 382 genes / 482 (Hutchison et al., 2006)	<i>H.influenzae</i> 670 genes/ ~1272 (Akerley et al. 2002)	<i>E.coli</i> 620 genes / 3746 (Gerdes et al. 2003) 234 genes / 2994 (Hashimoto et al. 2005)
<i>S.cerevisiae</i> 1105 genes/ 5916 (Giaever et al. 2002)	<i>C.elegans</i> 1722 genes/ 19427 (Kamath et al. 2003)	<i>S.aureus</i> 150 genes (Yi et al. 2001)	<i>S.pneumoniae</i> 110 genes (Thanassi et al. 2002)

### Comparative genomics

2 genomes 256 genes (Mushegian & Koonin 1996)	34 genomes 80 genes (Harris et al 2003)	100 genomes 60 genes (Koonin et al. 2003)	147 genomes 35 genes (Charlebois & Doolittle 2004)
---	---	---	--

## Number of genes in the minimal set depends on

### Experiments:

- life/environmental conditions of the organism during the experiment  
→ bacteria live in very good lab conditions

### Computational detection of sequence homology:

- parameters and tools to detect homologies  
→ there are genomes with more than 60% of genes with unknown function

Genes relevant to **environmental conditions** are missing

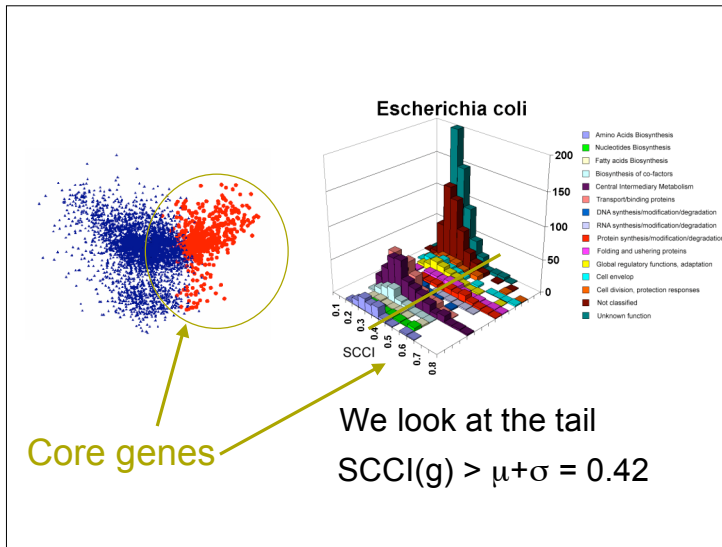
**Stress response genes** are missing

Genes with **uncharacterized functions** are missing

## Genes in minimal gene sets

$$SCCI(g) > \mu + \sigma$$

- Genes with uncharacterised function
- Genes dependent on specific environmental conditions
- Stress response genes
- Highly expressed genes (belonging to most species)
- Non-orthologous genes



Map of core genes of 27 organisms (based on 200 most biased genes)

AcI Bha Bhu Bth Bba Cdi Efa Eca Eco Hin Lpl Lia Mac Pmi Pih Pab Sty Sat Son Sfl Sag Sam Spa Sps Svy Vch Ype

**INFORMATION STORAGE AND PROCESSING**

**J Translation and associated functions**  
ribosomal proteins (including subunits) 49 05 48 34 11 49 49 41 45 46 50 53 39 51 47 49 47 46 48 44 52 46 51 52 22 53 51  
elongation factors 5 4 4 4 1 4 4 5 5 5 4 4 2 5 5 3 5 5 6 5 5 4 6 5 3 7 3  
initiation factors 2 2 1 1 3 1 1 1 1 1 1 2 2 2 3 1 1 1 1 2 2 2 2 2 2  
aminoacyl-transfer-RNA-synthetases 1 2 13 5 5 5 7 9 6 6 8 6 6 9 5 7 7 7 11 6 7 11 9 10 2  
polyribonucleotide nucleotidyltransferase 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
ribosome recycling/releasing/binding factors 1 1 1 1 1 1 2 1 1 1 2 1 1 1 2 1 1 1 1 2 1 1 1 1 1

**K Transcription**  
cold shock proteins 2 1 3 5 1 1 2 3 3 2 2 3 3 3 2 3 1 3 3 2 3 1 1 1 3 7  
RNA polymerase 3 1 4 1 1 3 3 4 5 3 5 5 2 3 4 4 3 3 3 5 6 4 4 5 5 4 4  
transcription antiterminator 1  
transcription terminator 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 1 1 1 1 1 1 1 2 1

**L DNA replication, recombination and repair**  
Bacterial nucleoid DNA-binding protein 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 2 1 1 1 2 1 1 1 2 2  
RNA helicase 1  
single-strand binding protein 1  
Recombination protein 1

**CELLULAR PROCESSES AND SIGNALING**

**D Cell division and chromosome partitioning**  
cell division proteins 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 1 1 1 1 1 1 1 1 1 1 1

**O Posttranslational modification, protein turnover, chaperones**  
chaperone proteins 3 3 3 2 3 4 3 3 2 3 3 3 2 3 5 3 3 3 3 1 2 3 2 4 5 3  
peptidyl-prolyl cis-trans isomerase 2 1 1 1 2 1 3 3 3 1 2 2 3 3 3 3 3 3 1 1 2 3 1 1 2 3  
thioesteron 1 3 1 1 1 2 1 1 1 2 2 1 1 1 1 1 1 1 1 2 1 2 1 2 1  
alkyl hydroperoxide reductase protein 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1  
trigger factor 1  
Clp protease 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1  
ribisophosphate pyrophosphokinase 1 2 1 1

**M Cell envelop biogenesis, outer membrane channel forming, conductance**  
outer membrane proteins 1 2 2 2 2  
lipoproteins 4 5 4 1 2 1 1 1 2 1 1 3 2 2 2 1 1 3 4 3 4 2 1 1 2 3 3

**N Cell mobility and secretion**  
secretory proteins 1 2 1 1 1 2 3 2 3 3 3 4 3 4 2 1 2 3 3 3 3 3 1 1 2 3  
flagella proteins 1 2 1 3  
membrane GTP-binding proteins 1 3

**P Inorganic ion transport and metabolism**  
superoxide dismutase 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1  
phosphate binding proteins 2 2 1 1 3 2 1 1 1 1 1 1 1 1 3 3 2 2 1 1 1 2 1 1 1 1 1  
metal-ion binding proteins 4 2 1 1 3 2 2 1 4 1 1 1 1 3 3 2 2 2 1 1 2 1 1 1 1 1 1

Metabolism

AcI Bha Bhu Bth Bba Cdi Efa Eca Eco Hin Lpl Lia Mac Pmi Pih Pab Sty Sat Son Sfl Sag Sam Spa Sps Svy Vch Ype

**C Energy production and conversion**  
hydrazinase 2 1  
dehydrogenases 7 8 5 10 7 7 6 5 6 5 6 5 6 14 6 5 6 6 6 4 11 5 2 3 5 5 3 6 6  
membrane-bound ATP-synthase 4 3 2 3 2 5 2 4 3 4 3 3 3 5 4 3 3 3 4 3 2 2 2 3 3 4 3 3  
succinyl-CoA synthetase 2 2 1 1 2 3 3 2 1 2 2 2 4 3 2 1 2 2 2 4 3 2 2 2 2 2  
ferredoxin 2 2 1  
other reductases 2 1 1 1 2 3 2 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
cytochrome oxidase 2 1  
cytochrome 1 1 1 2 1 1 3 3 1  
ferrous pyrophosphatase 1  
phosphate acetyltransferase 1  
ferrous acetyltransferase 1  
acetate kinase 1  
ferredoxin 1

**C Carbohydrate transport and metabolism**  
sucrase 1  
1,6-bisphosphate aldolase 1  
glyceraldehyde-3-phosphate dehydrogenase 1 1 1 2 1 1 1 1 1 1 2 1 1 1 1 1 1 1 2 2 1 2 1 1 1 1 1 1  
hexokinase 1  
transketolase 1 1 1 1 2 1 1 2 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
transporter (ABC and others) 3 2 2 2 2 2 2 5 1 8 9 1 1 1 5 6 6 4 9 13 12 9 3 2 3 2 3 2  
PTS system 1  
triphosphatase isomerase 1  
phosphoglycerate kinase/isomase 1 2 1 1 1 2 2 2 2 2 2 2 1 1 3 3 3 3 3 3 3 3 2 3 1 2 3 3  
pyruvate kinase 1  
6-phosphofruktokinase 1  
glucose-6-phosphate isomerase 1  
transaldolase 1  
6-phosphogluconate dehydrogenase (gdh) 1  
isomerase 1

**E Amino acids transport and metabolism**  
transporters 1 2 3 1 1 2 1 2 3 1 1 2 1 2 1 5 4 2 3 3 1 1 1 1 1 1 1 1  
glutamine synthetase 1  
serine hydroxymethyltransferase 2 1 2 1 1 1 1 1 2 1 2 2 1 1 4 1 2 1 1 1 1 1 1 1 1 1 1 1  
asparaginase 1  
ketol-acid reductoisomerase 1 1 1 3 2 1 2 3 2 3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
dehydrogenases 1 1 3 2 2 2 2 1  
synthase 1 1 1 1 2 1 1 2 1  
lysine 1

**F Nucleotide transport and metabolism**  
nucleoside diphosphate kinase 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 1 1 1 1 1 1 1 1 1 1 1  
other kinases 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 1 1 1 1 1 1 1 1 1 1 1 1  
nucleoside di(tri)phosphate reductase 1  
GMP reductase 1  
GMP synthase 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1  
CTP synthase 1  
adenylosuccinate synthetase 1  
transferase 1 2 1 1 1 1 2 1 1 1 1 1 1 1 1 2 2 3 1 1 1 1 1 1 1 1 1 1  
purine-nucleoside phosphorylase 2 1  
deoxyribose-phosphate aldolase 1

**H Coenzyme metabolism**  
S-adenosylmethionine synthetase 1

**I Lipid metabolism**  
acyl carrier protein 1 1 1 1 2 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
acetyl-CoA carboxylase 3 2 1 1 2 1 1 2  
beta-hydroxyacyl carrier protein synthase 1

Genes with specific metabolic functions are in the tail

**Photosynthesis metabolism : Synechocystis**  
Phycobilisome proteins  
Photosystem I and II  
Fructose-1,6-bisphosphate-aldolase

**Methan metabolism : Methanosarcina acetivorans**  
Methanol-5 hydroxybenzimidazolycobamidec methyltransferase  
Methyl coenzyme M reductase  
Methylocobamide methyltransferase isozyme M  
Carbonyl proteins  
Ack, Pta, cdhA

**Ferredoxin metabolism : Pyrococcus abyssi**  
Ferredoxin  
Ferredoxin oxidoreductase  
Keto-valine-ferredoxin oxidoreductase γ-chain

**Carbohydrates metabolism : Streptococcus mutans**  
Transport and metabolism of cellobiose, sucrose, beta-glucoside  
Metabolism of mannitol  
Genes for metabolism of glucose, fructose, mannose, maltose/maltodextrin



**Stress response genes** are in the tail

### Comparison with data from comparative genomics

Most represented functional classes of genes issued by comparing *M.genitalium* and *H.influenzae* (Mushegian and Koonin, 1996) **correspond to** most represented functional classes in functional genomic cores

Core genes expected to be essential but **missed** in (Mushegian&Koonin) :

**Transcription** : Sigma factors (rpo), termination factors (rho), chaperons (hsp90)

**Energy metabolism** : PTS proteins

**Translation** : no tRNA nucleotidyltransferase is found (consistently with comparative genomics)

### Comparison with experimental data

**difficult to make**

there are no a priori false positives nor false negatives

#### *E. coli*:

620 essential genes (Gerdes *et al.*, 2003)  
234 essential genes (Hashimoto *et al.*, 2005)

#### *E. coli* (Gerdes, 2003) :

620 essential genes over 3746 analyzed ones  
520 core genes: 62.5% are essential

Enolase (*eno*) is a core gene and it does not belong to the 620 genes claimed to be essential for *E. coli*

#### *E. coli* (Hashimoto *et al.*, 2005) :

234 essential genes, 1890 non-essential, 900 unknown behavior over 2994 analyzed ones (after genome minimization)

520 core genes : 129 essential, 278 non-essential,

53 unknown behavior,

63 deleted after minimization

Most are stress response genes

*B. subtilis* (Kobayashi et al., 2003) :

248 essential genes  
519 core genes : 126 essential

Most genes involved in Embden-Meyerhof-Parnas pathway are core genes in agreement with their unexpected essentiality for (Kobayashi et al. 2003)

## Collaborations and references

### Algorithm and microbial SCCI codon space :

- F.Képès, CNRS & génopole Evry
- A.Zinovyev, IHÉS & Institut Curie (Paris)

A. Carbone, A. Zinovyev, F. Képès, Codon adaptation index as a measure of dominating codon bias, *Bioinformatics*, **19**, 2005–2015, 2003.

A. Carbone, F. Képès, A. Zinovyev, Codon Bias Signatures, Organization of Microorganisms in Codon Space, and Lifestyle, *Molecular Biology and Evolution*, **22**, 547–561, 2004.

### Metabolic networks comparison :

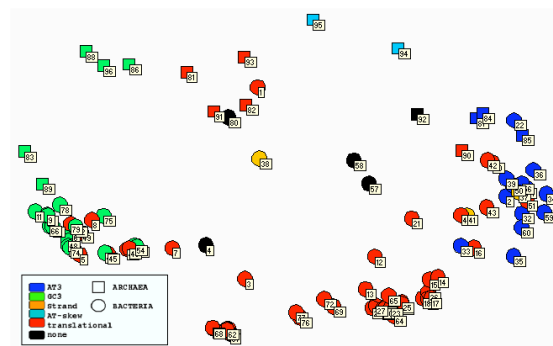
- D.Madden, IHÉS & IGI (USA)

A. Carbone, R. Madden, Insights on the Evolution of Metabolic Networks of Unicellular Translationally Biased Organisms from Transcriptomic Data and Sequence Analysis, *Journal of Molecular Evolution*, **59**, 1–25, 2005.

### Minimal gene sets :

A. Carbone, Computational prediction of genomic functional cores specific to different microbes, *Journal of Molecular Evolution*, 2006, in press.

## Bootstrapping information from translationally biased organisms



translationally biased organisms are everywhere

## Small genomes : *M.genitalium* and *B.aphidicola*

*Buchnera aphidicola* str Bp      504 coding genes  
498 genes homologous to *E.coli* genes

*Mycoplasma genitalium*      484 coding genes  
266 genes homologous to *E.coli* genes

189 genes are shared by *B.aphidicola* and *M.genitalium*

129 of these genes have high SCCI in *E.coli*

