

M2 - BIM

Comparative Genomics

Introduction

Alessandra Carbone
Université Pierre et Marie Curie
(Alessandra.Carbone@lip6.fr)

Resolution

- Low resolution
 - entry: complete genomes
 - example of event: **rearrangement**
- High resolution
 - entry: nucleotide sequences
 - example of event: **specific mutation**

3

What is Comparative Genomics?

Input: genomic sequences and their annotation

There are several different ways to compare genome at different **resolutions**

- **Complete genomes**
 - Alignments between genomes
 - Synteny (conservation of gene order)
 - Unusual regions
- **Groups of genes**
 - Gene families and species specific genes
 - Gene clustering based on functions
- **Genes: variations in gene sequences**
 - Codon usage, SNPs, indels, pseudogenes

2

Evolutionary events at the scale of complete genomes

- Rearrangements
- gene/region duplication
- gene/region loss
- Chromosome \leftrightarrow plasmid exchange
- Lateral gene transfer

4

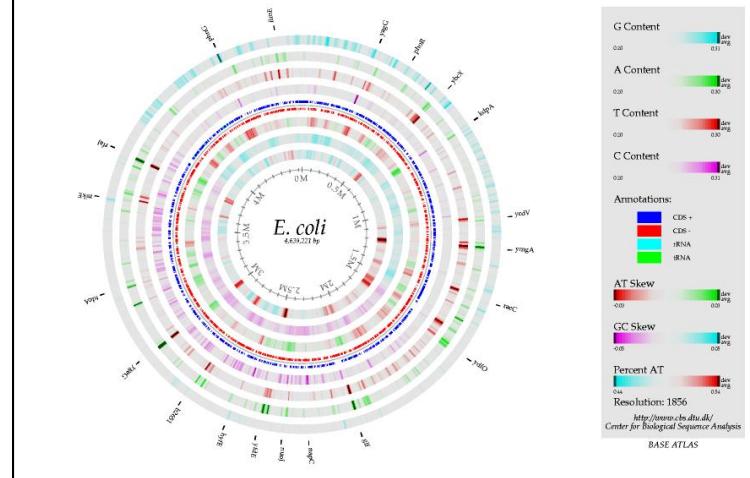
1

Genome analysis

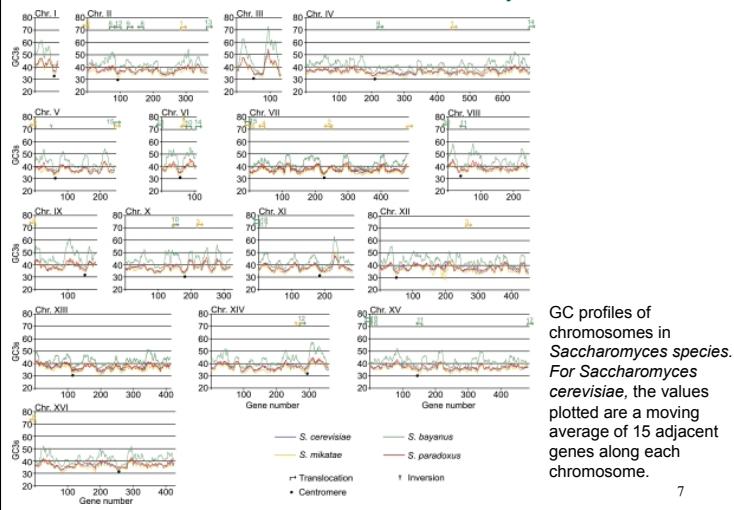
- Variation in
 - Genome size
 - GC content
 - Codon usage
 - Amino acid composition
 - Genome organisation
 - Single circular chromosomes
 - Linear chromosome + extra chromosomal elements

5

Analyse statistique à large échelle



GC content distribution in yeast chromosomes



7

Genomes as units of comparison

- Strains of the same species
- Very close species
- Very distant species
 - Orthologs list
 - Evolution of single genes
 - Evolution of organisms

8

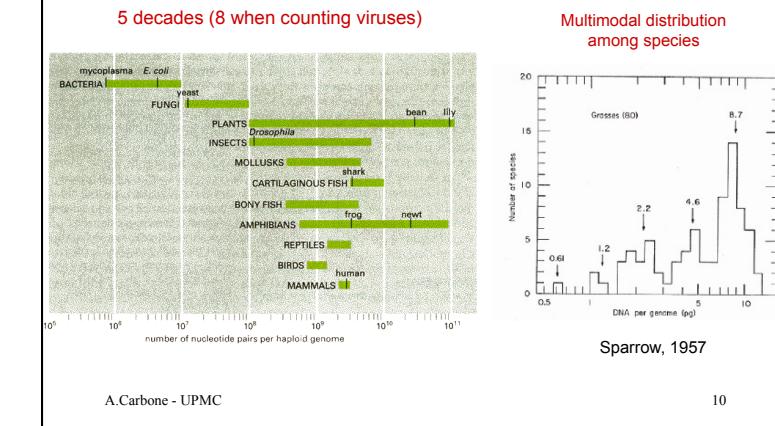
Identification of similarities/differences between genomes to understand:

- How 2 organisms evolved/co-evolved?
- Why certain bacteria produce diseases and others do not?
- How to identify therapeutic targets?
-

GC analysis paved the way to answer these questions

9

Genomes lengths



Famous genomes and their lengths

- *Haemophilus influenzae* (1.8 Mb)
 - Human pathogen, first sequenced genome (1995)
 - *Escherichia coli* (4.6 Mb)
 - Human pathogen and model organism (1997)
 - *Agrobacterium tumefaciens* (6 Mb)
 - Plant pathogen and biotechnological tool (2001)
-
- The genome of *Plasmodium falciparum*, the malaria parasite, sizes 23 Mb.
 - The human genome is around 150 times larger, the mouse genome > 100 times, and the drosophila genome > 5 times.

11

Questions on genome length:

How ancestral genomes could become longer and longer along evolution?

In what manner their genetic material differs?

B.Dujon lectures will address these questions

A. Carbone - UPMC

12

MUMmer, a tool for comparing entire genomes

Comparisons and analysis at the level of the nucleic acids or the proteins

- MUMmer (for the alignment of entire genomes)
-
- Comparative genomics of Parasites @ TIGR
- Microbial Genome Database (MDG) - Japan
- Comprehensive Microbial Resource page@TIGR
-

13

- MUMmer has been developed by Steven Salzberg's group at TIGR

- NAR (1999) 27:2369-2376
- Delcher AL, Phillippy A, Carlton J, Salzberg SL. Fast algorithms for large-scale genome alignment and comparison. *Nucleic Acids Res.* 2002 Jun 1;30(11):2478-83.
- Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL. Versatile and open software for comparing large genomes. *Genome Biol.* 2004;5(2):R12
- <http://mummer.sourceforge.net>

14

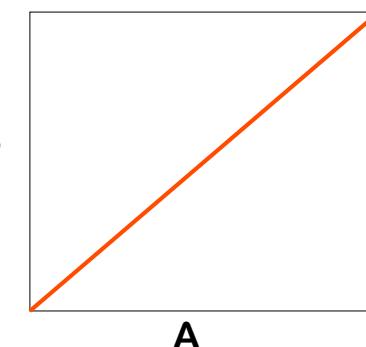
Blocks of genes as units of comparison

- Colinearity in gene order (synteny)
- Groups of conserved genes appearing with the same order in two different genomes: syntenic blocks or syntonic clusters
- Translocation: mouvement of **genomic segments** from a position to another in the genome.

15

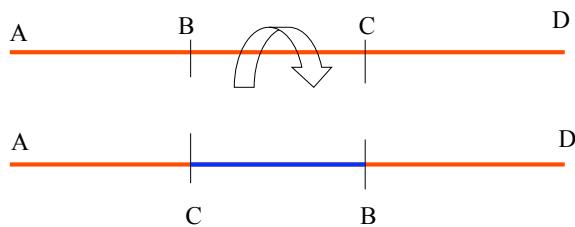
The two genomes case

If sequences were identical we would see:

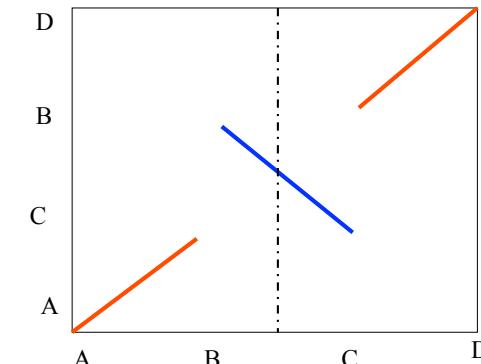


16

An inversion

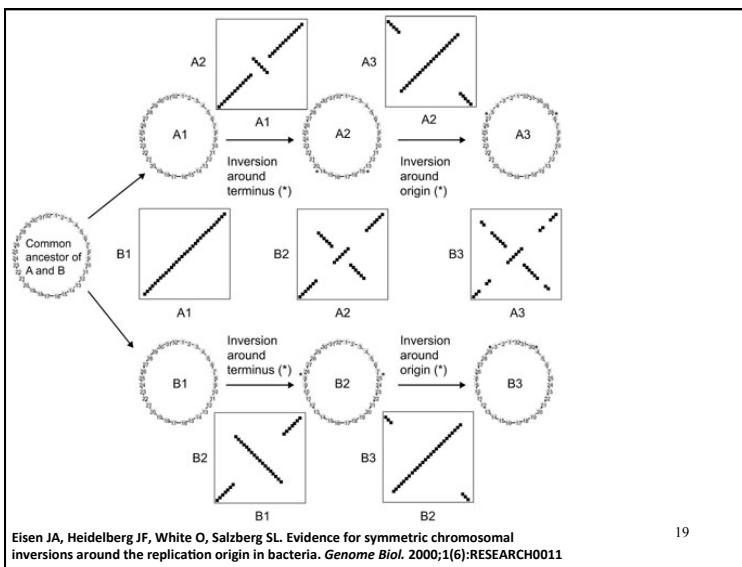


17

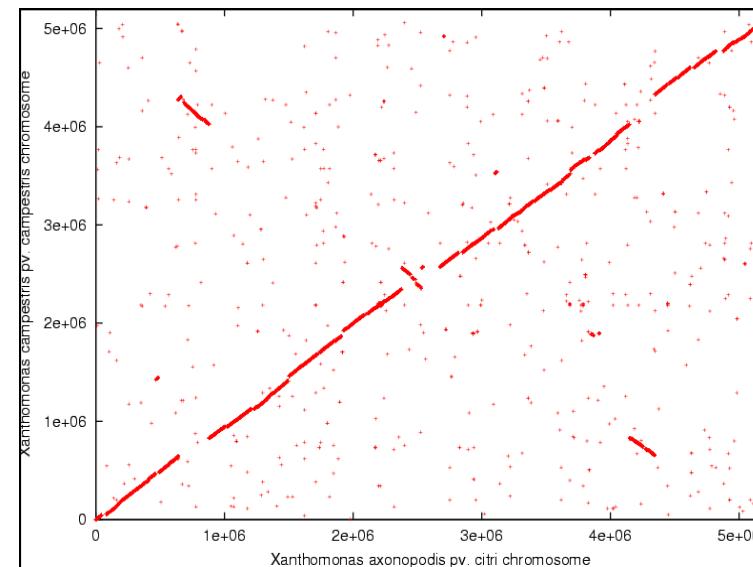


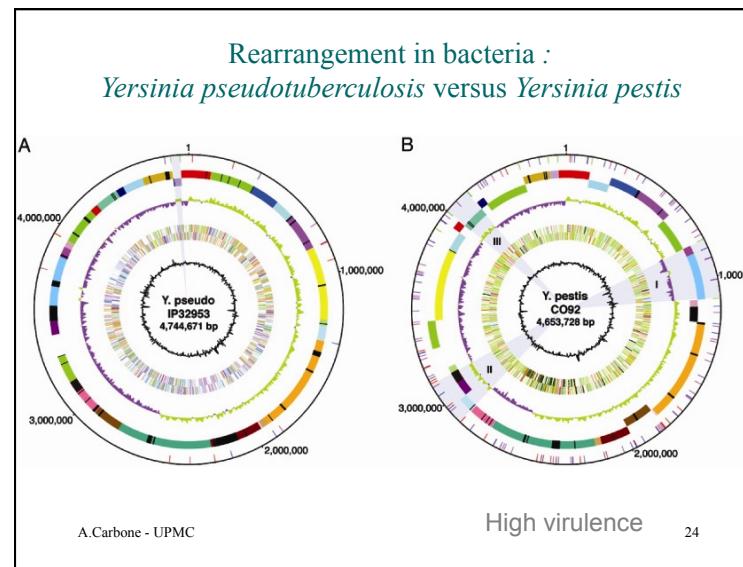
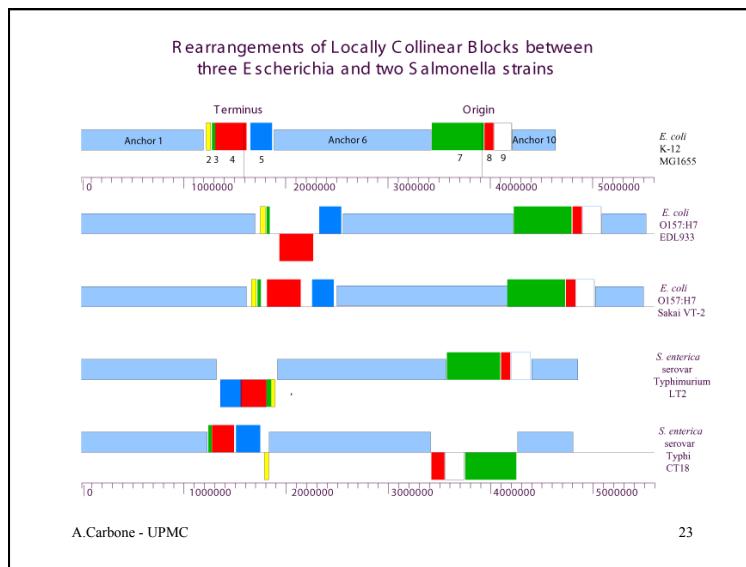
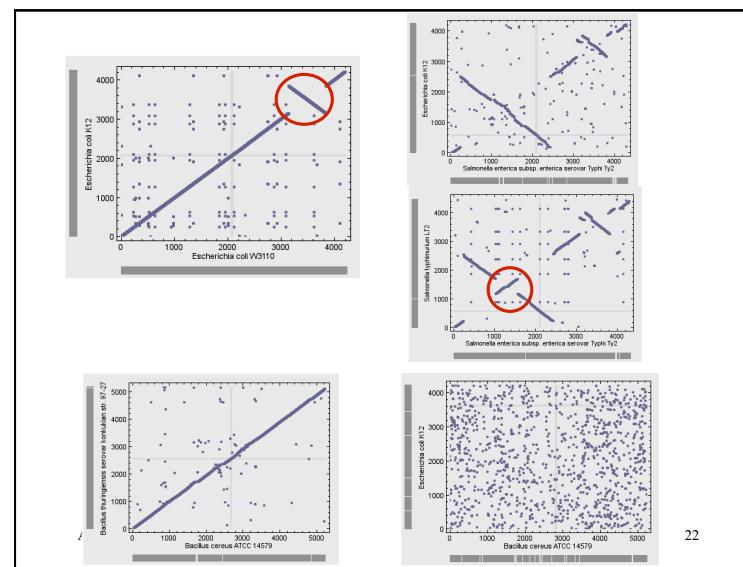
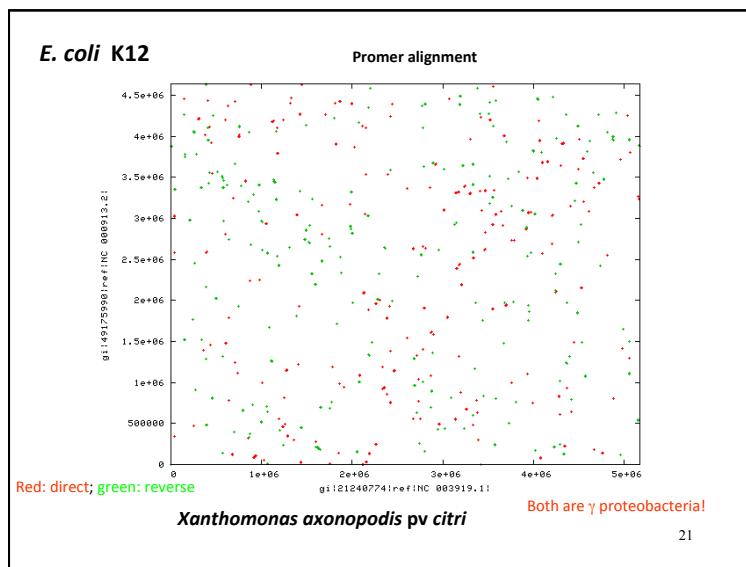
Such inversions seem to happen around the replication origin or the terminus

18



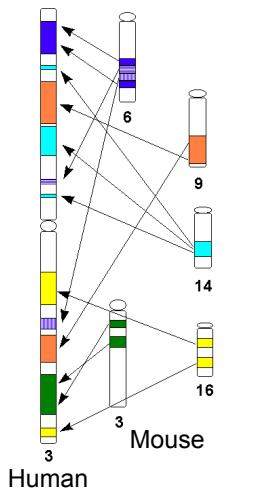
19





Human and mouse

- There is a meaningful quantity of genome rearrangement between human and mouse.
- Here we see a chart of the human chromosome 3.
- It contains sequences that are homologous to at least 5 chromosomes in the mouse.



A.Carbone - UPMC

An example : human and mouse

- The mouse has 2.1×10^9 bp vs 2.9×10^9 bp in human.
- About 95% of the genetic material is shared.
- 99% of shared genes on a total of 30,000.
- 300 homologous genes in the two species are involved into immunity, detoxification, smell and sex.

A.Carbone - UPMC

26

Genome rearrangements are rare events compared to **pointwise mutations**:

- 10 substitutions per generation of an organism
- 1 non fatal rearrangement every 5-10 millions years

The low frequency of rearrangement events allows to establish the existence of the evolutionary process. As a consequence, passing through the reconstruction of all rearrangements, we can reconstruct evolutionary hypothesis.

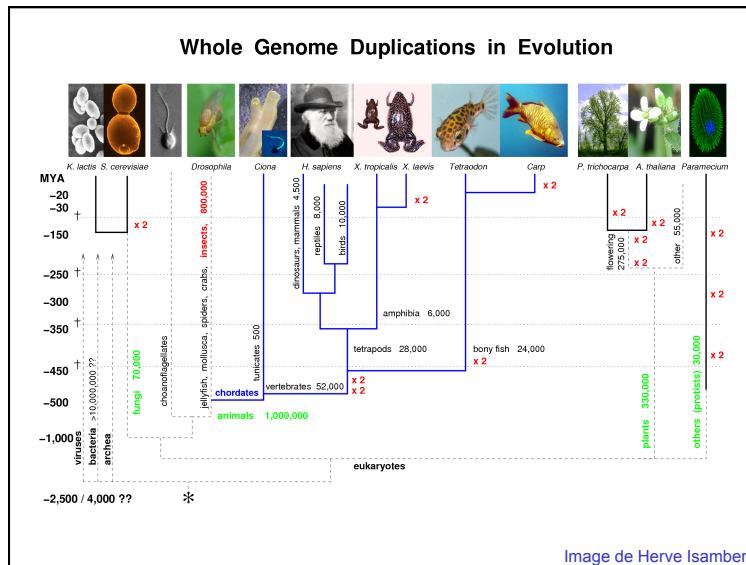
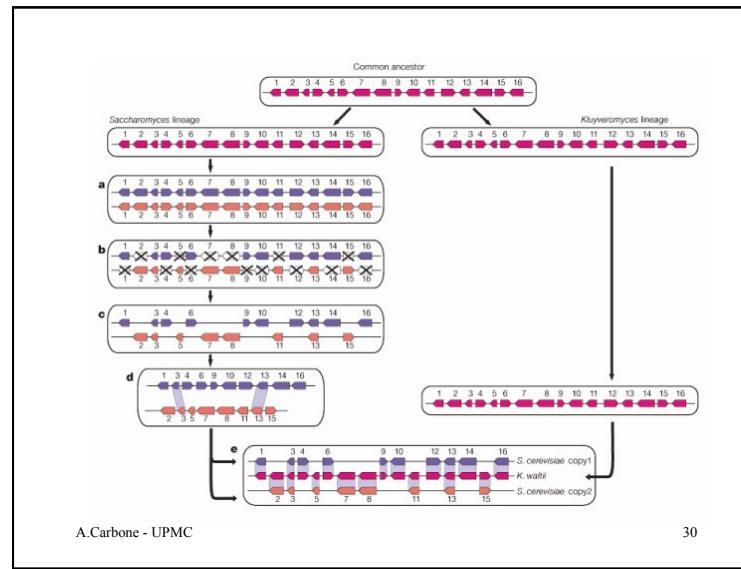
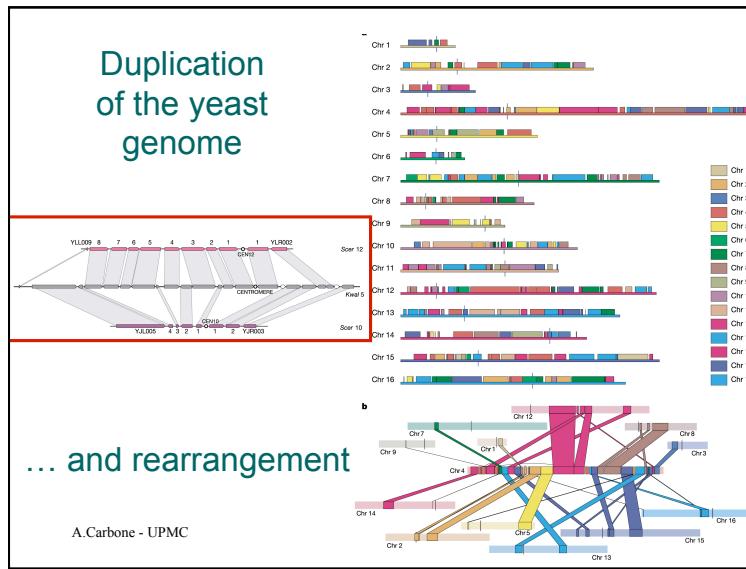
A.Carbone - UPMC

27

Search for synteny

- Considers regions in the two genomes showing important similarities in terms of
 - Sequence
 - Conservation of gene order
- Most probably it relies on the existence of a common ancestor.

28



- ## Genes as units of comparison
- Number
 - Content (sequence)
 - Localisation (position)
 - Gene order
 - Gene clustering (genes belonging to the same metabolic pathway belong to the same cluster)
- 32

How to search for genes in different species

- DNA databases are much larger than protein databases.
- Translation of a DNA sequence into a protein sequence induces an information loss.
- Protein sequences are more conserved than DNA sequences.

The translation of a DNA sequence into a protein sequence brings always better results!

A.Carbone - UPMC

33

NCBI protein–protein **BLAST**

Nucleotide Protein Translations Retrieve results for an RID

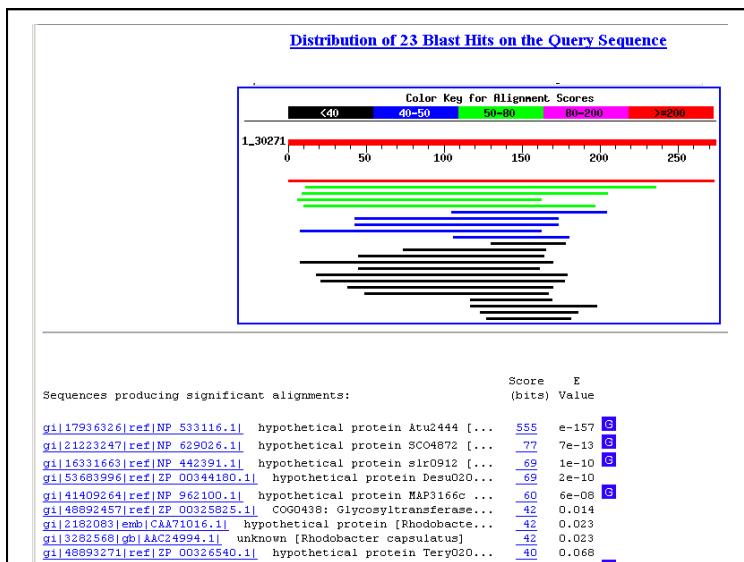
MNIEDRRFIFQELRSVEGYIDPPDALVFKA... [Search](#)

Set subsequence From: To:

Choose database nr

Do CD-Search

Now: **BLAST!** or [Reset query](#) [Reset all](#)



The notion of Reciprocal Best Hit

Take a sequence S_1 ,
BLAST S_1 and get the 1st BLAST choice S_2 (sequence with best E-value)
BLAST S_2 and get the 1st BLAST choice S_3 (sequence with best E-value)
If S_3 is S_2 then we say that S_2 and S_3 are reciprocal best hits (RBH)

Weaker notions of homology exist:

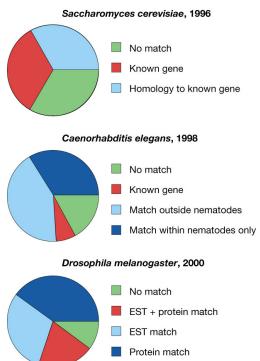
Take a sequence S_1 ,
BLAST S_1 and get the 1st BLAST choice S_2 (sequence with best E-value) if it is at
>30% sequence identity from S_1

A.Carbone - UPMC

36

Old paradigm: few genes, the same genes

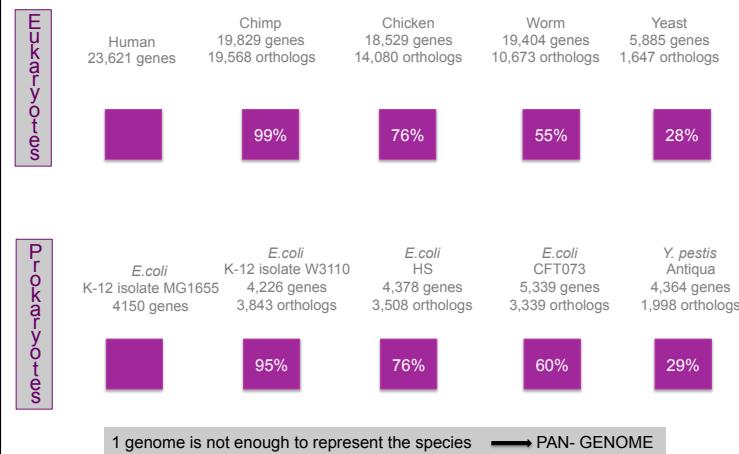
Organism	Year	Millions of bases sequenced	Total coverage (%)	Coverage of euchromatin (%)	Predicted number of genes	Number of genes per million bases sequenced
<i>Saccharomyces cerevisiae</i>	1996	12	93	100	5,800	483
<i>Ceenorhabditis elegans</i>	1998	97	99	100	19,099	197
<i>Drosophila melanogaster</i>	2000	116	64	97	13,601	117
<i>Arabidopsis thaliana</i>	2000	115	92	100	25,498	221
Human chromosome 21	2000	34	75	100	225	7
Human chromosome 22	1999	34	70	97	545	16
Human genome rough draft (public sequence)	2001	2,693	84	90	31,780	12
Human genome rough draft (Celera sequence)	2001	2,654	83	88–93	39,114	15



A.Carbone - UPMC

37

Intra-species variations



Comparison of several *E. coli* strains

Organised Genome Dynamics in the *Escherichia coli* Species
Results in Highly Diverse Adaptive Paths
PLoS Genetics, 2009

A.Carbone - UPMC

39

Strains	Host	Sample	Clinical condition	Phylogenetic group*	Extraintestinal mouse model phenotype* (Number of mice killed out of 10)	Genome sequence reference
K-12 MG1655	Human	Faeces	Commensal	A	NK (0)	[111]
K-12 W3110	Human	Faeces	Commensal	A	NK (0)	National Institute of Science and Technology
IAI1	Human	Faeces	Commensal	B1	NK (0)	This work
S5989	Human	Faeces	Diarrhoea (EPEC)	B1	K (10)	This work
<i>S. boydii</i> 4 227 (Sb 227)	Human	Faeces	Shigellosis	S1	ND ^b	[116]
<i>S. sonnei</i> 046 (Ss 046)	Human	Faeces	Shigellosis	SS	ND	[116]
<i>S. flexneri</i> 2a 301 (Sf 301)	Human	Faeces	Shigellosis	S3	ND	[117]
<i>S. flexneri</i> 2a 2457T (Sf 2457T)	Human	Faeces	Shigellosis	S3	NK (0)	[118]
<i>S. flexneri</i> 5b 8401 (Sf 8401)	Human	Faeces	Shigellosis	S3	ND	[119]
<i>S. dysenteriae</i> 1 197 (Sd 197)	Human	Faeces	Shigellosis	SD1	ND	[116]
O157:H7 EDL933	Human	Faeces	Diarrhoea (EHEC)	E	NK (1)	[120]
O157:H7 Sakai	Human	Faeces	Diarrhoea (EHEC)	E	NK (1)	[121]
UMN026	Human	Urine	Cystitis (ExPEC)	D	K (10)	This work
IAI39	Human	Urine	Pyelonephritis (ExPEC)	D	K (8)	This work
UT89	Human	Urine	Cystitis (ExPEC)	B2	K (10)	[122]
APEC_O1	Chicken	Lung	Colisepticemia (ExPEC)	B2	K (10)	[123]
S88	Human	Cerebro-spinal fluid (ExPEC)	New born meningitis	B2	K (10)	This work
CFT073	Human	Blood	Pyelonephritis (ExPEC)	B2	K (10)	[30]
ED1A	Human	Faeces	Healthy subject	B2	NK (0)	This work
S36	Human	Urine	Pyelonephritis (ExPEC)	B2	K (10)	[124]
<i>E. fergusonii</i>	Human	Faeces	Unknown	Outgroup	NK (1)	This work

The strains in bold correspond to the strains sequenced in this work.

*EAF (enteroaggregative *E. coli*), EHEC (enterohaemorrhagic *E. coli*), ExPEC (extraintestinal pathogenic *E. coli*).

^aThe *E. coli* and Shigella phylogenetic groups are as defined in [22] and [6], respectively.

^bNK killer; NK, Non Killer [32].

ND, not determined.

doi:10.1371/journal.pgen.1000344.t001

A.Carbone - UPMC

40

Chromosome features	<i>E. coli</i> K-12 MG1655		<i>E. coli</i> strains				<i>E. fergusonii</i> ATCC	
	55989	IAI1	ED1a	S88	IAI39	UMN026		
Genome Size (bp)	4 639 675	5 154 862	4 700 560	5 209 548	5 032 268	5 132 068	5 202 090	4 588 711
G+C content (%)	50.8	50.7	50.8	50.7	50.7	50.6	50.7	49.9
rRNA operons	7 (+55)	7 (+55)	7 (+55)	7 (+55)	7 (+55)	7 (+55)	7 (+55)	
tRNA genes	86	94	86	91	91	88	88	87
Total Protein-coding genes ^a	4306	4969	4491	5129	4859	4906	4918	4336
Pseudogenes ^b (nb)	81	79	51	95	90	80	45	22
Protein coding density ^c	85.7	87.4	87.6	86.2	87	86.1	87.8	84.7
Assigned function ^d (%)	80	74	77	74	77	76.5	77	
Conserved hypothetical (%)	12.5	23	21.5	23	22	20	22	20
Orphans (%)	7.5	3	1.5	3	1	2	1.5	3
IS-like genes (nb)	66	150	42	118	47	224	92	29
Phage-associated genes (nb)	231	406	201	657	507	393	429	235

^aThe number of protein-coding genes is given without the number of coding sequences annotated as artifactual genes (Supplementary Table 2A).
^bThe number of pseudogenes computed for each genome corresponds to the real number of genes that are pseudogenes; one pseudogene can be made of only one CDS (in this case the gene is partial compared to the wild type form in other *E. coli* strains) or of several CDSs (generally two or three CDSs corresponding to the different fragments of the wild type form in other *E. coli* strains). These lists of pseudogenes are available in Supplementary Table 1.
^cThe computed protein coding density takes into account the total length of protein genes excluding overlaps between genes, artifacts, and RNA genes.
^dProtein genes with assigned function include the total number of definitive and putative functional assignments.
doi:10.1371/journal.pgen.1000344.t002

D'autres informations manquent : ARN non codants par exemple

A.Carbone - UPMC

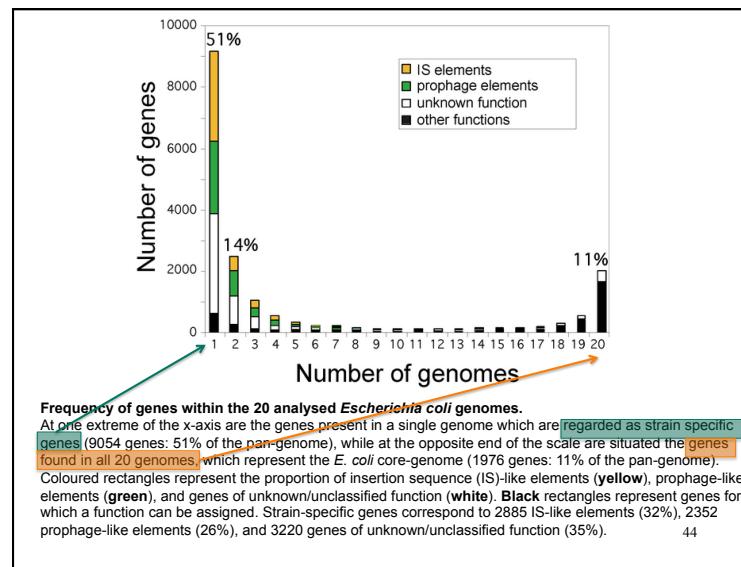
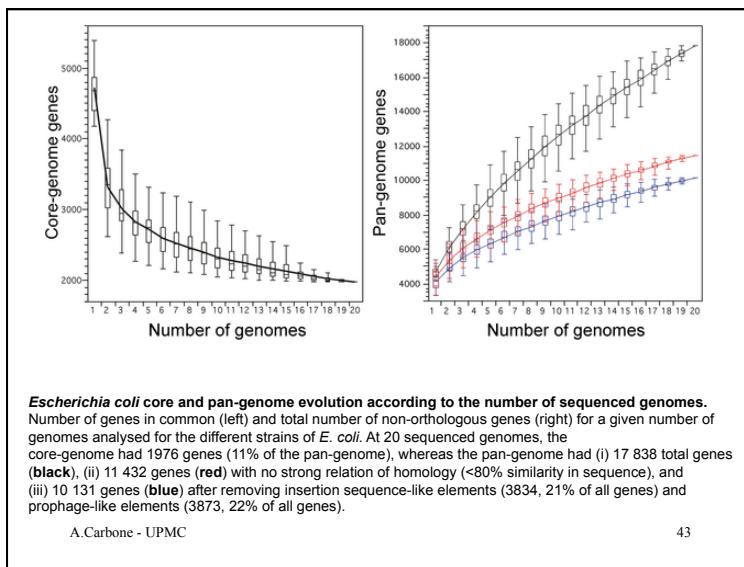
41

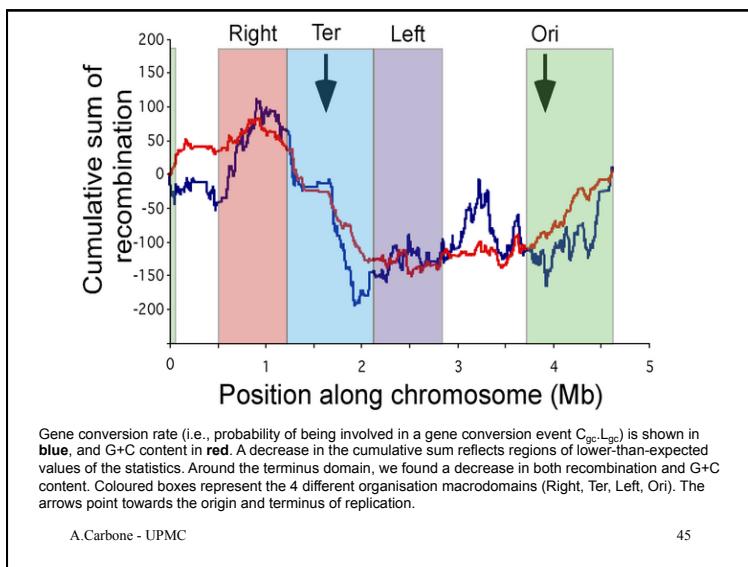
Plasmid features	<i>E. coli</i> strains			<i>E. fergusonii</i> ATCC	
	55989	ED1a	S88	UMN026	
Genome Size (bp)	72 482	119 594	133 853	122 301	33 809 55 150
G+C content (%)	46.1	49.2	49.3	50.5	42 48.5
Total Protein-coding genes ^a	100	150	144	149	49 54
Pseudogenes ^b (nb)	7	11	9	8	0 5
Protein coding density ^c	75.6	86.2	87	79.4	87.5 88.7
Assigned function ^d (%)	74	53	65	65.7	35.4 46.6
Orphans (%)	17	31.5	25.8	27.8	12.5 20.7
Hypothetical (%)	9	15.5	9.2	6.5	52.2 32.7
(S-)like genes (nb)	18	14	14	15	0 4

^aThe number of protein-coding genes is given without the number of coding sequences annotated as artifactual genes (Supplementary Table 2A).
^bThe number of pseudogenes computed for each genome corresponds to the real number of genes that are pseudogenes; one pseudogene can be made of only one CDS (in this case the gene is partial compared to the wild type form in other *E. coli* strains) or of several CDSs (generally two or three CDSs corresponding to the different fragments of the wild type form in other *E. coli* strains). These lists of pseudogenes are available in Supplementary Table 1.
^cThe computed protein coding density takes into account the total length of protein genes excluding overlaps between genes, artifacts, and RNA genes.
^dProtein genes with assigned function include the total number of definitive and putative functional assignments.
doi:10.1371/journal.pgen.1000344.t003

A.Carbone - UPMC

42

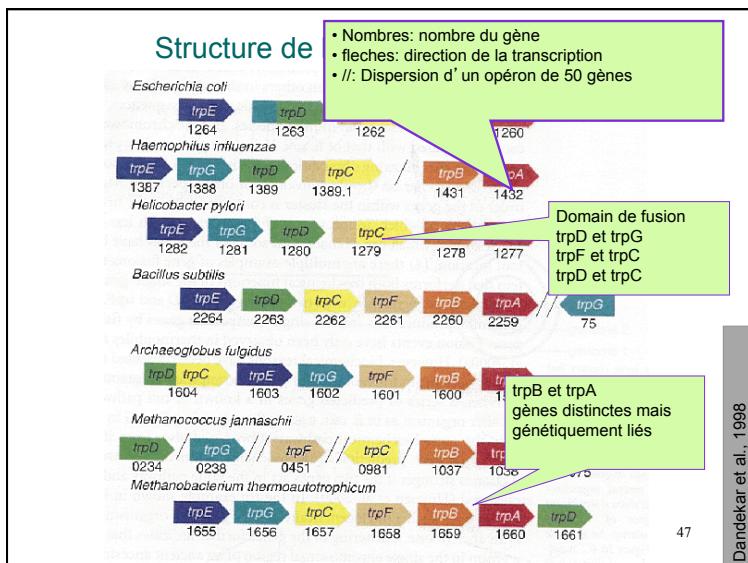




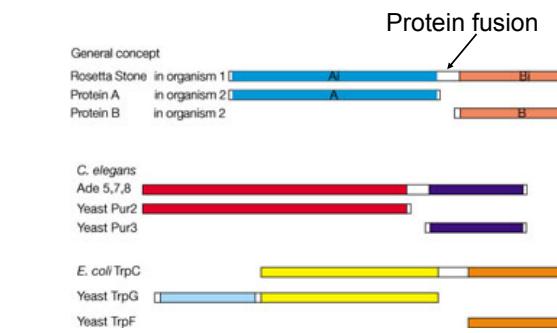
From gene order to gene function

- Gene order is very conserved for close species but it changes by rearrangement.
- By considering larger evolutionary distances, we lose the correspondence between gene order for orthologous genes.
- Groups of genes with similar biochemical functions have the tendency to remain localized.
 - Genes involved in tryptophane synthesis (trp genes) in *E. coli* and other prokaryotes

46

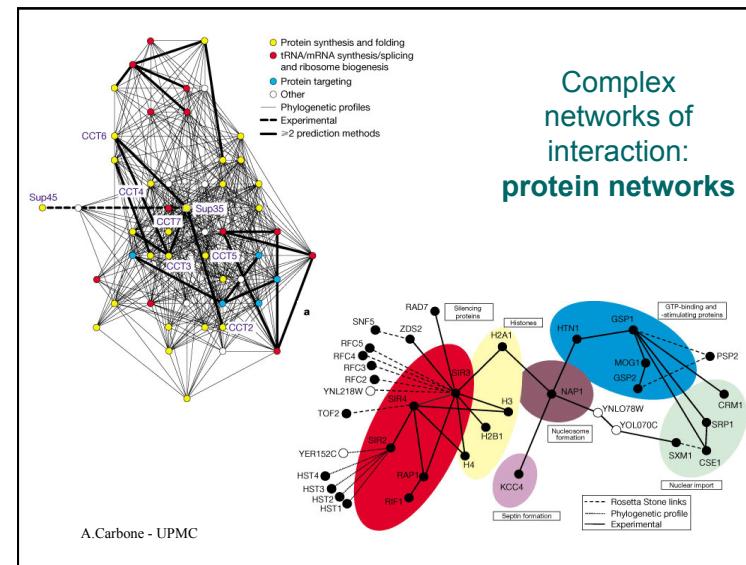
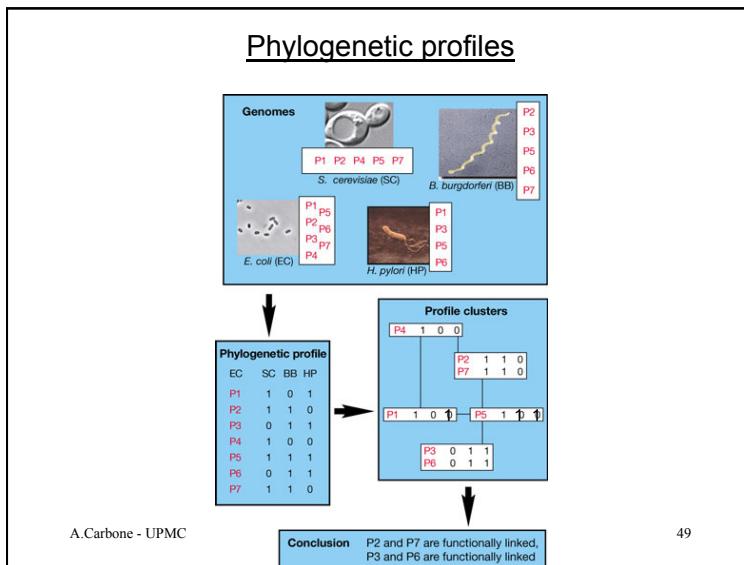


Search of genes in different species: criteria for the detection of their functional relation



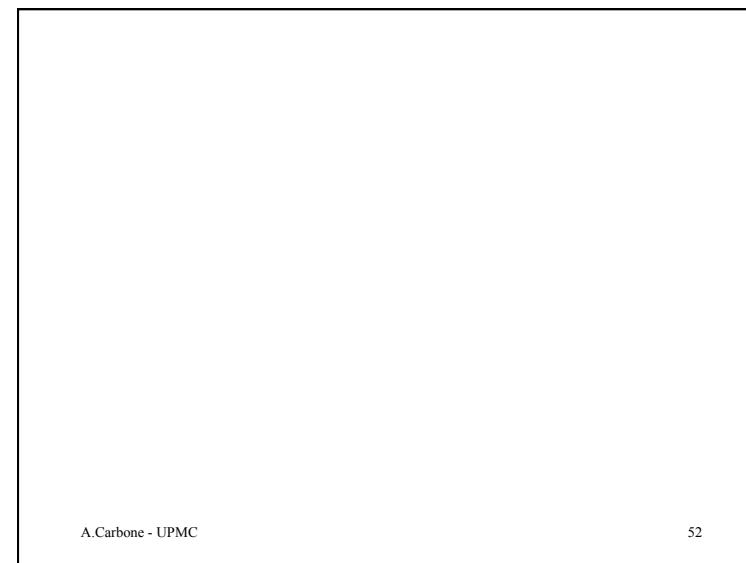
A. Carbone - UPMC

48



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
		<i>Ciona intestinalis</i>	4122	25
	Fungi	<i>Saccharomyces cerevisiae</i>	6607	23
Bacteria	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
	Stramenopiles	<i>Giardia theta</i>	24822	51
	Parabasalia	<i>Bigelowiella natans</i>	21000	47
	Proteobacteria	<i>Phaeodactylum tricornutum</i>	10408	20
	Elusimicrobia	<i>Trichomonas vaginalis</i>	59681	43
	Firmicutes	<i>Salmonella enterica</i>	4720	33
	Actinobacteria	<i>Helicobacter pylori</i>	1593	32
Archaea	Euryarchaeota	<i>Elusimicrobium minutum</i>	1548	35
		<i>Staphylococcus aureus</i>	2620	62
	Firmicutes	<i>Enterococcus faecalis</i>	3278	27
	Actinobacteria	<i>Streptococcus pneumoniae</i>	1990	23
	Cyanobacteria	<i>Mycobacterium tuberculosis T46</i>	4134	42
	Euryarchaeota	<i>Microcystis aeruginosa</i>	5356	50
	Korarchaeota	<i>Methanobrevibacter smithii</i>	1710	25
Thaumarchaeota	Korarchaeota	<i>Methanococcus maripaludis</i>	1807	31
		<i>Halobacterium salinarum</i>	2749	47
	Crenarchaeota	<i>Candidatus Korarchaeum cryptofilum</i>	1612	25
	Crenarchaeota	<i>Cenarchaeum symbiosum A</i>	2017	48
Crenarchaeota	Crenarchaeota	<i>Pyrobaculum oguniense TE7</i>	2869	38



Génomique comparative des phages

A. Carbone. Codon bias is a major factor explaining phage evolution in translationally biased hosts, *Journal of Molecular Evolution*, 66(3):210–23, 2008.

Synthetic biology

Bacteria and environment

Phages

Bacteria are the phage environment

A significant fraction of the prokaryotic community is infected by phages

The total number of viruses, which is much larger than the total prokaryotic abundance, varies strongly in different environments and it is correlated with bacterial abundance and activity

Codon bias is a major factor explaining phage evolution in transl biased hosts

116 phages of 11 translationally biased host bacteria :

actinobacteria

Mycobacterium smegmatis

proteobacteria gamma

Mycobacterium tuberculosis

Escherichia coli

Vibrio cholerae

Salmonella typhimurium

Bacillus subtilis

Listeria monocytogenes

Staphylococcus aureus

Lactococcus lactis

Streptococcus pyogenes

Chlamydophila caviae

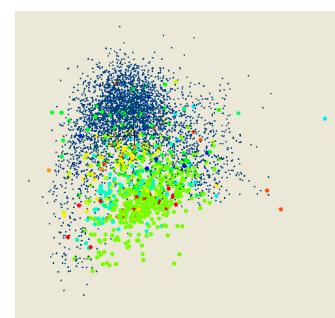
firmicutes bacillales

firmicutes lactobacillales

chlamydiales

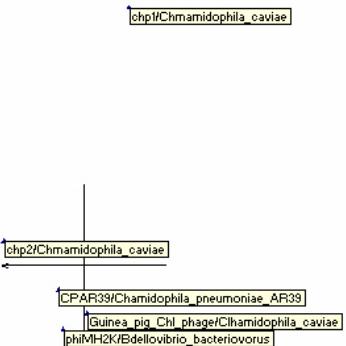
102 of these phages display at least one gene with high bias

Vibrio cholerae + 10 phages



Membrane proteins and
permease proteins («hydrophobic» aa)

Patterns of adaptation go beyond host species



Phage classification does not reflect host phylogeny

COGs: Classification phylogénétique des protéines codées dans les génomes complets

Screenshot of the COG (Clusters of Orthologous Groups) website. The page displays a table of proteins grouped by COG codes. The columns include Code, Name, Proteins, and Principal component analysis of genomes. The table lists several clusters, such as A (Archaea), O (Bacteria), M (Methanobacterium thermocatenulatum), R (Thermoplasma acidophilum), K (Pyrococcus horikoshii), Z (Pyrococcus abyssi), X (Sphaerotilus sphaerotilis), Q (Aequorea aequorea), and V (Thermotoga maritima). The COG website also includes links for "List of COGs", "Distribution", "Co-occurrences", "Phylogenetic patterns search", and "Functional categories".

62