



Pan- and core-Genomics



Center for Biological Sequence Analysis

Department of Systems Biology

Dave Ussery

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain





http://egg.umh.es/course.html

Course outline

	Monday, 12 December	Tuesday, 13 December	Wednesday, 14 December	Thursday, 15 December	Friday, 16 December
9:30-10:30	Introduction to Genomics E. coli O104 as an example	3 Generations of Genome Sequencing	Pan- and Core-genomics	Introduction to Metagenomics	Tools to compare bacterial genomes
10:30-10:45 (short Break)			Contraction of the second		and the second second
10:45-11:30	Sequences as biological information	BLAST matrices and atlases	<i>Veilonella</i> : halfway between Gram_positive & Gram neg.	Transcriptomics and metranscriptomics	Vibrio
11:30-12:00 (Coffee Break)					
12:00-13:00	Computer exercises: Introduction to UNIX	Computer exercises: Extract DNA, find rRNAs	Computer exercises: Gene-finding	Computer exercises: BLAST Matrix	Bifidobacterium Pseudomonas
13:00-15:00 (Lunch break)			8.382.762		
15:00-16:30	Computer exercises: GenomeAtlas construction	Computer exercises: Construction of 16S rRNA trees	Computer exercises: Amino acid and Codon Usage	Computer exercises: Pan- core-genome plots	E. coli, Bacillus, Yersinia
16:30-17:00 (Coffee break)					
17:00-18:30	Computer exercises: GenomeAtlas construction	Computer exercises: Construction of 16S rRNA trees	Computer exercises: Amino acid and Codon Usage	Computer exercises: Pan- core-genome plots	certificate presentation, closing

http://www.cbs.dtu.dk/staff/dave/genomics_course/Alicante2011.php

Proc. Natl. Acad. Sci. USA Vol. 93, pp. 10268–10273, September 1996 Evolution

1. What is a core-genome? What is a pan-genome? A minimal gene set for cellular life derived by comparison of complete bacterial genomes

ARCADY R. MUSHEGIAN AND EUGENE V. KOONIN*

National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894

Communicated by Clyde Hutchinson, University of North Carolina, Chapel Hill, NC, May 17, 1996 (received for review March 11, 1996)





Haemophilus influenza 86-028NP



1. What is a pan-genome?

E. coli pan-genome ~50,000 gene families



4144 proteins

~3000

E. coli gene _{families}

BMC Genomics 2009, **10:**385



14 December, 2011 6

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

Vol 464 4 March 2010 doi:10.1038/nature08821

nature

ARTICLES

1. What is a meta-genome?

A human gut microbial gene catalogue established by metagenomic sequencing

Junjie Qin¹*, Ruiqiang Li¹*, Jeroen Raes^{2,3}, Manimozhiyan Arumugam², Kristoffer Solvsten Burgdorf⁴, Chaysavanh Manichanh⁵, Trine Nielsen⁴, Nicolas Pons⁶, Florence Levenez⁶, Takuji Yamada², Daniel R. Mende², Junhua Li^{1,7}, Junming Xu¹, Shaochuan Li¹, Dongfang Li^{1,8}, Jianjun Cao¹, Bo Wang¹, Huiqing Liang¹, Huisong Zheng¹, Yinlong Xie^{1,7}, Julien Tap⁶, Patricia Lepage⁶, Marcelo Bertalan⁹, Jean-Michel Batto⁶, Torben Hansen⁴, Denis Le Paslier¹⁰, Allan Linneberg¹¹, H. Bjørn Nielsen⁹, Eric Pelletier¹⁰, Pierre Renault⁶, Thomas Sicheritz-Ponten⁹, Keith Turner¹², Hongmei Zhu¹, Chang Yu¹, Shengting Li¹, Min Jian¹, Yan Zhou¹, Yingrui Li¹, Xiuqing Zhang¹, Songgang Li¹, Nan Qin¹, Huanming Yang¹, Jian Wang¹, Søren Brunak⁹, Joel Doré⁶, Francisco Guarner⁵, Karsten Kristiansen¹³, Oluf Pedersen^{4,14}, Julian Parkhill¹², Jean Weissenbach¹⁰, MetaHIT Consortium[†], Peer Bork², S. Dusko Ehrlich⁶ & Jun Wang^{1,13}

To understand the impact of gut microbes on human health and well-being it is crucial to assess their genetic potential. Here we describe the Illumina-based metagenomic sequencing, assembly and characterization of 3.3 million non-redundant microbial genes, derived from 576.7 gigabases of sequence, from faecal samples of 124 European individuals. The gene set, ~150 times larger than the human gene complement, contains an overwhelming majority of the prevalent (more frequent) microbial genes of the cohort and probably includes a large proportion of the prevalent human intestinal microbial genes. The genes are largely shared among individuals of the cohort. Over 99% of the genes are bacterial, indicating that the entire cohort harbours between 1,000 and 1,150 prevalent bacterial species and each individual at least 160 such species, which are also largely shared. We define and describe the minimal gut metagenome and the minimal gut bacterial genome in terms of functions present in all individuals and most bacteria, respectively.

It has been estimated that the microbes in our bodies collectively make up to 100 trillion cells, tenfold the number of human cells, and suggested that they encode 100-fold more unique genes than our own genome¹. The majority of microbes reside in the gut, have a profound influence on human physiology and nutrition, and are crucial for human life^{2,3}. Furthermore, the gut microbes contribute to energy harvest from food, and changes of gut microbiome may be associated with bowel diseases or obesity^{4–8}.

individuals from the United States or Japan^{8,16,17}. To get a broader overview of the human gut microbial genes we used the Illumina Genome Analyser (GA) technology to carry out deep sequencing of total DNA from faecal samples of 124 European adults. We generated 576.7 Gb of sequence, almost 200 times more than in all previous studies, assembled it into contigs and predicted 3.3 million unique open reading frames (ORFs). This gene catalogue contains virtually all of the prevalent gut microbial genes in our cohort, provides a



The international MetaHIT (Metagenomics of the Human Intestinal Tract) project has published a gene catalogue of the <u>human gut microbiome</u> derived from 124 healthy, overweight and obese human adults, as well as inflammatory disease patients, from Denmark and Spain. The data provide the first insights into this gene set - over 150 times larger than the human gene complement - and permit the definition of both a minimal gut metagenome and a minimal gut bacterial genome. Credit: Roger Harris /Science Photo Library.

2. How is the core- and pan-genome calculated?

Proc Natl Acad Sci USA, 102:13950-13955 (2005).

Genome analysis of multiple pathogenic isolates of Streptococcus agalactiae: Implications for the microbial "pan-genome"

Herve Tettelina,b, Vega Masignanib,c, Michael J. Cieslewiczb,d,e, Claudio Donatic, Duccio Medinic, Naomi L. Warda,f, Samuel V. Angiuolia, Jonathan Crabtreea, Amanda L. Jonesg, A. Scott Durkina, Robert T. DeBoya, Tanja M. Davidsena, Marirosa Morac, Maria Scarsellic, Immaculada Margarit y Rosc, Jeremy D. Petersona, Christopher R. Hausera, Jaideep P. Sundarama, William C. Nelsona, Ramana Madupua, Lauren M. Brinkaca, Robert J. Dodsona, Mary J. Rosovitza, Steven A. Sullivana, Sean C. Daughertya, Daniel H. Hafta, Jeremy Selenguta, Michelle L. Gwinna, Liwei Zhoua, Nikhat Zafara, Hoda Khouria, Diana Radunea, George Dimitrova, Kisha Watkinsa, Kevin J. B. O'Connorh, Shannon Smithi, Teresa R. Utterbacki, Owen Whitea, Craig E. Rubensg, Guido Grandic, Lawrence C.

Madoffe,j, Dennis L. Kaspere,j, John L. Telfordc, Michael R. Wesselsd,e, Rino Rappuolic,k,I, and Claire M. Frasera,b,k,m



Fig. 2. GBS core genome. The number of shared genes is plotted as a function of the number *n* of strains sequentially added (see *Materials and Methods*). For each *n*, circles are the $8!/[(n - 1)! \cdot (8 - n)!]$ values obtained for the different strain combinations. Squares are the averages of such values. The continuous curve represents the least-squares fit of the function $F_c = \kappa_c \exp[-n/\tau_c] + \Omega$ (see Eq. **1** in *Supporting Text*) to data. The best fit was obtained with correlation $r^2 = 0.990$ for $\kappa_c = 610 \pm 38$, $\tau_c = 2.16 \pm 0.28$, and $\Omega = 1,806 \pm 16$. The extrapolated GBS core genome size Ω is shown as a dashed



Fig. 3. GBS pan-genome. The number of specific genes is plotted as a function of the number *n* of strains sequentially added (see *Materials and Methods*). For each *n*, circles are the $8!/[(n - 1)! \cdot (8 - n)!]$ values obtained for the different strain combinations; squares are the averages of such values. The blue curve is the least-squares fit of the function $F_s(n) = \kappa_s \exp[-n/\tau_s] + tg(\theta)$ (see Eq. 2 in *Supporting Text*) to the data. The best fit was obtained with correlation $r^2 = 0.995$ for $\kappa_s = 476 \pm 62$, $\tau_s = 1.51 \pm 0.15$, and $tg(\theta) = 33 \pm 3.5$. The extrapolated average number $tg(\theta)$ of strain-specific genes is shown as a









D.W. Ussery et al., *Computing for Comparative Microbial Genomics*, Computational Biology 8, DOI 10.1007/978-1-84800-255-5_14, © Springer-Verlag London Limited 2009





Fig. 12.5 The pan-genome and core genome for five different Proteobacterial genera. The *Salmonella* graph represents one species (*Salmonella enterica*), whereas the *E. coli/Shigella* figure contains both *E. coli* and four different *Shigella* species. The other graphs represent multiple species per genus. All graphs are drawn on the same scale

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

223



Figure 3 Core and pan genome plot of the *Pseudomonas* genomes. *Azotobacter* is *highlighted*. The number of pan genome for all strains used is 29,696 (*dashed line on top*) and the core genome is reduced to

443 (*dashed line below*). The *light purple color* indicates where *P. putida* is added, and the *light red column* indicates where *A. vinelandii* is added

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain



Microbial Ecology, 60:708-720, (2010).

64 E. coli genomes





Wednesday, 14 December, 2011

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

14 December, 2011 16





Figure 3 Pan- and core genome plot of *Bacteroides* and *Parabacteroides* genomes. The *blue line* (core) represents the conserved number of orthoMCL protein families. The *red line* (pan) indicates the cumulative number of orthoMCL protein families in the genomes. *Green bar* indicate the number of novel orthoMCL protein families in

the genome. The relative size of the core protein families to the total genome size (% Core) is based on the 1,085 protein families shared by *Bacteroides* and *Parabacteroides* (excluding *B. capillosus* and *B. pectinophilus*), On average, 27% of the proteins is shared in the core protein families





Microbial Ecology, in the press (2011)

DOI: 10.1007/s00248-011-9948-y



Figure 1 Pan- and core genome plots of the six analyzed genera. The genomes were analyzed in alphabetical order of species names



given in the pie chart are calculated by exclusion of classes R, S and X. Only values \geq 5% are shown

Figure 9 COG statistics for the genes found in the pan-genome (*left*) and core genome (*right*) of *Bifidobacterium* genomes. The key for the COG classes is explained below the pie charts. Percentages

<u>Microbial Ecology</u>, *in the press (2011)* DOI: 10.1007/s00248-011-9948-y

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain



Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

ď

'n

`х



Figure 11 COG statistics for the genes found in the pan-genome (*left*) and core genome (*right*) of the collection of genomes from all included organisms, divided into non-pathogenic isolates (probiotic,

fermentative and normal human gut flora) at the top and pathogenic isolates at the bottom

Microbial Ecology, in the press (2011)

DOI: 10.1007/s00248-011-9948-y

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

3. Who cares? What are pan-genomes good for?

Standard operating procedure for computing pangenome trees coli O157:H7 str. EC4196 Standards in Genomic Sciences (2010) 2:135-141 coli O157:H7 str. EC4113 E coli O157:H7 str. EC508 coli O157:H7 str. EC4501 coli O157:H7 str. EC4076 coli O157:H7 str. EC4115 coli O157:H7 str. EC4042 Standards in coli O157:H7 str. EC4486 coli O157:H7 str. EC869 **Genomic Sciences** coli O157:H7 str. EC4206 coli 0157:H7 str. FC4401 F coli O157:H7 str. EDL933 coli O157:H7 str. Sakai 00 100 coli O157:H7 str. TW14588 An Open Access Journal of the Genomic Standards Consortium coli O157:H7 str. LANL ECF coli O157:H7 str. LANL ECA coli O157:H7 str. EC4045 coli K12 str. W3110 100 coli K12 str. MG1655 100 . coli K12 str. DH1 . coli K12 str. BW2952 E. coli K12 DH10B E. coli K12ATCC8739 . coli B str. REL606 . coli BL21 (DE3Korea) . coli BL21 (DE3 AU) E coli BL21 (DE3 DOE) F coli HS E. coli SE11 coli IAI1 F 95 coli 55989 F. coli E24377A coli O26:H11 str. 11368 coli O103:H2 str. 12009 coli O111:H-str.11128 F coli O103 Oslo coli SMS-3-5 E coli UMN026 F. coli 53638 coli IAI39 coli UTI89 coli S88 F. coli CFT073 E E. coli SE15 F coli 536 coli 0127:H6 str. E2348/69 E coli ED1a F 99 . coli F11 F. coli APECO1 coli E110019 F F coli E22 coli B7A F coli 101-1 boydii CDC 3083-94 boydii Sb227 sonnei Ss046 flexneri 2a str. 2457T 70 flexneri 2a str. 301 100 flexneri 5 str. 8401 murium str. LT2 70 fergusonii ATCC 35469 -99 E. Е albertii TW07627 dysenteriae 1012 S. dysenteriae Sd197 Relati ve manhattan distance 0.20 0.15 0.10 0.05 0.00

Figure 3 Pan-genome clustering of *E. coli* (*black*) and related species (*colored*), based on the alignment of their variable gene content. The genomes now cluster according to species and a relatedness between

E. coli K12 derivatives (*green block*) and group B isolates (*orange block*) is visible Microbial Ecology, **60**:708-720, (2010).

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain

CENTERFO RBIOLOGI



Pangenome trees



Figure 5. The left panel show the pan-genome family tree for the 14 strains of *Staphylococcus aureus* completed as NCBI. Here the weights have been chosen according to the blue bars in Figure 1, i.e. the "stabilome" genes have been emphasized. In the right panel, the same data have been used, but weights are now chosen to emphasize the "mobilome" genes. In both cases ORFans have been discarded.

A. Jacobsen et al.





Figure 1 16S rRNA tree generated with the neighbor-joining method with 1,000 bootstrap resamplings. The tree shows the evolutionary relationships of *A. vinelandii* with the *Pseudomonas* genus and other Gammaproteobacteria based on their 16S rRNA sequences

Genomics of Prokaryotes Workshop



Universidad Miguel Hernández

Alicante, Spain

Wednesday, 14 December, 2011

DOI: 10.1007/s00248-011-9914-8

Microbial Ecology, in the press (2011)

14 December, 2011 26





Relative manhattan distance

Figure 2 Pan genome family tree. The tree shows the phylogenetic relationships based on the gene families found in the pan genome, excluding the families found in only one genome. Color coding is again based on Table 1

Microbial Ecology, in the press (2011)

DOI: 10.1007/s00248-011-9914-8

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain







Figure 4 BLASTMatrix shows the homology between the pairs of genomes compared (green) and homology within each genome itself (red) based on the 50/50 rule

<u>Microbial Ecology</u>, *in the press (2011)* DOI: 10.1007/s00248-011-9914-8

Genomics of Prokaryotes Workshop Universidad Miguel Hernández Alicante, Spain





Consider the following taxonomic groups:

Proteobacteria (phyla) Enterobacteriaceae (family) Escherichia (genus) E. coli (species)

Which <u>core</u>-genome is likely to be the largest?
Which <u>pan-genome</u> is likely to be the largest?

•How would you expect the addition of a low-quality unfinished genome to affect the size of the core- and pan-genome?