

Metagenomics



Smithsonian Institution - Comparative Genomics Course 10/6/2011

Jay Evans and Scott Cornman Bee Research Lab, USDA-ARS Beltsville, MD

- ❖ FIRST PRINCIPLES
- ❖ DATA COLLECTION
- ❖ OCEANS, TERMITES, MAMMALS
- ❖ BEE STORIES
- ❖ RETURN TO FUNCTION
- ❖ ANALYSIS



- ❖ **FIRST PRINCIPLES**
- ❖ **DATA COLLECTION**
- ❖ **OCEANS, TERMITES, MAMMALS**
- ❖ **BEE STORIES**
- ❖ **RETURN TO FUNCTION**
- ❖ **ANALYSIS**



Molecular biological access to the chemistry of unknown soil microbes: a new frontier for natural products

Jo Handelsman¹, Michelle R Rondon¹, Sean F Brady², Jon Clardy² and Robert M Goodman¹



Cultured soil microorganisms have provided a rich source of natural-product chemistry. Because only a tiny fraction of soil microbes from soil are readily cultured, soil might be the greatest untapped resource for novel chemistry. The concept of cloning the metagenome to access the collective genomes and the biosynthetic machinery of soil microflora is explored here.

“Tapping into this source should be a great, joint adventure for biologists and chemists”

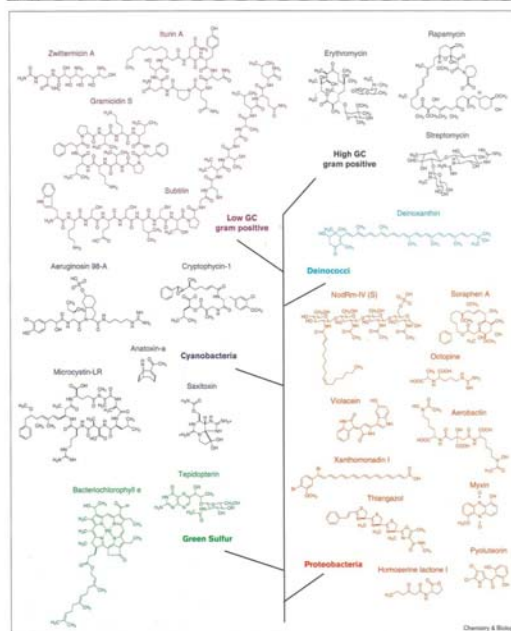
Figure 1



Morphological diversity typical of microorganisms cultured from soil on a broad spectrum medium, tryptic soy agar.

Chemistry and Biology, Oct 1998

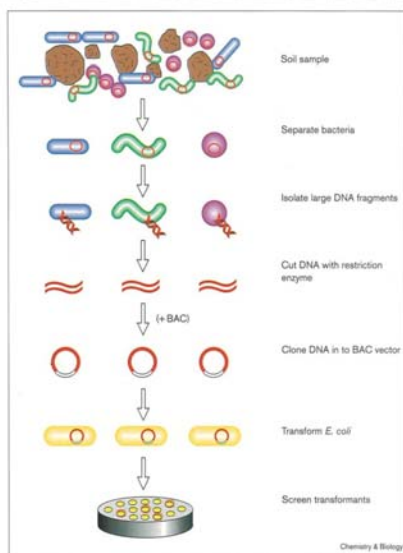
Figure 2



Examples of the chemical diversity in some of the major phyla of bacteria.

Figure 3

Cloning the metagenome is our process for isolating new pathways for the synthesis of bioactive molecules from noncultured soil microorganisms. DNA is extracted directly from soil, using gentle methods to preserve high-molecular-weight DNA. The DNA is cut using a restriction enzyme and cloned into a bacterial artificial chromosome (BAC), a vector which can carry large fragments of DNA in *E. coli*. The BAC clones are then screened for biological activity and for the production of novel natural products.



Chemistry and Biology, Oct 1998

"The excitement surrounding this new field lies in the vast diversity of unknown soil microflora and the chemical richness that they are thought to contain"

"The methodology has been made possible by advances in molecular biology and eukaryotic genomics, which have laid the groundwork for **cloning and functional analysis of the collective genomes** of soil microflora, which we term the metagenome of the soil"

Daughter of metagenomics....



Plus combinatorial chemistry, HTP proteomics, other advances in databases and thinking...



- ❖ **FIRST PRINCIPLES**
- ❖ **DATA COLLECTION**
- ❖ **OCEANS, TERMITES, MAMMALS**
- ❖ **BEE STORIES**
- ❖ **RETURN TO FUNCTION**
- ❖ **ANALYSIS**



Issues in design & analysis

Focus on deep sequencing

Does it work?

Metadata

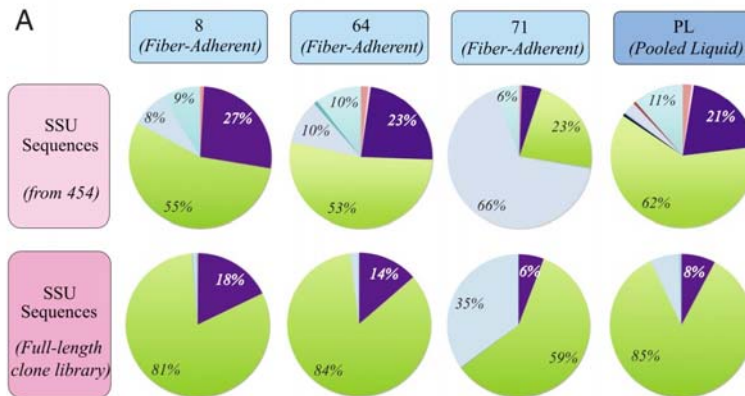
Assembly vs. mapping

Measuring diversity (rDNA) vs. function (enzymes)

OTU-based analysis

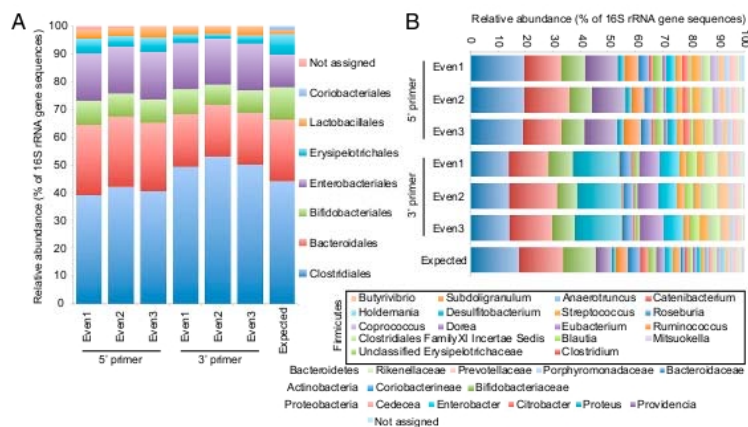
Gene-based analysis

NextGen sequencing has less bias than cloning



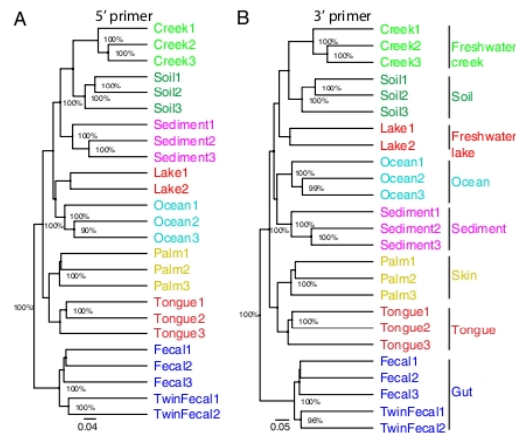
Bulc et al. PNAS 2008

Replication of community diversity estimates with Illumina short reads



Caporaso et al. PNAS 2010

Replication of community diversity estimates with Illumina short reads



Caporaso et al. PNAS 2010

Sequencing platform comparison

(per minimum unit: plate/lane/chip)

454: ~400 bp length, several hundred thousand reads

Illumina (GA/HiSeq): 120 bp, 30 million/150 million reads

ABI Solid: 75 bp max, total output \geq Illumina

-highest accuracy, fewer software options

Ion Torrent: 200 bp, 4-8 million reads

-fast

Sample preparation

Filtration

Extraction biases exist (use a consistent method)

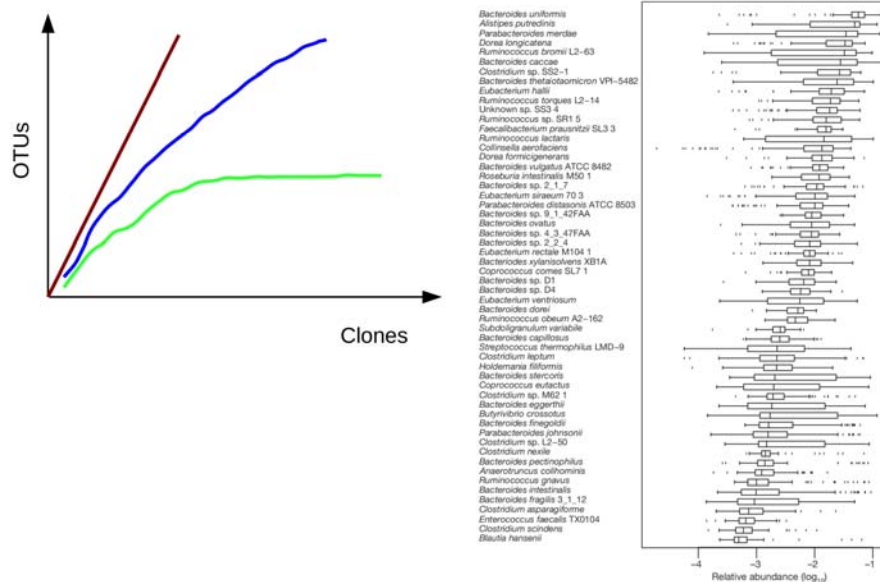
Amplicons vs. whole genome

- limitations of universal primers

Normalization of RNA (transcriptomics)

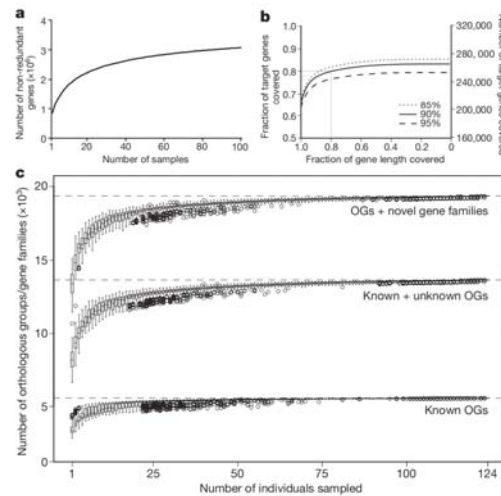
Library preparation for NGS can be difficult, talk with your sequencing center

Rarefaction curves estimate completeness



Qin et al., *Nature*, 2010

Rarefaction curves estimate completeness



Qin et al., *Nature*, 2010

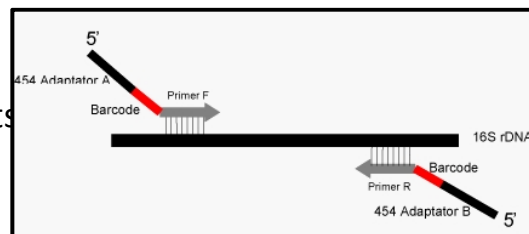
Multiplexing / Barcoding

Independent sequencing libraries run together

- identified by distinct adaptor primers
- divides basic sequencing run into multiple subsets

Applications:

- Biological replicates (key for statistical comparisons)
- Time series
- Multiple treatments



Collection of metadata

Table 1. Sampling Locations and Environmental Data

ID	Sample Location	Country	Date, mm/dd/yy	Time	Location	Sample Depth, m	Water Depth, m	T (°C) ^a	S ^b (ppt)	Size Fraction (µm)	Habitat Type	Chl <i>a</i> Sample Month (Annual Mean) ± SE mg/m ³	Good Sequences	
G000a	Sargasso Stations 13 and 11	Bermuda (UK)	07/26/03	1:00	31°32'5" N 63°35'42" W	5.0	>4,200	20.0	20.5	35.6	0.1-0.8	Open ocean	0.17 (0.09 ± 0.02)	644,551
G000b	Sargasso Stations 13 and 11	Bermuda (UK)	07/26/03	5:55	31°10'30" N 64°19'27" W	5.0	>4,200	20.0	20.5	36.7	0.22-0.8	Open ocean	0.17 (0.09 ± 0.02)	317,180
G000c	Sargasso Stations 3	Bermuda (UK)	02/25/03	13:00	32°09'30" N 64°00'36" W	5.0	>4,200	19.8	36.7	0.22-0.8	Open ocean	0.17 (0.09 ± 0.02)	368,835	
G000d	Sargasso Stations 13	Bermuda (UK)	07/25/03	17:00	31°32'5" N 63°35'42" W	5.0	>4,200	20.0	35.6	0.22-0.8	Open ocean	0.17 (0.09 ± 0.02)	332,240	
G001a	Hydrostation 5	Bermuda (UK)	05/15/03	11:40	32°10'00" N 64°30'00" W	5.0	>4,200	22.9	35.7	0.8-3.0	Open ocean	0.10 (0.10 ± 0.01)	142,852	
G001b	Hydrostation 5	Bermuda (UK)	04/14/03	11:40	32°10'00" N 64°30'00" W	5.0	>4,201	22.9	35.7	0.8-3.0	Open ocean	0.10 (0.10 ± 0.01)	90,705	
G001c	Hydrostation 5	Bermuda (UK)	05/15/03	11:40	32°10'00" N 64°30'00" W	5.0	>4,202	22.9	35.7	0.1-0.8	Open ocean	0.1 (0.1 ± 0.01)	92,551	
G002	Gulf of Maine	USA	08/21/03	6:32	42°30'11" N 67°14'24" W	1.0	106	18.7	29.7	0.1-0.8	Coastal	1.4 (1.17 ± 0.19)	131,540	
G003	Browns Bank, Gulf of Maine	Canada	08/21/03	11:50	42°51'10" N 66°13'2" W	1.0	119	11.7	29.9	0.1-0.8	Coastal	1.4 (1.12 ± 0.19)	61,605	
G004	Outside Halifax, Nova Scotia	Canada	08/22/03	5:25	44°26'14" N 64°54'40" W	2.0	142	17.1	28.1	0.1-0.8	Coastal	0.4 (0.78 ± 0.17)	52,709	
G005	Bedford Basin, Nova Scotia	Canada	08/22/03	19:21	44°51'25" N 63°58'14" W	1.0	64	15.0	33.2	0.1-0.8	Embayment	6 (6.76 ± 0.50)	61,131	
G006	Bay of Fundy, Nova Scotia	Canada	08/23/03	10:47	45°26'42" N 64°54'48" W	1.0	11	11.2		0.1-0.8	Littoral	2.8 (1.87 ± 0.18)	59,679	
G007	Northern Gulf of Maine	Canada	08/23/03	8:25	43°37'56" N 66°50'50" W	1.0	119	17.9	31.7	0.1-0.8	Coastal	1.4 (1.12 ± 0.19)	50,900	
G008	Newport Harbor, RI	USA	11/16/03	16:45	41°29'9" N 71°21'4" W	1.0	12	9.4	25.5	0.1-0.8	Coastal	2.2 (1.59 ± 0.17)	126,655	
G009	Block Island, NY	USA	11/17/03	10:30	41°52'38" N 71°39'8" W	1.0	32	11.0	51.0	0.1-0.8	Coastal	4.0 (2.72 ± 0.24)	79,325	
G010	Cape May, NJ	USA	11/18/03	9:50	39°56'24" N 74°41'6" W	1.0	10	12.0	31.0	0.1-0.8	Coastal	2.0 (2.75 ± 0.23)	70,304	
G011	Delaware Bay, NJ	USA	11/18/03	11:30	39°25'4" N 75°30'15" W	1.0	8	11.0		0.1-0.8	Littoral	4.8 (3.73 ± 1.05)	174,435	
G012	Chesapeake Bay, MD	USA	12/18/03	11:22	38°56'49" N 76°25'2" W	1.0	25	5.7	14.7	0.1-0.8	Fauna	21.0 (15.0 ± 1.01)	126,162	
G013	Off Nags Head, NC	USA	12/18/03	9:08	36°01'14" N 75°33'41" W	1.0	28	9.5		0.1-0.8	Coastal	3.0 (2.24 ± 0.24)	130,616	
G014	South of Charleston, SC	USA	12/20/03	17:12	32°50'25" N 79°15'50" W	1.0	31	18.8		0.1-0.8	Coastal	1.70 (1.02 ± 0.22)	128,885	
G015	Off Key West, FL	USA	01/08/04	6:25	24°29'18" N 83°4'12" W	2.0	47	25.3	36.0	0.1-0.8	Coastal	0.2 (0.27 ± 0.09)	127,362	
G016	Gulf of Mexico	USA	01/08/04	14:15	24°10'29" N 84°23'40" W	2.0	2,222	26.4	35.8	0.1-0.8	Coastal sea	0.16 (0.11 ± 0.01)	127,137	
G017	Tocatlán Channel	Mexico	01/08/04	13:47	20°11'21" N 85°24'49" W	2.0	4,511	27.6	35.8	0.1-0.8	Open ocean	0.13 (0.09 ± 0.01)	227,281	
G018	Rosario Bank	Honduras	01/10/04	8:12	18°21'12" N 80°47'5" W	2.0	4,470	27.4	35.4	0.1-0.8	Open ocean	0.14 (0.09 ± 0.01)	161,745	
G019	Northeast of Colón	Panama	01/12/04	9:05	10°42'58" N 80°15'18" W	2.0	3,196	27.7	35.4	0.1-0.8	Coastal	0.23 (0.15 ± 0.02)	135,328	
G020	Lake Gatun	Panama	01/13/04	10:24	9°52'12" N 79°50'10" W	2.0	4	26.5	0.06	0.1-0.8	Fresh water		296,350	
G021	Gulf of Panama	Panama	01/19/04	16:40	07°45'15" N 79°51'28" W	2.0	16	27.6	30.7	0.1-0.8	Coastal	0.50 (0.75 ± 0.22)	151,795	
G022	250 miles from Panama City	Panama	01/20/04	16:50	6°20'14" N 80°54'14" W	2.0	2,411	28.3	31.3	0.1-0.8	Open ocean	0.23 (0.28 ± 0.02)	121,662	
G023	30 miles from Cocos Island	Costa Rica	01/21/04	15:00	5°38'24" N 80°33'25" W	2.0	1,129	28.7	37.6	0.1-0.8	Open ocean	0.07 (0.19 ± 0.05)	155,051	
G024	Dirty Rock, Cocos Island	Costa Rica	01/26/04	10:51	5°55'10" N 80°5'16" W	1.1	30	28.3	31.4	0.8-3.0	Hinging reef	0.11 (0.19 ± 0.01)	120,671	
G025	134 miles NE of Galapagos	Ecuador	02/01/04	16:10	1°15'51" N 90°17'42" W	2.0	2,376	27.8	32.6	0.1-0.8	Open ocean	0.22 (0.28 ± 0.02)	102,708	
G027	Devil's Crown, Floreana	Ecuador	02/04/04	11:51	1°12'58" S 90°25'22" W	2.0	2.3	25.5	34.0	0.1-0.8	Coastal	0.40 (0.38 ± 0.03)	222,003	
G028	Coastal Floreana	Ecuador	02/04/04	15:47	1°13'1" S 90°19'11" W	2.0	156	25.0		0.1-0.8	Coastal	0.35 (0.35 ± 0.02)	189,032	
G029	North James Bay, Santa Fe	Ecuador	02/08/04	18:02	0°12'0" S 90°50'7" W	2.0	12	26.2	34.5	0.1-0.8	Coastal	0.10 (0.39 ± 0.03)	131,529	
G030	Warm seep, Boca Raton	Florida	02/08/04	11:54	0°19'20" N 81°58'0" W	19.0	19	26.9		0.1-0.8	Warm seep		359,152	
G031	Upwelling, Forastaria	Ecuador	02/10/04	14:40	0°19'4" S 91°29'6" W	12.0	15	18.0		0.1-0.8	Coastal upwelling	0.57 (0.39 ± 0.03)	436,401	
G032	Mangrove, Isabela	Ecuador	02/11/04	11:20	0°25'08" S 91°16'10" W	6.0	0.67	26.4		0.1-0.8	Mangrove		146,018	
G033	Punta Concordia Lagoon, Pichincha	Ecuador	02/15/04	13:35	1°35'42" S 92°25'45" W	0.2	0.33	37.6	40	0.1-0.8	Hypersaline		602,255	

Submission of metadata

Specialized Structured Comments

1. MIGS/MIMS/MIENS

Minimum information checklists have been developed by the [Genomic Standards Consortium](#) (GSC) as a means of reporting core descriptive information about the environment from which an organism(s) was collected. Core descriptors include information about the origins of the nucleic acid sequence (genome), its environment (eg. latitude and longitude, date and time of sampling, habitat) and sequence processing (sequencing and assembly methods).

Three different metadata lists have been developed to describe genomic, metagenomic, and environmental sequences:

- o MIGS - Minimum Information About a Genome Sequence
- o MIMS - Minimum Information About a Metagenome Sequence
- o MIENS - Minimum Information About an Environmental Sequence

The tag-value pairs that are included for each submission type can be validated for compliance with the GSC recommended list. The recommended lists of core descriptors that should be included for each of these sequence types can be found [here](#).

Validation tools within Sequin and tbl2asn will report if structured comments include all of the GSC recommended compliant core descriptors. Submissions that include all of the compliant tags will have a Keyword Included within the GenBank flatfile:

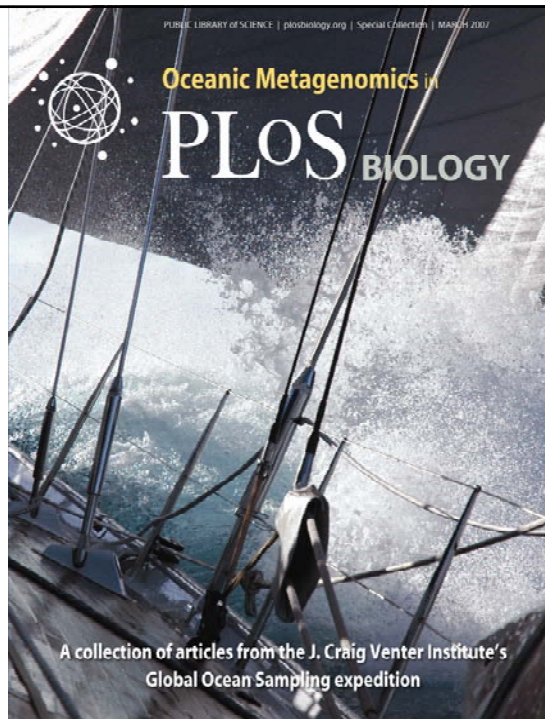
KEYWORD GSC:MIGS:2.1

Structured comments that are not compliant based on the GSC guidelines can still be included within GenBank submissions - they just will not include the keyword.

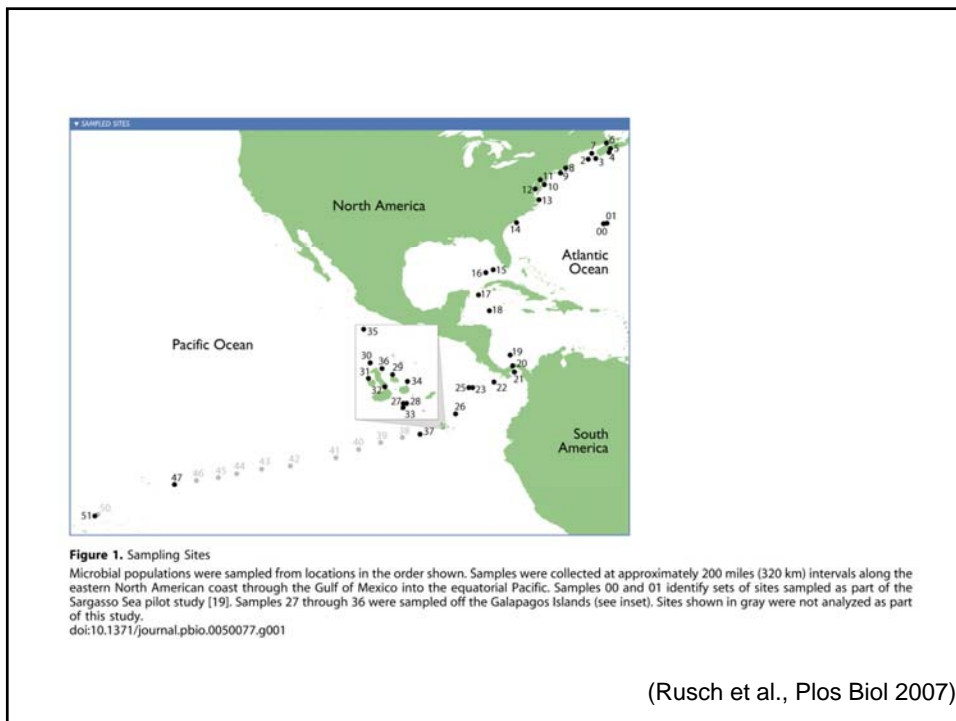
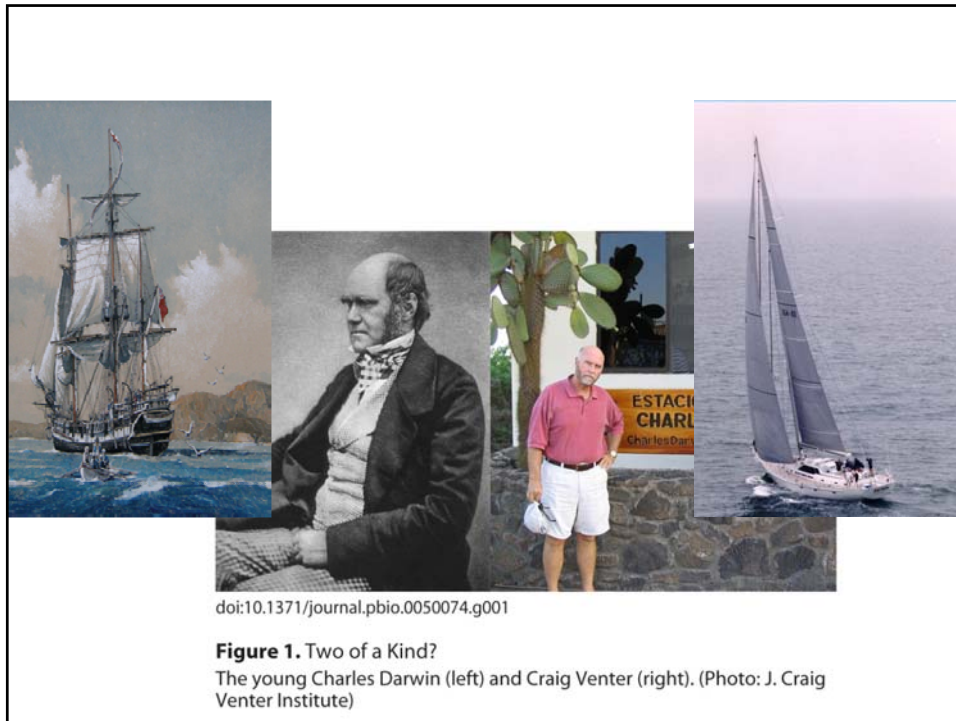
In order for this validation to occur, you will need to include within the first column in your table a tag that defines the prefix and suffix for the start and end tags within the structured comment, for example:

StructuredCommentPrefix MIGS-Data

- ❖ FIRST PRINCIPLES
- ❖ DATA COLLECTION
- ❖ OCEANS, TERMITES, MAMMALS
- ❖ BEE STORIES
- ❖ RETURN TO FUNCTION
- ❖ ANALYSIS



Various authors, 2007



(Rusch et al., Plos Biol 2007)

Table 1. Some Major Methods for Studying Individual Microbes Found in the Environment

Method	Summary	Comments
Microscopy	Microbial phenotypes can be studied by making them more visible. In conjunction with other methods, such as staining, microscopy can also be used to count taxa and make inferences about biological processes.	The appearance of microbes is not a reliable indicator of what type of microbe one is looking at.
Culturing	Single cells of a particular microbial type are grown in isolation from other organisms. This can be done in liquid or solid growth media.	This is the best way to learn about the biology of a particular organism. However, many microbes are uncultured (i.e., have never been grown in the lab in isolation from other organisms) and may be unculturable (i.e., may not be able to grow without other organisms).
rRNA-PCR	The key aspects of this method are the following: (a) all cell-based organisms possess the same rRNA genes (albeit with different underlying sequences); (b) PCR is used to make billions of copies of basically each and every rRNA gene present in a sample; this amplifies the rRNA signal relative to the noise of thousands of other genes present in each organism's DNA; (c) sequencing and phylogenetic analysis places rRNA genes on the rRNA tree of life; the position on the tree is used to infer what type of organism (a.k.a. phylotype) the gene came from; and (d) the numbers of each microbe type are estimated from the number of times the same rRNA gene is seen.	This method revolutionized microbiology in the 1980s by allowing the types and numbers of microbes present in a sample to be rapidly characterized. However, there are some biases in the process that make it not perfect for all aspects of typing and counting.
Shotgun genome sequencing of cultured species	The DNA from an organism is isolated and broken into small fragments, and then portions of these fragments are sequenced, usually with the aid of sequencing machines. The fragments are then assembled into larger pieces by looking for overlaps in the sequence each possesses. The complete genome can be determined by filling in gaps between the larger pieces.	This has now been applied to over 1,000 microbes, as well as some multicellular species, and has provided a much deeper understanding of the biology and evolution of life. One limitation is that each genome sequence is usually a snapshot of one or a few individuals.
Metagenomics	DNA is directly isolated from an environmental sample and then sequenced. One approach to doing this is to select particular pieces of interest (e.g., those containing interesting rRNA genes) and sequence them. An alternative is ESS, which is shotgun genome sequencing as described above, but applied to an environmental sample with multiple organisms, rather than to a single cultured organism.	This method allows one to sample the genomes of microbes without culturing them. It can be used both for typing and counting taxa and for making predictions of their biological functions.

doi:10.1371/journal.pbio.0050082.t001

(Rusch et al., Plos Biol 2007)

Table 2. Methods of Binning

Method	Description	Comments
Genome assembly	Identify regions of overlap between different fragments from the same organism to build larger contiguous pieces (contigs).	Getting deep enough sampling for this to work is very expensive except for low diversity systems or for very abundant taxa.
Reference genome alignment	Identify ESS fragments or contigs that are very similar to already assembled sections of the genome of single microbial types.	(a) One of the most effective ways to sort through ESS data, if the reference genome is very closely related to an organism in the sample; (b) the reason why more reference genomes are needed; (c) does not handle regions present in uncultured organisms but not in the reference.
Phylogenetic analysis	Build evolutionary trees of genes encoded by ESS fragments or contigs. Assign fragments or contigs to taxonomic groups based on nearest neighbor(s) in trees.	(a) Very powerful, but level of resolution depends on whether fragments encode useful phylogenetic markers and on how well sampled the database is for the neighbor analysis; (b) would work much better if more genomes were available from across the tree of life.
Word frequency and nucleotide composition analysis	Measure word frequency and composition of each fragment. Group by clustering algorithms or principal component analysis.	(a) Has the potential to work because organisms sometimes have "signatures" of word frequencies that are found throughout the genome and are different between species; (b) very challenging for small fragments.
Population genetics	Build alignments of fragments or contigs with similarity to each other (but not as much as needed for assembly). Examine haplotype structure, predicted effective population size, and synonymous and non synonymous substitution patterns.	May be most useful as a way of subdividing bins created by other methods.

Note that some methods can be applied to ESS fragments or to bins identified by other methods.
doi:10.1371/journal.pbio.0050082.t002

(Rusch et al., Plos Biol 2007)

Expanding the Protein Family Universe

Table 1. The Complete Dataset Consisted of Sequences from NCBI-nr, ENS, TGI-EST, PG, and GOS, for a Total of 28,610,944 Sequences

Dataset	Source	Number of Amino Acid Sequences	Mean Sequence Length	Brief Description
NCBI-nr	NCBI	2,317,995	339	Consists of protein sequences submitted to SWISS-PROT, PDB, PIR, and PRF, and also predicted proteins from both finished and unfinished genomes in GenBank, EMBL, and DDBJ.
PG ORFs	NCBI	3,049,695	160	ORFs identified from 222 prokaryotic genome projects. Organisms are listed in Protocol S1.
TGI-EST ORFs	TIGR Gene Index	5,458,820	119	ORFs identified from 72 datasets in which each dataset consists of EST assemblies. Organisms are listed in Protocol S1.
ENS	Ensembl	361,668	466	Sequences from 12 species, including human, mouse, rat, chimp, zebrafish, fruit fly, mosquito, honey bee, dog, two species of puffer fish, chicken, and worm.
GOS ORFs	J. Craig Venter Institute	17,422,766	134	ORFs identified from an assembly of 7.7 million reads. These reads include both the reads from the <i>Sorcerer II</i> GOS Expedition and the reads from the earlier Sargasso Sea study. Also included are 36,318 ORFs identified from an assembly of sequences collected from the viral size (< 0.1 μ m) fraction of one sample.

doi:10.1371/journal.pbio.0050016.t001

(Rusch et al., Plos Biol 2007)

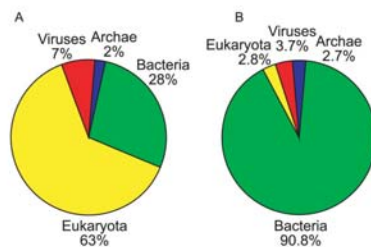


Figure 1. Proportion of Sequences for Each Kingdom
 (A) The combined set of NCBI-nr, PG, TGI-EST, and ENS has 3,167,979 sequences. The eukaryotes account for the largest portion and is more than twice the bacterial fraction.
 (B) Predicted kingdom proportion of sequences in GOS. Out of the 5,654,638 GOS sequences, 5,058,757 are assigned kingdoms using a BLAST-based scheme. The bacterial kingdom forms by far the largest fraction in the GOS set.
 doi:10.1371/journal.pbio.0050016.g001

Table 7. Taxonomic Makeup of GOS Samples Based on 16S Data from Shotgun Sequencing

Phylum or Class	Fraction ^a
Alpha Proteobacteria	0.32
Unclassified Proteobacteria	0.155
Gamma Proteobacteria	0.132
Bacteroidetes	0.13
Cyanobacteria	0.079
Firmicutes	0.075
Actinobacteria	0.046
Marine Group A	0.022
Beta Proteobacteria	0.017
OP11	0.008
Unclassified Bacteria	0.008
Delta Proteobacteria	0.005
Planctomycetes	0.002
Epsilon Proteobacteria	0.001

(Rusch et al., Plos Biol 2007)

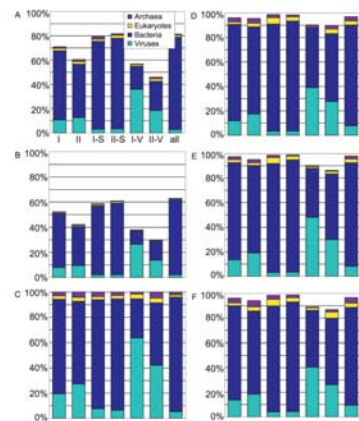


Figure 16. GOS-Only Clusters Are Enriched for Sequences of Viral Origin Independently of the Kingdom Assignment Method Employed

For each panel, clusters are as in Figure 4. For (A–C), a kingdom is assigned to each neighboring ORF within each cluster set; the percentage of all neighboring ORFs with a given kingdom assignment is plotted. In (A) and (D), a kingdom is assigned to a neighboring ORF by a majority vote of the top four BLAST matches to a protein in NCBI-nr (Materials and Methods). In (B) and (E), a kingdom is assigned if all eight highest-scoring BLAST matches agree in kingdom. In (C) and (F), all ORFs on a scaffold are assigned the same kingdom by voting among all ORFs with BLAST matches to NCBI-nr on that scaffold (Materials and Methods). In all graphs, only clusters with at least one assignable neighbor are considered. When compared to the size-matched controls, in all cases the GOS-only clusters show enrichment for viral sequences.

doi:10.1371/journal.pbio.0050016.g016

(Rusch et al., Plos Biol 2007)

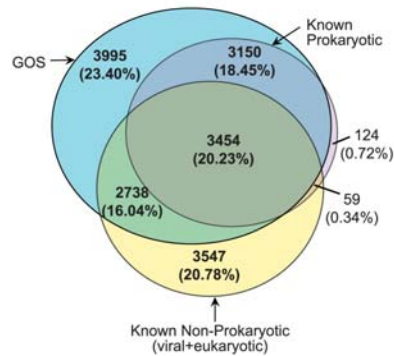


Figure 3. Venn Diagram Showing Breakdown of the 17,067 Medium and Large Clusters by Three Categories—GOS, Known Prokaryotic, and Known Nonprokaryotic

doi:10.1371/journal.pbio.0050016.g003

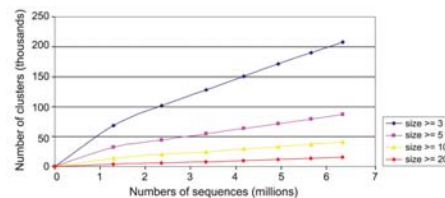
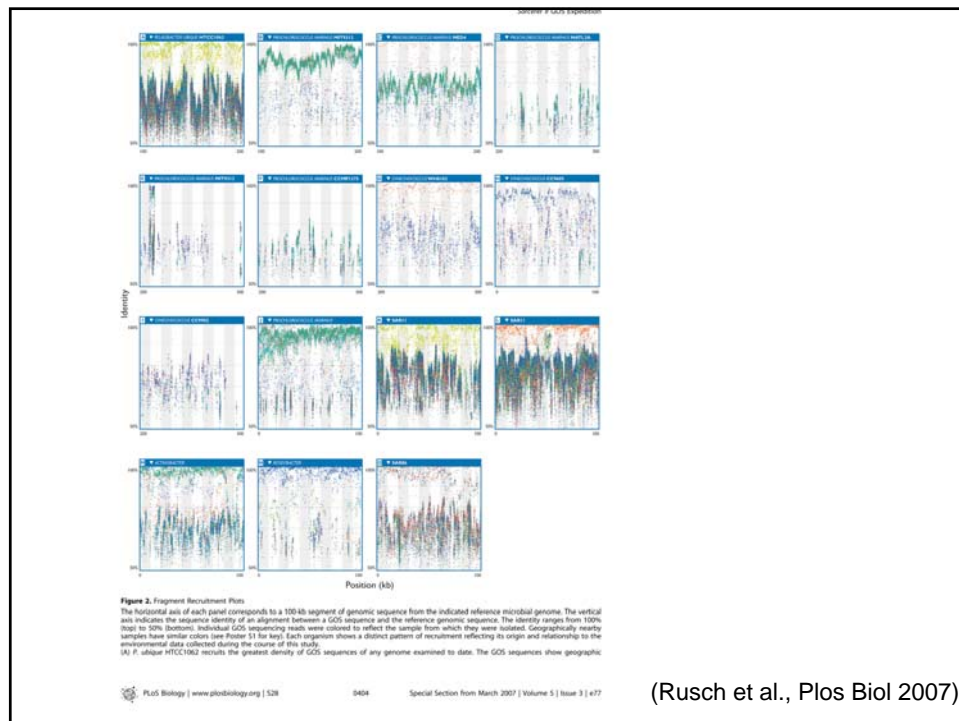
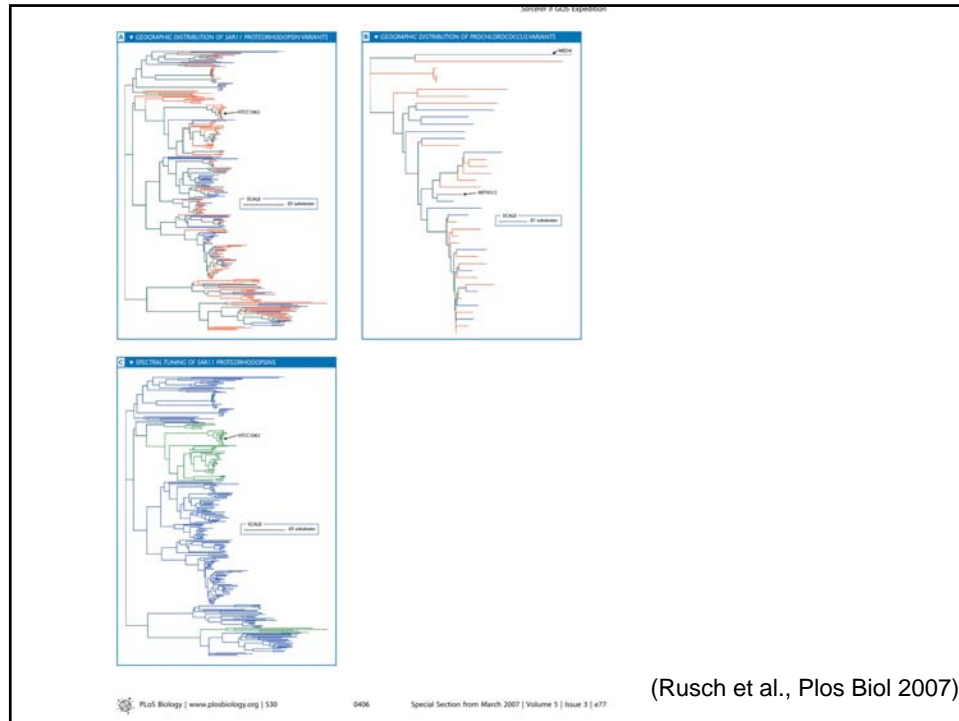


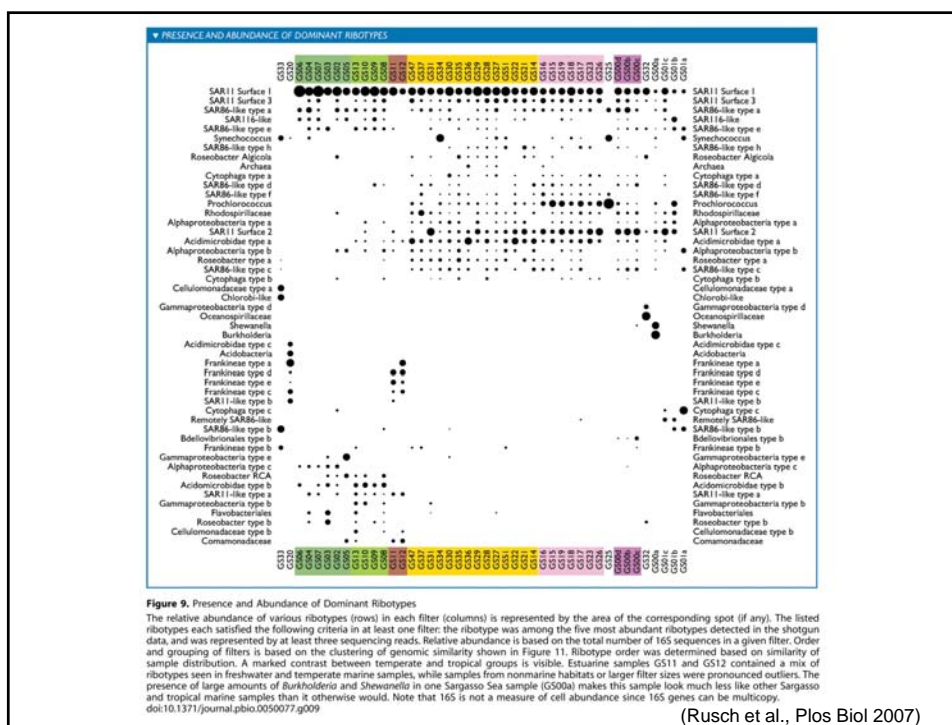
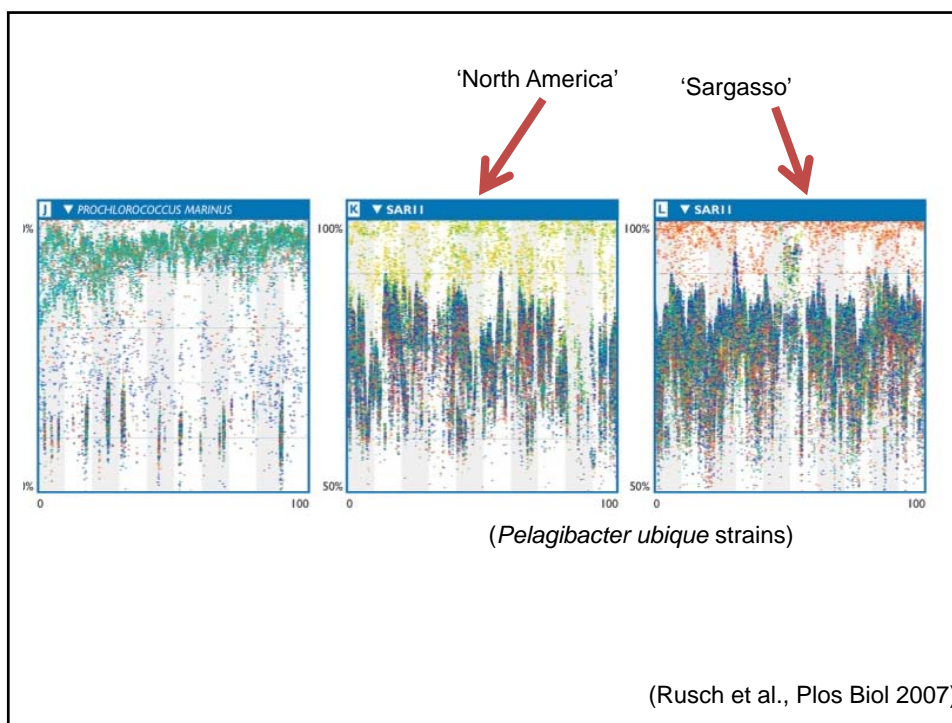
Figure 2. Rate of Discovery of Clusters as (Nonredundant) Sequences Are Added

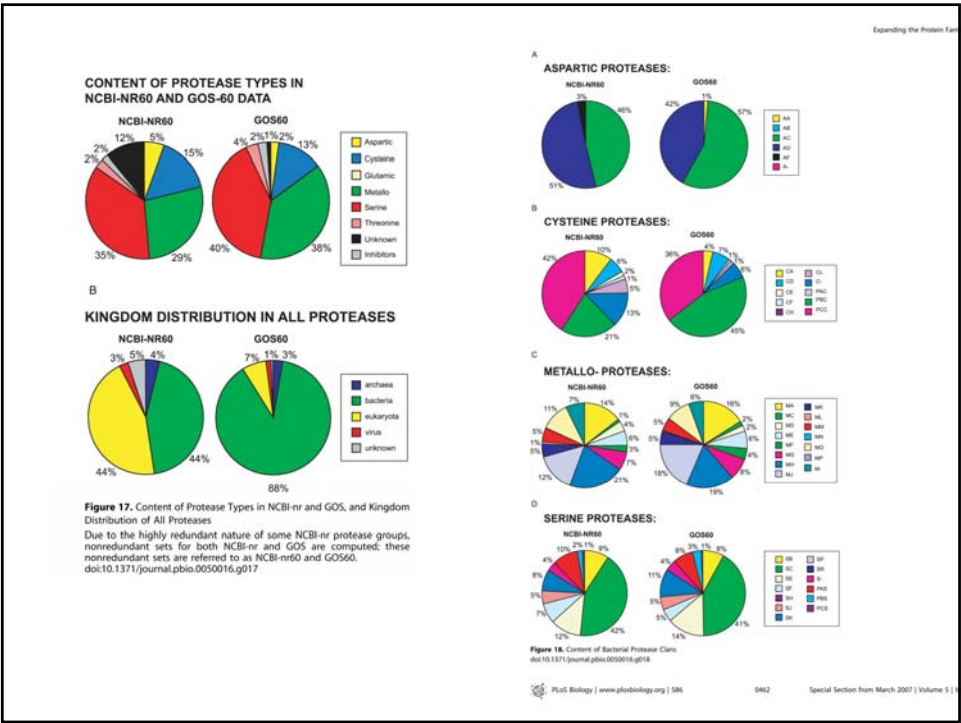
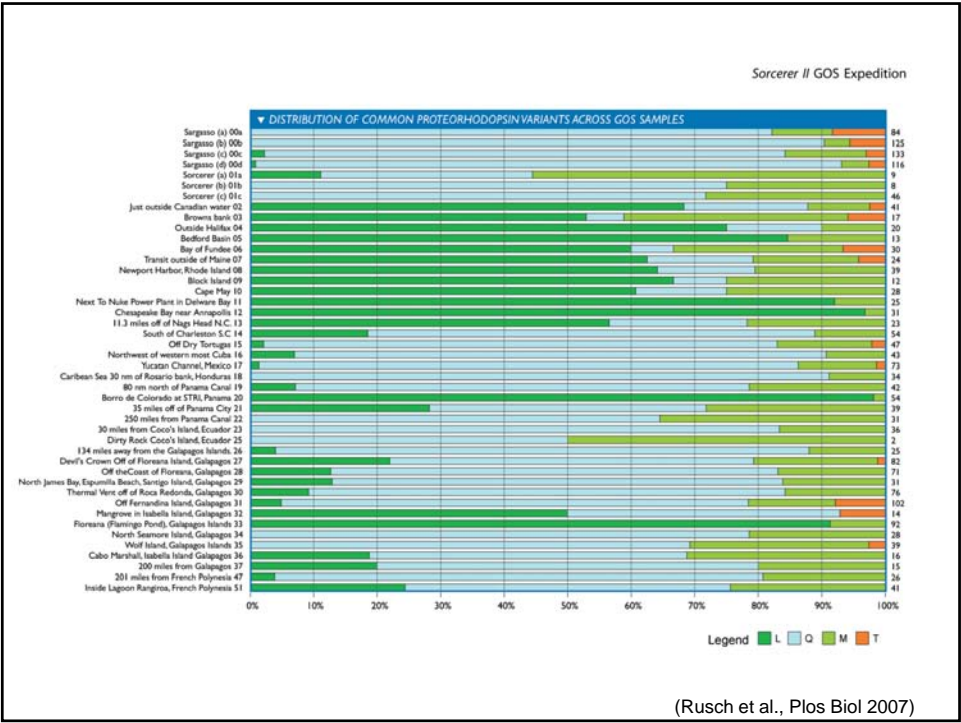
The x-axis denotes the number of sequences (in millions) and the y-axis denotes the number of clusters (in thousands). Seven datasets with increasing numbers of (nonredundant) sequences are chosen as described in the text. The blue curve shows the number of core sets of size ≥ 3 for the seven datasets. Curves for core set sizes ≥ 5 , ≥ 10 , and ≥ 20 are also shown. Linear regression gives slopes 0.027 ($R^2 = 0.999$), 0.011 ($R^2 = 0.999$), 0.0053 ($R^2 = 0.999$), and 0.0024 ($R^2 = 0.996$) for size ≥ 3 , size ≥ 5 , size ≥ 10 , and size ≥ 20 , respectively.

doi:10.1371/journal.pbio.0050016.g002

(Rusch et al., Plos Biol







LETTERS

(Warnecke et al., *Nature* 2007)

Metagenomic and functional analysis of hindgut microbiota of a wood-feeding higher termite

Falk Warnecke^{1*}, Peter Luginbühl^{2*}, Natalia Ivanova¹, Majid Ghassemin², Toby H. Richardson^{2†}, Justin T. Stege², Michelle Cayouette², Alice C. McHardy^{3†}, Gordana Djordjevic², Nahla Aboushadi², Rotem Sorek¹, Susannah G. Tringe¹, Mircea Podar¹, Hector Garcia Martin¹, Victor Kunin¹, Daniel Dalevi¹, Julita Madejska¹, Edward Kirton¹, Darren Platt¹, Ernest Szeto¹, Asaf Salamov¹, Kerrie Barry¹, Natalia Mikhailova¹, Nikos C. Kyrpides¹, Eric G. Matson¹, Elizabeth A. Ottesen⁶, Xinning Zhang², Myriam Hernández², Catalina Murillo², Luis G. Acosta², Isidore Rigoutsos¹, Giselle Tamayo², Brian D. Green², Cathy Chang², Edward M. Rubin¹, Eric J. Mathur^{2†}, Dan E. Robertson², Philip Hugenholtz¹ & Jared R. Leadbetter^{2*}

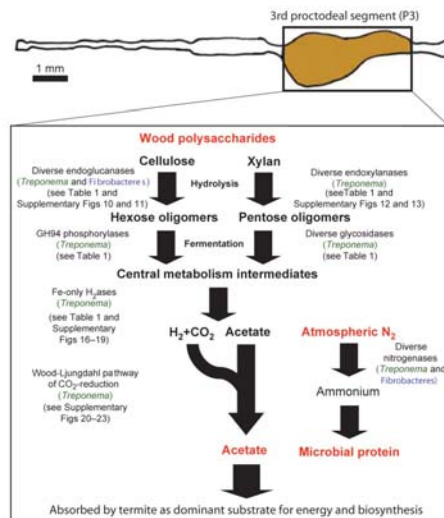
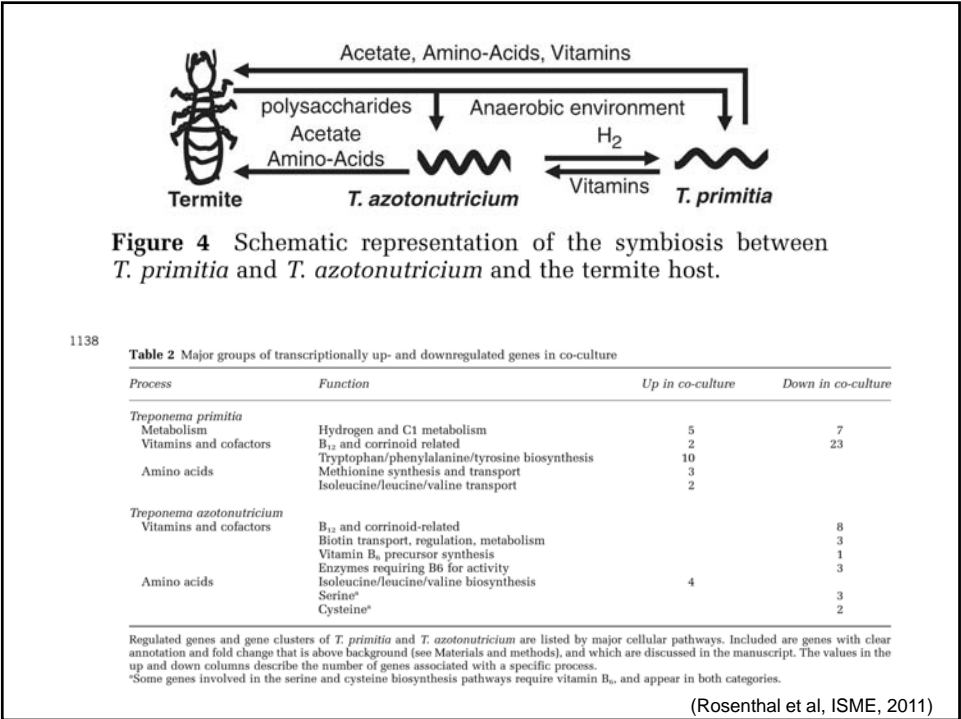
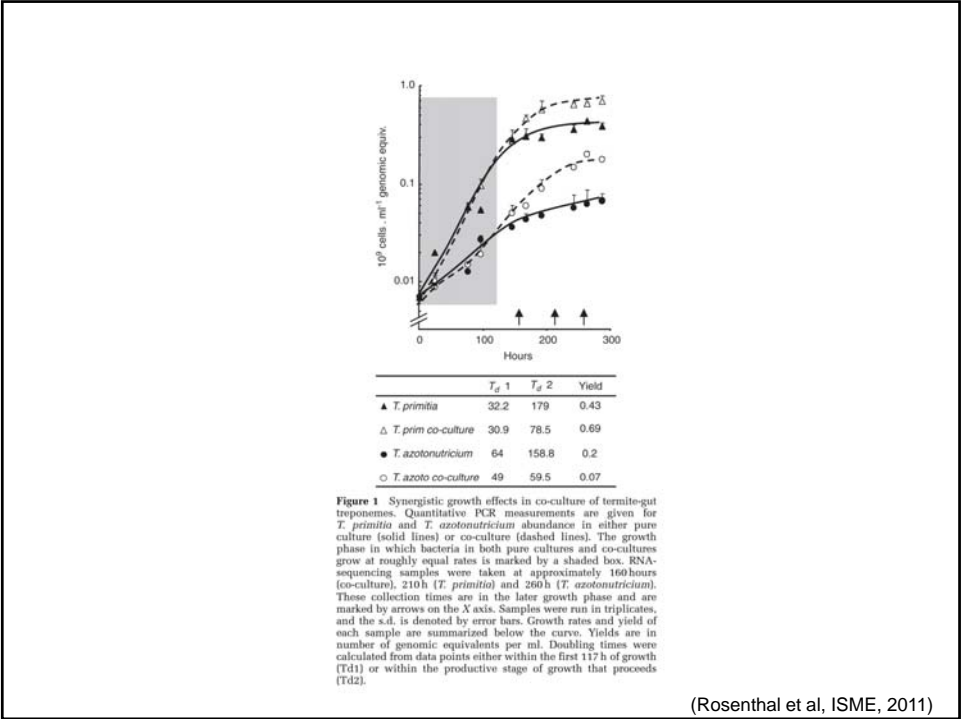


Figure 2 | Model of nutritional symbiosis-relevant metabolism by *Nasutitermes* P3 luminal bacteria. Wood is triturated by the insect's mandibles into small particles and predigested by poorly studied upstream processes before transit into the P3 compartment. The P3 lumen is dominated by diverse species of *Treponema* (Spirochaetes) and Fibrobacteres. There was no evidence for methanogenesis or lignin degradation in the metagenomic data set.

fishing Group

561
(Warnecke et al., *Nature* 2007)



Mammals are metagenomic in that they are composed of not only their own gene complements but also those of all of their associated microbes. To understand the coevolution of the mammals and their indigenous microbial communities, we conducted a network-based analysis of bacterial 16S ribosomal RNA gene sequences from the fecal microbiota of humans and 59 other mammalian species living in two zoos and in the wild. The results indicate that host diet and phylogeny both influence bacterial diversity, which increases from carnivory to omnivory to herbivory; that bacterial communities codiversified with their hosts; and that the gut microbiota of humans living a modern life-style is typical of omnivorous primates.

REPORTS



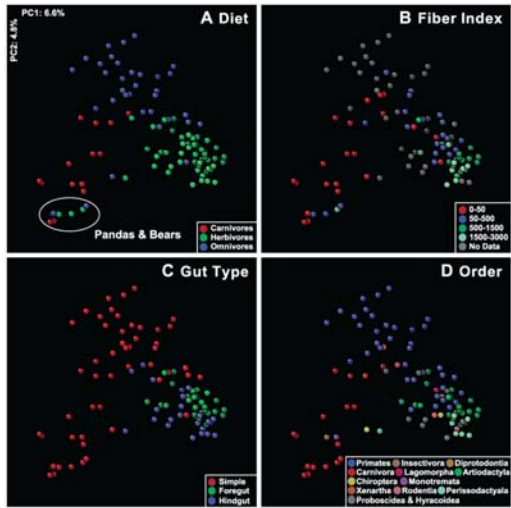
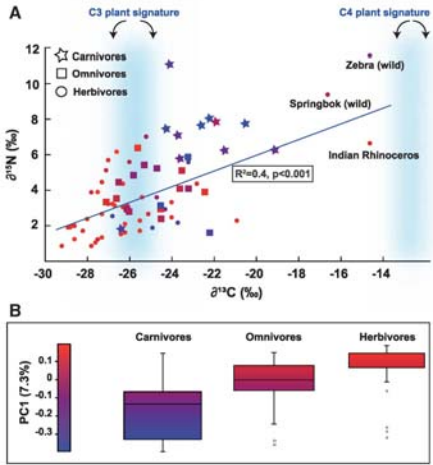


Fig. 2. Mammalian fecal bacterial communities clustered using principal coordinates analysis (PCoA) of the UniFrac metric matrix. PC1 and PC2 are plotted on x and y axes. Each circle corresponds to a fecal sample colored according to (A) diet, (B) diet fiber index, (C) gut morphology/physiology, and (D) host taxonomic order. The same data (samples) are shown in each panel. The percentage of the variation explained by the plotted principal coordinates is indicated on the axes.

(Ley et al., Science, 2008)

Fig. 3. Markers of trophic level mapped onto the variance in fecal microbial community diversity. (A) Stable isotope values for C and N plotted for each fecal sample, presented according to diet group. Symbols are colored according to their PC1 value; PC1 is the first principal coordinate of the PCoA of the un-weighted UniFrac metric. $\delta^{13}\text{C}$ ranges for C3 and C4 plants [per mil (‰)] are highlighted in blue. R^2 is for $\delta^{13}\text{C}$ versus $\delta^{15}\text{N}$. (B) Box plots are shown for the three diet groups (central line is the mean; box outline equals 1 SD; the bar denotes 2 SD; circles are outliers). The majority of fecal $\delta^{13}\text{C}$ values are intermediate between the average for C4 plants (−12.5‰) and C3 plants (−26.7‰).



(Ley et al., Science, 2008)

ARTICLES

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 Nielsen⁶, 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 Tap⁹, Patricia Lepage⁶, Marcelo Bertalan¹⁰, Jean-Michel Batto¹⁰, Torben Hansen¹⁰, Denis Le Paslier¹⁰, Allan Linneberg¹¹, H. Bjørn Nielsen¹¹, Eric Pelletier¹⁰, Pierre Renault¹⁰, Thomas Sicheritz-Ponten¹¹, 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 Brunak¹³, Joel Doré⁶, Francisco Guarner³, Karsten Kristiansen¹³, Oluf Pedersen^{13,14}, 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.

Vol 464|4 March 2010| doi:10.1038/nature08821

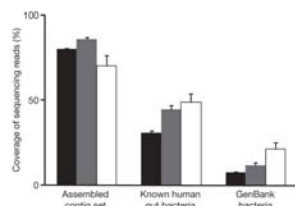


Figure 1 | Coverage of human gut microbiome. The three human microbial sequencing read sets—Illumina GA reads generated from 124 individuals in this study (black; $n = 124$), Roche/454 reads from 18 human twins and their mothers (grey; $n = 18$) and Sanger reads from 13 Japanese individuals (white; $n = 13$)—were aligned to each of the reference sequence sets. Mean values \pm s.e.m. are plotted.

AA

NATURE | Vol 464 | 4 March 2010

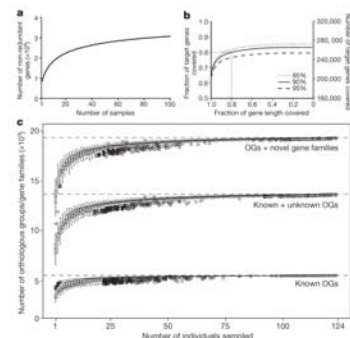
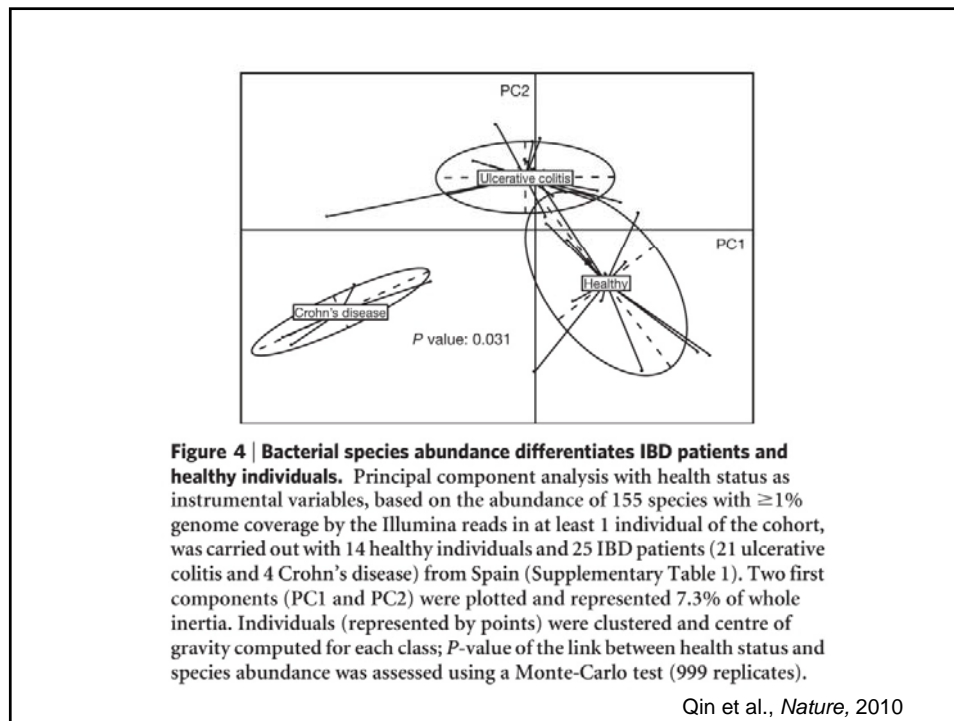
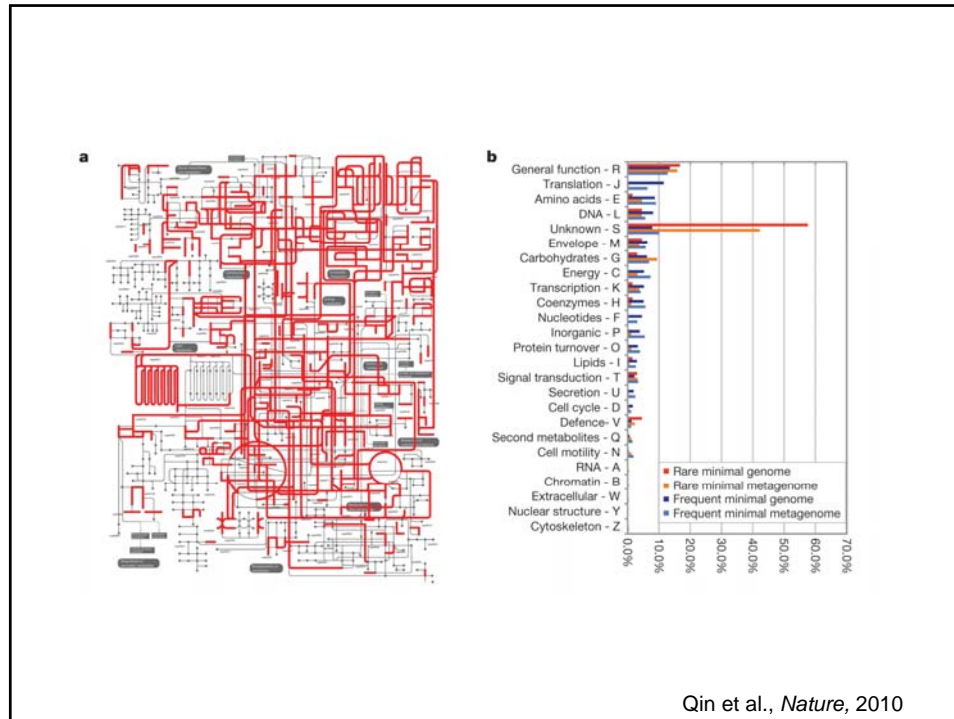


Figure 2 | Predicted ORFs in the human gut microbiome. **a**, Number of unique genes as a function of the extent of sequencing. The gene accumulation curve corresponds to the S_{obs} (Mao Tau) values (number of observed genes), calculated using EstimateS³¹ (version 8.2.0) on randomly chosen 100 samples (due to memory limitation). **b**, Coverage of genes from 89 frequent gut microbial species (Supplementary Table 12). **c**, Number of functions captured by number of samples investigated, based on known (well characterized) orthologous groups (OGs; bottom), known plus unknown orthologous groups (including, for example, putative, predicted, conserved hypothetical functions; middle) and orthologous groups plus novel gene families (>20 proteins) recovered from the metagenome (top). Boxes denote the interquartile range (IQR) between the first and third quartiles (25th and 75th percentiles, respectively) and the line inside denotes the median. Whiskers denote the lowest and highest values within 1.5 times IQR from the first and third quartiles, respectively. Circles denote outliers beyond the whiskers.

Qin et al., *Nature*, 2010



MICROBIAL ECOLOGY

Human gut microbes associated with obesity

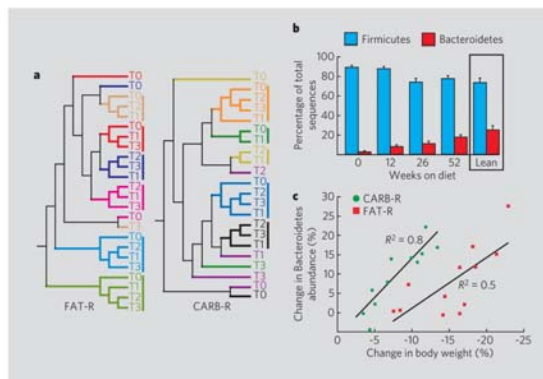


Figure 1 | Correlation between body-weight loss and gut microbial ecology. **a**, Clustering of 16S ribosomal RNA gene sequence libraries of faecal microbiota for each person (in different colours) and time point in diet therapy (T0, baseline; T1, 12 weeks; T2, 26 weeks; T3, 52 weeks) in the two diet-treatment groups (fat restricted, FAT-R; carbohydrate restricted, CARB-R), based on UniFrac analysis of the 18,348-sequence phylogenetic tree. **b**, Relative abundance of Bacteroidetes and Firmicutes. For each time point, values from all available samples were averaged (n was 11 or 12 per time point). Lean-subject controls include four stool samples from two people taken 1 year apart, plus three other stool samples. Mean values \pm s.e. are plotted. **c**, Change in relative abundance of Bacteroidetes in subjects with weight loss above a threshold of 2% weight loss for the CARB-R diet and 6% for the FAT-R diet.

Ley et al., NATURE|Vol 444|21/28 December 2006

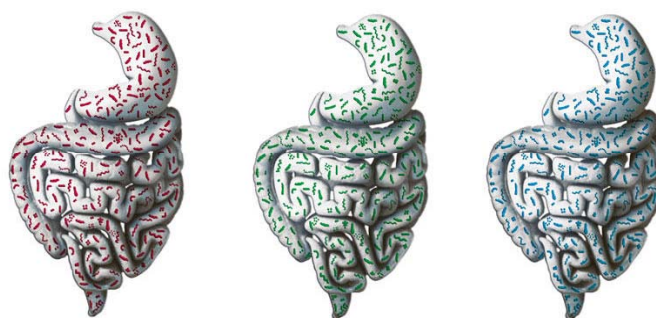
BRIEF COMMUNICATIONS

ARTICLE

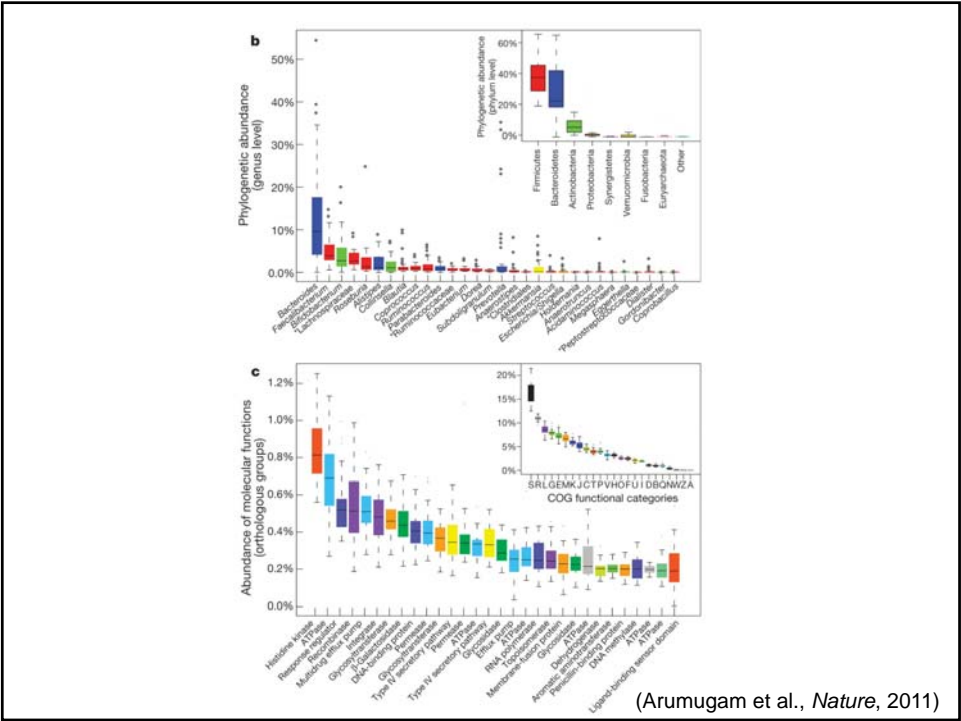
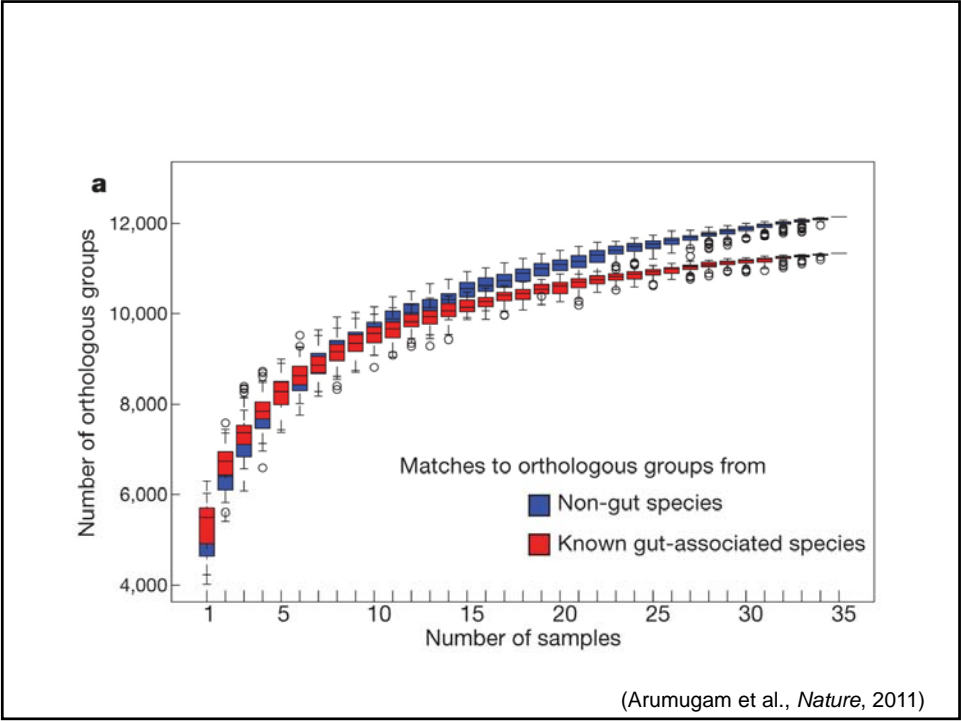
doi:10.1038/nature09944

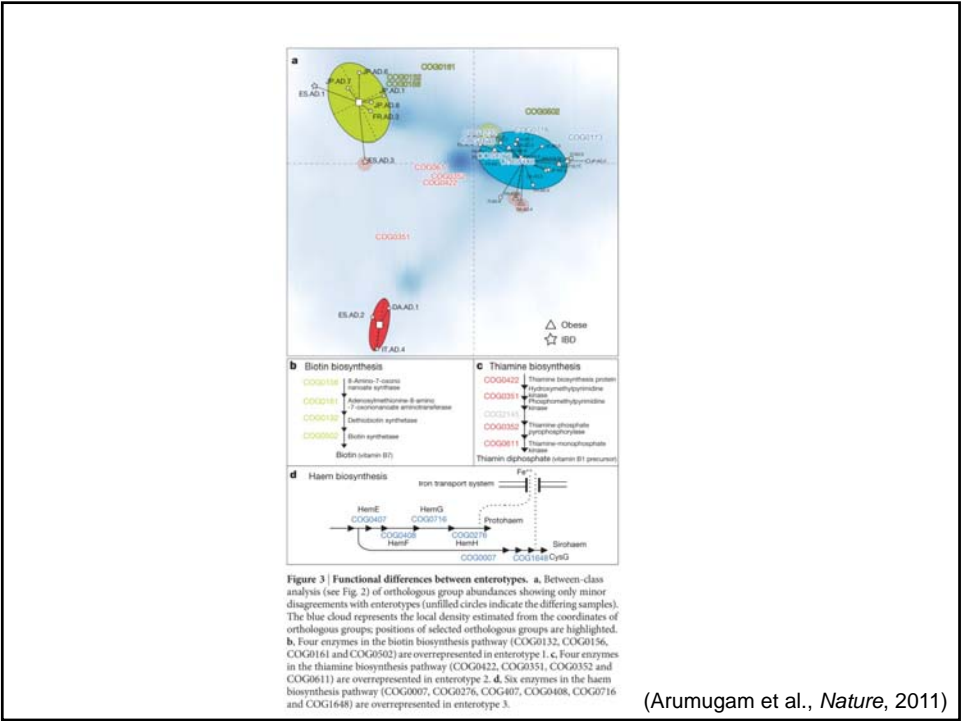
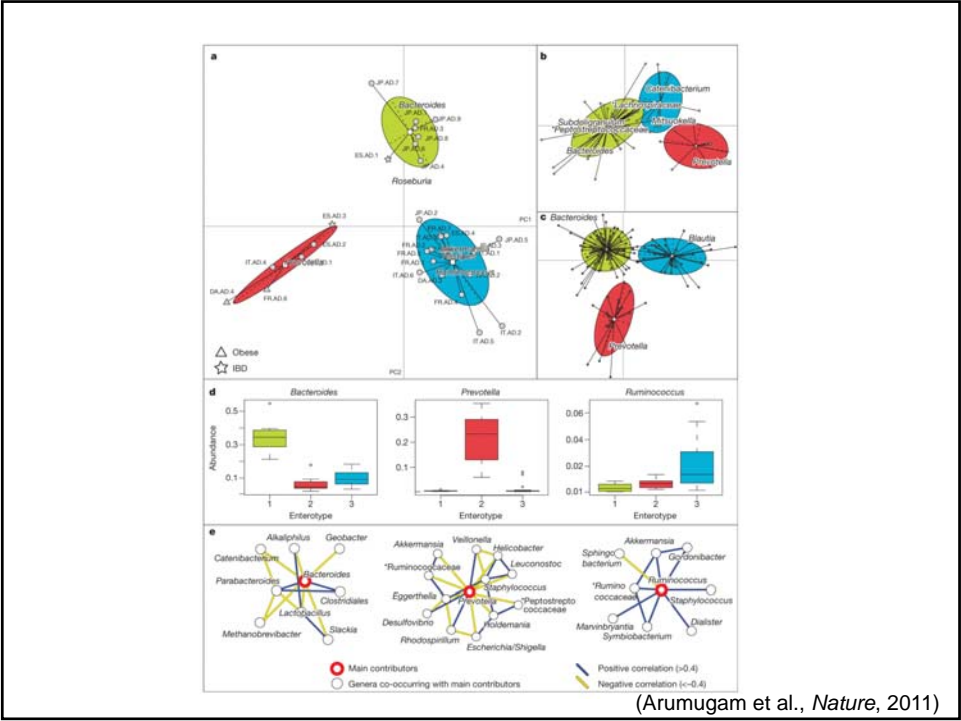
Enterotypes of the human gut microbiome

Manimozhayan Arumugam^{1*}, Jeroen Raes^{1,2*}, Eric Pelletier^{3,4,5}, Denis Le Paslier^{3,4,5}, Takuji Yamada¹, Daniel R. Mende¹, Gabriel R. Fernandes^{1,6}, Julien Tap^{1,7}, Thomas Bruns^{3,4,5}, Jean-Michel Batto⁷, Marcelo Bertalan⁸, Natalia Borrue⁹, Francesc Casellas¹⁰, Leyden Fernandez¹⁰, Laurent Gautier⁸, Torben Hansen^{11,12}, Masahira Hattori¹³, Tetsuya Hayashi¹⁴, Michiel Kleerebezem¹⁵, Ken Kurokawa¹⁶, Marion Leclerc⁷, Florence Levenez⁷, Chaysavanh Manichanh⁹, H. Bjørn Nielsen⁹, Trine Nielsen¹¹, Nicolas Pons⁷, Julie Poulain⁷, Junjie Qin¹⁷, Thomas Sicheritz-Ponten^{8,18}, Sebastian Tims¹³, David Torrents^{10,19}, Edgardo Ugarte¹, Erwin G. Zoetendal¹⁵, Jun Wang^{17,20}, Francisco Guarner⁹, Oluf Pedersen^{11,21,22,23}, Willem M. de Vos^{15,24}, Søren Brunak⁸, Joel Doré², MetaHIT Consortium†, Jean Weissenbach^{3,4,5}, S. Dusko Ehrlich⁷ & Peer Bork^{1,25}



(Arumugam et al., Nature, 2011)





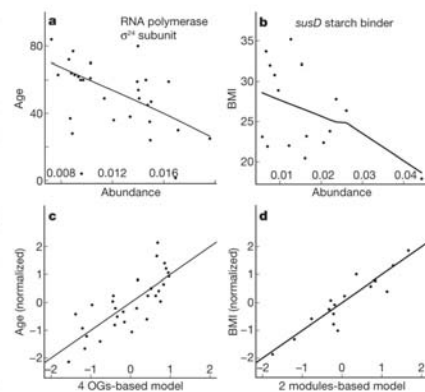


Figure 4 | Correlations with host properties. **a**, Pairwise correlation of RNA polymerase facultative σ^{28} subunit (COG1595) with age ($P = 0.03$, $\rho = -0.59$). **b**, Pairwise correlation of SusD, a family of proteins that bind glycan molecules before they are transported into the cell, and BMI ($P = 0.27$, $\rho = -0.29$, weak correlation). **c**, Multiple orthologous groups (OGs) (COG0085, COG0086, COG0438 and COG0739; see Supplementary Table 18) significantly correlating with age when combined into a linear model (see Supplementary Methods section 13 and ref. 40 for details; $P = 2.75 \times 10^{-5}$, adjusted $R^2 = 0.57$). **d**, Two modules, ATPase complex and ectoine biosynthesis (M00051), significantly correlating with BMI when combined into a linear model ($P = 6.786 \times 10^{-6}$, adjusted $R^2 = 0.82$).

(Arumugam et al., *Nature*, 2011)

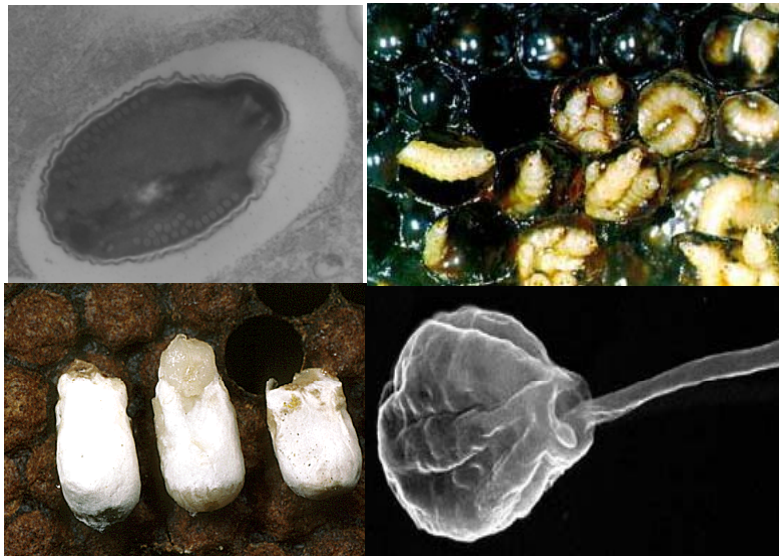
- ❖ FIRST PRINCIPLES
- ❖ DATA COLLECTION
- ❖ OCEANS, TERMITES, MAMMALS
- ❖ BEE STORIES
- ❖ RETURN TO FUNCTION
- ❖ ANALYSIS



HONEY BEE PARASITES



AND MORE HONEY BEE PARASITES



The New York Times
nytimes.com

PRINTER-FRIENDLY FORMAT
SPONSORED BY CHI

February 27, 2007

Honeybees Vanish, Leaving Keepers in Peril

By [ALEXEI BARRIONUEVO](#)

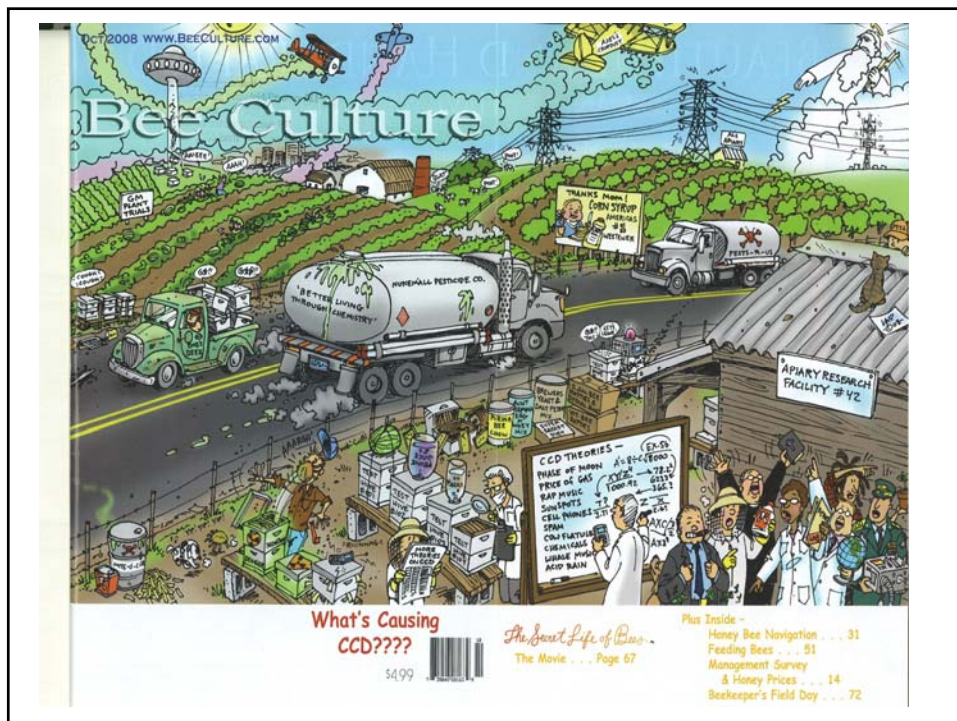
VISALIA, Calif., Feb. 23 — David Bradshaw has endured countless stings during his life as a beekeeper, but he got the shock of his career when he opened his boxes last month and found half of his 100 million bees missing.

In 24 states throughout the country, beekeepers have gone through similar shocks as their bees have been disappearing inexplicably at an alarming rate, threatening not only their livelihoods but also the production of numerous crops, including California almonds, one of the nation's



CCD TRAITS

'RAPID' WORKER LOSS
NO DEAD BODIES
EARLY SPRING
PATCHY IN SPACE/TIME



HORIZONTAL TRANSMISSION + STRESS?



A Metagenomic Survey of Microbes in Honey Bee Colony Collapse Disorder

Diana L. Cox-Foster,¹ Sean Conlan,² Edward C. Holmes,^{3,4} Gustavo Palacios,² Jay D. Evans,⁵ Nancy A. Moran,⁶ Phenix-Lan Quan,² Thomas Brieese,² Mady Hornig,² David M. Geiser,⁷ Vince Martinson,⁸ Dennis vanEngelsdorp,^{1,9} Abby L. Kalkstein,¹ Andrew Drysdale,² Jeffrey Hui,² Junhui Zhai,² Liwang Cui,¹ Stephen K. Hutchison,¹⁰ Jan Fredrik Simons,¹⁰ Michael Egholm,¹⁰ Jeffery S. Pettis,⁵ W. Ian Lipkin^{2*}

In colony collapse disorder (CCD), honey bee colonies inexplicably lose their workers. CCD has resulted in a loss of 50 to 90% of colonies in beekeeping operations across the United States. The observation that irradiated combs from affected colonies can be repopulated with naive bees suggests that infection may contribute to CCD. We used an unbiased metagenomic approach to survey microflora in CCD hives, normal hives, and imported royal jelly. Candidate pathogens were screened for significance of association with CCD by the examination of samples collected from several sites over a period of 3 years. One organism, Israeli acute paralysis virus of bees, was strongly correlated with CCD.

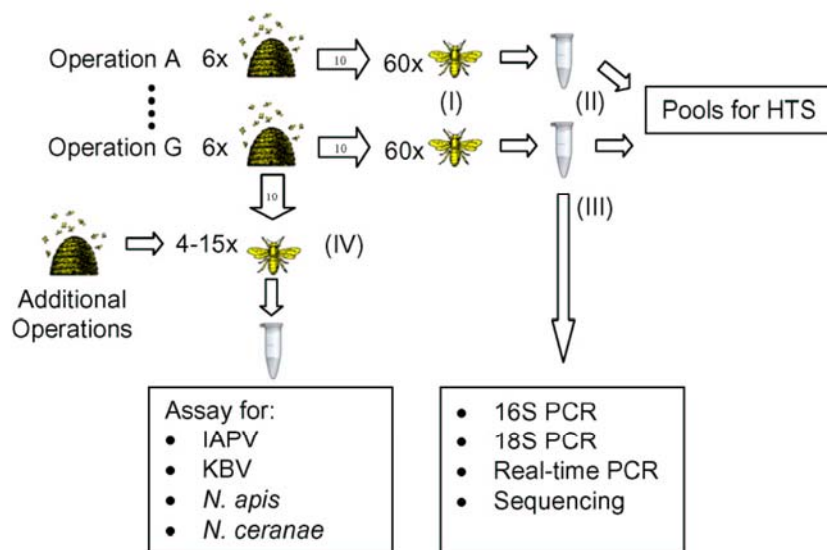


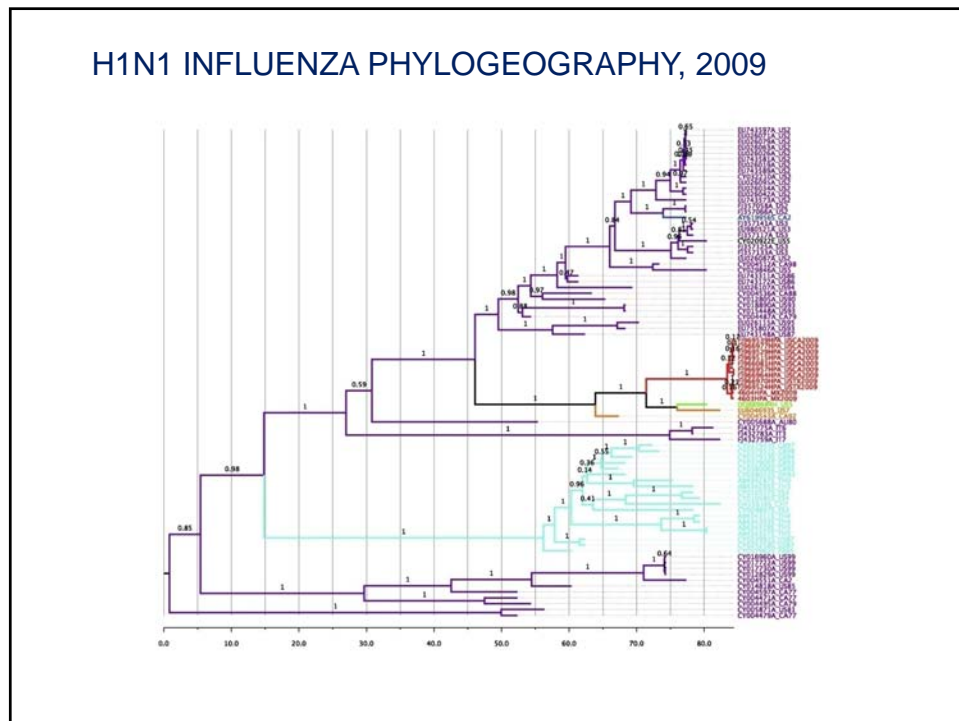
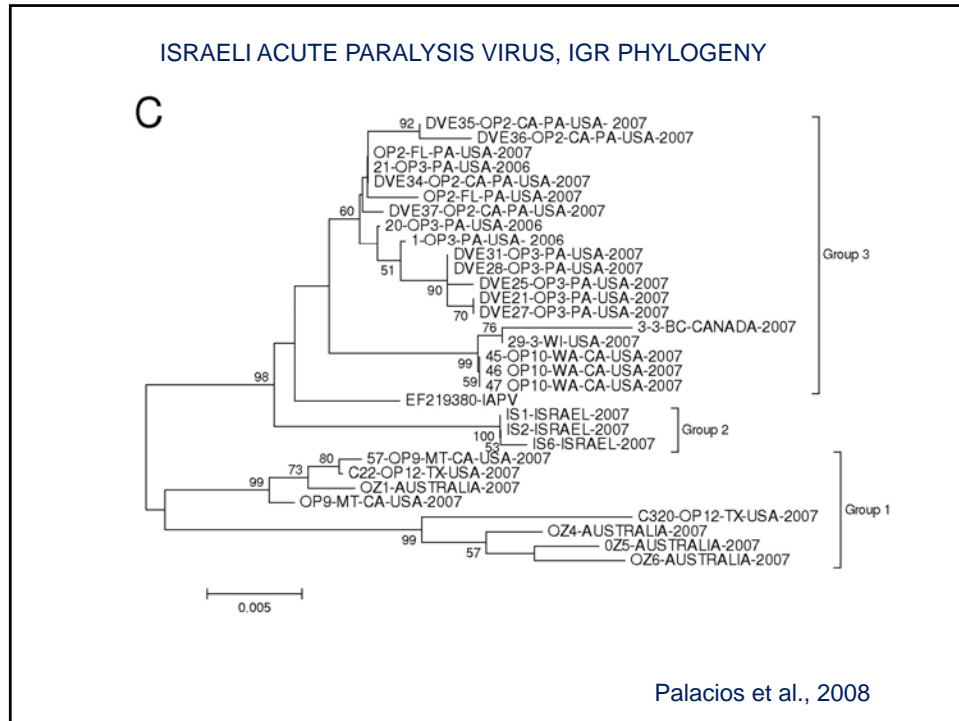
Table 1. Closest sequenced relatives identified through BLAST analysis of the high-throughput sequence data.

Kingdom	Taxon (rank)	Organism
Bacteria	Firmicutes (phylum)	<i>Lactobacillus</i> sp.*†
		Uncultured Firmicutes†
Bacteria	Actinobacteria (class)	<i>Bifidobacterium</i> sp.*
Bacteria	Alphaproteobacteria (class)	<i>Bartonella</i> sp.*†
		<i>Gluconacetobacter</i> sp.*†
Bacteria	Betaproteobacteria (class)	<i>Simonsiella</i> sp.*†
Bacteria	Gammaproteobacteria (class)	Two uncultured species*†
Fungus	Entomophthorales (order)	<i>Pandora delphacis</i>
Fungus	Mucorales (order)	<i>Mucor</i> spp.
Fungus/microsporidian	Nosematidae (family)	<i>Nosema ceranae</i>
Fungus/microsporidian	Nosematidae (family)	<i>Nosema apis</i>
Eukaryota	Trypanosomatidae (family)	<i>Leishmania/Leptomonas</i> sp.
Metazoan	Varroidae (family)	<i>Varroa destructor</i>
Virus	(Unclassified)	CBPV‡
Virus	<i>Iflavirus</i> (genus)	SBV
Virus	<i>Iflavirus</i> (genus)	DWV‡
Virus	Dicistroviridae (family)	BQCV
Virus	Dicistroviridae (family)	KBV‡
Virus	Dicistroviridae (family)	ABPV
Virus	Dicistroviridae (family)	IAPV of bees‡

*Found by Jeyaprakash *et al.* (10). †Found by Babendreier *et al.* (9). ‡Indicates viruses not yet classified by the International Committee on the Taxonomy of Viruses but that exhibit the key features of the indicated taxon.

IDENTIFICATION OF CANDIDATES TO PURSUE

Agent	Number of positive samples <i>n</i> (% positive of samples tested)			Positive Predictive Value (%)	Sensitivity (%)	Specificity (%)
	CCD (<i>n</i> = 30)	non-CCD (<i>n</i> = 21)	Total (<i>n</i> = 51)			
IAPV	25 (83.3%)	1 (4.8%)	26 (51.0%)	96.1	83.3	95.2
KBV	30 (100%)	16 (76.2%)	46 (90.2%)	65.2	100	23.8
<i>N. apis</i>	27 (90%)	10 (47.6%)	37 (72.5%)	73.0	90.0	52.4
<i>N. ceranae</i>	30 (100%)	17 (80.9%)	47 (92.1%)	63.8	100	19.0
All 4 agents	23 (76.7%)	0 (0%)	23 (45.0%)	100	76.7	100

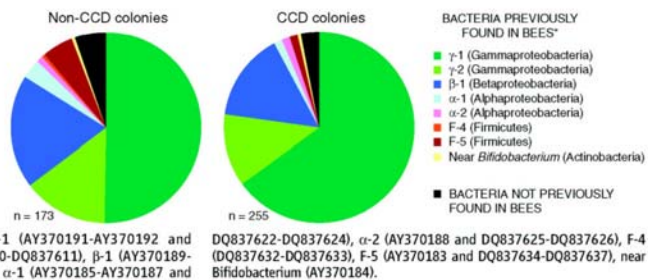


1) Identification of Candidates to Pursue

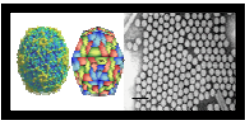
Agent	Number of positive samples <i>n</i> (% positive of samples tested)			Positive Predictive Value (%)	Sensitivity (%)	Specificity (%)
	CCD (<i>n</i> = 30)	non-CCD (<i>n</i> = 21)	Total (<i>n</i> = 51)			
IAPV	25 (83.3%)	1 (4.8%)	26 (51.0%)	96.1	83.3	95.2
KBV	30 (100%)	16 (76.2%)	46 (90.2%)	65.2	100	23.8
<i>N. apis</i>	27 (90%)	10 (47.6%)	37 (72.5%)	73.0	90.0	52.4
<i>N. ceranae</i>	30 (100%)	17 (80.9%)	47 (92.1%)	63.8	100	19.0
All 4 agents	23 (76.7%)	0 (0%)	23 (45.0%)	100	76.7	100

(Cox-Foster et al., *Science*, 2007)

Fig. 1. Summary of bacterial groups from *A. mellifera* derived from colonies categorized as non-CCD and CCD. For both categories, the top BLAST hit for over 96% of sequences from 16S rRNA clones was a sequence obtained in previous studies on bacterial associates of *A. mellifera*. Asterisk indicates that the bacteria were categorized according to the cluster designations of Babendreier et al. (9) [except for the *Bifidobacterium*-like sequence of Jeyaprakash et al. (10)]. *n*, number of sequences. GenBank accession numbers corresponding to the categories are: γ -1 (AY370191-AY370192 and DQ837602-DQ837609), γ -2 (DQ837610-DQ837611), β -1 (AY370189-AY370190 and DQ837616-DQ837621), α -1 (AY370185-AY370187 and

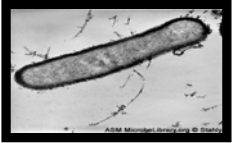


(Cox-Foster et al., *Science*, 2007)




Viruses
 ABPV
 BQCV
 CBPV
 DWV
 IAPV
 KBV
 SBV
 VDV


Beenomics 2011



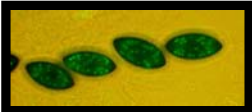
*American foulbrood
bacterium*



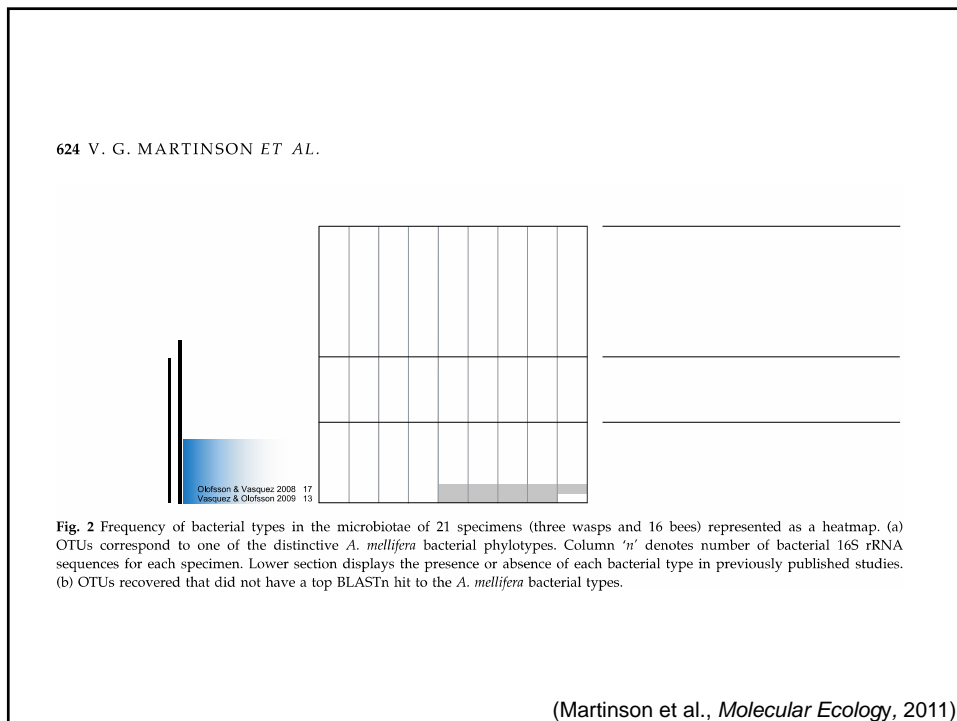
Chalkbrood fungus
9/9/11

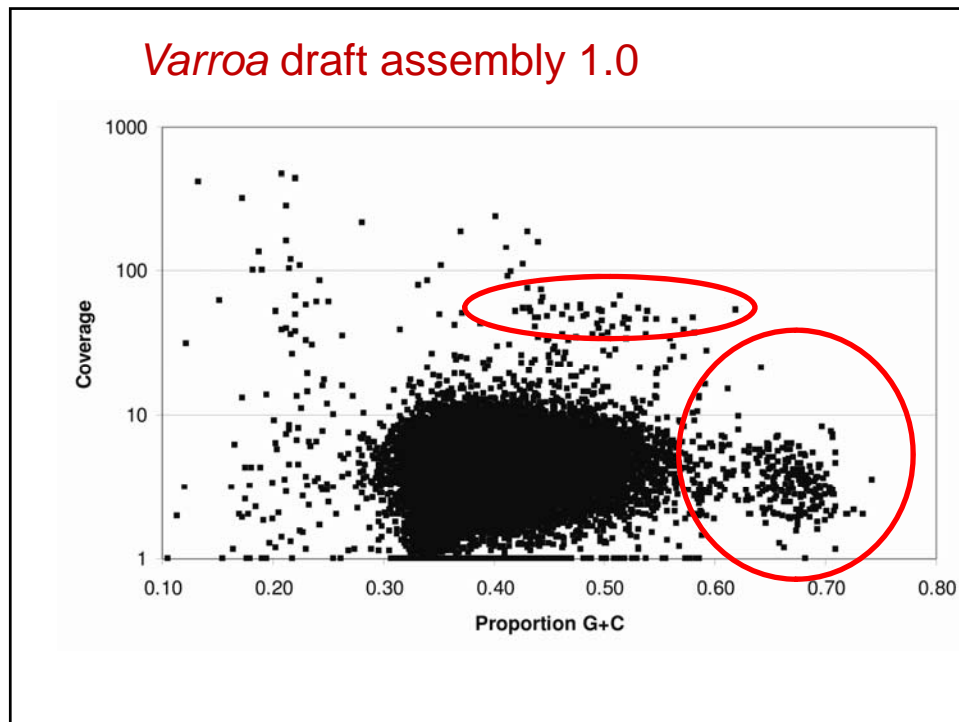


Varroa mites



Nosema ceranae
Nosema apis



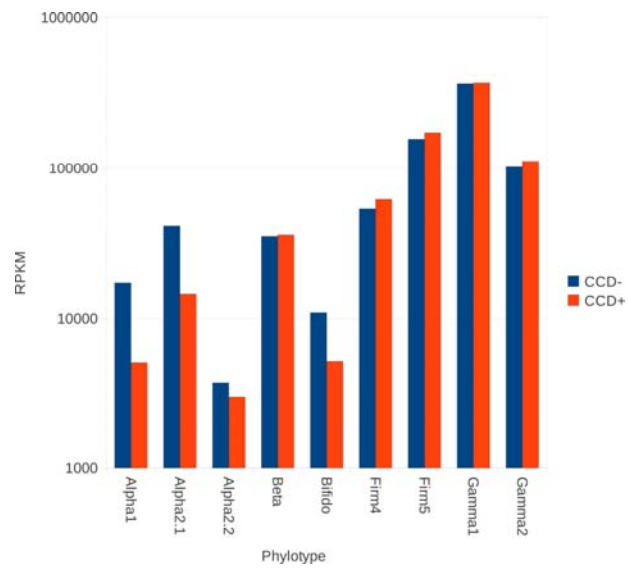


ILLUMINA CCD +/- survey

Healthy (n=63) and collapsed (n=61) colonies sampled in 2007 from eastern and western U.S.

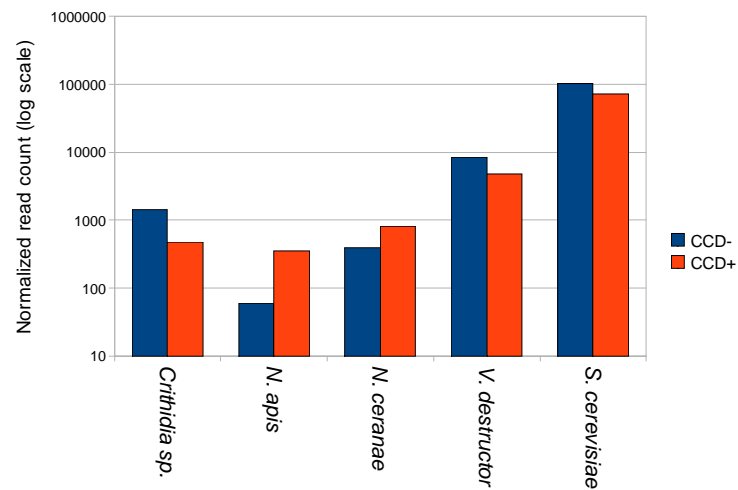
- 8 workers each, random-primed cDNA libraries
- Illumina RNA single-end (CCD-) and paired-end sequencing (CCD+), one lane each
- >19 million reads for CCD-, >20 million paired reads for CCD+
- Assembly with Velvet, mapping with Bowtie

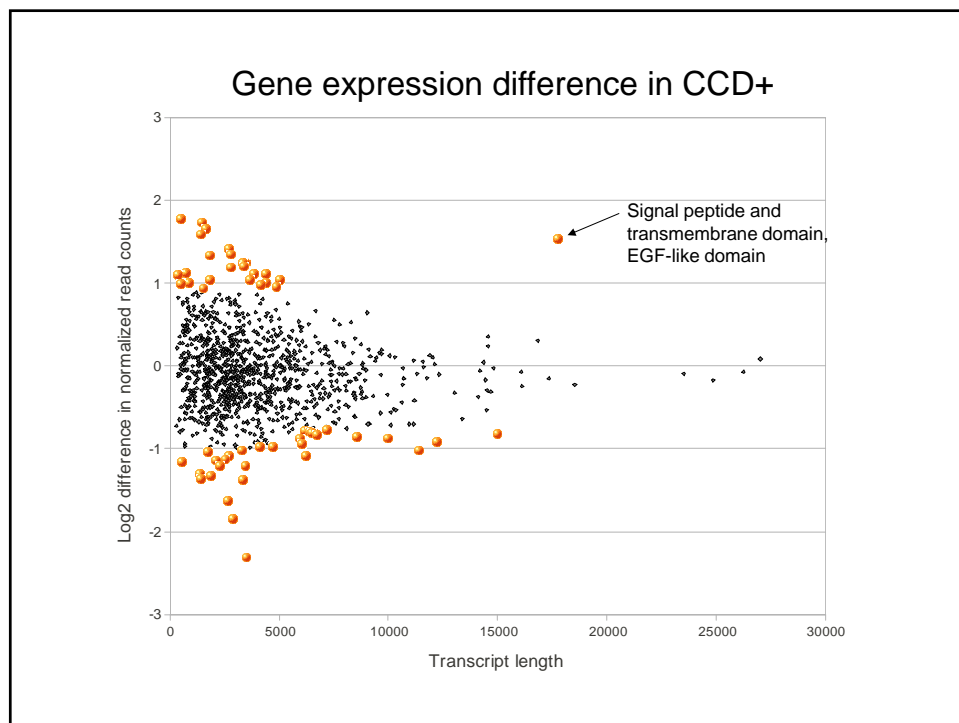
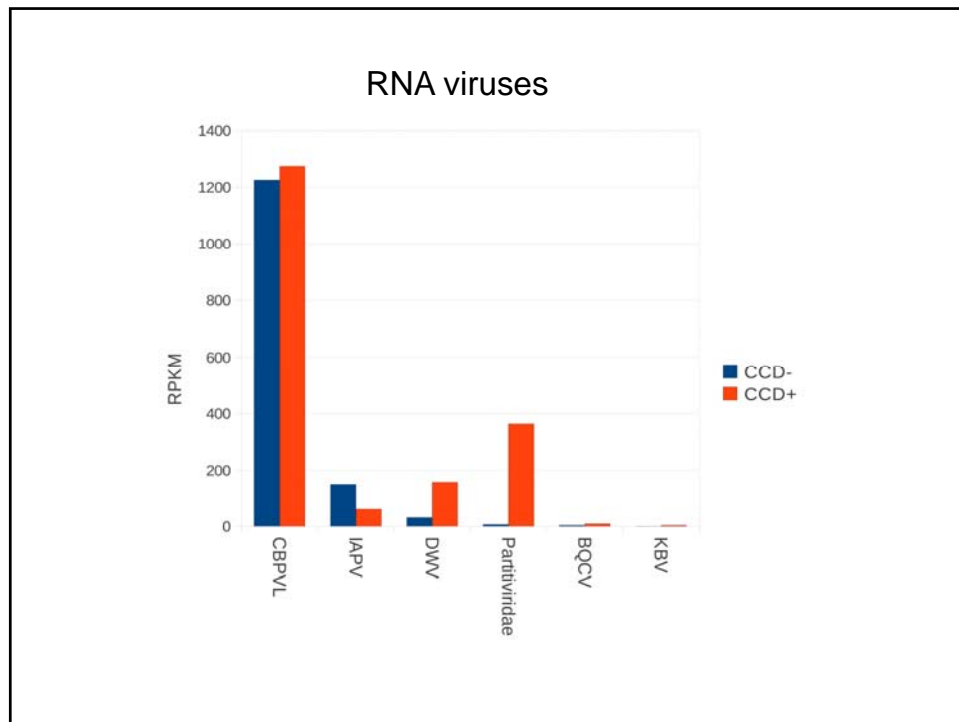
Major honey bee gut bacteria



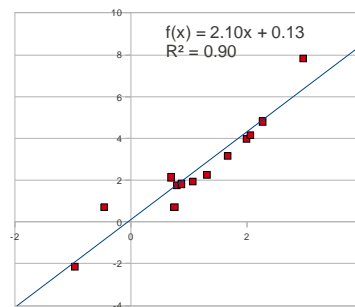
OTUs from Martinson et al. 2011 Molecular Ecology 20:619

Eukaryotic associates





Gene	Log2Diff	Homology	Class
GB11923	1.8	HSP20-like	Chaperone
GB10737	1.9	Apidermin2	Cuticle
GB10674	2.1	Apidermin3	Cuticle
GB30202	2.3	Apidermin1	Cuticle
GB19612	1.5	Paxillin	Focal adhesion-associated adaptor
GB17331	1.8	Peritrophic matrix	Peritrophic matrix
GB16446	-1.5	RBM5	RNA binding protein
GB13519	-2.1	RAP55	mRNA silencing/degradation



Gene	Log2Diff	Homology	Class
GB11923	1.8	HSP20-like	Chaperone
GB10737	1.9	Apidermin2	Cuticle
GB10674	2.1	Apidermin3	Cuticle
GB30202	2.3	Apidermin1	Cuticle
GB19612	1.5	Paxillin	Focal adhesion-associated adaptor
GB17331	1.8	Peritrophic matrix	Peritrophic matrix
GB16446	-1.5	RBM5	RNA binding protein
GB13519	-2.1	RAP55	mRNA silencing/degradation

Evidence of age structure (older bees missing in CCD), consistent with CCD phenotype?

Gene	Log2Diff	Homology	Class
GB11923	1.8	HSP20-like	Chaperone
GB10737	1.9	Apidermin2	Cuticle
GB10674	2.1	Apidermin3	Cuticle
GB30202	2.3	Apidermin1	Cuticle
GB19612	1.5	Paxillin	Focal adhesion-associated adaptor
GB17331	1.8	Peritrophic matrix	Peritrophic matrix
GB16446	-1.5	RBM5	RNA binding protein
GB13519	-2.1	RAP55	mRNA silencing/degradation

Related to higher viral loads in CCD+? RBM5 promotes apoptosis, RAP55 silences mRNA within cytoplasmic P bodies. Both functions are potential responses to viral stress.

Conclusions

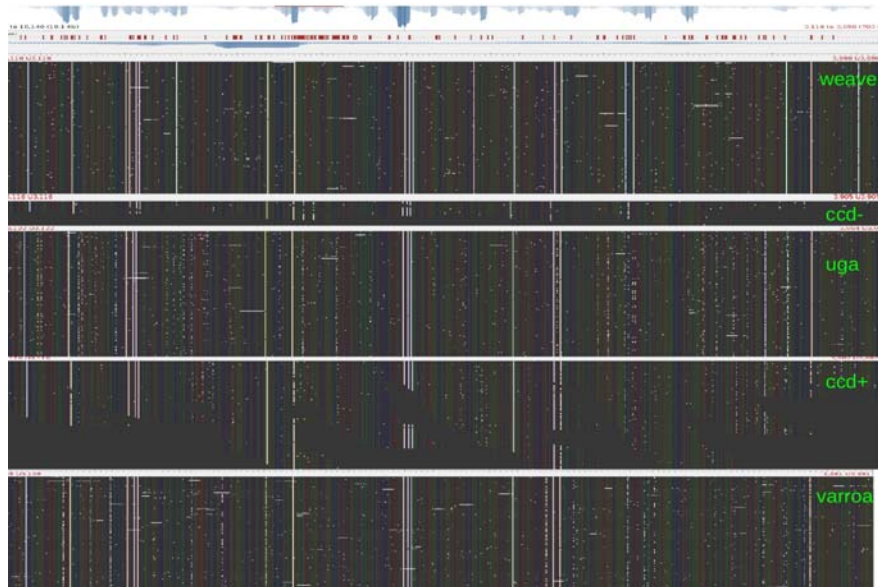
CCD survey:

- Bacterial signature
- Known viruses increased, novel viruses found
- *Nosema* increased
- No strong immune/detox signal

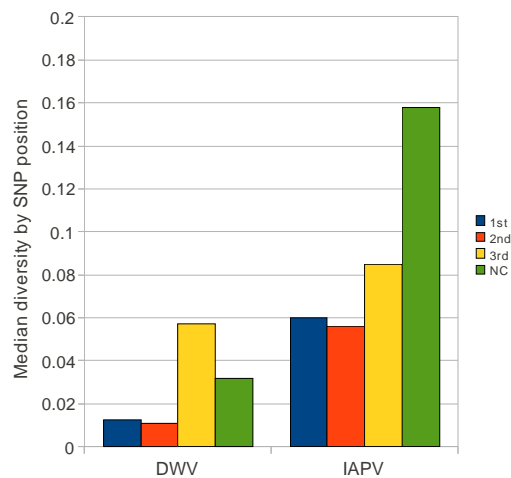
Viral polymorphism

Quantify genetic variation among viral strains of the same species

- Virulence genotypes
- Recombination, chimeras
- Population structure
- Host-specific genotypes



Lower noncoding polymorphism in DWV



Deformed wing virus

	Median per-site Fst				
bp					
hfb	0.002				
varroa	0.001	0.002			
ugasick	0.001	0.001	0.002		
weaver-	0.001	0.001	0.001	0.001	
ccd+	0.012	0.008	0.013	0.010	0.011

Deformed wing virus

	Median per-site Fst				
bp					
hfb	0.002				
varroa	0.001	0.002			
ugasick	0.001	0.001	0.002		
weaver-	0.001	0.001	0.001	0.001	
ccd+	0.012	0.008	0.013	0.010	0.011

Israel acute paralysis virus

	Median per-site Fst	
bp		
ccd-	0.83542	
ccd+	0.17019	0.3333



- ❖ FIRST PRINCIPLES
- ❖ DATA COLLECTION
- ❖ OCEANS, TERMITES, MAMMALS
- ❖ BEE STORIES
- ❖ RETURN TO FUNCTION
- ❖ ANALYSIS



Social evolution in multispecies biofilms

Sara Mitri^{1,2,3}, João B. Xavier², and Kevin R. Foster^{1,2,3}

¹Department of Zoology, University of Oxford, Oxford OX1 3PS, United Kingdom; ²Oxford Centre for Integrative Systems Biology, Oxford University, Oxford OX1 3QU, United Kingdom; and ³Program in Computational Biology, Memorial Sloan-Kettering Cancer Center, New York, NY 10065

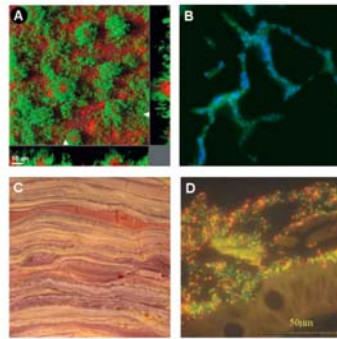


Fig. 1. Microbial diversity: examples of natural microbial communities. (A) A two-species bacterial biofilm cultivated in the laboratory in which one strain evolves to increase its exploitation of the other. Adapted by permission from Macmillan Publishers Ltd: *Nature* (78), copyright 2007. (B) A two-strain bacterial aggregate detected on a bean leaf surface (magnification 500 \times) [Appl Environ Microbiol (2005) 71(9):5484–5493, 10.126/AEM.71.9.5484–5493.2005. Reproduced with permission from the American Society for Microbiology] (79). (C) Stromatolite fossil that is ~2 billion y old. Modern stromatolites consist of multilayered sheets of microorganisms, and are a good example of very diverse, yet spatially structured microbial communities (copyright Merv Feick, <http://www.Indianafossils.com>). (D) The detection of two of the species present in a bacterial biofilm covering the intestinal mucosa of a self-limiting colitis patient, imaged using triple-color fluorescence in situ hybridization [J Clin Microbiol (2005) 43(7):3380–3389, 10.1128/JCM.43.7.3380–3389.2005. Reproduced with permission from the American Society for Microbiology] (80).

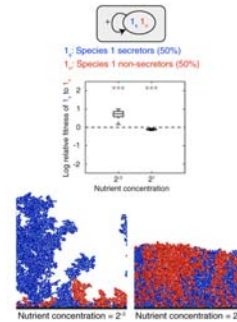


Fig. 2. Secretors and nonsecretors of a single species. Equal proportions of two strains of the same species are inoculated and left to grow to a fixed total biomass. Strain 1_s secretes a product that benefits both strains. 1_n does not secrete the product. Product secretion incurs a cost of 30% of the cells' growth rate. Boxplots show log relative fitness (Materials and Methods) of secreting to nonsecreting cells [$\log(x(1, w(1, j)))$] in 40 replicates with high and low nutrient concentrations. The dashed line shows the level at which the two phenotypes are equally fit. Asterisks indicate the significance of the difference between secretor and nonsecretor fitness, *** $P < 0.001$; ** $P < 0.01$; * $P < 0.05$; ns, not significant. Below each boxplot is an image generated using the simulation from one of the 40 simulations that was closest to the median in the boxplot. It is shown that secretors can outcompete nonsecretors when the two phenotypes are well segregated, whereas they are at a disadvantage under conditions leading to high mixing.

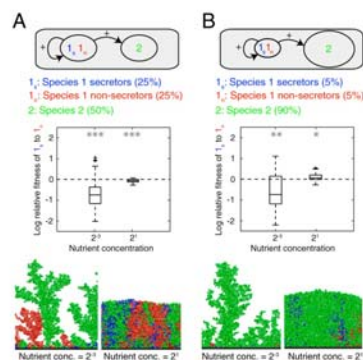


Fig. 3. Ecological competition with a second species. A second species is added to the competition between secretors and nonsecretors (Fig. 2). This second species is intended to also approximate the effects of a mixture of many species (Box 1). Species 1 is equally divided into secretor and nonsecretor strains, whereas species 2 represents either 50% (A) or 90% of the cells inoculated (B). All cells are then left to grow to a fixed total biomass. Strain 1_s secretes a product that benefits both strains of its own species, as well as species 2. 1_n and species 2 do not secrete any products. Product secretion incurs a cost of 30% of the cells' growth rate. See Fig. 2 legend for explanations on data representation. It is shown that when cells are highly segregated, secretor cells lose their advantage (compared with Fig. 2, Bottom Left), independently of the two proportions of species 2. At high levels of mixing, however, secretors can outcompete nonsecretors when there is a high proportion of species 2 cells. The image (B, Bottom Right) shows the social insulation effect discussed in the text.

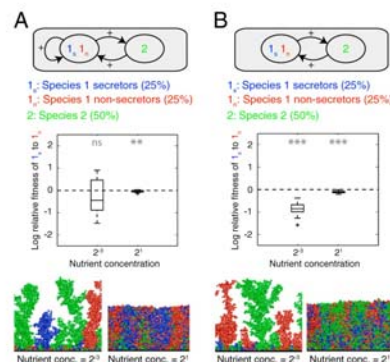
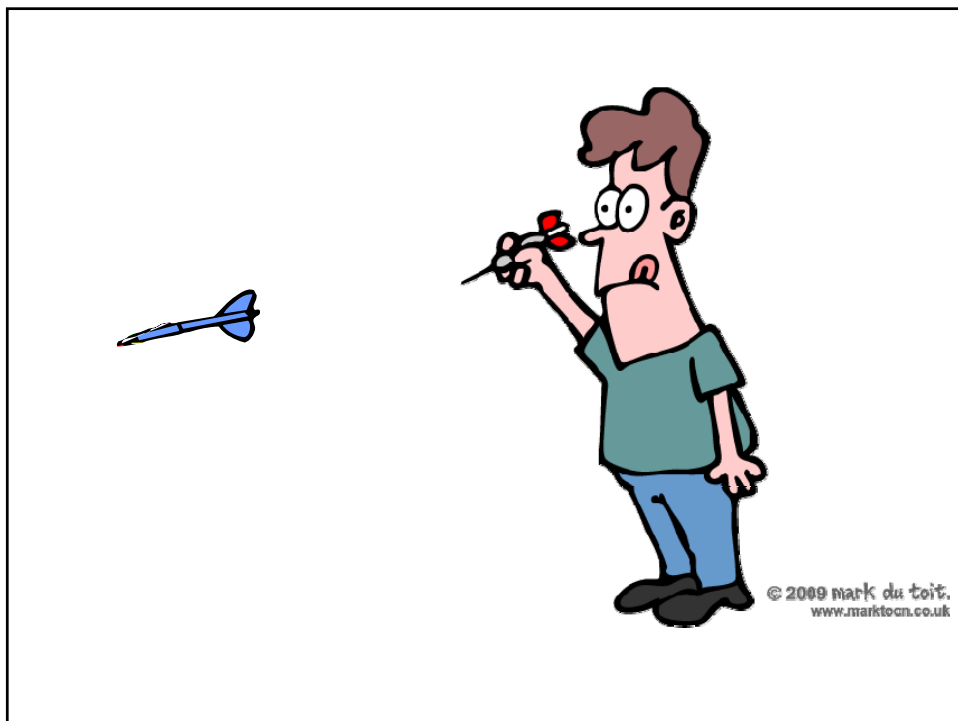
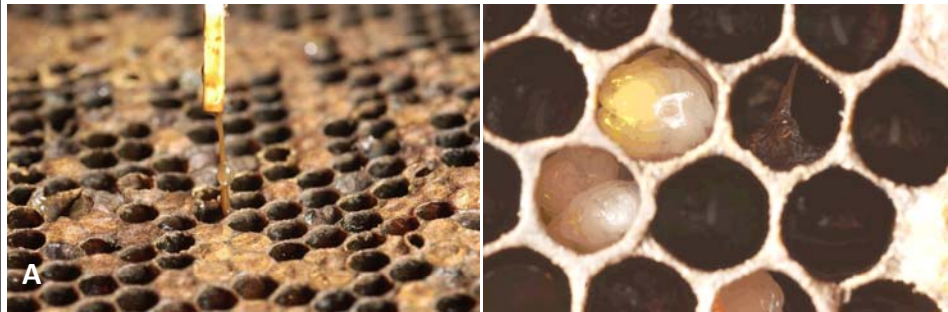


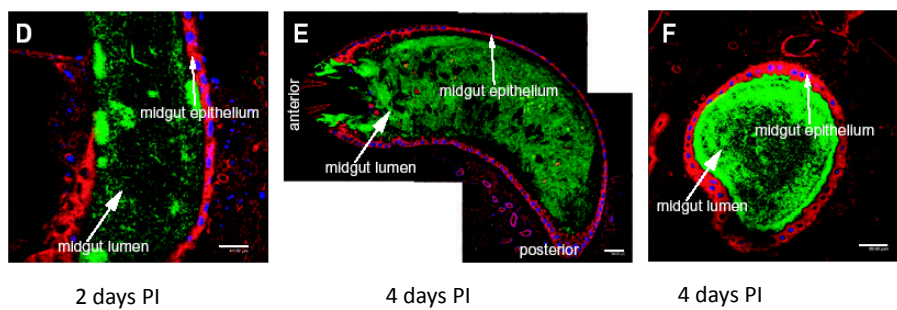
Fig. 4. Multispecies mutualism. Species 2 now secretes a product that is beneficial to species 1, resulting in a mutualism between the two species. Species 1 is equally divided into secretor and nonsecretor strains, whereas species 1 and 2 are inoculated in equal proportions and left to grow to a fixed total biomass. Strain 1_s secretes a product that either benefits both strains of its own species, as well as species 2 (A), or species 2 only (B). Product secretion by 1_s incurs a cost of 30% of the cells' growth rate. In turn, species 2 secretes a cost-free product that benefits species 1. 1_n does not secrete any products. See Fig. 2 legend for explanations on data representation. It is shown that secretor cells do not have a clear advantage over nonsecretors in any one of the four conditions considered here. This result is because mixing is important for the benefits of the two secreting strains to be shared, but is detrimental because it allows nonsecretors to grow faster than secretors, thereby undermining the mutualistic interaction.



Paenibacillus larvae (American foulbrood disease)



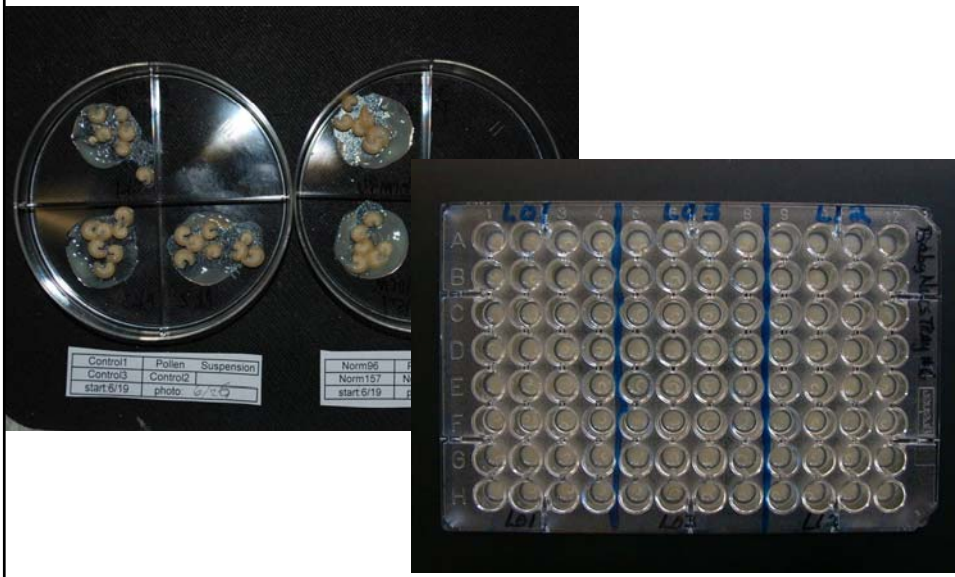
Stephen Pernal, AgCanada



YUE ET AL. ENV MICROB. 2008

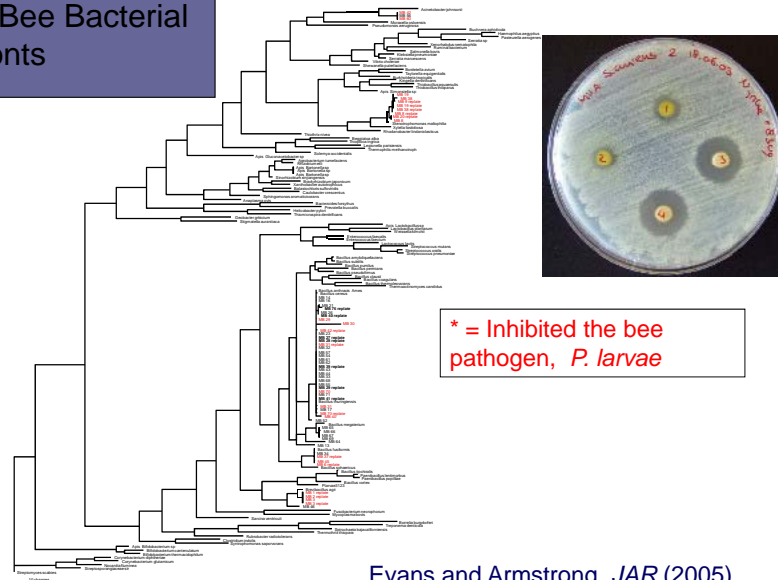


Challenge Experiments

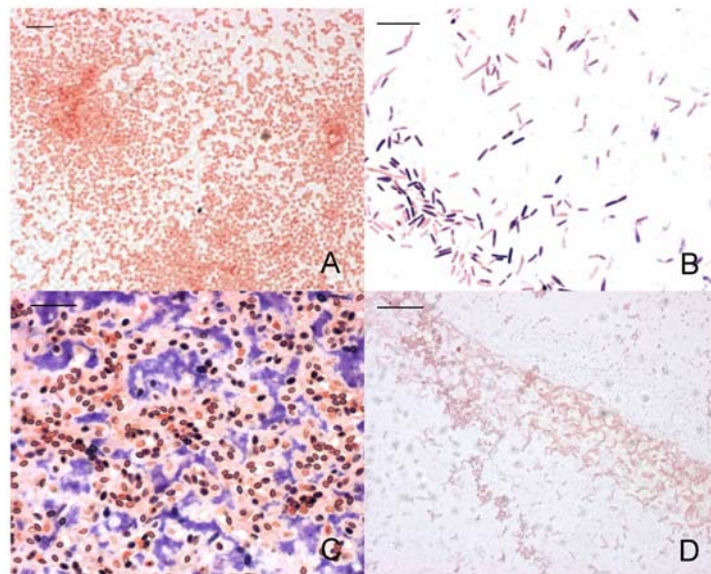




16S Placement of Honey Bee Bacterial Symbionts



Evans and Armstrong, *JAR* (2005)
BMC Ecology, 6:4 (2006)



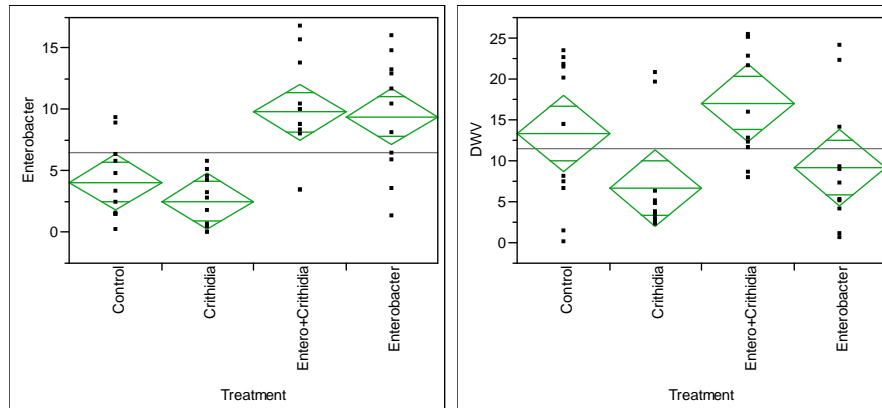
Evans and Armstrong, BMC Ecology, 2006



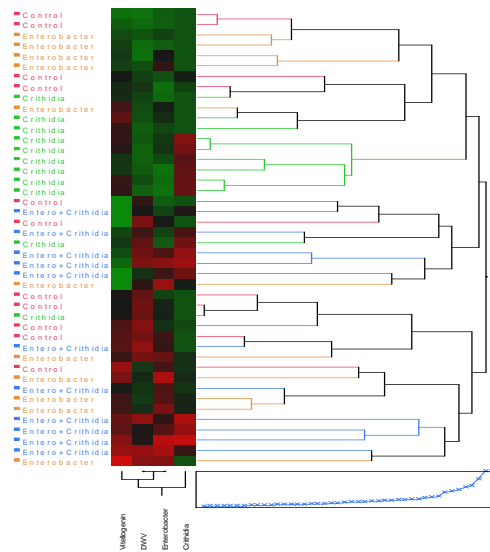
Bacteria x Pathogen Interactions



Controlled Inoculations, Bacteria Plus Trypanosome



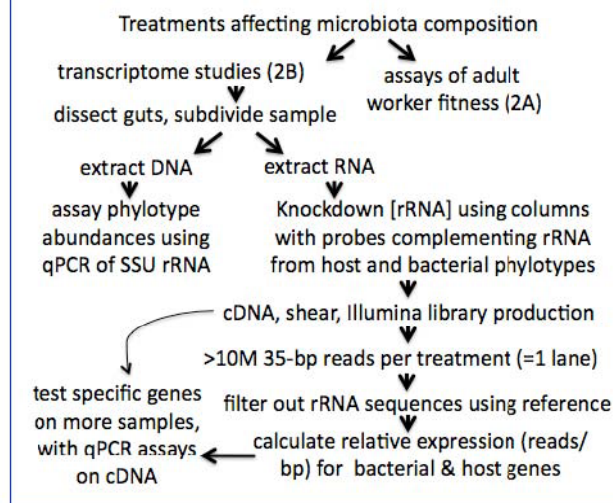
Controlled Inoculations, Bacteria Plus Trypanosome



The screenshot shows the NSF website header with the logo and tagline "WHERE DISCOVERIES BEGIN". A search bar is visible on the right. The navigation menu includes links for HOME, FUNDING, AWARDS, DISCOVERIES, NEWS, PUBLICATIONS, STATISTICS, and ABOUT. The "Awards" section is highlighted, showing an abstract for award #1046153. The abstract title is "Dimensions: Genomics, functional roles, and diversity of the symbiotic gut microbiotae of honey bees and bumble bees" by Moran/Evans/Winfree, 2011-2015.

Moran/Evans/Winfree 2011-2015

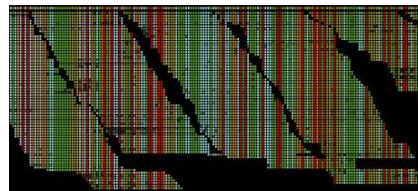
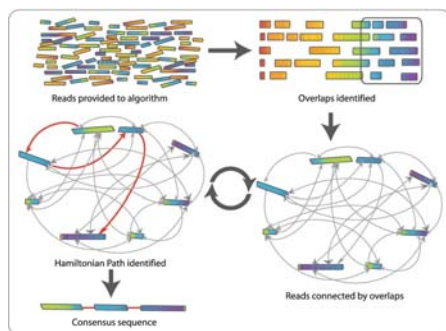
Fig. 4. Functional studies of *A. mellifera* microbiota: host fitness and transcriptome



- ❖ FIRST PRINCIPLES
- ❖ DATA COLLECTION
- ❖ OCEANS, TERMITES, MAMMALS
- ❖ BEE STORIES
- ❖ RETURN TO FUNCTION
- ❖ ANALYSIS



Assembly vs. mapping



Assembly: every read aligned to each other and then resolved to optimal 'path'. Computationally intensive. Outputs contigs.

Mapping: reads aligned one at a time to an existing reference. Computationally easy. Outputs an alignment file.

OTUs

Operational, *a priori* way to describe the number of taxonomic groups

Usually based on rDNA sequence (16S)

Clustering at arbitrary %ID (97-98% typical for bacteria, 95% for virus). Different programs give different clusters.

Comparison of OTUs

Phylogenetic tree

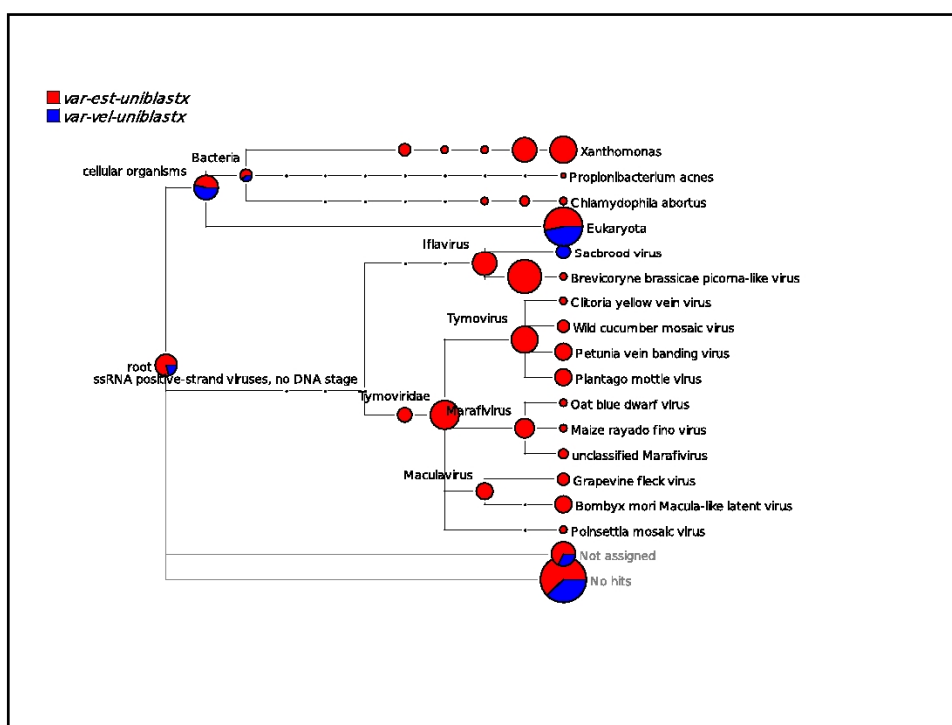
- MEGAN

Ribosomal classifier, align to reference database

- RDP Classifier

Nucleotide composition model

- PhyloPythia, PhymmBL



	A	B	C	D	E	F	G	H	I	J	K	L	
1	KOPR - KOPR Micro Biome Report - Africa (version 2) / March 2018												
2	Fermisearch Technology, Inc. Report based on metatranscriptomic sequencing and Benign's Manual												
3	Query File: readseqdiffs.txt												
4	Runned Date: Thu Feb 15 14:57:03 EST 2017												
5	Confidence threshold 95%												
6	Species threshold 95%												
7	NODE_16451_length_191_cov_1.895264-213	Bacteria	99%	Fusobacteriia	14%	Fusobacteriia	14%	Fusobacteriia	14%	Lysitrichaceae	12%	Bifidobacteriia	
8	NODE_20699_length_148_cov_1.555267-211	Bacteria	98%	Proteobacteria	49%	Alphaproteobacteria	20%	Rhodocycla	12%	Flavobacteriia	12%	Flavobacteriia	
9	NODE_17470_length_148_cov_1.750000-207	Bacteria	100%	Firmicutes	39%	Clostridia	24%	Clostridiales	20%	Lachnospiraceae	14%	Sporobacteriia	
10	NODE_1722_length_146_cov_1.214226-168	Bacteria	100%	Proteobacteria	36%	Bacteroidia	16%	Thiotrichales	16%	Prophymonadaceae	16%	Thiospirillum	
11	NODE_16450_length_143_cov_2.405004-165	Bacteria	97%	Firmicutes	5%	Thermococci	5%	Thermococci	5%	Thermococci	4%	Geitingeria	
12	NODE_15386_length_119_cov_2.268008-141	Bacteria	95%	Bacteroidia	28%	Bacteroidia_incertae_sedis	11%	Moraxillum	11%				
13	NODE_15571_length_115_cov_1.632174-137	Bacteria	99%	Bacteroidia	19%	Sphingobacteriia	10%	Sphingobacteriia	10%	Phosphoribacteriaceae	8%	Saibacteriia	
14	NODE_15790_length_114_cov_1.428928-135	Bacteria	96%	Bacteroidia	24%	Bacteroidia_incertae_sedis	6%	Moraxillum					
15	NODE_1184_length_110_cov_1.514454-132	Bacteria	97%	Actinobacteria	27%	Actinobacteria	27%	Actinobacteriia	26%	Actinomycetaceae	23%	Mycobacteriia	
16	NODE_11863_length_110_cov_26.299999-132	Bacteria	95%	Bacteroidia	30%	Flavobacteriia	21%	Fusobacteriia	21%	Corynebacterium	11%	Lysitrichaceae	
17	NODE_14391_length_106_cov_1.220846-128	Bacteria	100%	Proteobacteria	34%	Alphaproteobacteria	13%	Acetivibromonadaceae	13%	Acetivibromonadaceae	13%	Flavobacteriia	
18	NODE_3674_length_103_cov_2.291262-125	Bacteria	99%	Proteobacteria	29%	Betaproteobacteriia	10%	Betaproteobacteriia_incertae_sedis	2%	Thiospirillum			
19	NODE_15511_length_102_cov_1.872459-124	Bacteria	99%	Firmicutes	36%	Bacilli	21%	Bacillales	6%	Bacillus_incertae_sedis	20%	Saibacteriia	
20	NODE_11552_length_98_cov_2.356816-119	Bacteria	96%	Bacteroidia	22%	Sphingobacteriia	14%	Sphingobacteriia	14%	Phosphoribacteriaceae	4%	Thiospirillum	
21	NODE_13474_length_83_cov_1.232024-107	Bacteria	100%	Proteobacteria	99%	Alphaproteobacteria	96%	Rhodocyclales	96%	Acetivibromonadaceae	96%	Kosakia	
22	NODE_20281_length_84_cov_1.070000-106	Bacteria	97%	Proteobacteria	48%	Deinoproteobacteriia	14%	Deinoproteobacteriia	10%	Deinoproteobacteriia	10%	Brightwellia	
23	NODE_20712_length_83_cov_1.540265-105	Bacteria	97%	Chloroflexi	12%	Chloroflexi	12%	Chloroflexi	12%	Chloroflexi	12%	Chloroflexi	
24	NODE_4458_length_80_cov_2.184144-104	Bacteria	97%	Fusobacteriia	8%	Fusobacteriia	8%	Fusobacteriia	8%	Lysitrichaceae		Staphylococcus	
25	NODE_20402_length_80_cov_1.687500-102	Bacteria	98%	Proteobacteria	49%	Alphaproteobacteria	11%	Rhodocyclales	3%	Anaplastomonadaceae	3%	Anaplastomonadaceae	
26	NODE_11368_length_80_cov_3.547600-102	Bacteria	98%	Bacteroidia	44%	Sphingobacteriia	21%	Sphingobacteriia	21%	Phosphoribacteriaceae		Woodwardia	
27	NODE_18401_length_72_cov_1.772155-101	Bacteria	99%	Bacteroidia	38%	Sphingobacteriia	18%	Sphingobacteriia	18%	Phosphoribacteriaceae	11%	Fabaceae	
28	NODE_1201_length_503_cov_31.886787-95	Bacteria	100%	Cyanobacteriia	22%	Cyanobacteriia	22%	Chlorococcales	22%	Chlorococcales	22%	Chlorococcales	
29	NODE_2654_length_318_cov_262.897876-95	Bacteria	100%	Bacteroidia	34%	Fusobacteriia	22%	Fusobacteriia	22%	Corynebacteriia	11%	Flavobacteriia	
30	NODE_1440_length_205_cov_8.930009-241	Bacteria	96%	Bacteroidia	30%	Sphingobacteriia	23%	Sphingobacteriia	23%	Cytophagaceae	19%	Pseudomonadaceae	
31	NODE_2655_length_102_cov_1.540265-105	Bacteria	99%	Bacteroidia	26%	Bacteroidia	14%	Bacteroidia	14%	Phosphoribacteriaceae	14%	Pseudomonadaceae	
32	NODE_4448_length_123_cov_1.875610-159	Bacteria	96%	Firmicutes	41%	Clostridia	33%	Clostridiales	33%	Syntrophomonadaceae	14%	Syntrophomonadaceae	
33	NODE_1371_length_110_cov_3.421404-158	Bacteria	99%	Firmicutes	26%	Clostridia	26%	Clostridiales	26%	Phosphoribacteriaceae	14%	Phosphoribacteriaceae	
34	NODE_2934_length_109_cov_30.963303-149	Bacteria	99%	Firmicutes	68%	Bacilli	52%	Lactobacillales	48%	Acetivibromonadaceae	48%	Abiotritina	
35	NODE_2113_length_104_cov_2.29926-140	Bacteria	98%	Firmicutes	31%	Clostridia	27%	Clostridiales	22%	Verrucomonadaceae		Anaplastomonadaceae	
36	NODE_2248_length_102_cov_2.7048-138	Bacteria	99%	Firmicutes	42%	Bacilli	18%	Lactobacillales	15%	Acetivibromonadaceae	15%	Abiotritina	
37	NODE_4343_length_100_cov_2.20000-136	Bacteria	100%	Proteobacteria	71%	Gammaproteobacteriia	30%	Acetivibromonadaceae	12%	Syntrophomonadaceae	12%	Ruminococcaceae	
38	NODE_4343_length_95_cov_1.903303-136	Bacteria	100%	Proteobacteria	11%	Betaproteobacteriia	11%	Betaproteobacteriia	11%	Phosphoribacteriaceae	11%	Phosphoribacteriaceae	

Species richness

Alpha: the diversity of species at one site/habitat

Beta: how distinct different sites/habitats are from each other

Typically use some kind of index that considers how much of each species is present, rather than just total number of species

Testing difference between samples

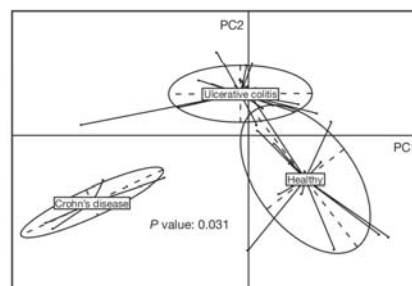
Need a metric/index

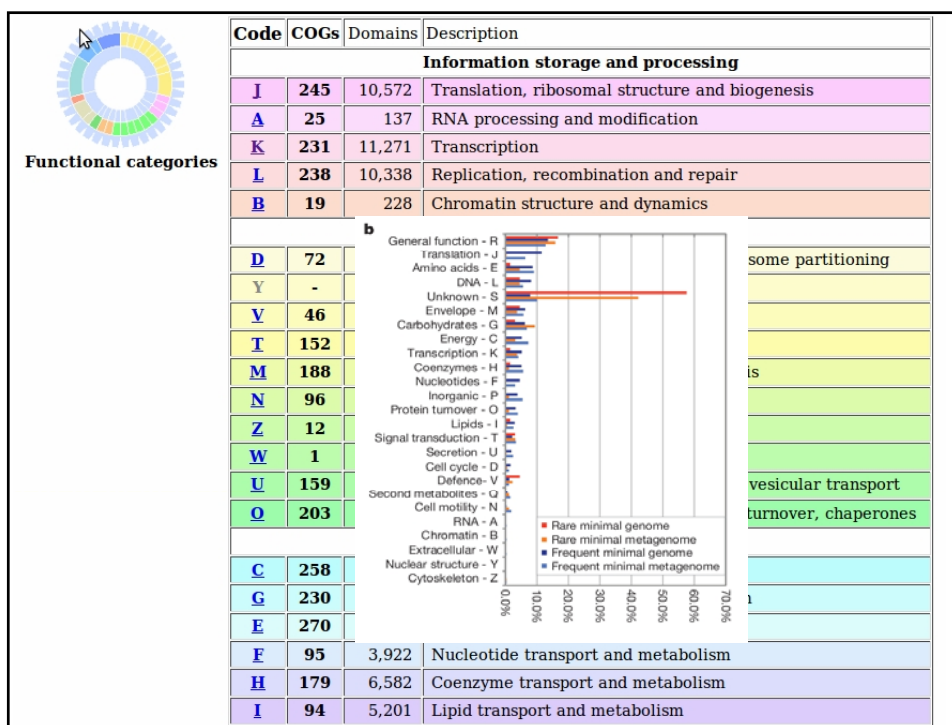
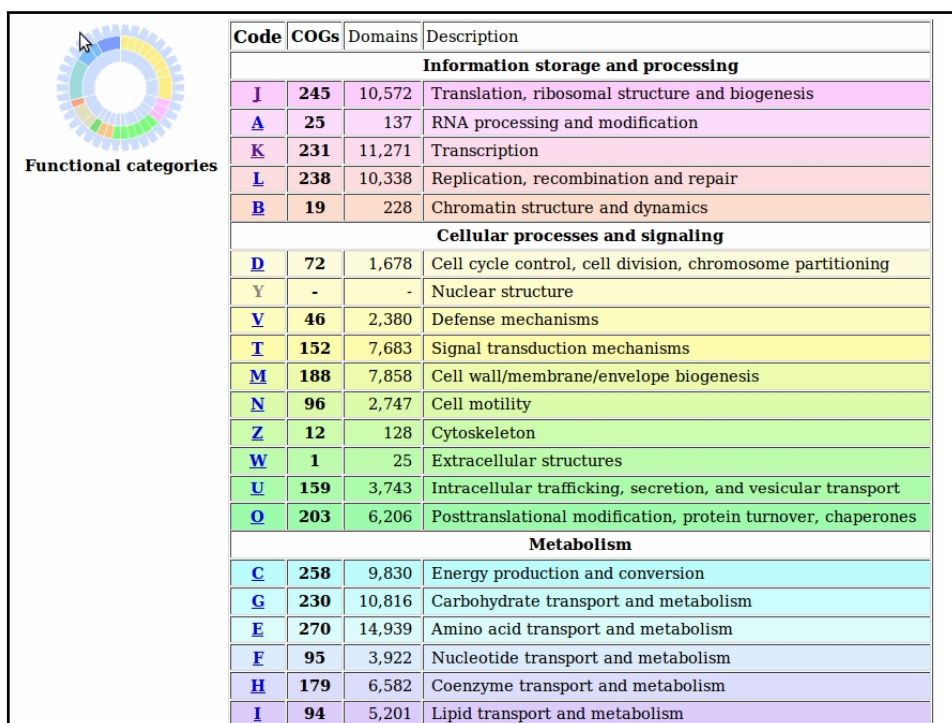
- Overlap/relative abundance of OTUs

Clustering/PCA/NMDS

Some programs:

- Estimates
- Unifrac





Gene ontology terms

Tree Browser

Filter tree view

Filter by ontology: **All**
 biological process
 cellular component
 molecular function

Filter Gene Product Counts: **All**
 Data source: **All**
 Species: **All**
 Arabidopsis thaliana
 Aspergillus fumig...
 Aspergillus niger

View Options: Tree view ☒ Full ☐ Compact
 Set filters
 Remove all filters

all : all [536514 gene products]

- ☐ GO:0008150 : biological process [405884 gene products]
- ☐ GO:0022810 : biological adhesion [6597 gene products]
- ☐ GO:0065007 : biological regulation [91014 gene products]
- ☒ **GO:0009758 : carbohydrate utilization [9 gene products]**
- ☐ GO:0043610 : regulation of carbohydrate utilization [5 gene products]
- ☐ GO:0015976 : carbon utilization [248 gene products]
- ☐ GO:0001906 : cell killing [1080 gene products]
- ☐ GO:0006283 : cell proliferation [7343 gene products]
- ☐ GO:0071840 : cellular component organization or biogenesis [44285 gene products]
- ☐ GO:0009987 : cellular process [242433 gene products]
- ☐ GO:0016265 : death [10051 gene products]
- ☐ GO:0032022 : developmental process [38990 gene products]
- ☐ GO:0051234 : establishment of localization [52184 gene products]
- ☐ GO:0040007 : growth [10464 gene products]
- ☐ GO:0002376 : immune system process [9394 gene products]
- ☐ GO:0051179 : localization [58555 gene products]
- ☐ GO:0040011 : locomotion [9713 gene products]
- ☐ GO:0008160 : metabolic process [500160 gene products]

Actions:
 Last action Opened
 GO:0009758
 Graphical View
 Permalink
 Download...
 OBO
 PDF-XML
 Graphviz dot

Domain/ontology mapping programs

Pfam/Interpro scan

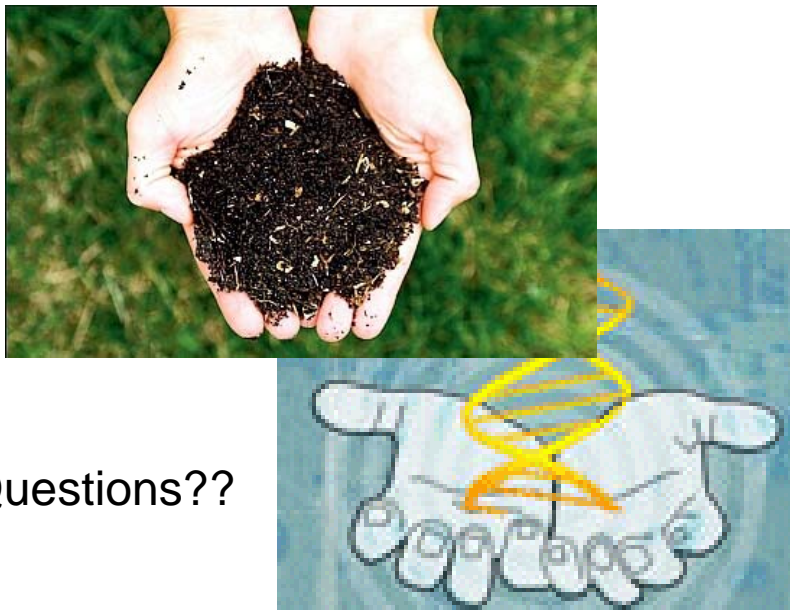
Blast2GO

KEGG Mapper

eggNOG

Read mapping to a reference

- Bowtie, BWA, Stampy, Novoalign
- Competitive mapping
 - Matches uniquely?
- Parameters for allowing match (% mismatch, indels)
- Measure abundance of taxa, expression of genes independent of assembled contigs



Questions??