



**HAL**  
open science

## Enterotypes of the human gut microbiome

Manimozhiyan Arumugam, Jeroen Raes, Eric Pelletier, Denis Le Paslier, Takuji Yamada, Daniel Mende, Gabriel Fernandes, Julien Tap, Thomas Bruls, Jean-Michel Batto, et al.

► **To cite this version:**

Manimozhiyan Arumugam, Jeroen Raes, Eric Pelletier, Denis Le Paslier, Takuji Yamada, et al.. Enterotypes of the human gut microbiome. *Nature*, 2013, 473 (7346), pp.174 - 180. <10.1038/nature09944>. <hal-03649169>

**HAL Id: hal-03649169**

**<https://hal.science/hal-03649169v1>**

Submitted on 22 Apr 2022

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Published in final edited form as:

*Nature*. 2011 May 12; 473(7346): 174–180. doi:10.1038/nature09944.

## Enterotypes of the human gut microbiome

Manimozhiyan Arumugam<sup>1,\*</sup>, Jeroen Raes<sup>1,2,\*</sup>, Eric Pelletier<sup>3,4,5</sup>, Denis Le Paslier<sup>3,4,5</sup>, Takuji Yamada<sup>1</sup>, Daniel R. Mende<sup>1</sup>, Gabriel R. Fernandes<sup>1,6</sup>, Julien Tap<sup>1,7</sup>, Thomas Bruls<sup>3,4,5</sup>, Jean-Michel Batto<sup>7</sup>, Marcelo Bertalan<sup>8</sup>, Natalia Borrueal<sup>9</sup>, Francesc Casellas<sup>9</sup>, Leyden Fernandez<sup>10</sup>, Laurent Gautier<sup>8</sup>, Torben Hansen<sup>11</sup>, Masahira Hattori<sup>12</sup>, Tetsuya Hayashi<sup>13</sup>, Michiel Kleerebezem<sup>14</sup>, Ken Kurokawa<sup>15</sup>, Marion Leclerc<sup>7</sup>, Florence Levenez<sup>7</sup>, Chaysavanh Manichanh<sup>9</sup>, H. Bjørn Nielsen<sup>8</sup>, Trine Nielsen<sup>11</sup>, Nicolas Pons<sup>7</sup>, Julie Poulain<sup>3</sup>, Junjie Qin<sup>16</sup>, Thomas Sicheritz-Ponten<sup>8</sup>, Sebastian Tims<sup>14</sup>, David Torrents<sup>10,17</sup>, Edgardo Ugarte<sup>3</sup>, Erwin G. Zoetendal<sup>14</sup>, Jun Wang<sup>16,18</sup>, Francisco Guarner<sup>9</sup>, Oluf Pedersen<sup>11,19</sup>, Willem M. de Vos<sup>14,20</sup>, Søren Brunak<sup>8</sup>, Joel Doré<sup>7</sup>, MetaHIT Consortium<sup>†</sup>, Jean Weissenbach<sup>3,4,5</sup>, S. Dusko Ehrlich<sup>7,#</sup>, and Peer Bork<sup>1,21,#</sup>

<sup>1</sup>European Molecular Biology Laboratory, Meyerhofstrasse 1, 69117 Heidelberg, Germany. <sup>2</sup>VIB —Vrije Universiteit Brussel, 1050 Brussels, Belgium. <sup>3</sup>Commissariat à l'Energie Atomique, Genoscope, 91000 Evry, France. <sup>4</sup>Centre National de la Recherche Scientifique, UMR8030, 91000 Evry, France <sup>5</sup>Université d'Evry Val d'Essonne 91000 Evry, France <sup>6</sup>Department of Biochemistry and Immunology, Universidade Federal de Minas Gerais, Av. Antônio Carlos 6627, 31270-901 Belo Horizonte, MG, Brazil. <sup>7</sup>Institut National de la Recherche Agronomique, 78350 Jouy en Josas, France. <sup>8</sup>Center for Biological Sequence Analysis, Technical University of Denmark, DK-2800 Kongens Lyngby, Denmark. <sup>9</sup>Digestive System Research Unit, University Hospital Vall d'Hebron, Ciberehd, 08035 Barcelona, Spain. <sup>10</sup>Barcelona Supercomputing Center, Jordi Girona 31, 08034 Barcelona, Spain <sup>11</sup>Hagedorn Research Institute, 2820 Gentofte, Denmark. <sup>12</sup>Computational Biology Laboratory Bld, The University of Tokyo Kashiwa Campus, Kashiwa-no-ha 5-1-5, Kashiwa, Chiba, 277-8561, Japan <sup>13</sup>Division of Bioenvironmental Science, Frontier Science Research Center, University of Miyazaki, 5200 Kiyotake, Miyazaki 889-1692, Japan. <sup>14</sup>Laboratory of Microbiology, Wageningen University, 6710BA Ede, The Netherlands. <sup>15</sup>Tokyo Institute of Technology, Graduate School of Bioscience and Biotechnology, Department of Biological Information, 4259 Nagatsuta-cho, Midori-ku, Yokohama-shi, Kanagawa Pref. 226-8501, Japan <sup>16</sup>BGI-Shenzhen, Shenzhen 518083, China. <sup>17</sup>Institució Catalana de Recerca i Estudis Avançats (ICREA), Pg. Lluís Companys 23, 08010 Barcelona, Spain <sup>18</sup>Department of Biology, University of Copenhagen, DK-2200 Copenhagen, Denmark. <sup>19</sup>Institute of Biomedical Sciences, University of Copenhagen, Denmark. <sup>20</sup>University of Helsinki, FI-00014 Helsinki, Finland. <sup>21</sup>Max Delbrück Centre for Molecular Medicine, D-13092 Berlin, Germany.

### Abstract

Our knowledge on species and function composition of the human gut microbiome is rapidly increasing, but it is still based on very few cohorts and little is known about their variation across the world. Combining 22 newly sequenced fecal metagenomes of individuals from 4 countries with previously published datasets, we identified three robust clusters (enterotypes hereafter) that are not nation or continent-specific. We confirmed the enterotypes also in two published, larger cohorts suggesting that intestinal microbiota variation is generally stratified, not continuous. This further indicates the existence of a limited number of well-balanced host-microbial symbiotic states that might respond differently to diet and drug intake. The enterotypes are mostly driven by species composition, but abundant molecular functions are not necessarily provided by abundant species, highlighting the importance of a functional analysis for a community understanding. While individual host properties such as body mass index, age, or gender cannot explain the observed enterotypes, data-driven marker genes or functional modules can be identified for each

of these host properties. For example, twelve genes significantly correlate with age and three functional modules with the body mass index, hinting at a diagnostic potential of microbial markers.

## Introduction

Various studies of the human intestinal tract microbiome, based on the 16S ribosomal RNA-encoding gene, reported species diversity within and between individuals<sup>1-3</sup> and first metagenomics studies characterized the functional repertoire of the microbiomes of several American<sup>4-5</sup> and Japanese<sup>6</sup> individuals. Although a general consensus about the phylum level composition in the human gut is emerging<sup>1,3,7</sup>, the variation in species composition<sup>1-2</sup> and gene pools<sup>5,8</sup> within the human population is less clear. Furthermore, it is unknown whether inter-individual variation manifests itself as a continuum of different community compositions or whether individual gut microbiota congregate around some preferred, balanced and stable community compositions that can be classified. Studying such questions is complicated by the complexity of sampling, DNA preparation, processing, sequencing and analysis protocols<sup>9</sup> as well as by varying physiological, nutritional and environmental conditions. To analyze the feasibility of comparative metagenomics of the human gut across cohorts and protocols and to obtain first insights in commonalities and differences between gut microbiomes across different populations, we Sanger-sequenced 22 European metagenomes from Danish, French, Italian and Spanish individuals that were selected for diversity (Supplementary Notes Section 1), and combined them with existing Sanger (13 Japanese<sup>6</sup>, 2 American<sup>4</sup>) and 454 (2 American<sup>5</sup>) gut datasets – totaling 39 individuals.

## Global phylogenetic and functional variation of intestinal metagenomes

The vast majority of sequences in the newly sequenced 22 European samples belong to bacteria – only 0.14% of the reads could be classified as human contamination, all other eukaryotes together only comprised 0.5%, archaea 0.8% and viruses up to 5.8% (see Supplementary Notes Section 2.1 for details).

To investigate the phylogenetic composition of the 39 samples from 6 nationalities, we mapped metagenomic reads, using DNA sequence homology, to 1511 reference genomes (Supplementary Table 3) including 379 publicly available human microbiome genomes generated through the NIH Human Microbiome Project<sup>10</sup> and the European MetaHIT consortium<sup>11</sup> (Supplementary Methods Section 4.1). To consistently estimate the functional composition of the samples, we annotated the predicted genes from the metagenomes using eggNOG<sup>12</sup> orthologous groups (Supplementary Methods Section 6.2). We ensured that comparative analysis using these procedures was not biased by dataset origin, sample preparation, sequencing technology and quality filtering (see Supplementary Notes Section 1). We also investigated whether the relatively low and somewhat arbitrary amounts of sequence per sample (between 53-295 Mb) bias our results: we assigned habitat information to 1368 out of the 1511 reference genomes, distinguished between orthologous groups from gut and non-gut species and conclude that our dataset captures most of the functions from gut species even though functions from ‘non-gut’ species still accumulated with each additional sample (Fig. 1a; see Supplementary Notes Section 1.3).

We then characterized the phylogenetic variation across samples at the genus and phylum levels, and functional variation at gene and functional class levels. As infants are known to have very heterogeneous, unstable and distinctive microbiota<sup>6,13</sup>, we excluded the four respective Japanese samples from the analysis. Using calibrated similarity cutoffs (Supplementary Figure 1), on average, 52.8% of the fragments in each sample could be robustly assigned to a genus in our reference genome set (ranging from 22% to 80.5%), and

80% could be assigned to a phylum (ranging from 64.9% to 91%) implying that the trends observed (Fig. 1b) represent a large fraction of the metagenome.

The phylogenetic composition of the newly sequenced samples confirms that the Firmicutes and Bacteroidetes phyla constitute the vast majority of the dominant human gut microbiota<sup>7</sup> (Fig. 1b, inset). *Bacteroides* was the most abundant but also most variable genus across samples (Fig. 1b; Supplementary Notes Section 2.2), agreeing with previous observations<sup>6,14</sup>. Our function identification protocol led to a high functional assignment rate: 63.5% of all predicted genes in the Sanger-sequenced samples analyzed (41% of all predicted genes in two samples obtained by pyrosequencing; Supplementary Table 5) can be assigned to orthologous groups (OGs), and OG abundance patterns again agree with previous observations<sup>6,15</sup> (e.g. histidine kinases make up the largest group; Fig 1c; Supplementary Notes Section 2.3).

## Highly abundant functions from low-abundance microbes

Microbes in the human gut undergo selective pressure from the host as well as from microbial competitors. This typically leads to a homeostasis of the ecosystem in which some species occur in high and many in low abundance<sup>16</sup> (the “long-tail” effect, as seen in Fig. 1b), with some low-abundance species, like methanogens<sup>17</sup>, performing specialized functions beneficial to the host. Metagenomics enables us to study the presence of abundant functions shared by several low-abundance species, which could shed light on their survival strategies in the human gut. In the samples analyzed here, the most abundant molecular functions generally trace back to the most dominant species. However, we identified some abundant orthologous groups that are contributed primarily by low abundance genera (see Supplementary Figure 2, Supplementary Table 6 and Supplementary Notes Section 3). For example, low abundance *Escherichia* contribute over 90% of two abundant proteins associated with bacterial pilus assembly, FimA (COG3539) and PapC (COG3188), found in one individual (IT-AD-5). Pili enable the microbes to colonize the epithelium of specific host organs; they help microbes to stay longer in the human intestinal tract by binding to the human mucus or mannose sugars present on intestinal surface structures<sup>18</sup>. They are also key components in the transfer of plasmids between bacteria through conjugation, often leading to exchange of protective functions such as antibiotic resistance<sup>18</sup>. Pili can thus provide multiple benefits to these low-abundance microbes in their efforts to survive and persist in the human gut. This example illustrates that abundant species or genera cannot reveal the entire functional complexity of the gut microbiota. More reference genomes will facilitate better taxonomic assignments from samples and thus the detection of more low abundance species. However, there is not much room for as yet undetected, abundant genera. Even with our limited genus assignment rate of 52.8% of all reads, we estimate that we miss another 30.7% of the already classified genera due to our strict assignment criteria (Supplementary Figure 1), i.e. only 16.5% of all reads are likely to belong to hitherto unknown genera.

## Robust clustering of samples across nations: Identification of enterotypes

To get an overview of the species variation we used phylogenetic profile similarities obtained by mapping metagenomic reads to the 1511 reference genomes (Fig. 2a, see Supplementary Methods Section 4.1). We excluded the two American Sanger-sequenced samples<sup>4</sup> from further analysis because of an unusual, very low fraction of *Bacteroidetes*, and suspected technical artifacts<sup>19</sup>. Multidimensional cluster analysis and Principal Component Analysis (PCA) revealed that the remaining 33 samples formed three distinct clusters which we designate enterotypes (see Supplementary Notes Section 4.1, Supplementary Figure 3a and Supplementary Table 8). Each of these three enterotypes are

identifiable by the variation in the levels of one of three genera: *Bacteroides* (enterotype 1), *Prevotella* (enterotype 2) and *Ruminococcus* (enterotype 3; Fig. 2a and 2d), which was reproduced using independent array-based HITChip<sup>20</sup> data in a subset of 22 European samples (Supplementary Figure 4 and Supplementary Notes Section 4.5). The same analysis on two larger published gut microbiome datasets of different origins (16S pyrosequencing data from 154 American individuals<sup>5</sup> and Illumina-based metagenomics data from 85 Danish individuals<sup>8</sup>, Supplementary Methods Section 5) shows that these datasets could also be represented best by three clusters (Supplementary Figure 3b and c, Supplementary Table 9 and Supplementary Table 10). Two of these are also driven by *Bacteroides* and *Prevotella*, while the third cluster is mostly driven by related groups of the order Clostridiales, *Blautia* and unclassified Lachnospiraceae in the 16S rDNA and Illumina data, respectively (Fig. 2b and 2c). This can be explained by a different reference data set in case of 16S rDNA data, different mapping behavior of short reads in case of the Illumina data or current taxonomic uncertainties in the Lachnospiraceae and Ruminococcaceae clades (see Supplementary Notes Section 4.2). The differences might also hint at community subpopulations within this enterotype, which might only be detectable with substantially more samples. Correlation analysis of the Sanger data revealed that abundances of each of the three discriminating genera strongly correlate (that is they co-occur or avoid each other) with those of other genera (Fig. 2d; see Supplementary Methods Section 11), suggesting that the enterotypes are in fact driven by groups of species that together contribute to the preferred community compositions.

We further demonstrate the robustness of the enterotypes using two distinct statistical concepts. First we used the silhouette coefficient<sup>21</sup> to validate that the three clusters are superior to clusterings obtained from various randomizations of the genus profile data, suggesting a potential role for the interactions between co-occurring genera (see Supplementary Figure 5 and Supplementary Notes Section 4.3). Second we used supervised learning and cross validation to establish that these clusters have non-random characteristics that can be modeled and subsequently used to classify new samples (learning on clusters from randomized genus profiles led to considerably worse classification performance; see Supplementary Figure 6 and Supplementary Notes Section 4.4). These consistent results suggest that enterotypes will be identifiable in human gut metagenomes also from larger cohorts.

We then clustered the 33 samples using a purely functional metric: the abundance of the assigned orthologous groups (Fig. 3a). Remarkably, this clustering also showed a similar grouping of the samples with only minor differences (5 samples placed in different clusters compared to Fig. 2a) indicating that function and species composition roughly coincide with some exceptions such as Spanish sample ES-AD-3 whose genus composition belongs to enterotype 2 while its functional composition is similar to members of enterotype 1. This individual has high levels of phage-related genes compared to the other samples (see Supplementary Figure 7), hinting at partial temporal variability and dynamics of the microbiota, and perhaps indicating phage or virus bursts.

The robustness and predictability of the enterotypes in different cohorts and at multiple phylogenetic and functional levels suggests that they are the result of well-balanced, defined microbial community compositions of which only a limited number exist across individuals. These enterotypes are not as sharply delimited as, for example, human blood groups; they are rather densely populated areas in a multidimensional space of community composition. They are nevertheless likely to characterize individuals, in line with previous reports that gut microbiota is rather stable in individuals and can even be restored after perturbation<sup>22-25</sup>.

## Phylogenetic and functional variation between enterotypes

To determine the phylogenetic and functional basis of the enterotypes, we investigated in detail their differences in composition at the phylum, genus, gene and pathway level as well as correlations in abundance of co-occurring genera (Figs. 2 and 3; also see Supplementary Methods Sections 10, 11 and 12). Enterotype 1, containing 8 samples, is enriched in *Bacteroides* ( $p < 0.01$ ; Supplementary Figure 8), which co-occurs, for example, with *Parabacteroides* (see Supplementary Table 11 for enriched genera and Fig. 2e for correlation networks of co-occurring genera in each enterotype). The drivers of this enterotype seem to derive energy primarily from carbohydrates and proteins through fermentation, since these closely related genera have a very broad saccharolytic potential<sup>26</sup> and since genes encoding enzymes involved in the degradation of these substrates (galactosidases, hexosaminidases, proteases) along with glycolysis and pentose phosphate pathways are enriched in this enterotype (see Supplementary Table 12 and Supplementary Table 13). Enterotype 2 contains 6 samples and is enriched in *Prevotella* ( $p < 0.01$ ; Supplementary Figure 9) and the co-occurring *Desulfovibrio*, who can act in synergy to degrade mucin glycoproteins present in the mucosal layer of the gut: *Prevotella* is a known mucin-degrader and *Desulfovibrio* could enhance the rate-limiting mucin desulfation step by removing the sulfate<sup>27</sup>. Enterotype 3 is the most frequent one and is enriched in *Ruminococcus* ( $p < 0.01$ ; Supplementary Figure 10) as well as co-occurring *Akkermansia*, both known to comprise species able to degrade mucins<sup>28</sup>. It is also enriched in membrane transporters, mostly of sugars, suggesting the efficient binding of mucin and its subsequent hydrolysis as well as uptake of the resulting simple sugars by these genera. The enriched genera suggest that enterotypes employ different routes to generate energy from fermentable substrates available in the colon reminiscent of a potential specialization in ecological niches or guilds. In addition to the conversion of complex carbohydrates into absorbable substrates, the gut microbiota is also beneficial to the human host by producing vitamins. Although all the vitamin metabolism pathways are represented in all samples, enterotypes 1 and 2 were enriched in biosynthesis of different vitamins: biotin (Fig. 3b), riboflavin, pantothenate and ascorbate in the former, and thiamine (Fig. 3c) and folate in the latter. These phylogenetic and functional differences among enterotypes thus reflect different combinations of microbial trophic chains with a likely impact on the synergistic inter-relations with the human hosts.

## Functional biomarkers for host properties

Enterotypes do not seem to differ in functional richness (Supplementary Figure 11), and virtually none of several measured host properties, namely nationality, gender, age or body mass index (BMI), significantly correlates with the enterotypes (with the exception of enterotype 1 which is enriched in Japanese individuals). However, some strong correlations do occur between host properties and particular functions, at the genes or module level (a module is a part of a pathway that is functionally tightly interconnected, see Supplementary Methods Sections 6.3, 13 and Supplementary Notes Section 6). The only significant correlation between a host property and a taxonomic group is a negative one between age and the abundance of an unknown Clostridiales genus ( $p < 0.02$ ) containing three obligate anaerobes (Supplementary Figure 12a; see Supplementary Notes Section 6.2). It should be noted that age is not constant across the nationalities (in our dataset, Italians are relatively old and Japanese young), but that individuals did not stratify by nationality, suggesting that this is not a confounding factor. Our data did not reveal any correlation between BMI and the Firmicutes/Bacteroidetes ratio and we thus cannot contribute to the ongoing debate on the relation between this ratio and obesity<sup>29-30</sup>.

In contrast to the little phylogenetic signal, we found several significant functional correlations with each of the host properties studied (after correcting for multiple testing to

avoid artifacts; see Supplementary Methods Section 13), suggesting that metagenomics-derived functional biomarkers might be more robust than phylogenetic ones. For example, the abundance of 10 orthologous groups (OGs) varies more between than within nationalities (Supplementary Table 14) although overall, the functional composition in total was remarkably similar among the nations (also with respect to the functional core; see Supplementary Figure 13). For gender, we find five functional modules and one OG that significantly correlate ( $p < 0.05$ ; e.g., enriched aspartate biosynthesis modules in males; see Supplementary Table 16). In addition, twelve OGs significantly correlate with age (Supplementary Table 17). For instance, starch degradation enzymes such as glycosidases and glucan phosphorylases increase with age (which could be a reaction to decreased efficiency of host breakdown of dietary carbohydrates with age<sup>31</sup>) and so does the secA preprotein translocase (Supplementary Figure 14). Conversely, an OG coding for the facultative sigma-24 subunit of RNA polymerase, which drives expression under various stress responses and is linked to intestinal survival<sup>32</sup>, decreases with age (Fig. 4a). One explanation for this could be the reduced need for stress response in the gut due to the age-associated decline in host immune response<sup>33</sup> (immunosenescence). Our analyses also identified three marker modules that correlate strongly with the hosts' BMI (Supplementary Table 19, Supplementary Figure 14), two of which are ATPase complexes, supporting the link found between the gut microbiota's capacity for energy harvest and host's obesity<sup>34</sup>. Interestingly, functional markers found by a data-driven approach (derived from the metagenomes without previous knowledge) gave much stronger correlations than genes for which a link would be expected (e.g. SusC/SusD, involved in starch utilization<sup>26</sup>; Fig. 4b). Linear models combining the abundance of only a few functional modules correlate even better with host properties (Fig 4c,d). It should be noted that given the possibility of many confounding variables due to the heterogeneity and size of our cohort, these observations will need to be substantiated using larger, independent cohorts in the future. Furthermore, patterns in metagenomics data can (partly) reflect indirect factors<sup>9</sup> such as genome size<sup>35</sup> (the smaller the average genome size of a sample, the higher would be the relative fraction of single copy genes therein), which does not matter for diagnostics though.

While individual host properties don't explain the enterotypes, the latter might be driven by a complex mixture of functional properties, by host immune modulation or by hitherto unexplored physiological conditions such as transit time or pH of luminal contents. Furthermore, the three major enterotypes could be triggered by the three distinct pathways for hydrogen disposal<sup>36</sup> (Supplementary Notes Section 6.4). Indeed, despite their low abundance, *Methanobrevibacter* (a methanogen) and *Desulfovibrio* (a known sulfate-reducer) are enriched in enterotypes 3 and 1, respectively.

Taken together, we have demonstrated the existence of enterotypes in the human gut microbiome and have identified three of them that vary in species and functional composition using data that spans several nations and continents. As our current data do not reveal which environmental or even genetic factors are causing the clustering, and as fecal samples are not representative of the entire intestine, we anticipate that the enterotypes introduced here will be refined with deeper and broader analysis of individuals' microbiomes. Presumably, enterotypes are not limited to humans but also occur in animals. Their future investigations might well reveal novel facets of the human and animal symbiotic biology and lead to the discovery of the microbial properties correlated with the health status of individuals. We anticipate that they might allow classification of human groups that respond differently to diet or drug intake. The enterotypes appear complex, are probably not driven by nutritional habits and cannot simply be explained by host properties such as age or BMI, although there are functional markers such as genes or modules that correlate remarkably well with individual features. The latter might be utilizable for diagnostic and perhaps even prognostic tools for numerous human disorders, for instance

colorectal cancer and obesity-linked co-morbidities such as metabolic syndrome, diabetes and cardio-vascular pathologies.

## Methods summary

### Sample collection

Human fecal samples from European individuals were collected and frozen immediately, and DNA was purified as described previously<sup>37</sup>.

### Sequencing

Random shotgun DNA libraries of 3kb were Sanger-sequenced using standard protocols established at Genoscope.

### Sequence processing

Cloning vector, sequencing primers and low quality bases were end-trimmed from raw Sanger reads, and possible human DNA sequences were removed. Reads were processed by the SMASH comparative metagenomics pipeline<sup>38</sup> for assembly and gene prediction.

### Phylogenetic annotation

Phylogenetic annotation of samples was performed by (1) aligning reads (Sanger/Illumina) against a database of 1511 reference genomes (listed in Supplementary Table 3) or (2) classifying 16S rDNA reads using RDP classifier<sup>39</sup>. Genus and phylum abundance was estimated after normalizing for genome size for the former, and for 16S gene copy number for the latter.

### Functional annotation

Genes were functionally annotated using BLASTP against eggNOG (v2) and KEGG (v50) databases. Protein abundances were estimated after normalizing for protein length. Functional abundance profiles at eggnog-, KEGG orthologous group-, functional module- and pathway-level were created.

### Clustering and classification

Samples were clustered using Jensen-Shannon distance and partitioning around medoid (PAM) clustering. Optimal number of clusters was estimated using Calinski-Harabasz (CH) index. We used the silhouette validation technique for assessing the robustness of clusters. Additionally, within a cross-validation scheme, we trained predictive decision tree models on clusters obtained using the same clustering method and evaluated the classification of hold-out samples by accuracy, average precision and average precision gain.

### Statistics

Correlations between metadata and feature abundances were computed as described previously<sup>40</sup>, based on multiple-testing corrected pairwise Spearman correlation analysis and stepwise regression for multi-feature model building. For categorical metadata and enterotype comparisons, samples were pooled into bins (male/female, obese/lean, one enterotype/rest, specific nationality/rest etc.) and significant features were identified using Fisher's exact test with multiple testing correction of p-values.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Footnotes

#Correspondence and requests for materials should be addressed to P.B. (bork@embl.de) or S.D.E. (dusko.ehrlich@jouy.inra.fr).

\*These authors contributed equally

†Lists of authors and affiliations of the additional MetaHIT members appear at the end of the paper.

**Author Contributions** All authors are members of the Metagenomics of the Human Intestinal Tract (MetaHIT) Consortium. Ju.W., F.G., O.P., W.M.V., S.B., J.D., Je.W., S.D.E. and P.B. managed the project. N.B., F.C., T.H., C.M., and T. N. performed clinical analyses. M.L. and F.L. performed DNA extraction. E.P., D.L.P., T.B., J.P. and E.U. performed DNA sequencing. M.A., J.R., S.D.E. and P.B. designed the analyses. M.A., J.R., T.Y., D.R.M., G.R.F., J.T., J.M.B., M.B., L.F., L.G., M.K., H.B.N., N.P., J.Q., T.S-P., S.T., D.T., E.G.Z., S.D.E. and P.B. performed the analyses. M.A., J.R., P.B. and S.D.E. wrote the manuscript. M.H., T.H., K.K. and the MetaHIT Consortium members contributed to the design and execution of the study.

Supplementary Information is linked to the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

Full Methods and any associated references are available in the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Author Information** Informed consent was obtained from the 22 European subjects. Sample collection and experiments were approved by the following ethics committees: MetaHIT (Danish) – ethical committee of the Capital Region of Denmark; MetaHIT (Spanish) – CEIC, Hospital Vall d’Hebron; MicroObes – Ethical Committee for Studies with Human Subjects of Cochin Hospital in Paris, France; MicroAge – Joint Ethical Committee of the University of Camerino. Raw Sanger read data from the European fecal metagenomes have been deposited to NCBI Trace Archive with the following project ids: MH6 (33049), MH13 (33053), MH12 (33055), MH30 (33057), CD1 (33059), CD2 (33061), UC4(33113), UC6(33063), NO1 (33305), NO3 (33307), NO4 (33309), NO8 (33311), OB2 (33313), OB1 (38231), OB6 (38233), OB8 (45929), A (63073), B(63075), C (63077), D (63079), E (63081), G (63083). Contigs, genes and annotations are available to download from [http://www.bork.embl.de/Docu/Arumugam\\_et\\_al\\_2011/](http://www.bork.embl.de/Docu/Arumugam_et_al_2011/).

The authors declare no competing financial interests.

**Additional MetaHIT Consortium members** María Antolín<sup>1</sup>, François Artiguenave<sup>2</sup>, Hervé M. Blottiere<sup>3</sup>, Mathieu Almeida<sup>3</sup>, Carlos Cara<sup>4</sup>, Christian Chervaux<sup>5</sup>, Antonella Cultrone<sup>3</sup>, Christine Delorme<sup>3</sup>, Gérard Denariáz<sup>5</sup>, Rozenn Dervyn<sup>3</sup>, Konrad U. Foerstner<sup>6,7</sup>, Carsten Friss<sup>8</sup>, Maarten van de Guchte<sup>3</sup>, Eric Guedon<sup>3</sup>, Florence Haimet<sup>3</sup>, Wolfgang Huber<sup>6</sup>, Alexandre Jamet<sup>3</sup>, Catherine Juste<sup>3</sup>, Ghalia Kaci<sup>3</sup>, Jan Knol<sup>5</sup>, Omar Lakhdari<sup>3</sup>, Severine Layec<sup>3</sup>, Karine Le Roux<sup>3</sup>, Emmanuelle Maguin<sup>3</sup>, Raquel Melo Minardi<sup>2</sup>, Jean Muller<sup>9,10</sup>, Raish Oozeer<sup>5</sup>, Julian Parkhill<sup>11</sup>, Pierre Renault<sup>3</sup>, Maria Rescigno<sup>12</sup>, Nicolas Sanchez<sup>3</sup>, Shinichi Sunagawa<sup>6</sup>, Antonio Torrejon<sup>1</sup>, Keith Turner<sup>11</sup>, Gaetana Vandemeulebrouck<sup>3</sup>, Encarna Varela<sup>1</sup>, Yohanan Winogradsky<sup>3</sup>, Georg Zeller<sup>6</sup>

<sup>1</sup>Digestive System Research Unit, University Hospital Vall d’Hebron, Ciberehd, Barcelona, Spain.

<sup>2</sup>Commissariat à l’Energie Atomique, Genoscope, 91000 Evry, France.

<sup>3</sup>Institut National de la Recherche Agronomique, 78350 Jouy en Josas, France.

<sup>4</sup>UCB Pharma SA, 28046 Madrid, Spain.

<sup>5</sup>Danone Research, 91120 Palaiseau, France.

<sup>6</sup>European Molecular Biology Laboratory, Meyerhofstrasse 1, 69117 Heidelberg, Germany.

<sup>7</sup>Darmstadt, Germany.

<sup>8</sup>Center for Biological Sequence Analysis, Technical University of Denmark, DK-2800 Kongens Lyngby, Denmark.

<sup>9</sup>Institute of Genetics and Molecular and Cellular Biology, CNRS, INSERM, University of Strasbourg.

<sup>10</sup>Genetic Diagnostics Laboratory, CHU Strasbourg Nouvel Hôpital Civil, Strasbourg, France

<sup>11</sup>The Wellcome Trust Sanger Institute, Hinxton, Cambridge CB10 1SA, UK.

<sup>12</sup>Istituto Europeo di Oncologia, 20100 Milan, Italy.

## Acknowledgments

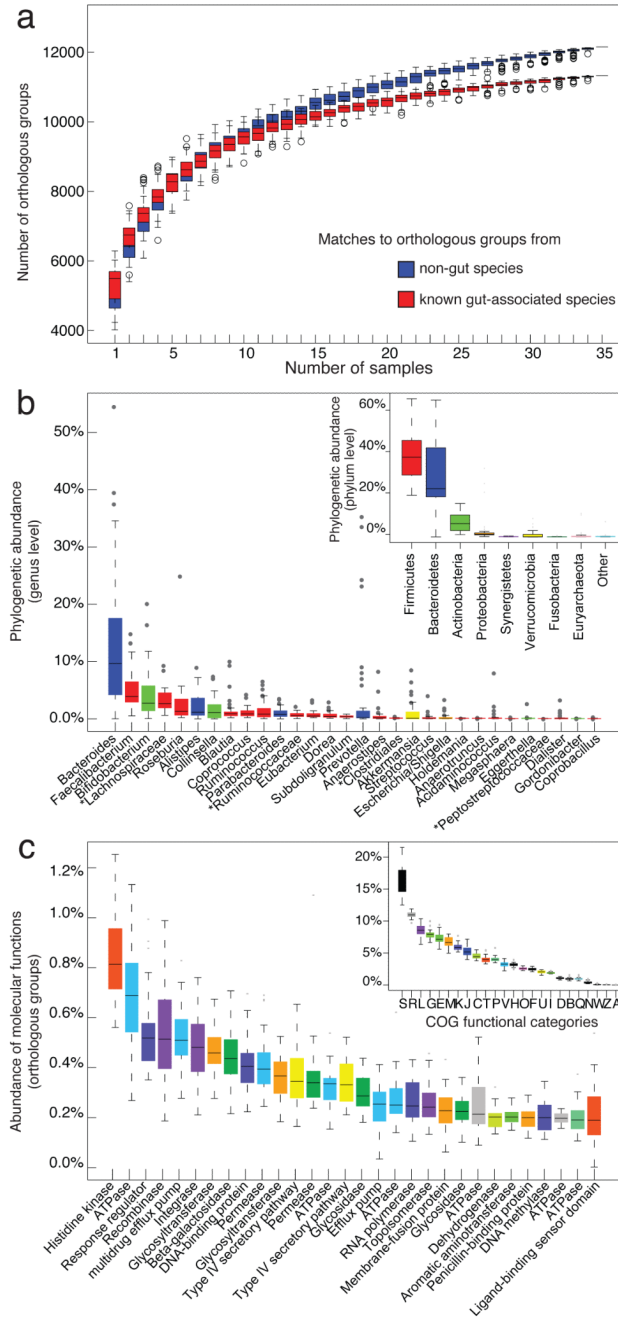
The authors are grateful to Christopher Creevey, Gwen Falony and members of the Bork group at EMBL for helpful discussions and assistance. We thank the EMBL IT core facility and Yan Yuan for managing the high-performance computing resources. The research leading to these results has received funding from the European Community's Seventh Framework Programme (FP7/2007-2013): MetaHIT, grant agreement HEALTH-F4-2007-201052 and from EMBL. Obese/non-obese volunteers for MicroObes study were recruited from the SU.VI.MAX cohort study coordinated by S. Hercberg, and metagenome sequencing was funded by ANR; volunteers for MicroAge study were recruited from the CROWNALIFE cohort study coordinated by A. Cresci, and metagenome sequencing was funded by GenoScope. Ciberehd is funded by the Instituto de Salud Carlos III (Spain). The study was supported by grants from the Lundbeck Foundation Centre for Applied Medical Genomics in Personalized Disease Prediction, Prevention and Care (LuCAMP). JR is supported by the IWIOB and the Odysseus programme of the Fund for Scientific Research Flanders (FWO). BGI was partially funded by the International Science and Technology Cooperation Project in China (0806). We are thankful to the Human Microbiome Project for generating the reference genomes from human gut microbes and the International Human Microbiome Consortium for stimulating discussions and the exchange of data.

## References

1. Eckburg PB, et al. Diversity of the Human Intestinal Microbial Flora. *Science*. 2005; 308:1635–1638. doi:10.1126/science.1110591. [PubMed: 15831718]
2. Hayashi H, Sakamoto M, Benno Y. Phylogenetic Analysis of the Human Gut Microbiota Using 16S rDNA Clone Libraries and Strictly Anaerobic Culture-Based Methods. *MICROBIOLOGY and IMMUNOLOGY*. 2002; 46:535. [PubMed: 12363017]
3. Lay C, et al. Colonic Microbiota Signatures across Five Northern European Countries. *Appl. Environ. Microbiol.* 2005; 71:4153–4155. doi:10.1128/aem.71.7.4153-4155.2005. [PubMed: 16000838]
4. Gill SR, et al. Metagenomic analysis of the human distal gut microbiome. *Science*. 2006; 312:1355–1359. [PubMed: 16741115]
5. Turnbaugh PJ, et al. A core gut microbiome in obese and lean twins. *Nature*. 2009; 457:480–484. [PubMed: 19043404]
6. Kurokawa K, et al. Comparative metagenomics revealed commonly enriched gene sets in human gut microbiomes. *DNA Res*. 2007; 14:169–181. [PubMed: 17916580]
7. Zoetendal EG, Rajilic-Stojanovic M, de Vos WM. High-throughput diversity and functionality analysis of the gastrointestinal tract microbiota. *Gut*. 2008; 57:1605–1615. [PubMed: 18941009]
8. Qin J, et al. A human gut microbial gene catalogue established by metagenomic sequencing. *Nature*. 2010; 464:59–65. [PubMed: 20203603]
9. Raes J, Bork P. Molecular eco-systems biology: towards an understanding of community function. *Nat Rev Microbiol*. 2008; 6:693–699. [PubMed: 18587409]
10. Nelson KE, et al. A catalog of reference genomes from the human microbiome. *Science*. 2010; 328:994–999. doi:10.1126/science.1183605. [PubMed: 20489017]
11. MetaHIT Consortium. MetaHIT draft bacterial genomes at the Sanger Institute. <http://www.sanger.ac.uk/resources/downloads/bacteria/metahit/>
12. Muller J, et al. eggNOG v2.0: extending the evolutionary genealogy of genes with enhanced non-supervised orthologous groups, species and functional annotations. *Nucl. Acids Res*. 2010; 38:D190–195. doi:10.1093/nar/gkp951. [PubMed: 19900971]
13. Palmer C, Bik EM, Digiulio DB, Relman DA, Brown PO. Development of the Human Infant Intestinal Microbiota. *PLoS Biol*. 2007; 5:e177. [PubMed: 17594176]
14. Tap J, et al. Towards the human intestinal microbiota phylogenetic core. *Environmental Microbiology*. 2009; 11:2574–2584. [PubMed: 19601958]
15. Jensen LJ, et al. STRING 8--a global view on proteins and their functional interactions in 630 organisms. *Nucl. Acids Res*. 2009; 37:D412–416. doi:10.1093/nar/gkn760. [PubMed: 18940858]
16. Dethlefsen L, Huse S, Sogin ML, Relman DA. The Pervasive Effects of an Antibiotic on the Human Gut Microbiota, as Revealed by Deep 16S rRNA Sequencing. *PLoS Biol*. 2008; 6:e280. [PubMed: 19018661]
17. Walker A. Say hello to our little friends. *Nat Rev Micro*. 2007; 5:572.

18. Krogfelt KA. Bacterial adhesion: genetics, biogenesis, and role in pathogenesis of fimbrial adhesins of *Escherichia coli*. *Rev Infect Dis*. 1991; 13:721–735. [PubMed: 1681580]
19. Salonen A, et al. Comparative analysis of fecal DNA extraction methods with phylogenetic microarray: effective recovery of bacterial and archaeal DNA using mechanical cell lysis. *J Microbiol Methods*. 2010; 81:127–134. doi:S0167-7012(10)00066-7 [pii] 10.1016/j.mimet.2010.02.007. [PubMed: 20171997]
20. Rajilic-Stojanovic M, et al. Development and application of the human intestinal tract chip, a phylogenetic microarray: analysis of universally conserved phylotypes in the abundant microbiota of young and elderly adults. *Environ Microbiol*. 2009
21. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*. 1987; 20:53–65.
22. Vanhoutte T, Huys G, Brandt E. d. Swings J. Temporal stability analysis of the microbiota in human feces by denaturing gradient gel electrophoresis using universal and group-specific 16S rRNA gene primers. *FEMS Microbiology Ecology*. 2004; 48:437–446. [PubMed: 19712312]
23. Tannock GW, et al. Analysis of the Fecal Microflora of Human Subjects Consuming a Probiotic Product Containing *Lactobacillus rhamnosus* DR20. *Appl. Environ. Microbiol*. 2000; 66:2578–2588. doi:10.1128/aem.66.6.2578-2588.2000. [PubMed: 10831441]
24. Seksik P, et al. Alterations of the dominant faecal bacterial groups in patients with Crohn's disease of the colon. *Gut*. 2003; 52:237–242. doi:10.1136/gut.52.2.237. [PubMed: 12524406]
25. Costello EK, et al. Bacterial Community Variation in Human Body Habitats Across Space and Time. *Science*. 2009; 326:1694–1697. doi:10.1126/science.1177486. [PubMed: 19892944]
26. Martens EC, Koropatkin NM, Smith TJ, Gordon JI. Complex glycan catabolism by the human gut microbiota: the Bacteroidetes Sus-like paradigm. *J Biol Chem*. 2009; 284:24673–24677. [PubMed: 19553672]
27. Wright DP, Rosendale DI, Robertson AM. Prevootella enzymes involved in mucin oligosaccharide degradation and evidence for a small operon of genes expressed during growth on mucin. *FEMS Microbiology Letters*. 2000; 190:73–79. doi:10.1111/j.1574-6968.2000.tb09265.x. [PubMed: 10981693]
28. Derrien M, Vaughan EE, Plugge CM, de Vos WM. *Akkermansia muciniphila* gen. nov., sp. nov., a human intestinal mucin-degrading bacterium. *Int J Syst Evol Microbiol*. 2004; 54:1469–1476. doi:10.1099/ijs.0.02873-0. [PubMed: 15388697]
29. Ley RE, Turnbaugh PJ, Klein S, Gordon JI. Microbial ecology: human gut microbes associated with obesity. *Nature*. 2006; 444:1022–1023. [PubMed: 17183309]
30. Schwartz A, et al. Microbiota and SCFA in Lean and Overweight Healthy Subjects. *Obesity*. 2009; 18:190. [PubMed: 19498350]
31. Woodmansey EJ. Intestinal bacteria and ageing. *J Appl Microbiol*. 2007; 102:1178–1186. doi:JAM3400 [pii] 10.1111/j.1365-2672.2007.03400.x. [PubMed: 17448153]
32. Kovacicova G, Skorupski K. The alternative sigma factor sigma(E) plays an important role in intestinal survival and virulence in *Vibrio cholerae*. *Infect Immun*. 2002; 70:5355–5362. [PubMed: 12228259]
33. Fujihashi K, Kiyono H. Mucosal immunosenescence: new developments and vaccines to control infectious diseases. *Trends Immunol*. 2009; 30:334–343. [PubMed: 19540811]
34. Turnbaugh PJ, et al. An obesity-associated gut microbiome with increased capacity for energy harvest. *Nature*. 2006; 444:1027–1031. [PubMed: 17183312]
35. Raes J, Korb J, Lercher MJ, von Mering C, Bork P. Prediction of effective genome size in metagenomic samples. *Genome Biol*. 2007; 8:R10. [PubMed: 17224063]
36. Gibson GR, et al. Alternative pathways for hydrogen disposal during fermentation in the human colon. *Gut*. 1990; 31:679–683. [PubMed: 2379871]
37. Godon JJ, Zumstein E, Dabert P, Habouzit F, Moletta R. Molecular microbial diversity of an anaerobic digester as determined by small-subunit rDNA sequence analysis. *Appl Environ Microbiol*. 1997; 63:2802–2813. [PubMed: 9212428]
38. Arumugam M, Harrington ED, Foerster KU, Raes J, Bork P. SmashCommunity: a metagenomic annotation and analysis tool. *Bioinformatics*. 2010; 26:2977–2978. doi:10.1093/bioinformatics/btq536. [PubMed: 20959381]

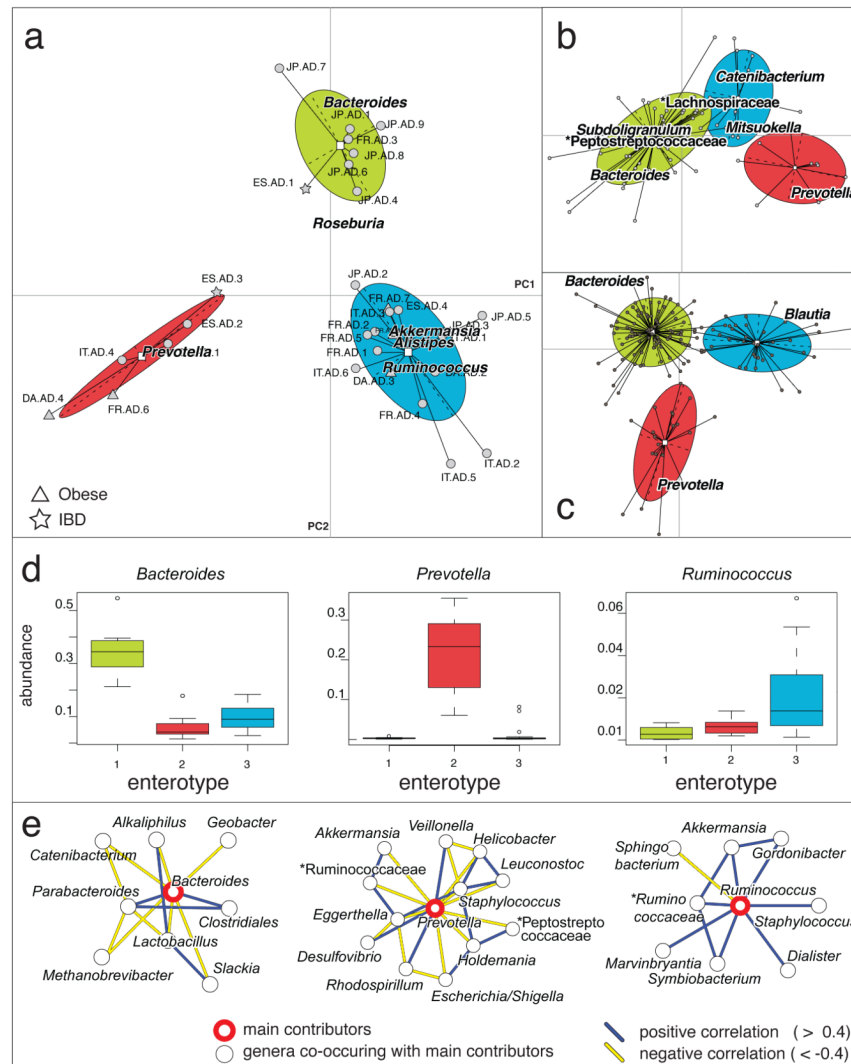
39. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian classifier for rapid assignment of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol.* 2007; 73:5261–5267. [PubMed: 17586664]
40. Gianoulis TA, et al. Quantifying environmental adaptation of metabolic pathways in metagenomics. *Proc Natl Acad Sci U S A.* 2009; 106:1374–1379. [PubMed: 19164758]



**Fig. 1. Functional and phylogenetic profiles of human gut microbiome**

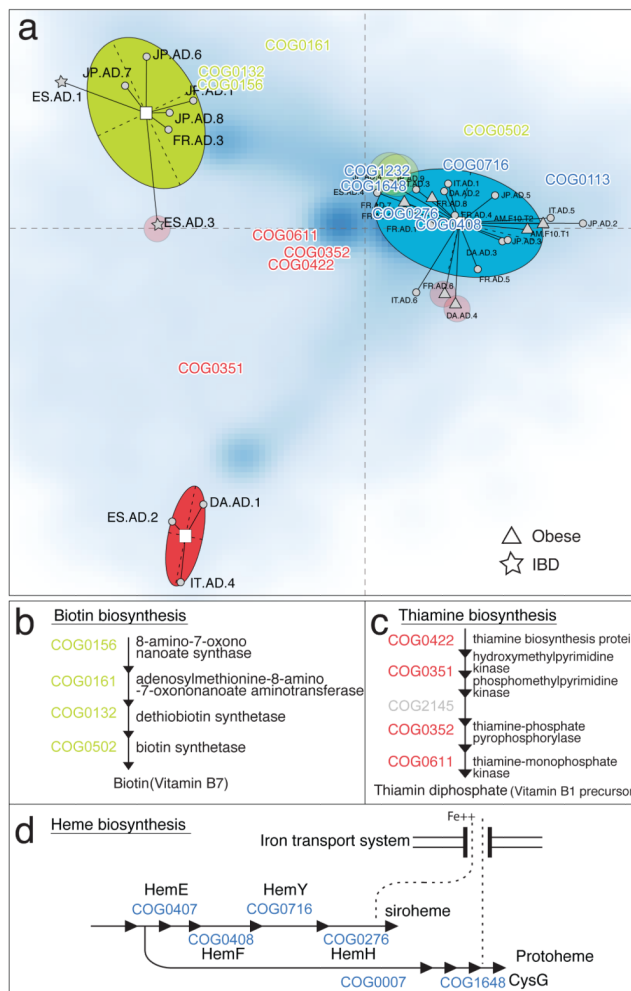
**(a)** Simulation of the detection of distinct orthologous groups (OGs) when increasing the number of individuals (samples). Complete genomes were classified by habitat-information and the OGs divided into those that occur in known gut-species (red) and those that have not yet associated to gut (blue). The former are close to saturation when sampling 35 individuals (excluding infants) whereas functions from non-gut (probably rare and transient) species are not. **(b)** Genus abundance variation box plot for the 30 most abundant genera as determined by read abundance. Genera are colored by their respective phylum (see inset for color key). Inset: phylum abundance box plot. Genus and phylum level abundances were measured

using reference genome based mapping with 85% and 65% sequence similarity cutoffs. Unclassified genera under a higher rank are marked by asterisks. **(c)** Orthologous group (OG) abundance variation box plot for the 30 most abundant OGs as determined by assignment to eggNOG<sup>12</sup>. OGs are colored by their respective functional category (see inset for color key). Inset: abundance box plot of 24 functional categories.



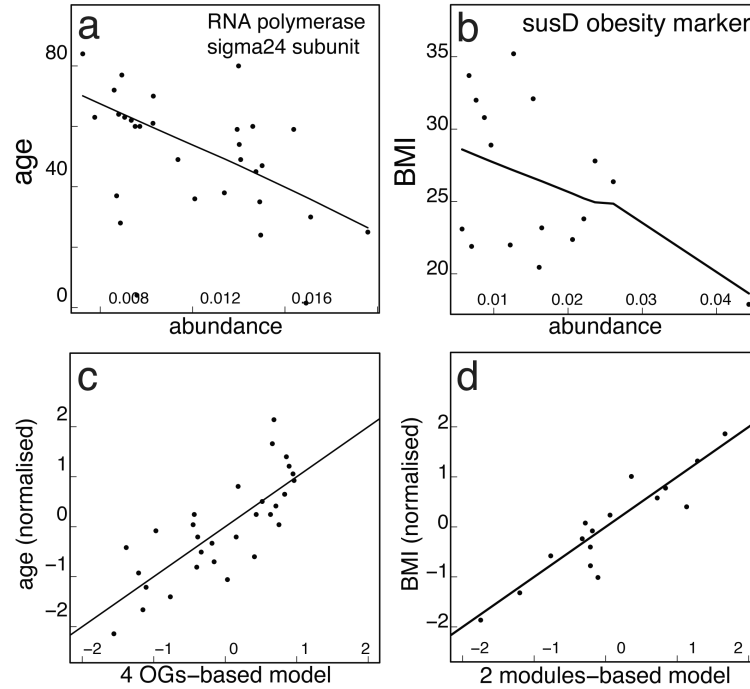
**Fig. 2. Phylogenetic differences between enterotypes**

Between class analysis, which visualizes results from Principal Component Analysis and clustering, of the genus compositions of (a) 33 Sanger metagenomes estimated by mapping the metagenome reads to 1511 reference genome sequences using an 85% similarity threshold, (b) Danish subset containing 85 metagenomes from a published Illumina dataset<sup>8</sup> and (c) 154 pyrosequencing-based 16S sequences<sup>5</sup> reveal three robust clusters that we call enterotypes. Two principal components are plotted using the ade4 package in R with each sample represented by a filled circle. The center of gravity for each cluster is marked by a rectangle and the colored ellipse covers 67% of the samples belonging to the cluster. (d) Abundances of the main contributors of each enterotype from the Sanger metagenomes. (e) Co-occurrence networks of the three enterotypes from the Sanger metagenomes. Unclassified genera under a higher rank are marked by asterisks in (b) and (e).



**Fig. 3. Functional differences between enterotypes**

(a) Between class analysis (see Fig. 2) of orthologous group (OG) abundances showing only minor disagreements with enterotypes (transparent circles indicate the differing samples). The blue cloud represents the local density estimated from the coordinates of OGs; positions of selected OGs are highlighted. (b) Four enzymes in the biotin biosynthesis pathway (COG0132, COG0156, COG0161 and COG0502) are overrepresented in enterotype 1. (c) Four enzymes in the thiamine biosynthesis pathway (COG0422, COG0351, COG0352 and COG0611) are overrepresented in enterotype 2. (d) Six enzymes in the heme biosynthesis pathway (COG0007, COG0276, COG407, COG0408, COG0716 and COG1648) are overrepresented in enterotype 3.



**Fig. 4. Correlations with host properties**

(a) Pairwise correlation of RNA polymerase facultative sigma24 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 body mass index ( $p=0.27$ ,  $\rho=-0.29$ , weak correlation). (c) Multiple 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;  $p=2.75e-05$ , adjusted  $R^2=0.57$ ). (d) Two modules, ATPase complex and ectoine biosynthesis (M00051), significantly correlating with BMI when combined into a linear model ( $p=6.786e-06$ , adjusted  $R^2=0.82$ ).