

The Dutch Microbiome Project defines factors that shape the healthy gut microbiome

Rinse Weersma (✉ r.k.weersma@umcg.nl)

University of Groningen and University Medical Center Groningen

Ranko Gacesa

University of Groningen and University Medical Center Groningen

Alexander Kurilshikov

University Medical Center Groningen <https://orcid.org/0000-0003-2541-5627>

Arnau Vich Vila

<https://orcid.org/0000-0003-4691-5583>

Trishla Sinha

University of Groningen and University Medical Center Groningen

Marjolein Klaassen

University of Groningen and University Medical Center Groningen

Laura Bolte

University of Groningen and University Medical Center Groningen

Sergio Andreu-Sanchez

University Medical Center Groningen <https://orcid.org/0000-0002-3503-9971>

Lianmin Chen

University Medical Center Groningen <https://orcid.org/0000-0003-0660-3518>

Valerie Collij

UMCG

Shixian Hu

UMCG

Jackie Dekens

UMCG

Virissa Lenters

University Medical Centre Utrecht

Johannes Björk

University of Notre Dame <https://orcid.org/0000-0001-9768-1946>

J. Casper Swarte

University of Groningen and University Medical Center Groningen <https://orcid.org/0000-0002-3709-9193>

Morris Swertz

University of Groningen <https://orcid.org/0000-0002-0979-3401>

B. H. Jansen

University of Groningen and University Medical Center Groningen

Jody Gelderloos-Arends

University of Groningen and University Medical Center Groningen

Marten Hofker

University of Groningen and University Medical Center Groningen

Roel Vermeulen

Utrecht University

Serena Sanna

University Medical Center Groningen <https://orcid.org/0000-0002-3768-1749>

Hermie Harmsen

University of Groningen and University Medical Center Groningen

Cisca Wijmenga

University Medical Centre Groningen <https://orcid.org/0000-0002-5635-1614>

Jingyuan Fu

University Medical Center Groningen <https://orcid.org/0000-0001-5578-1236>

Alexandra Zhernakova

University Medical Center Groningen <https://orcid.org/0000-0002-4574-0841>

Biological Sciences - Article

Keywords: Exposome, Three-generational Cohort, Environment, Cohousing, Heritability, Exposure

Posted Date: December 18th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-117376/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **The Dutch Microbiome Project defines factors that shape the healthy**
2 **gut microbiome**

3 R. Gacesa^{1,2*}, A. Kurilshikov^{2*}, A. Vich Vila^{1,2}, T. Sinha², M.A.Y. Klaassen^{1,2}, L.A. Bolte^{1,2}, S.
4 Andreu-Sánchez^{2,3}, L. Chen^{2,3}, V. Collij^{1,2}, S. Hu^{1,2}, J.A.M. Dekens^{2,4}, V.C. Lenters⁵, J.R. Björk^{1,2},
5 J.C. Swarte^{1,2}, M.A. Swertz^{2,6}, B.H. Jansen^{1,2}, J. Gelderloos-Arends², M. Hofker^{3¶}, R.C.H.
6 Vermeulen^{5,7}, S. Sanna^{2,8}, H.J.M. Harmsen^{9#}, C. Wijmenga^{2#}, J. Fu^{2,3#^}, A. Zhernakova^{2#^}, R.K.
7 Weersma^{1#^}

8 * shared first authors: RG and AK

9 # shared last authors: HH, CW, JF, AZ, RW

10 ^ corresponding authors: JF, AZ, RW

11 ¶ deceased

12
13
14 **Author affiliations**

15 ¹ University of Groningen and University Medical Center Groningen, Department of
16 Gastroenterology and Hepatology, Groningen, the Netherlands

17 ² University of Groningen and University Medical Center Groningen, Department of
18 Genetics, Groningen, The Netherlands

19 ³ Department of Pediatrics, University of Groningen and University Medical Center
20 Groningen, Groningen, The Netherlands

21 ⁴ University Medical Center Groningen, Center for Development and Innovation

22 ⁵ University Medical Centre Utrecht, Julius Center for Health Sciences and Primary Care,
23 3584 CG Utrecht, The Netherlands

24 ⁶: University of Groningen and University Medical Center Groningen, Genomics Coordination
25 Center, Groningen, The Netherlands

26 ⁷ Utrecht University, Institute for Risk Assessment Sciences (IRAS), Department of Population
27 Health Sciences, 3584 CM Utrecht, The Netherlands

28 ⁸ Institute for Genetic and Biomedical Research (IRGB), National Research Council (CNR),
29 Cagliari, Italy

30 ⁹ Department of Medical Microbiology and Infection prevention: University of Groningen
31 and University Medical Center Groningen, Groningen, The Netherlands

32 **Abstract**

33 The gut microbiome is associated with diverse diseases, but the universal signature of an
34 (un)healthy microbiome remains elusive and there is a need to understand how genetics,
35 exposome, lifestyle and diet shape the microbiome in health and disease. To fill this gap, we
36 profiled bacterial composition, function, antibiotic resistance and virulence factors in the
37 gut microbiomes of 8,208 Dutch individuals from a three-generational cohort comprising
38 2,756 families. We then correlated this to 241 host and environmental factors, including
39 physical and mental health, medication use, diet, socioeconomic factors and childhood and
40 current exposome. We identify that the microbiome is primarily shaped by environment and
41 cohousing. Only ~13% of taxa are heritable, which are enriched with highly prevalent and
42 health-associated bacteria. By identifying 2,856 associations between microbiome and
43 health, we find that seemingly unrelated diseases share a common signature that is
44 independent of comorbidities. Furthermore, we identify 7,519 associations between
45 microbiome features and diet, socioeconomics and early life and current exposome, of
46 which numerous early-life and current factors are particularly linked to the microbiome.
47 Overall, this study provides a comprehensive overview of gut microbiome and the
48 underlying impact of heritability and exposures that will facilitate future development of
49 microbiome-targeted therapies.

50 **Introduction**

51 Alterations in gut microbiota composition and function are associated with a broad range of
52 human health disorders, from gastrointestinal (GI) and metabolic diseases to mental
53 disorders¹⁻³. The influence of the gut bacteria and microbial pathways on host metabolism
54 and immunity, together with the fact that the microbiota can be modified, have led to
55 heightened interest in developing microbiota-targeted therapies^{4,5}. In this context, however,
56 the characteristics of a “healthy” microbiome remain largely unclear, as does the extent to
57 which the gut microbiome is driven by intrinsic and external factors that might be amenable
58 to microbiota-targeted therapies.

59 The capacity to define a “healthy” microbiome has been hampered by the large variation in
60 microbiome composition between individuals and the large variability in the processing of
61 faecal samples between studies. Population-based studies, including our previous study on
62 ~1500 Dutch volunteers⁶, have shown that this variation is partially accounted for by diet,
63 lifestyle, host genetics and environmental factors, including early-life exposures⁶⁻⁸.
64 However, while the number of microbiome studies has increased exponentially, in-depth
65 integrative analyses of large standardized, well-phenotyped cohorts remain scarce, even
66 though this integrative perspective is essential to disentangle meaningful host–microbiota
67 associations and identify potential targets for microbiota-directed interventions.

68 To address these issues, we initiated the Dutch Microbiome Project (DMP) and analysed the
69 gut microbiome composition and functionality of 8,208 individuals from the Lifelines study –
70 a three-generational population cohort from a geographic region of approximately 11,400
71 km² in the Northern Netherlands⁹. By exploring 241 host characteristics, we determined the
72 impact of shared household, socioeconomic status, lifestyle and environmental exposures
73 on the gut microbiota, which together explain 12.9% of the inter-individual variation in gut
74 microbiota composition and 16.3% of the variation in microbiome functionality. Based on
75 family relations, we estimated the heritability of the gut microbiota, which varies widely
76 between different species (max $H^2 = 0.3$). We also define the core microbiota and keystone
77 species in the taxonomic networks of this Dutch population. After correcting for the impact
78 of technical factors, age, sex, BMI and Bristol stool scale of faecal samples, we identify clear
79 patterns of microbial taxa and pathway associations between multiple unrelated groups of
80 diseases and with self-perceived health, which allows us to define features of a healthy or
81 unhealthy microbiome and to predict the health status independently of comorbidities.
82 Finally, we identified intrinsic and external factors, including environmental pressures (i.e.
83 exposome), that influence the (un)healthy microbiome. Overall, this study provides a
84 comprehensive large-scale exploration of the gut microbiome that will facilitate future
85 investigations of microbiome-targeted therapies for improvement of long-term human
86 health.

87 **Methods**

88 ***Population cohort and metadata collection***

89 The Lifelines Dutch Microbiome Project (DMP) cohort has been developed as a part of the
90 Lifelines cohort study. Lifelines is a multi-disciplinary prospective population-based cohort
91 study using a unique three-generation design to examine the health and health-related
92 behaviours of 167,729 people living in the North of the Netherlands. It employs a broad
93 range of investigative procedures to assess the biomedical, socio-demographic, behavioural,
94 physical and psychological factors that contribute to health and disease in the general
95 population, with a special focus on multi-morbidity and complex genetics^{9,10}. A total of
96 8,719 distinct fresh frozen faecal and blood samples were collected from Lifelines
97 participants in 2015 and 2016 to form the DMP cohort. Whole-genome shotgun sequencing
98 was performed on one aliquot from each of 8,534 faecal samples, of which 8,208 were
99 retained for downstream analysis after stringent quality control. Metadata information
100 collected from the participants was grouped into the following categories: family structure,
101 diseases, gastrointestinal (GI) complaints, general health score, medication use,
102 anthropometrics, birth-related factors, reported childhood (< 16 years) exposures, current
103 exposome (air pollutants, greenspace, urbanicity, pets and smoking), socioeconomic
104 characteristics and diet (Supplementary table 2d).

105 ***Informed consent***

106 The Lifelines study was approved by the medical ethical committee from the University
107 Medical Center Groningen (METc number: 2017/152). Additional written consent was
108 signed by all DMP participants or their parents or legal representatives (for children aged
109 under 18).

110 ***Metadata***

111 Metadata was collected by questionnaires and curated as described previously¹⁰ and below.
112 We included 241 phenotypes from a broad range of categories, including socioeconomic
113 factors, self-reported diseases and medications, quality of life, mental health, education and
114 employment, nutrition, smoking, stress and childhood environmental factors.
115 Questionnaires were developed and processed by the Lifelines cohort study¹⁰ as described
116 at www.lifelines.nl. Additional in-depth data curation and acquisition was performed to
117 assess dietary intake, air pollution and environmental exposures, medication use and gut
118 health, as described below.

119 ***Diet***

120 Habitual diet was assessed through a semi-quantitative Food Frequency Questionnaire (FFQ)
121 collected 4 years prior to DMP faecal sampling¹⁰. The FFQ was designed and validated by the

122 division of Human Nutrition of Wageningen University using standardized methods¹¹. It
123 assesses how often a food was consumed over the previous month on a scale ranging from
124 'never' to '6-7 days per week', along with the usual amount taken. The average daily
125 nutrient intake was calculated using the Dutch Food Composition database (NEVO, RIVM)
126 and a mono- and disaccharide-specific food composition table¹², resulting in the generation
127 of data on 21 dietary factors. Energy adjustment was performed by means of the nutrient
128 density method¹³. Published dietary scores and inter-nutrient ratios were calculated as
129 indicators of dietary quality and composition^{14,15}.

130 To validate the assumed stability of FFQs across time^{12,16}, we used questionnaires from
131 128,501 Lifelines participants to study diet consistency between the baseline questionnaire,
132 collected 4 years prior to this study, and a second smaller-scope nutrient-specific
133 questionnaire collected concurrently with DMP faecal sampling. Sixty-five dietary questions,
134 reflecting consumption of major food categories such as fruits, vegetables, fish, meat,
135 bread, grains and sweets, as well as special dietary regimes (e.g. vegan or macrobiotic diet),
136 were compared between the first and second time point. The majority of dietary items were
137 available for >44,000 individuals at both timepoints (Supplementary table 8). To quantify the
138 consistency of answers and the degree of change in the 4-year period, we calculated the
139 consistency (proportion of identical answers) and Euclidean distance between all
140 participants who answered a question the same way at baseline and follow-up and the
141 distance between participants who answered differently at follow-up (Supplementary table
142 8).

143 ***Exposome***

144 Elements of the exposome, neighbourhood urbanicity and income were assessed for the
145 participant's home address at the time of faecal sampling. Exposure to two air pollutants,
146 particulate matter with aerodynamic diameter $\leq 2.5 \mu\text{m}$ (PM_{2.5}) and nitrogen dioxide (NO₂),
147 was assigned based on land use regression models developed in the European Study of
148 Cohorts for Air Pollution Effects (ESCAPE) project^{17,18}. These estimates are based on
149 measurement data from 2009 and reflect long-term ambient air pollution exposures¹⁹.

150 Greenspace was assigned using the Normalized Difference Vegetation Index (NDVI), which
151 reflects the average density of green vegetation within a 100-meter circular buffer around
152 the participant's residential address. The NDVI was derived from a LANDSAT 5 (TM) satellite
153 image taken in 2016 and captures the density of green vegetation at a spatial resolution of
154 30x30 m based on land surface reflectance of visible (red) and near-infrared parts of
155 spectrum.

156 Neighbourhood urbanicity was assigned based on a five-category scale of surrounding
157 address density developed by Statistics Netherlands (1 = very urban, ≥ 2500 addresses per

158 km² to 5 = very rural, < 500 addresses per km², data for 2015)²⁰. Neighbourhood income was
159 considered a proxy of neighbourhood socioeconomic position and defined as the proportion
160 of persons with low (< 40th percentile) income (Statistics Netherlands, data for 2015)²⁰.

161 ***Stool characteristics, diseases and medication***

162 Participants recorded a bowel movement diary, Bristol stool scale, daily medication use and
163 GI symptoms daily for seven days in the week of stool sample collection, and these records
164 were used to extract information on drug use, stool characteristics, stool frequency and GI
165 symptoms during the week of stool collection. The validated ROME III questionnaires²¹ were
166 used to characterize functional GI disorders, and participants were classified as having either
167 no functional GI diseases or irritable bowel syndrome, functional diarrhoea, functional
168 constipation or functional bloating. Information about the presence of other diseases was
169 self-reported and collected using Lifelines questionnaires. Diseases were grouped into 11
170 disease categories. The presence of cancer was grouped into a separate category defined as
171 “any cancer”, independent of cancer type. Non-alcoholic fatty liver disease fibrosis score²²
172 and fatty liver index²³ were calculated from the anthropometrics and blood measurements,
173 as described previously^{22,23}. Diseases with < 20 cases were excluded from further analysis.
174 Self-reported medications were grouped into categories based on Anatomical Therapeutic
175 Chemical classification (ATC codes) at the most-specific ATC level (5-digit ATC code if
176 possible). ATC categories with < 20 users were grouped into a higher level (4-digit or 3-digit)
177 ATC class, and categories with < 20 individuals that could not be grouped according to ATC
178 classification were excluded from further analysis. In total, 62 drug groups were included
179 (Supplementary table 2a).

180 ***Sample collection, DNA extraction and sequencing***

181 Faecal sample collection was performed by participants at home. Participants were asked to
182 freeze stool samples within 15 min of stool production. The frozen samples were collected
183 by Lifelines personnel, transported to the Lifelines biorepository on dry ice and stored at -
184 80°C until DNA extraction. Microbial DNA was isolated with the QIAamp Fast DNA Stool Mini
185 Kit (Qiagen, Germany), according to the manufacturer’s instructions, using the QIAcube
186 (Qiagen) automated sample preparation system. Metagenomic sequencing was performed
187 at Novogene, China using the Illumina HiSeq 2000 platform to generate approximately 8 Gb
188 of 150 bp paired-end reads per sample (mean 7.9 gb, st.dev 1.2 gb).

189 ***Profiling microbiome composition and function***

190 Metagenomes were profiled consistent with previous data analysis of 1000IBD²⁴ and
191 Lifelines-DEEP¹⁰ cohorts, as follows. KneadData tools (v0.5.1)²⁵ were used to process
192 metagenomic reads (in fastq format) by trimming the reads to PHRED quality 30 and

193 removing Illumina adapters. Following trimming, the KneadData integrated Bowtie2 tool
194 (v2.3.4.1)²⁶ was used to remove reads that aligned to the human genome (GRCh37/hg19).

195 Taxonomic composition of metagenomes was profiled by MetaPhlAn2 tool (v2.7.2)²⁷ using
196 the MetaPhlAn database of marker genes mpa_v20_m200. Profiling of genes encoding
197 microbial biochemical pathways was performed using the HUMAnN2 pipeline (v0.11.1)²⁸
198 integrated with the DIAMOND alignment tool (v0.8.22)²⁸, UniRef90 protein database
199 (v0.1.1) and ChocoPhlAn pan-genome database (v0.1.1). As a final quality control step,
200 samples with unrealistic microbiome composition (eukaryotic or viral abundance > 25% of
201 total microbiome content or total read depth < 10 million) were excluded, leaving 8,208
202 samples for further analyses. Analyses were performed using locally installed tools and
203 databases on CentOS (release 6.9) on the high-performance computing infrastructure
204 available at our institution and using the MOLGENIS data platform²⁹.

205 In total, we detected 1,253 taxa (4 kingdoms, 21 phyla, 35 classes, 62 orders, 128 families,
206 270 genera and 733 species) and 564 pathways in at least one of the samples in the quality
207 controlled-dataset. To deal with sparse microbial data in the downstream analysis, we
208 focused on bacterial and archaeal species/pathways with a mean relative abundance >
209 0.01% that are present in at least 5% of participants. This yielded 257 taxa (6 phyla, 11
210 classes, 15 orders, 30 families, 59 genera and 136 species) and 277 pathways. Together,
211 these microbial features accounted for 97.86% and 87.82% of the average taxonomic and
212 functional compositions, respectively.

213 Based on the abundance profiles of taxa that passed the filtering process, alpha diversity, as
214 measured by richness and Shannon entropy, was calculated at family-, genus- and species-
215 level using *specnumber* and the function *diversity*, respectively, in R package *vegan* (v.3.6.1).
216 Rarefaction and extrapolation (R/E) sampling curves for estimation of total richness of
217 species and genera in the population were constructed using a sample size-based
218 interpolation/extrapolation algorithm implemented in the *iNEXT* package for R³⁰.

219 ***Profiling of bacterial virulence factors and antibiotic resistance genes***

220 Metagenomes were searched for bacterial virulence factors (VFs) using the shortBRED
221 toolkit (v0.9.5)³¹ and the database of Virulence Factors of Pathogenic Bacteria (VFDB) core
222 dataset of DNA sequences (downloaded on 01/11/2018)³². The shortBRED tool
223 *shortbred_identify.py* (v0.9.5) was used to identify unique markers for VFs, with the
224 UniRef90 database (downloaded on 01/11/2018) used as negative control, and the
225 *shortbred_quantify.py* tool (v0.9.5) was used to perform a quantification of these markers in
226 metagenomes. Quantification of antibiotic resistance genes was performed using shortBRED
227 tool *shortbred_quantify.py* (v0.9.5), with markers generated using *shortbred_identify.py*
228 (v0.9.5) on the CARD database of bacterial antibiotic resistance genes (downloaded

229 01/11/2018)³³, with the UniRef90 database used as negative control. This identified 190 VFs
230 and 303 antibiotic resistance gene families (ARs), of which 47 VFs and 98 ARs were present
231 in at least 5% of participants with relative abundance > 0.01%. These accounted for 95.22%
232 of VF composition and 98.08% of AR composition, respectively.

233 ***Estimation of heritability of microbiome***

234 To accommodate covariate effects, we estimated the heritability of microbiome features
235 using a variance components model implemented in the software POLY (v.0.5.1)^{34,35}. We
236 first considered a base model in which variance is partitioned into a polygenic component,
237 V_g , shared between individuals that is proportional to their kinship coefficient, and an
238 environmental component, V_e , that is unique to each individual. Thus, if Y is the measured
239 trait, the variance of Y is $\text{Var}(Y) = V_g + V_e$, and its broad heritability is $H^2 = V_g / (V_g + V_e)$. Of note,
240 this measure reflects the overall impact of genes on the phenotype, thus including all
241 potential models of action, as opposed to narrow heritability, which only includes additive
242 effects. After fitting this base model, we also considered two refined models that included i)
243 an additional variance component to model a shared environment for people belonging to
244 the same family, V_h , and ii) a covariate with values 0/1 that distinguishes family members
245 currently living in the same house from those who do not. We compared the significance of
246 these two models to the basic model using a likelihood ratio test and found the base model
247 to be the most appropriate fit for all microbiome traits.

248 We restricted the analysis of heritability to the relative abundances of the 242 microbial
249 taxa present in at least 250 individuals and focused on 4,745 individuals in 2,756 families in
250 which at least two individuals had available microbiome data. In total, the analysis included
251 2,756 parent-child pairs, 530 sibling pairs and 815 pairs with second-degree or more distant
252 relationships. All models were adjusted for age, sex, BMI, read depth and stool frequency,
253 and values of relative abundances of taxa were transformed using the centred log-ratio (clr)
254 transformation³⁶. Benjamini-Hochberg correction was used to control the multiple testing
255 false discovery rate (FDR), and results with $FDR < 0.1$ were considered significant.

256 ***Estimation of the effect of co-housing on microbiome***

257 We estimated the effect of co-housing on microbiome composition, function, ARs and VFs
258 by comparing beta-diversities of the microbiomes of co-housing study participants (1,710
259 unrelated pairs, 285 parent-child pairs and 144 sibling pairs) to those of participants who
260 did not share housing (2,000 unrelated pairs, 301 parent-child pairs and 299 sibling pairs).
261 Microbiome distance was calculated for every pair using Bray-Curtis dissimilarity, and mean
262 dissimilarities within groups were compared using Mann-Whitney U tests using the
263 Benjamini-Hochberg correction to control multiple testing FDR. Results were considered
264 significant at $FDR < 0.05$.

265 ***Calculation of microbiome–phenotype associations***

266 The microbiome composition variance explained by phenotypes was calculated by
267 permutational multivariate analysis of variance using distance matrices, implemented in the
268 *adonis* function for R package *vegan* (v.2.4-6), using 20,000 permutations and a Bray-Curtis
269 distance matrix calculated using relative abundances of microbial species. A separate
270 analysis was performed to calculate the microbiome functional potential explained by
271 phenotypes using equivalent methodology. The functional dissimilarity matrix was
272 calculated using the Bray-Curtis dissimilarity index calculated on the relative abundances of
273 MetaCyc microbial biochemical pathways.

274 Prior to association analysis of phenotypes and microbiome features, the microbiome data
275 was transformed using the clr transformation. The geometric mean for clr transformation of
276 relative abundances of taxa was calculated on species-level and applied to higher levels. The
277 associations between phenotypes and microbial features (microbial taxa, MetaCyc
278 functional pathways, CARD and VFDB entities) were calculated using linear regression,
279 adjusting for age, sex and BMI of the individual along with Bristol stool scale of the faecal
280 sample and technical factors (DNA concentration, sequencing read depth, sequencing batch
281 and sampling season). Benjamini-Hochberg correction was used to control for multiple
282 testing with the number of tests equal to the number of tested feature–phenotype pairs.
283 Results were considered significant at $FDR < 0.05$.

284 ***Quantification of microbiome–disease–drug interactions***

285 To disentangle interactions between the gut microbiome, medication and diseases, we
286 explored the effect of a selection of drugs and the common diseases for which these drugs
287 are used: functional GI disorders and proton pump inhibitors, type 2 diabetes and
288 antidiabetic drugs, and depression and selective serotonin reuptake inhibitors. We extracted
289 subsets of participants with the disease and those who report using disease-associated
290 medication and used these subsets to construct multivariate models including both
291 medication use and disease, additionally adjusting for other factors as described above in
292 *Calculation of microbiome-phenotype associations*.

293 ***Definition of core microbiome and prediction of keystone microbiome features***

294 To identify core microbial species and pathways, we used a bootstrapping-based selection
295 approach. We randomly sampled 1% to 100% of the samples of the cohort a hundred times
296 and calculated the standard deviation of the presence rate of each microbial
297 species/pathway at different sampling percentages. Microbial features with a presence rate
298 of more than 95% of samples were defined as the core microbiome, and we identified 9
299 core microbial species and 143 core microbial pathways (Supplementary table 1a,b).

300 To analyse microbiome community structure, we constructed microbial species and
301 pathway co-abundance networks using SparCC, as previously published^{37,38}. Relative
302 abundances of taxa were converted to estimated read counts by multiplying abundance
303 percentages by total sequenced reads per sample after quality control. For pathway
304 analysis, the read counts (RPKM) from HUMAnN2 were directly used for SparCC. Significant
305 co-abundance was controlled at FDR 0.05 level using 100 permutations. In each
306 permutation, the abundance of each microbial feature was randomly shuffled across
307 samples. In this way, we obtained 6,473 species and 55,407 pathway co-abundances at FDR
308 < 0.05. Features that ranked in the top 20% in the number of network connections (node
309 degree) were considered keystone species or pathways, resulting in 28 keystone species and
310 53 keystone pathways (Supplementary table 1e).

311 ***Identification of microbiome clusters***

312 To identify microbial clusters and assess the presence of gut enterotypes in our cohort, we
313 performed the partitioning around the medoid method on the relative abundances of
314 microbial species and used the Calinski-Harabasz index to select the optimal number of
315 clusters, as previously published in a study of gut enterotypes³⁹. Enrichment of phenotypes
316 in each cluster was assessed by logistic regression in R.

317 ***Calculation of microbiome signatures predictive of diseases and health***

318 We calculated the microbial signatures predictive of the 36 most common ($N_{\text{cases}} > 100$)
319 diseases in our dataset. In addition, we defined a “healthy” phenotype as an absence of any
320 self-reported disease. Using this definition, 2,937 (36%) out of 8,208 individuals were
321 defined as “healthy”.

322 To build prediction models for common diseases, the dataset was randomly split into
323 training (90%) and test (10%) sets. Next, we performed elastic net L1/L2 regularized
324 regression (R package *glmnet* v.4.0) on the training set, using Shannon diversity, clr-
325 transformed microbial taxa, clr-transformed MetaCyc bacterial pathways and age, sex and
326 BMI as fixed covariates (not penalized in the models). The model for each disease was
327 calculated independently using five-fold cross-validation to select the optimal lambda
328 penalization factor (at L1/L2 mixing parameter alpha fixed at 0.5). The lambda with minimal
329 cross-validation error was used in the downstream analysis.

330 In total, we defined three probabilistic models: a “null” signature that only includes effects
331 of general covariates (age, sex and BMI), a “microbiome” signature that includes all selected
332 microbiome features and a “combined” signature that includes both the effects of
333 microbiome features and general covariates.

334 ***Data availability***

335 Raw sequencing data and corresponding metadata that support the findings of this study
336 are available on request from the corresponding author R.K.W. The participant metadata
337 are not publicly available as they contain information that could compromise research
338 participant privacy/consent.

339 The authors declare that all other data supporting the findings of this study are available
340 within the paper and its supplementary information files.

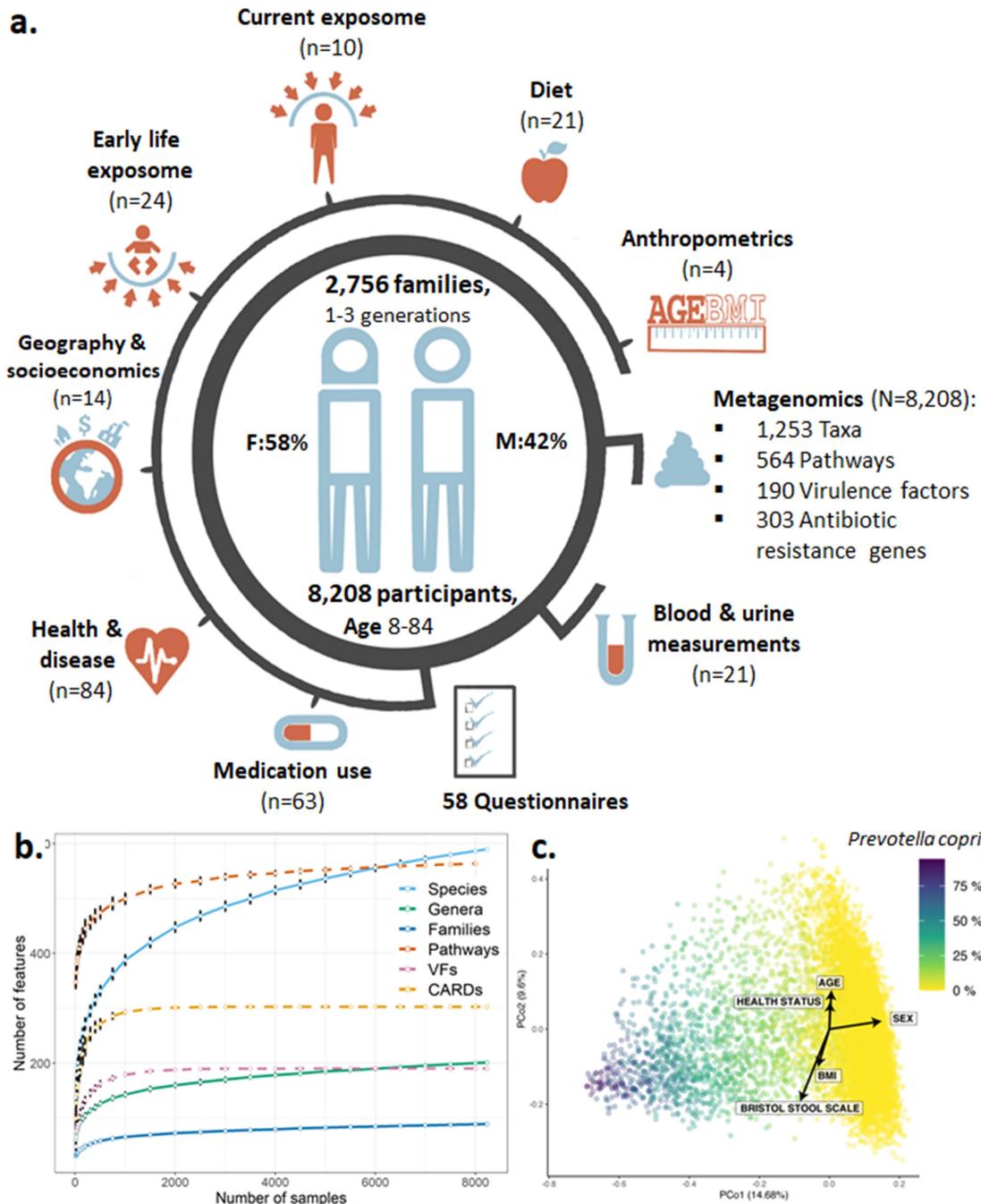
341

342 **Results**

343 ***Overall picture of gut microbial composition***

344 In this study we characterized the composition and function of the gut microbiota of 8,208
345 individuals from the Northern Netherlands. Participants were of a wide age range (8-84
346 years), 57.4% were female and 4,745 individuals clustered into 2,756 families. The
347 participants were largely (99.5%) of Dutch European ancestry (Fig.1a, Supplementary table
348 2d).

349 A total of 1,253 taxa (4 kingdoms, 21 phyla, 35 classes, 62 orders, 128 families, 270 genera
350 and 733 species) and 564 metabolic pathways were present in the dataset, of which 257
351 taxa and 277 pathways had a relative abundance higher than 0.01 and were present in > 5%
352 of individuals. Our sample size allowed us to capture ~100% of the estimated microbial
353 functional potential (including biochemical pathways, virulence factors and antibiotic
354 resistance genes) and microbial taxa at level of genus or higher. By subsampling the cohort,
355 we estimated that the presence rates of these microbial features become stable when 40%
356 or more of the cohort is sampled (~3,300 samples). However, we also observed that the
357 number of microbial species discovered continued to increase with increasing sample size,
358 with the total number of species in the population estimated to be 612 (standard deviation
359 = 20) at 25,000 samples, suggesting that the sampled population contains rare undiscovered
360 microbial species (Fig. 1b, Supplementary fig. 1). Gut microbiota composition was found to
361 be highly variable across the population, with the relative abundance of the phylum
362 Bacteroidetes, for example, ranging from 5% to > 95% (Supplementary fig. 2c). However, the
363 most abundant microbial pathways were largely stable across the population
364 (Supplementary fig. 2d).



365

366 **Figure 1: Summary of the Dutch Microbiome Project**

367 **a**, Graphical summary of the cohort and overview of available metadata (n = number of variables
 368 collected, N = sample size). **b**, Number of microbial features discovered in relation to sample size.
 369 Error bars denote the standard deviation of 100 resamplings. **c**, Biplot of principal coordinate analysis
 370 visualising the beta-diversity of the microbiome data. Colour indicates relative abundance of
 371 *Prevotella copri*. Arrows indicate the influence of self-reported health, anthropometrics and faecal
 372 sample metadata. VFs: bacterial virulence factors, CARD: antibiotic resistance gene families.

373

374 **Core microbes and pathobionts are keystone species in the gut ecosystem**

375 To pinpoint bacterial species and pathways potentially critical for the organisation and
376 maintenance of the gut ecosystem, we investigated 8,208 samples in our cohort for bacteria
377 present in > 95% of individuals (*core microbes*) and for bacteria that form the central nodes
378 in bacterial co-abundance networks (*keystone features*)^{38,40}. We identified 9 core species
379 (*Subdoligranulum sp.*, *Alistipes onderdonkii*, *A. putredinis*, *A. shahii*, *Bacteroides uniformis*, *B.*
380 *vulgatus*, *Eubacterium rectale*, *Faecalibacterium prausnitzii* and *Oscillibacter sp.*) that are
381 highly consistent with those found in previous studies of UK, US and European and non-
382 western populations (Supplementary Table 1a)^{6,41-44}. We also identified 28 species and 53
383 pathways as potential keystone features defined by more than 109 and 337 significant co-
384 abundances, respectively (false discovery rate (FDR) < 0.05).

385 We observed that 5 of the 9 core microbial species (*A. putredinis*, *A. shahii*, *F. prausnitzii*,
386 *Oscillibacter sp.* and *Subdoligranulum sp.*) are also keystone species, implying that these 5
387 microbes are not only highly prevalent, they play central roles in the gut microbiome
388 ecosystem in the Dutch population. For example, *F. prausnitzii*, a major butyrate producer
389 that is depleted in many chronic diseases^{45,46}, had a significant co-abundance with the
390 majority of *Bacteroidetes* and *Bifidobacterium* species (Supplementary table 1e). In addition
391 to these core species, we identified potential keystone species with low prevalence in the
392 population (prevalence ≤ 0.1), including *Ruminococcus gnavus* and multiple species from
393 genus *Clostridium*. These keystone bacteria were further observed to be positively
394 associated with multiple diseases in the current study (Supplementary table 3b), as also
395 seen in previously published studies^{47,48}.

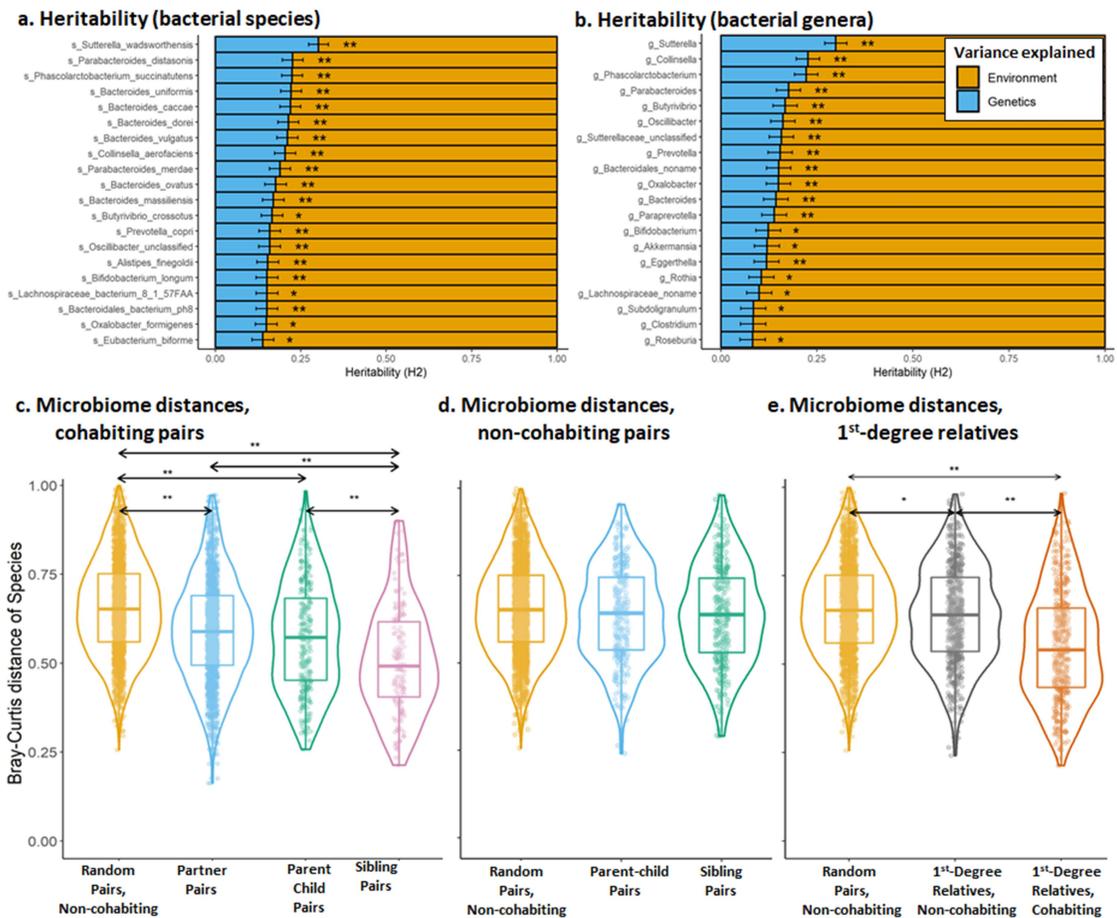
396 **Prevotella copri defines microbiome clusters**

397 To evaluate if gut microbiomes in our cohort show distinct clusters, we examined the data
398 using principal coordinate analysis (PCoA) and identified that the first principal coordinate is
399 largely driven by *Prevotella copri* (rSpearman=0.68, $P=3.6 \times 10^{-180}$, Fig. 1c, Supplementary
400 table 9). This bacterium is bimodally distributed in our cohort and defines two robust
401 clusters based on its presence or absence (Supplementary fig. 3a,b). It was previously
402 suggested that a gut microbiota with a low abundance of *Prevotella* is associated with
403 higher incidence of irritable bowel syndrome (IBS)⁴⁹, and we also observe that the cluster
404 with a high abundance of *P. copri* has a lower risk of IBS (OR = 0.72, 95% CI 0.86-0.61) and is
405 positively associated with general health (OR = 1.24, 95% CI 1.40-1.11, FDR < 0.05,
406 Supplementary fig. 3c). While previous studies have reported distinct enterotypes
407 dominated by *Bacteroides*, *Prevotella* and *Ruminococcaceae*^{39,50}, we only observed two such
408 clusters in our cohort, possibly because our cohort is ethnically uniform and comes from a
409 constrained geographic area.

410 Unlike microbiome composition, the PCoA of functional potential was not dominated by any
411 single pathway, and the top features explaining variance were found to be queuosine
412 biosynthesis, peptidoglycan biosynthesis and L-isoleucine biosynthesis pathways
413 (Supplementary table 9).

414 ***Microbiome is largely determined by cohabitation, but core species are heritable***

415 We next explored the relative contributions of family structure, co-housing and other
416 exposome factors in shaping the gut microbiome. We utilised the multi-generational family
417 structure of our cohort to estimate the heritability of microbial taxa and identified 31
418 heritable taxa (12.8% of the tested taxa) at FDR < 0.1 (Fig. 2a,b, Supplementary table 5). The
419 highest heritable values were observed for *Sutterella wadsworthensis* ($H^2 = 0.3$) and
420 multiple species from genera *Bacteroides*, *Collinsella* and *Phascolarctobacterium* ($H^2 \sim 0.21$).
421 While comprising a relatively minor part of the microbiome, the heritable taxa were found
422 to be enriched within the set of core microbes (3/9 core taxa were heritable, Chi-squared
423 test p-value < 1.0e-6) and to have significantly higher prevalence and abundance in the
424 population compared to non-heritable taxa (Mann-Whitney U test p-values < 1.0e-6). In
425 addition, 11/31 heritable taxa were positively associated with health, whereas only 3/31
426 heritable taxa were negatively associated with health (Supplementary table 5).



427

428 **Figure 2: Heritability and impact of cohabitation on the gut microbiome**

429 **a**, Top heritable species. **b**, Top heritable genera. ** taxa significantly heritable at FDR < 0.1.
 430 * taxa with nominally significant heritability (p-value < 0.05). Panels **c**, **d** and **e** show pairwise
 431 microbiome distance comparisons. Bray-Curtis dissimilarities calculated using microbial
 432 species of groups of **c**) random, non-cohabiting pairs compared to cohabiting partners,
 433 parent-child pairs and sibling pairs, **d**) random pairs compared to non-cohabiting parent-
 434 child and sibling pairs and **e**) random pairs compared to non-cohabiting first-degree relatives
 435 and cohabiting first-degree relatives. Significantly different groups: ** for FDR < 1.0e-5 and *
 436 for FDR < 0.05.

437 We further investigated the effect of cohabitation versus heritability by comparing the
 438 microbiome composition of first-degree relatives (parent-child pairs and siblings) living
 439 together to those of related pairs living separately and unrelated cohabiting participants
 440 (cohabiting partners) and to those of people living separately (randomly selected unrelated
 441 participants). The microbiomes of cohabiting pairs were more similar than the microbiomes
 442 of participants living separately, regardless of the relatedness of these pairs, with parent-
 443 child pairs, sibling pairs and unrelated partners all having microbiomes significantly closer to
 444 each other than randomly sampled non-cohabiting pairs (Fig. 2c, p-values < 1.0e-5).

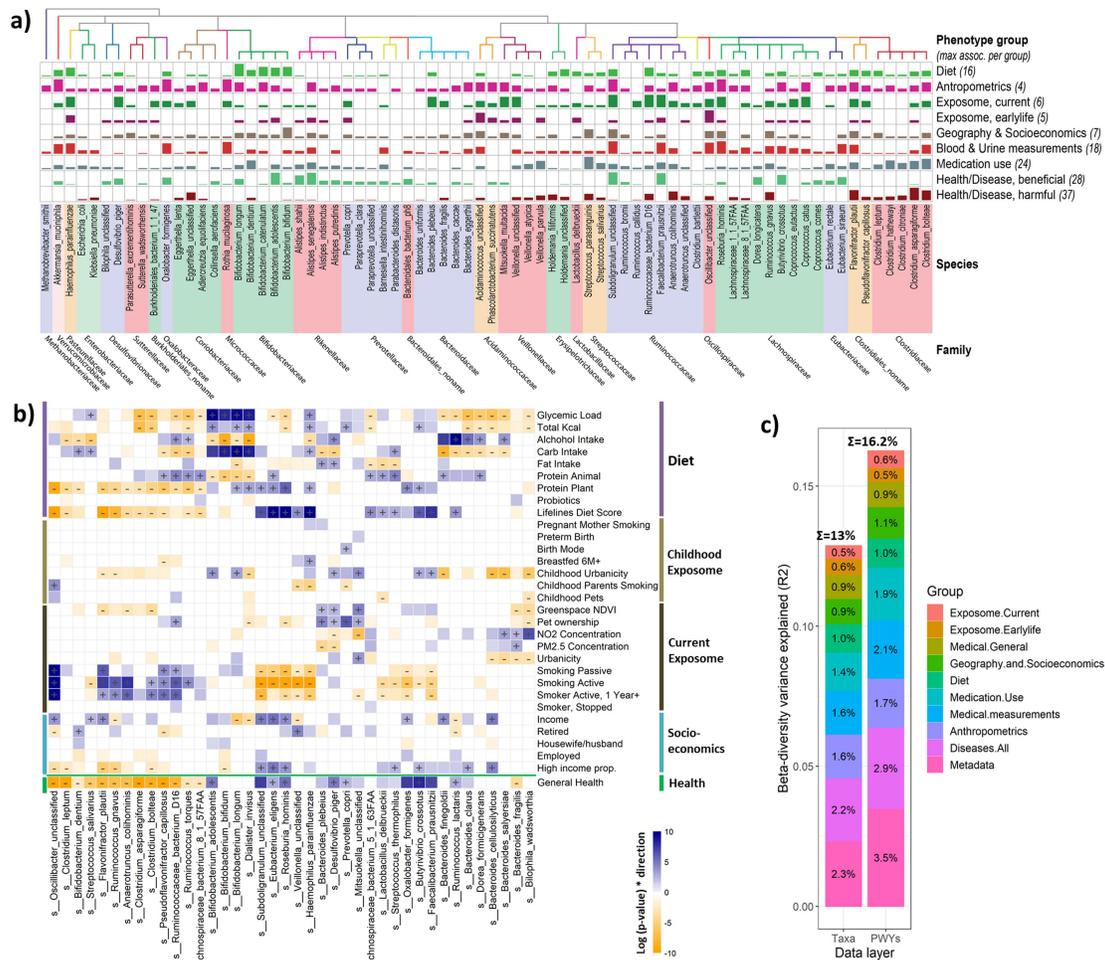
445 Microbiomes of cohabiting parent–child and sibling pairs were significantly more similar
446 than equivalent pairs living separately, while related pairs living separately were not
447 significantly different from random unrelated pairs (Fig. 2d,e). We further observed similar
448 patterns in the compositions of microbial pathways, virulence factors and antibiotic
449 resistance genes (Supplementary fig. 4a-c).

450 These results indicate that whole-microbiome composition is significantly influenced by
451 cohabitation, with genetics playing a less important role. However, a minor proportion of
452 the gut microbiome, e.g. species from genera *Sutterella* and *Collinsella*, is significantly
453 heritable.

454 **Overview of microbiome–phenotype associations**

455 We then explored the associations of microbial features to 241 individual measures
456 including technical factors, anthropometrics, early-life and current exposome, diet, self-
457 reported diseases and medication use, medical measurements and socioeconomics (Fig. 3a).
458 These phenotypes explained 12.9% of microbiome taxonomic composition and 16.3% of
459 microbiome function, with the largest contribution coming from technical factors, stool
460 characteristics, diseases, medication use and anthropometrics (Fig. 3c, Supplementary
461 tables 4a-c).

462 After correcting for technical factors, we observed 4,530 associations of phenotypes with
463 taxa, 5,224 with pathways, 1,848 with antibiotic resistance genes and 385 with virulence
464 factors (Supplementary table 3g, Supplementary fig. 5). Individually, the largest number of
465 associations were observed for keystone and core taxa, including *Flavonifractor plautii*, *F.*
466 *prausnitzii*, *Alistipes senegalensis* and species of genera *Clostridium* and *Subdogranulum*
467 (Supplementary tables 3b,h). Below, we describe a selection of association results.
468 Supplementary tables 3a-h provide a complete overview.



469

470 **Figure 3: Microbiome–phenotype associations**

471 **a**, Selected study-wide-significant associations (FDR < 0.05) per phenotype group, clustered
 472 by taxonomy. Bar height indicates the number of associations relative to the maximal
 473 number of associations for the phenotype group. **b**, Microbiome–phenotype associations for
 474 diet, childhood and current exposure and socioeconomics in comparison to healthy
 475 microbiome signature. Microbial species are clustered by association p-value using
 476 hierarchical clustering and coloured by direction of association. Study-wide significant
 477 associations (FDR < 0.05) are marked with + or -. Coloured associations without a mark
 478 indicate nominally significant associations (p-value < 0.05). **c**, The variance in microbiome
 479 composition and function explained by phenotype groups.

480 Age, sex and BMI ranked among the top phenotypes in our analysis of interindividual
 481 variation of beta-diversity, explaining 0.6%, 0.53% and 0.32%, respectively, as did individual
 482 microbiome feature associations (Supplementary tables 3a-e). Bristol stool scale explained
 483 the largest proportion of beta diversity for any single factor in our cohort (R² = 0.77%, FDR =
 484 0.012), and the sampling season also explained significant proportion of variance (R² =
 485 0.36%, FDR = 0.012, Supplementary table 4a), both highlighting the importance of assessing

486 faecal sample consistency and collection time-frame effects in microbiome studies. Based
487 on these results, we included corrections for age, sex, BMI, Bristol stool scale and sampling
488 season into all our association models alongside the corrections for technical factors
489 described above (Supplementary tables 3a-e).

490 ***Definition of healthy and disease-associated gut microbiome***

491 To define the microbiome signatures of health and disease, we associated microbiome
492 features to self-reported health and 81 diseases that had at least 20 cases in the cohort,
493 including GI and hepatological, cardiovascular and metabolic, mental and neurological,
494 pulmonary, and other disorders. We identified 1,206 significant health/disease associations
495 with bacterial taxa, 1,182 with microbial pathways, 390 with antibiotic resistance gene
496 families and 76 with bacterial virulence factors (FDR < 0.05, Supplementary tables 3b-f).
497 Different diseases showed varying numbers of associated microbes, with the strongest
498 signatures observed for cardiovascular and metabolic disorders, such as non-alcoholic fatty
499 liver disease and Type 2 Diabetes (T2D), and for GI disorders including Inflammatory Bowel
500 Disease (IBD) and IBS (Supplementary table 3f). We observed consistent microbiome-
501 disease patterns across the majority of diseases (Fig. 4), allowing us to pinpoint microbiome
502 signatures shared between different/independent diseases as well as features that define a
503 healthy (i.e. *absence of disease*) microbiome.

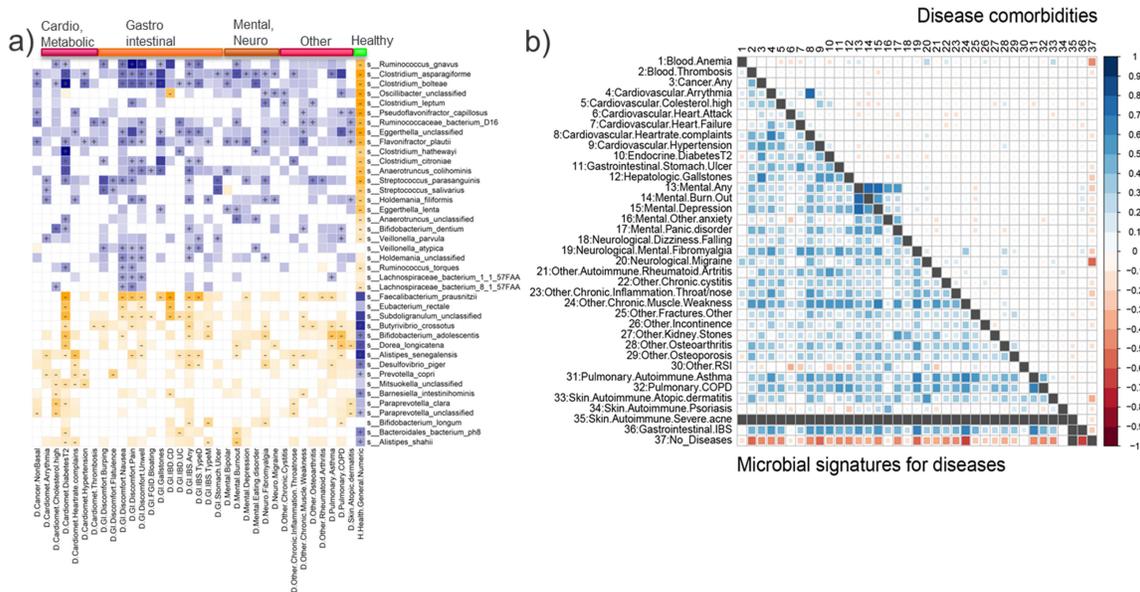
504 The shared microbiome signatures of disease (described in further detail in the
505 **Supplementary Discussion**) mainly consisted of increases in species from *Anaerotruncus*,
506 *Ruminococcus*, *Bacteroides*, *Holdemania*, *Flavonifractor*, *Eggerthella* and *Clostridium* genera
507 and decreases in *Faecalibacterium*, *Bifidobacterium*, *Butyrivibrio*, *Subdoligranulum*,
508 *Oxalobacter*, *Eubacterium* and *Roseburia*. Gut microbiome pathways shared across diseases
509 with different aetiologies mainly consisted of increases in biosynthesis of L-ornithine,
510 ubiquinol and menaquinol, enterobacterial common antigen, Kdo-2-lipid-A and
511 molybdenum cofactor and decreases in biosynthesis of amino acids, deoxyribonucleosides
512 and nucleotides, anaerobic energy metabolism and fermentation to short chain fatty acids
513 (mainly butanoate). Furthermore, virulence factors were increased in some diseases,
514 including T2D and GI disorders, with the largest effect observed for bacterial adherence and
515 iron-uptake factors (VF036, VF0228, VF0236, VF0404 and VF0394).

516 To further validate the commonality of the microbiome signatures for diseases, we
517 constructed L1/L2 regularized regression prediction models for the 36 most common
518 diseases with > 100 cases each in the cohort using a randomly selected 90% of our data as a
519 training set and 10% as a test set.

520 The microbiome-driven predictive signature of health (*lack of any reported disease*) showed
521 a prediction area under the curve (AUC) of 0.62 on the training set and of 0.57 on the test

522 set, which was comparable with the AUC for prediction of general anthropomorphic
523 parameters (age, sex and BMI; AUC of 0.58 on both training/test sets). Combining
524 microbiome and anthropometric data resulted in AUCs of 0.64 and 0.61 (training/test set).
525 L1/L2 regression selected 73 microbial features for the “health/disease” model fit, of which
526 the majority were microbial taxa rather than pathways (Supplementary table 6a). These
527 features were found to be largely consistent with the healthy microbiome signature we had
528 identified in the association analysis (Fig. 4), with 22/31 species selected by the model also
529 associated with > 5 diseases in our analyses. The microbiome signatures for 29 out of 36
530 tested diseases showed AUCs > 0.55 on the test set, with the highest predictive power
531 observed for T2D (AUC = 0.81), gallstones (AUC = 0.66) and kidney stones (AUC = 0.64). No
532 predictive power (AUC ≤ 0.5) was observed for severe acne, undefined anaemia, non-hip
533 fractures and reports of high cholesterol (Supplementary table 6b). Additionally,
534 microbiome predictions of diseases showed high correlation despite the low disease co-
535 morbidities in our cohort (Fig. 4b). Finally, we calculated the recently developed Gut
536 Microbiome Health Index (GMHI)⁵¹ for our data and identified a significant difference in
537 GMHI between healthy and unhealthy individuals (p-value = 1.1e-11, balanced accuracy
538 0.61, Supplementary fig. 6). Our microbiome signature of health showed high overlap with
539 signatures reported in the GMHI study⁵¹, with 43/50 GMHI signals replicating across studies
540 at genus or species level (Supplementary table 1). In addition to the observed consistency in
541 families *Clostridiaceae*, *Ruminococcaceae* and *Lachnospiraceae* and genera *Bacteroides*,
542 *Bilophila*, *Coprococcus*, *Faecalibacterium*, *Ruminococcus* and *Sutterella*, we also discovered
543 55 novel microbiome–health associations not identified in the GMHI study, including for
544 species from *Butyrivibrio* and *Akkermansia* and *Prevotella* genera (Supplementary tables
545 3b,10).

546



547

548 **Figure 4: Microbiome signatures of health and diseases.**

549 **a**, Heatmap of microbial species associated with categories of diseases and health status.
 550 Diseases are sorted and labelled by disease type. Microbial species are clustered by
 551 association p-value (indicated by colour intensity) using hierarchical clustering. Associations
 552 are coloured by direction of effect (blue = positive, orange = negative), with associations
 553 significant at study-wise FDR < 0.05 marked with +/- for positive/negative correlations.
 554 Coloured associations without a label indicate nominally significant associations (p-value <
 555 0.05, no multiple testing correction). **b**, Comparison of correlations between signatures
 556 predictive for diseases (lower triangle) and comorbidities of these diseases in the cohort
 557 (upper triangle).

558 **The disease-associated microbiome is a combination of disease and medication signatures**

559 We observed a large overlap between disease-associated microbiome patterns and
 560 microbiome associations with medication use (Supplementary fig. 7, Supplementary tables
 561 3a-g), with the largest effects observed for proton pump inhibitors (PPIs), antibiotics,
 562 biguanide antidiabetics, osmotic laxatives and intestinal anti-inflammatory agents (84, 56,
 563 47 and 32 associations with microbial taxa, respectively, at FDR < 0.05). To disentangle
 564 medication use from the diseases for which these drugs are prescribed, we constructed
 565 microbiome-drug-disease multivariate models for diseases where drug use is strongly
 566 correlated and/or specific for the disease: antidiabetics in T2D, selective serotonin reuptake
 567 inhibitors in depression and PPIs in functional GI disorders/IBS. We observed that both drug
 568 use and presence of the disease are independently and significantly associated with
 569 microbiome features in these models (Supplementary table 7), with consistent effect
 570 directions and strength. This observation may indicate that the unhealthy gut microbiome
 571 signature reflects both the disease and the medication used to treat it.

572 ***Childhood environment is associated to healthy microbiome in adult life***

573 Since it is known that the first 2–3 years of life are crucial for microbiome development, we
574 examined the influence of early-life (age < 4 years) factors on the adult microbiome. We
575 identified 106 associations with taxa, 30 with pathways, 22 with antibiotic resistance genes
576 and 2 with virulence factors (FDR < 0.05), with only minimal effects observed for birth mode,
577 breastfeeding and preterm birth (Fig. 3b, Supplementary fig. 8). Childhood living
578 environment (defined on a scale from 1 = rural to 5 = highly urban) was significantly
579 associated with the adult microbiome (54, 8 and 7 associations with taxa, pathways and
580 CARDs, respectively, at FDR < 0.05). A rural childhood environment was associated with an
581 increase in various species of bacteria, including *P. copri*, *F. prausnitzii*, *Rothia mucilaginosa*
582 and various species of *Bifidobacterium* and *Mitsuokella* genera, abundances of which were
583 also positively associated with increased general health. In contrast, gut abundances of
584 species from genera *Bacteroides*, *Alistipes* and *Bilophila* were reduced in participants with
585 urban childhood environments.

586 Furthermore, parental smoking was associated with the microbiome composition of their
587 children (15, 9 and 4 associations with taxa, pathways and CARDs, respectively, at FDR <
588 0.05). Here we observed associations between parental smoking and decreased abundances
589 of *F. prausnitzii* and *P. copri*, consistent with observations for current smokers. Finally, pet
590 ownership during childhood was associated with the adult microbiome (7 associations at
591 FDR < 0.05), with decreases in *Alistipes finegoldii*, *Lactobacillus delbrueckii* and species from
592 *Dialister* and *Bilophila* genera observed in participants who had pets as children.

593 We observed minimal effects for birth mode, breastfeeding and preterm birth (Fig. 3b,
594 Supplementary fig. 8). These included an association of C-section with increases in *P. copri* in
595 adults and of higher birth weight with increases in *Bacteroides eggerthii* and *Butyrivibrio*
596 *crossotus*.

597 ***Urbanicity, smoking, pollutants and pets are linked to health-associated bacteria***

598 Next, we studied environmental factors at time of sampling. We identified significant shared
599 positive association patterns between healthy microbiome signatures, pet ownership, rural
600 living environment and greenspace surface area in the living environment (Fig. 3b), including
601 increases in *P. copri*, *Bacteroides plebeius*, *Desulfovibrio piger* and species of genus
602 *Mitsuokella* and decreases in *Bacteroides fragilis* and *Bilophila wadsworthia*. These
603 associations were in contrast to the associations seen for increased measurements of NO₂
604 and small particulate matter pollutants, which showed microbiome signatures in the
605 opposite direction and negative association with health (Supplementary fig. 9).

606 Smoking phenotypes, including current active and passive smoking and a history of past
607 smoking, were among the strongest phenotypes associated with microbiome composition in

608 our study. Active smoking was associated with 41 species and 84 pathways (FDR < 0.05, Fig.
609 3b), with 60% of these factors also associated with previous smoking, suggesting a long-
610 lasting signature of smoking. Intriguingly, 15 of these were further associated with passive
611 smoking, highlighting the need to consider passive smoking in disease risk models. The
612 directionality of associations was consistent across the three smoking phenotypes and
613 shared direction with microbiome signatures of diseases, including decreased abundance of
614 the butyrate-producing bacteria *Roseburia hominis*, *F. prausnitzii*, *Coprococcus catus* and
615 *Subdoligranulum spp*; decreased abundance of the facultative oral bacteria *Veillonella spp*
616 and *Haemophilus parainfluenzae*; and increased abundance of several species within order
617 Clostridiales, including *Flavonifractor plautii*, *Pseudoflavonifractor capillosus* and
618 *Oscillibacter spp* (Supplementary fig. 9).

619 **High dietary quality is linked to health-associated microbes**

620 We identified 378 associations between 20 dietary factors, which were found to be stable
621 over a 5-year period between food frequency questionnaire (FFQ) collection and faecal
622 sampling (Supplementary table 8, Supplementary Discussion), and 82 species
623 (Supplementary table 3f). The Lifelines diet score (LLDS), a diet quality score based on
624 international nutrition literature¹⁴, showed the highest number of associations (79
625 associations with taxa, 44 with pathways, 20 with antibiotic resistance genes and 8 with
626 virulence factors at FDR < 0.05), followed by total alcohol intake, glycemic load, protein diet
627 score (reflecting quantity and source of protein) and total carbohydrate intake. The LLDS
628 and protein intake scores, including animal protein, showed association patterns that
629 overlapped with signals between microbiome features and increased general health, such as
630 decreases in *Clostridia* species, increases in *Butyrivibrio* and *Roseburia* genera and pathways
631 involved in ubiquinol and menaquinol synthesis. In contrast, total dietary carbohydrate
632 intake and glycemic load showed the opposite associations (Fig. 3b, Supplementary fig. 10).

633 **Socio-economic factors are associated with microbiome composition**

634 Several measurements related to the socio-economic status of the participants were
635 available for our cohort. This included information about the number of working hours,
636 income per month, religion and neighbourhood characteristics. In total, we observed 220
637 significant associations (FDR < 0.05). Of these, 72 were between bacterial abundances and
638 monthly salary, with higher income associated with a healthy microbiome signature, such as
639 an increase in *F. prausnitzii*, *Akkermansia municipihila* and *Bifidobacterium adolescentis*. In
640 contrast, the “unhealthy microbiota” *Clostridium bolteae*, *Ruminococcus gnavus* and
641 *Streptococcus parasanguinis* were more abundant in participants with lower incomes.
642 Similar patterns were observed when analysing the differences between people living in
643 high-income neighbourhoods versus low-income neighbourhoods. In our cohort, a high

644 income showed low but significant correlations with neighbourhood greenspace area, rural
645 living environment and higher LLDS (Spearman correlations 0.22, 0.17 and 0.07,
646 respectively; correlation test FDRs < 1.0e-6), all of which share microbiome patterns
647 observed in high income participants, implying that microbiome–income association is likely
648 a combination of multiple factors, including a healthier diet and lifestyle and a less urban
649 living environment.

650

651 Discussion

652 Microbiome studies have associated the gut microbiome to diverse chronic and acute
653 diseases¹⁻³, to the functioning of the immune system⁵² and to drug response⁵³, highlighting
654 that the gut microbiome is an essential factor in maintenance of human health. While
655 multiple microbiome-targeting therapies are currently being investigated in clinical trials⁵⁴,
656 the therapeutic potential of microbiome-targeting interventions is still confounded by a lack
657 of consensus in the definition of a healthy or unhealthy microbiome and by a limited
658 understanding of how heritability, exposome, lifestyle and diet shape microbiome
659 signatures of health and disease. To bridge this gap, we analysed the microbiomes of 8,208
660 individuals from an extensively phenotyped three-generation cohort, allowing us to pinpoint
661 core and keystone microbiome features and assess the effect of technical variables,
662 anthropometrics, heritability, cohousing, diseases, diet and environmental exposures on the
663 microbiome.

664 In addition to confirming known microbiome associations with age, sex and BMI^{43,55}, our
665 results highlight the importance of considering often-omitted confounders related to the
666 stool samples (stool consistency), sampling season and sample processing (such as DNA
667 concentration or sequencing batch). This is especially important when studying diseases
668 where age, sex, BMI and faecal consistency are often associated with the disease.

669 Our observation that the microbiome is primarily associated to cohabitation and
670 environment rather than relatedness corroborates previous studies that identified limited
671 overall microbiome heritability⁵⁶ and divergence in microbiomes of twins who stopped
672 cohousing⁵⁷. However, we do identify that the core component of the gut microbiome is
673 significantly heritable, whereas it is the disease-associated microbes that are largely
674 influenced by environment. These results suggest that the bacteria with low heritability that
675 are enriched in diseases, e.g. species from genera *Clostridia*, *Flavonifracter* and *Veillonella*,
676 might be more susceptible to microbiome-altering therapies than more heritable bacteria
677 such as those from genera *Sutterella*, *Collinsella* and *Bacteroides*.

678 By comparing associations between microbiome, health and diverse diseases, we identified
679 a common signal for dysbiosis in the gut (Fig. 4) that was largely consistent with a previous
680 study.⁵¹ The existence of shared dysbiosis has considerable implications for microbiome
681 research and development of microbiota-targeting diagnostics and therapies. It implies that
682 the gut microbiome is a biomarker of general health, which is supported by our prediction
683 models and by previous studies^{51,58}. However, this also complicates microbiome-based
684 diagnosis of individual diseases because single-disease models might be confounded by
685 signals shared across unrelated diseases, and testing such models for specificity in mixed-
686 disease cohorts will be an important step before clinical implementation. The shared

687 microbiome signature is also exciting because it suggests that microbiome-targeting
688 interventions could improve overall human health. This is supported by our observations
689 that lifestyle factors generally considered healthy, such as adherence to current dietary
690 recommendations and no smoking, associate with similar microbiome patterns to those
691 associated with general health. While microbiome–drug interactions are well described *in*
692 *vitro*⁵⁹, and characterized *in vivo* for antibiotics, PPIs⁶⁰ and antidiabetics⁶¹, our results
693 suggest that the general microbiome dysbiosis is a combination of drug- and disease-effects,
694 implying that many currently understudied drugs, such as SSRIs, have a negative effect on
695 the gut microbiome. This also highlights the importance of controlling for medication use.

696 By linking (un)healthy microbiome patterns to childhood and current exposome, diet and
697 socioeconomics, we observed that a healthier diet¹⁴, childhood and current exposure to
698 rural environment and pets, exposure to greenspace and higher income share signals with
699 healthy microbiome patterns. These observations support the microbiome diversity
700 hypothesis (also known as the hygiene hypothesis) – a postulate that reduction in exposure
701 to microbiota contributes to an increase in the frequency of autoimmune and allergic
702 diseases^{62,63}. Notably, while the classic hygiene hypothesis focuses on pathogens and the
703 impact of early-life exposures, our results suggest that exposures in adulthood also
704 contribute to (un)healthy microbiome patterns, implying that the environment shapes the
705 microbiome throughout human life and, as such, microbiome-targeted therapies could be
706 effective throughout an individual’s life. Furthermore, while we identified negative
707 correlations of diet scores, pets and rural environment with opportunistic pathogens, such
708 as *Clostridia* species, we also observed positive correlation with commensals such as
709 butyrate producers from genera *Bacteroides*, *Alistipes* and *Faecalibacterium*, implying that
710 exposure not only to pathogens but also to commensals from the environment plays an
711 important role in establishing a healthy gut ecosystem.

712 We also observed that smoking, a high-carbohydrate diet and exposure to NO₂ and small
713 particulate matter (PM_{2.5}) are positively correlated with disease-linked bacterial species
714 from genera *Clostridia* and *Ruminococcus*. While air pollutants were previously associated to
715 GI diseases in humans⁶⁴, and have been shown to affect the gut microbiota of mice, the
716 effects of air pollutants on the human gut microbiota are still unclear⁶⁵. Our results suggest
717 that air pollutants negatively impact the human gut microbiota and might increase risk for
718 development of GI diseases by contributing to general microbiome dysbiosis.

719 In addition to identifying links between current exposome and microbiome, we also found
720 that childhood exposures to smoking, pets and rural environment are associated with the
721 subsequent adult microbiome. While the effect sizes for these associations were lower than
722 those of equivalent current exposures measured alongside faecal sampling, the effect
723 directions and patterns were consistent, suggesting that environmental exposures can have

724 a long-lasting effect on the gut microbiome. This is further supported by our observation
725 that smokers who stopped smoking still showed microbiome associations similar to current
726 smokers, albeit with lower effect sizes. Notably, we observed limited associations between
727 microbiome and exposures during and shortly after birth, including breastfeeding, birth
728 mode and preterm birth. These exposures are known to shape infant microbiota^{66,67}, but
729 our findings indicate that their effects on the adult microbiome are indiscernible, possibly
730 because the selective pressures induced by these exposures end early in human life and
731 these early effects are superseded over time by other selective pressures such as diet,
732 exposure to microbiota of family members and environment.

733 Notably, while we measured 241 phenotypes from a broad set of categories, we could only
734 explain ~15% of the variation in microbiome composition and function between individuals.
735 While this level of explained variance is consistent with those found in previous large-scale
736 studies of European and American populations^{6,56,66}, it implies that the gut microbiome is
737 highly individual and that our current understanding of the factors that shape it is still
738 limited. A potential explanation for the existence of this “missing variance” is that the
739 microbiome composition and function are a result of an individual's history of lifestyle and
740 exposures, and cross-sectional measurement is thus insufficient to fully explain it. This is
741 supported by our observation that early-life exposures are associated with microbiome in
742 adult age and that cohousing participants have significantly closer microbiomes than non-
743 cohabiting individuals, regardless of relatedness. Future quantification of this missing
744 variance, potentially by long-time-frame longitudinal studies, will play a critical part in future
745 development of microbiome-targeting diagnostics and therapies.

746 **Conclusion**

747 We generated and analysed the largest, multi-generational gut microbiome cohort to date
748 that has been collected and profiled in a highly standardized manner, and linked it to
749 extensive phenotype data. We defined and described a gut dysbiosis shared across diverse
750 diseases and identified novel links between this dysbiosis and heritability, childhood and
751 current exposome, lifestyle and socioeconomics. This study demonstrates the power of
752 large-scale, well-phenotyped cohorts for dissecting the links between gut microbiome,
753 health, genetics and environment and provides a rich resource for future studies for
754 microbiome-directed interventions.

755

756

757 **Acknowledgements**

758 We would like to acknowledge and thank the late Marten Hofker who had the great wisdom
759 and vision to initiate the Lifelines DAG3/Dutch Microbiome Project.

760 The authors wish to acknowledge the services of the Lifelines Cohort Study, the contributing
761 research centres delivering data to Lifelines and all the study participants. The Lifelines
762 Biobank initiative has been made possible by subsidy from the Dutch Ministry of Health,
763 Welfare and Sport; the Dutch Ministry of Economic Affairs; the University Medical Center
764 Groningen (UMCG the Netherlands); the University of Groningen and the Northern
765 Provinces of the Netherlands. We would like to thank the Center for Information Technology
766 of the University of Groningen (RUG) for their support and for providing access to the
767 Peregrine high performance computing cluster and the Genomic Coordination Center
768 (UMCG and RUG) for their support and for providing access to Calculon and Boxy high
769 performance computing clusters, and the MOLGENIS team for data management and
770 analysis support. Metagenomics library preparation and sequencing was done at Novogene.
771 We also thank K. Mc Intyre for English and content editing.

772 **Keywords**

773 Gut microbiome, healthy microbiome, microbiome heritability, exposome, population
774 health

775 **Author contribution**

776 RG designed and implemented the metagenomic data analysis pipelines, analysed
777 metagenomic data, performed heritability analysis and drafted the manuscript. AK designed
778 the prediction models and implemented statistical methods for association analyses and
779 assisted in drafting of the manuscript. AVV, LC, VC, SH, MAYK, SAS, JB, LAB, VL, TS, MH and
780 SS assisted in other statistical analyses, interpretation of data and drafting of the
781 manuscript. MS provided data stewardship and analysis infrastructure. BHJ, JAMD and JGA
782 collected data, assisted in study planning and critically reviewed the manuscript. SS
783 supervised and coordinated heritability analysis. RCHV provided the air pollution data and
784 supervised the air pollution analysis. HJMH, SZ, RKW, JF and CW conceived, coordinated and
785 supported the study. All authors critically revised and approved the manuscript.

786 **Funding**

787 Sequencing of the cohort was funded by a grant from the CardioVasculair Onderzoek
788 Nederland grant (CVON 2012-03) to MH, JF and AZ. RG, HH and RKW are supported by the
789 collaborative TIMID project (LSHM18057-SGF) financed by the PPP allowance made
790 available by Top Sector Life Sciences & Health to Samenwerkende Gezondheidsfondsen
791 (SGF) to stimulate public-private partnerships and co-financing by health foundations that

792 are part of the SGF. RKW is supported by the Seerave Foundation and the Dutch Digestive
793 Foundation (16-14). AZ is supported by European Research Council (ERC) Starting Grant
794 715772, Netherlands Organization for Scientific Research (NWO) VIDI grant 016.178.056,
795 CVON grant 2018-27, and NWO Gravitation grant ExposomeNL 024.004.017. JF is supported
796 by NWO Gravitation grant Netherlands Organ-on-Chip Initiative (024.003.001) and CVON
797 grant 2018-27. CW is further supported by an ERC advanced grant (ERC-671274) and an
798 NWO Spinoza award (NWO SPI 92-266). LC is supported by a joint fellowship from the
799 University Medical Center Groningen and China Scholarship Council (CSC201708320268) and
800 a Foundation De Cock-Hadders grant (20:20-13). MS is supported by Netherlands
801 Organization for Scientific Research (NWO) VIDI grant 016 and EUCAN-connect, a project
802 funded by European Commission H2020 grant 824989.

803 **Competing interests declaration**

804 Authors declare no conflict of interest.

805 **Scripts and availability of data:**

806 Scripts used for data analysis can be found at:

807 [https://github.com/GRONINGEN-MICROBIOME-CENTRE/Groningen-](https://github.com/GRONINGEN-MICROBIOME-CENTRE/Groningen-Microbiome/tree/master/Projects/DMP)
808 [Microbiome/tree/master/Projects/DMP](https://github.com/GRONINGEN-MICROBIOME-CENTRE/Groningen-Microbiome/tree/master/Projects/DMP)

809 **References**

- 810 1. Lynch, S. V. & Pedersen, O. The Human Intestinal Microbiome in Health and Disease.
811 *N. Engl. J. Med.* **375**, 2369–2379 (2016).
- 812 2. Liang, D., Leung, R. K.-K., Guan, W. & Au, W. W. Involvement of gut microbiome in
813 human health and disease: brief overview, knowledge gaps and research opportunities. *Gut*
814 *Pathog.* **10**, 3 (2018).
- 815 3. Duvallet, C., Gibbons, S. M., Gurry, T., Irizarry, R. A. & Alm, E. J. Meta-analysis of gut
816 microbiome studies identifies disease-specific and shared responses. *Nat. Commun.* **8**, 1784
817 (2017).
- 818 4. Zmora, N., Soffer, E. & Elinav, E. Transforming medicine with the microbiome. *Sci.*
819 *Transl. Med.* **11**, (2019).
- 820 5. Zmora, N., Zeevi, D., Korem, T., Segal, E. & Elinav, E. Taking it Personally:
821 Personalized Utilization of the Human Microbiome in Health and Disease. *Cell Host Microbe*
822 **19**, 12–20 (2016).
- 823 6. Zhernakova, A. *et al.* Population-based metagenomics analysis reveals markers for
824 gut microbiome composition and diversity. *Science* **352**, 565–569 (2016).

- 825 7. Gaulke, C. A. & Sharpton, T. J. The influence of ethnicity and geography on human
826 gut microbiome composition. *Nat. Med.* **24**, 1495–1496 (2018).
- 827 8. Vatanen, T. *et al.* Genomic variation and strain-specific functional adaptation in the
828 human gut microbiome during early life. *Nat. Microbiol.* **4**, 470–479 (2019).
- 829 9. Scholtens, S. *et al.* Cohort Profile: LifeLines, a three-generation cohort study and
830 biobank. *Int. J. Epidemiol.* **44**, 1172–1180 (2015).
- 831 10. Tigchelaar, E. F. *et al.* Cohort profile: LifeLines DEEP, a prospective, general
832 population cohort study in the northern Netherlands: study design and baseline
833 characteristics. *BMJ Open* **5**, (2015).
- 834 11. Siebelink, E., Geelen, A. & de Vries, J. H. M. Self-reported energy intake by FFQ
835 compared with actual energy intake to maintain body weight in 516 adults. *Br. J. Nutr.* **106**,
836 274–281 (2011).
- 837 12. Brouwer-Brolsma, E. M. *et al.* A National Dietary Assessment Reference Database
838 (NDARD) for the Dutch Population: Rationale behind the Design. *Nutrients* **9**, (2017).
- 839 13. Willett, W. C. *et al.* Reproducibility and validity of a semiquantitative food frequency
840 questionnaire. *Am. J. Epidemiol.* **122**, 51–65 (1985).
- 841 14. Vinke, P. C. *et al.* Development of the food-based Lifelines Diet Score (LLDS) and its
842 application in 129,369 Lifelines participants. *Eur. J. Clin. Nutr.* **72**, 1111–1119 (2018).
- 843 15. G, M. *et al.* A Protein Diet Score, Including Plant and Animal Protein, Investigating
844 the Association with HbA1c and eGFR-The PREVIEW Project. *Nutrients* vol. 9
845 <https://pubmed.ncbi.nlm.nih.gov/28714926/> (2017).
- 846 16. Leeming, E. R., Johnson, A. J., Spector, T. D. & Le Roy, C. I. Effect of Diet on the Gut
847 Microbiota: Rethinking Intervention Duration. *Nutrients* **11**, (2019).
- 848 17. Eeftens, M. *et al.* Development of Land Use Regression Models for PM2.5, PM2.5
849 Absorbance, PM10 and PMcoarse in 20 European Study Areas; Results of the ESCAPE
850 Project. <https://pubs.acs.org/doi/full/10.1021/es301948k> (2012) doi:10.1021/es301948k.
- 851 18. Development of NO2 and NOx land use regression models for estimating air
852 pollution exposure in 36 study areas in Europe – The ESCAPE project. *Atmos. Environ.* **72**,
853 10–23 (2013).
- 854 19. Eeftens, M. *et al.* Stability of measured and modelled spatial contrasts in NO2 over
855 time. *Occup. Environ. Med.* **68**, 765–770 (2011).
- 856 20. StatLine. <https://opendata.cbs.nl/#/CBS/en/>.

- 857 21. Ford, A. C. *et al.* Validation of the Rome III Criteria for the Diagnosis of Irritable Bowel
858 Syndrome in Secondary Care. *Gastroenterology* **145**, 1262-1270.e1 (2013).
- 859 22. Angulo, P. *et al.* The NAFLD fibrosis score: A noninvasive system that identifies liver
860 fibrosis in patients with NAFLD. *Hepatology* **45**, 846–854 (2007).
- 861 23. Bedogni, G. *et al.* The Fatty Liver Index: a simple and accurate predictor of hepatic
862 steatosis in the general population. *BMC Gastroenterol.* **6**, 33 (2006).
- 863 24. Imhann, F. *et al.* The 1000IBD project: multi-omics data of 1000 inflammatory bowel
864 disease patients; data release 1. *BMC Gastroenterol.* **19**, 5 (2019).
- 865 25. Mclver, L. J. *et al.* bioBakery: a meta’omic analysis environment. *Bioinformatics* **34**,
866 1235–1237 (2018).
- 867 26. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat.*
868 *Methods* **9**, 357–359 (2012).
- 869 27. Truong, D. T. *et al.* MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nat.*
870 *Methods* **12**, 902–903 (2015).
- 871 28. Franzosa, E. A. *et al.* Species-level functional profiling of metagenomes and
872 metatranscriptomes. *Nat. Methods* **15**, 962–968 (2018).
- 873 29. Swertz, M.A. *et al.* The MOLGENIS toolkit: rapid prototyping of biosoftware at the
874 push of a button. *BMC Bioinformatics* **11** S12 (2010).
- 875 30. Hsieh, T. C., Ma, K. H. & Chao, A. iNEXT: an R package for rarefaction and
876 extrapolation of species diversity (Hill numbers). *Methods Ecol. Evol.* **7**, 1451–1456 (2016).
- 877 31. Kaminski, J. *et al.* High-Specificity Targeted Functional Profiling in Microbial
878 Communities with ShortBRED. *PLoS Comput. Biol.* **11**, (2015).
- 879 32. Chen, L. *et al.* VFDB: a reference database for bacterial virulence factors. *Nucleic*
880 *Acids Res.* **33**, D325-328 (2005).
- 881 33. Jia, B. *et al.* CARD 2017: expansion and model-centric curation of the comprehensive
882 antibiotic resistance database. *Nucleic Acids Res.* **45**, D566–D573 (2017).
- 883 34. Pilia, G. *et al.* Heritability of Cardiovascular and Personality Traits in 6,148 Sardinians.
884 *PLoS Genet.* **2**, (2006).
- 885 35. Chen, W.-M. & Abecasis, G. R. Estimating the power of variance component linkage
886 analysis in large pedigrees. *Genet. Epidemiol.* **30**, 471–484 (2006).
- 887 36. Pincus, R. Aitchison, J.: The Statistical Analysis of Compositional Data. Chapman and
888 Hall, London & New York 1986, XII, 416 pp. *Biom. J.* **30**, 794–794 (1988).

- 889 37. Friedman, J. & Alm, E. J. Inferring Correlation Networks from Genomic Survey Data.
890 *PLOS Comput. Biol.* **8**, e1002687 (2012).
- 891 38. Chen, L. *et al.* Gut microbial co-abundance networks show specificity in inflammatory
892 bowel disease and obesity. *Nat. Commun.* **11**, 1–12 (2020).
- 893 39. M, A. *et al.* Enterotypes of the human gut microbiome. *Nature* vol. 473
894 <https://pubmed.ncbi.nlm.nih.gov/21508958/> (2011).
- 895 40. Berry, D. & Widder, S. Deciphering microbial interactions and detecting keystone
896 species with co-occurrence networks. *Front. Microbiol.* **5**, 219 (2014).
- 897 41. J, Q. *et al.* A human gut microbial gene catalogue established by metagenomic
898 sequencing. *Nature* vol. 464 <https://pubmed.ncbi.nlm.nih.gov/20203603/> (2010).
- 899 42. Turnbaugh, P. J. *et al.* The Human Microbiome Project. *Nature* **449**, 804–810 (2007).
- 900 43. T, Y. *et al.* Human gut microbiome viewed across age and geography. *Nature* vol. 486
901 <https://pubmed.ncbi.nlm.nih.gov/22699611/> (2012).
- 902 44. Jk, G. *et al.* Human genetics shape the gut microbiome. *Cell* vol. 159
903 <https://pubmed.ncbi.nlm.nih.gov/25417156/> (2014).
- 904 45. *Faecalibacterium prausnitzii* and human intestinal health. *Curr. Opin. Microbiol.* **16**,
905 255–261 (2013).
- 906 46. Eppinga, H. *et al.* Similar Depletion of Protective *Faecalibacterium prausnitzii* in
907 Psoriasis and Inflammatory Bowel Disease, but not in Hidradenitis Suppurativa. *J. Crohns*
908 *Colitis* **10**, 1067–1075 (2016).
- 909 47. R, H. *et al.* Metagenome-wide association study of the alterations in the intestinal
910 microbiome composition of ankylosing spondylitis patients and the effect of traditional and
911 herbal treatment. *Journal of medical microbiology* vol. 69
912 <https://pubmed.ncbi.nlm.nih.gov/31778109/> (2020).
- 913 48. Kandeel, W. A. *et al.* Impact of *Clostridium* Bacteria in Children with Autism Spectrum
914 Disorder and Their Anthropometric Measurements. *J. Mol. Neurosci.* **70**, 897–907 (2020).
- 915 49. Tap, J. *et al.* Identification of an Intestinal Microbiota Signature Associated With
916 Severity of Irritable Bowel Syndrome. *Gastroenterology* **152**, 111-123.e8 (2017).
- 917 50. Vieira-Silva, S. *et al.* Statin therapy is associated with lower prevalence of gut
918 microbiota dysbiosis. *Nature* **581**, 310–315 (2020).
- 919 51. Gupta, V. K. *et al.* A predictive index for health status using species-level gut
920 microbiome profiling. *Nat. Commun.* **11**, 1–16 (2020).

- 921 52. Thaiss, C. A., Zmora, N., Levy, M. & Elinav, E. The microbiome and innate immunity.
922 *Nature* **535**, 65–74 (2016).
- 923 53. Precision Medicine Goes Microscopic: Engineering the Microbiome to Improve Drug
924 Outcomes. *Cell Host Microbe* **26**, 22–34 (2019).
- 925 54. Garber, K. First microbiome-based drug clears phase III, in clinical trial turnaround.
926 *Nature Reviews Drug Discovery* vol. 19 655–656 [https://www.nature.com/articles/d41573-](https://www.nature.com/articles/d41573-020-00163-4)
927 [020-00163-4](https://www.nature.com/articles/d41573-020-00163-4) (2020).
- 928 55. Castaner, O. *et al.* The Gut Microbiome Profile in Obesity: A Systematic Review.
929 *International Journal of Endocrinology* vol. 2018 e4095789
930 <https://www.hindawi.com/journals/ije/2018/4095789/> (2018).
- 931 56. Rothschild, D. *et al.* Environment dominates over host genetics in shaping human gut
932 microbiota. *Nature* **555**, 210–215 (2018).
- 933 57. Shotgun Metagenomics of 250 Adult Twins Reveals Genetic and Environmental
934 Impacts on the Gut Microbiome. *Cell Syst.* **3**, 572-584.e3 (2016).
- 935 58. Oh, M. & Zhang, L. DeepMicro: deep representation learning for disease prediction
936 based on microbiome data. *Sci. Rep.* **10**, 1–9 (2020).
- 937 59. M, Z., M, Z.-K., R, W. & Al, G. Mapping human microbiome drug metabolism by gut
938 bacteria and their genes. *Nature* vol. 570 <https://pubmed.ncbi.nlm.nih.gov/31158845/>
939 (2019).
- 940 60. Imhann, F. *et al.* Proton pump inhibitors affect the gut microbiome. *Gut* **65**, 740–748
941 (2016).
- 942 61. Forslund, K. *et al.* Disentangling type 2 diabetes and metformin treatment signatures
943 in the human gut microbiota. *Nature* **528**, 262–266 (2015).
- 944 62. Bach, J.-F. The hygiene hypothesis in autoimmunity: the role of pathogens and
945 commensals. *Nat. Rev. Immunol.* **18**, 105–120 (2018).
- 946 63. Scudellari, M. News Feature: Cleaning up the hygiene hypothesis. *Proc. Natl. Acad.*
947 *Sci.* **114**, 1433–1436 (2017).
- 948 64. Salim, S. Y., Kaplan, G. G. & Madsen, K. L. Air pollution effects on the gut microbiota.
949 *Gut Microbes* (2013) doi:10.4161/gmic.27251.
- 950 65. Impact of air quality on the gastrointestinal microbiome: A review. *Environ. Res.* **186**,
951 109485 (2020).
- 952 66. Manor, O. *et al.* Health and disease markers correlate with gut microbiome
953 composition across thousands of people. *Nat. Commun.* **11**, 1–12 (2020).

Figures

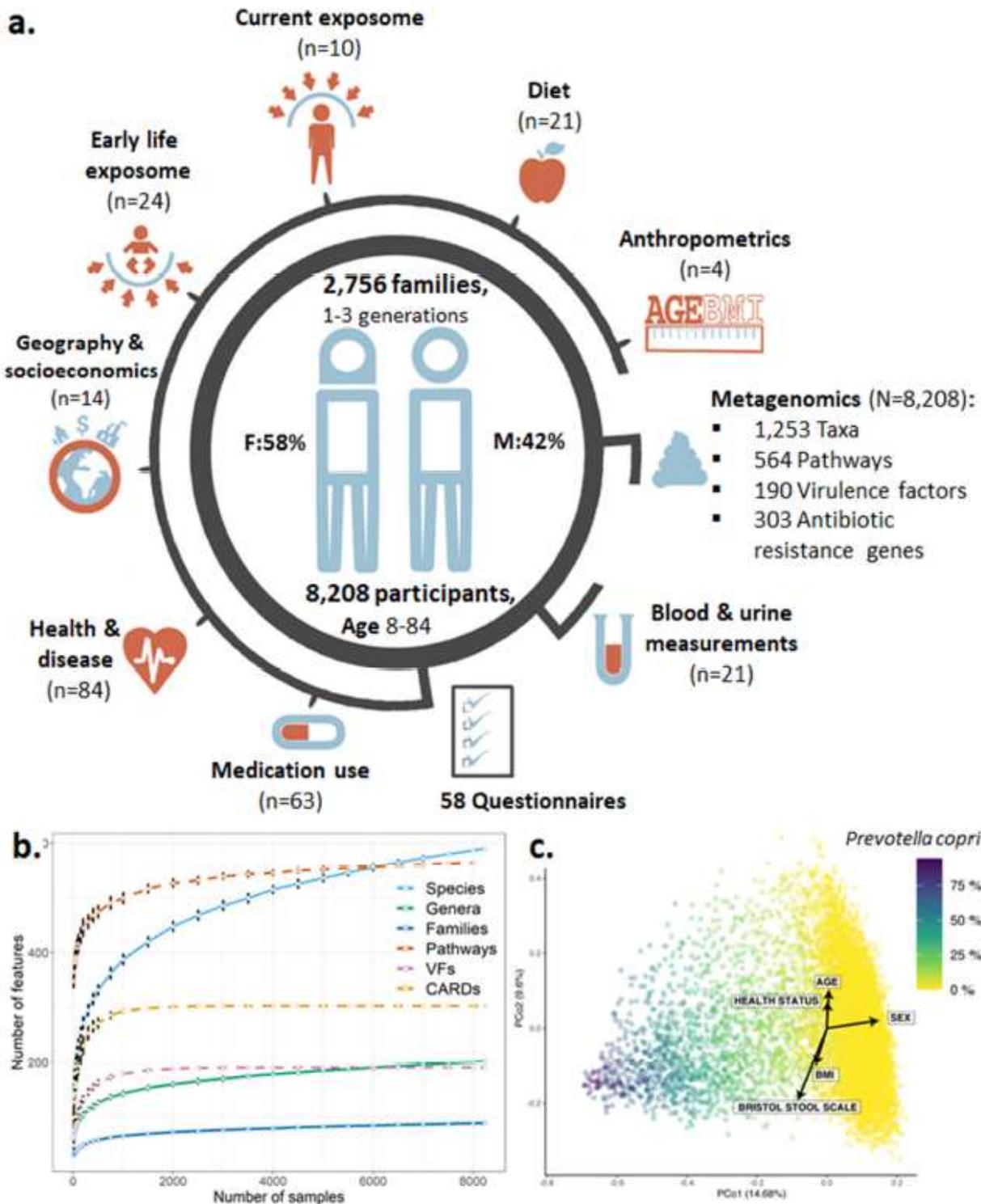


Figure 1

Summary of the Dutch Microbiome Project a, Graphical summary of the cohort and overview of available metadata (n = number of variables collected, N = sample size). b, Number of microbial features discovered in relation to sample size. Error bars denote the standard deviation of 100 resamplings. c,

Biplot of principal coordinate analysis visualising the beta-diversity of the microbiome data. Colour indicates relative abundance of *Prevotella copri*. Arrows indicate the influence of self-reported health, anthropometrics and faecal sample metadata. VFs: bacterial virulence factors, CARD: antibiotic resistance gene families.

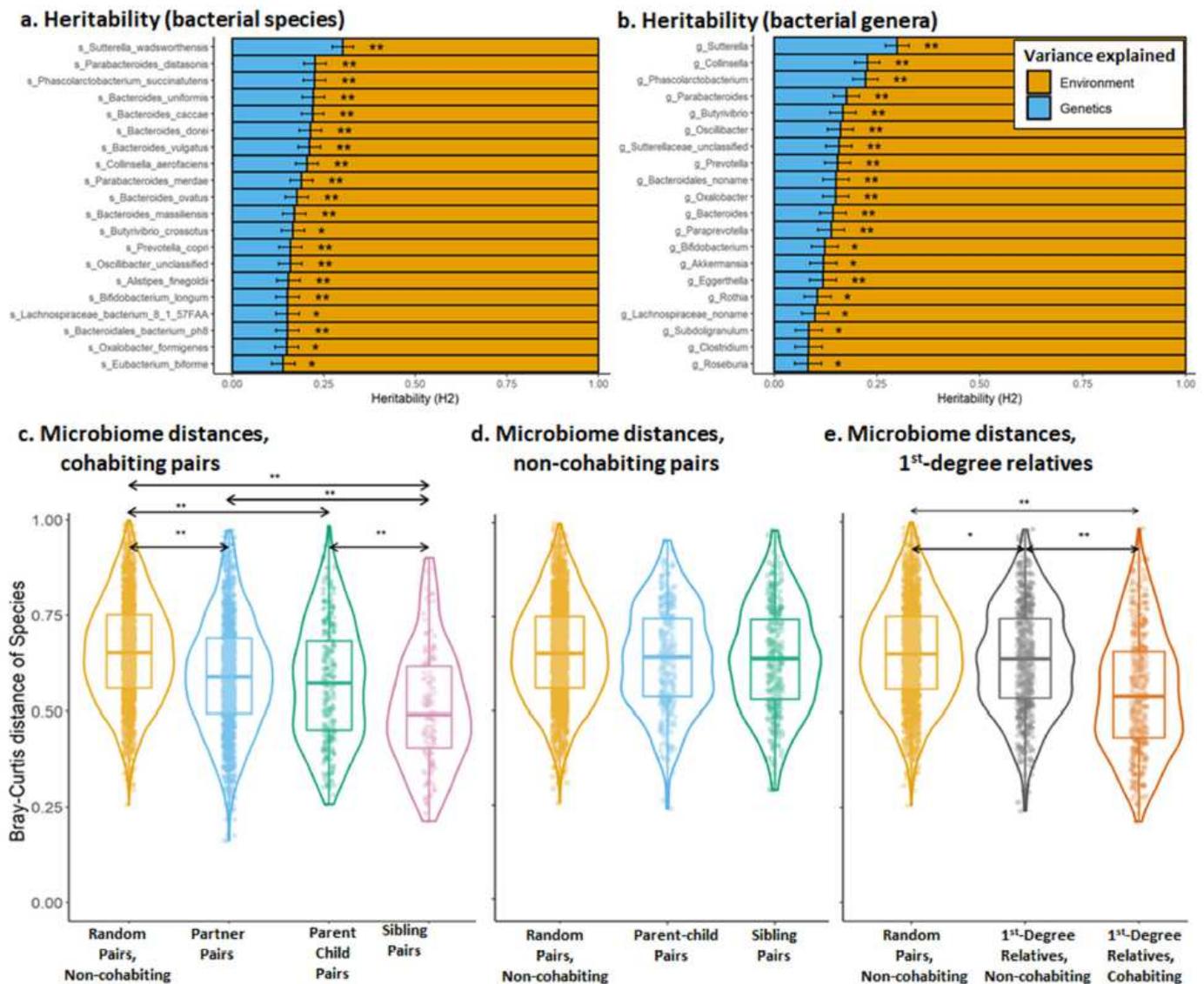


Figure 2

Heritability and impact of cohabitation on the gut microbiome a, Top heritable species. b, Top heritable genera. ** taxa significantly heritable at FDR < 0.1. * taxa with nominally significant heritability (p-value < 0.05). Panels c, d and e show pairwise microbiome distance comparisons. Bray-Curtis dissimilarities calculated using microbial species of groups of c) random, non-cohabiting pairs compared to cohabiting partners, parent-child pairs and sibling pairs, d) random pairs compared to non-cohabiting parent child and sibling pairs and e) random pairs compared to non-cohabiting first-degree relatives and cohabiting first-degree relatives. Significantly different groups: ** for FDR < 1.0e-5 and * for FDR < 0.05.

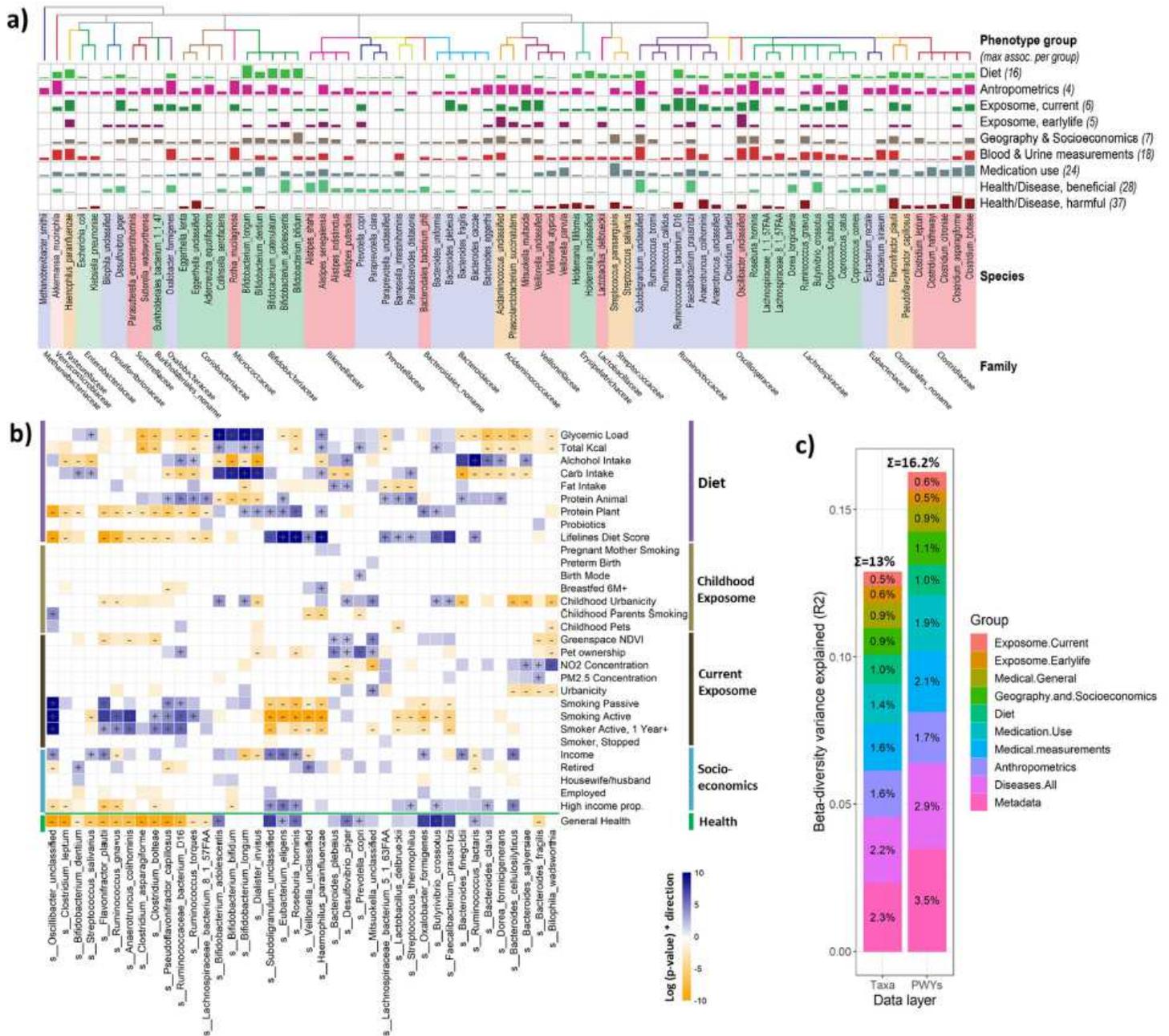


Figure 3

Microbiome-phenotype associations a, Selected study-wide-significant associations (FDR < 0.05) per phenotype group, clustered by taxonomy. Bar height indicates the number of associations relative to the maximal number of associations for the phenotype group. b, Microbiome-phenotype associations for diet, childhood exposome and current exposome and socioeconomics in comparison to healthy microbiome signature. Microbial species are clustered by association p-value using hierarchical clustering and coloured by direction of association. Study-wide significant associations (FDR < 0.05) are marked with + or -. Coloured associations without a mark indicate nominally significant associations (p-value < 0.05). c, The variance in microbiome composition and function explained by phenotype groups.

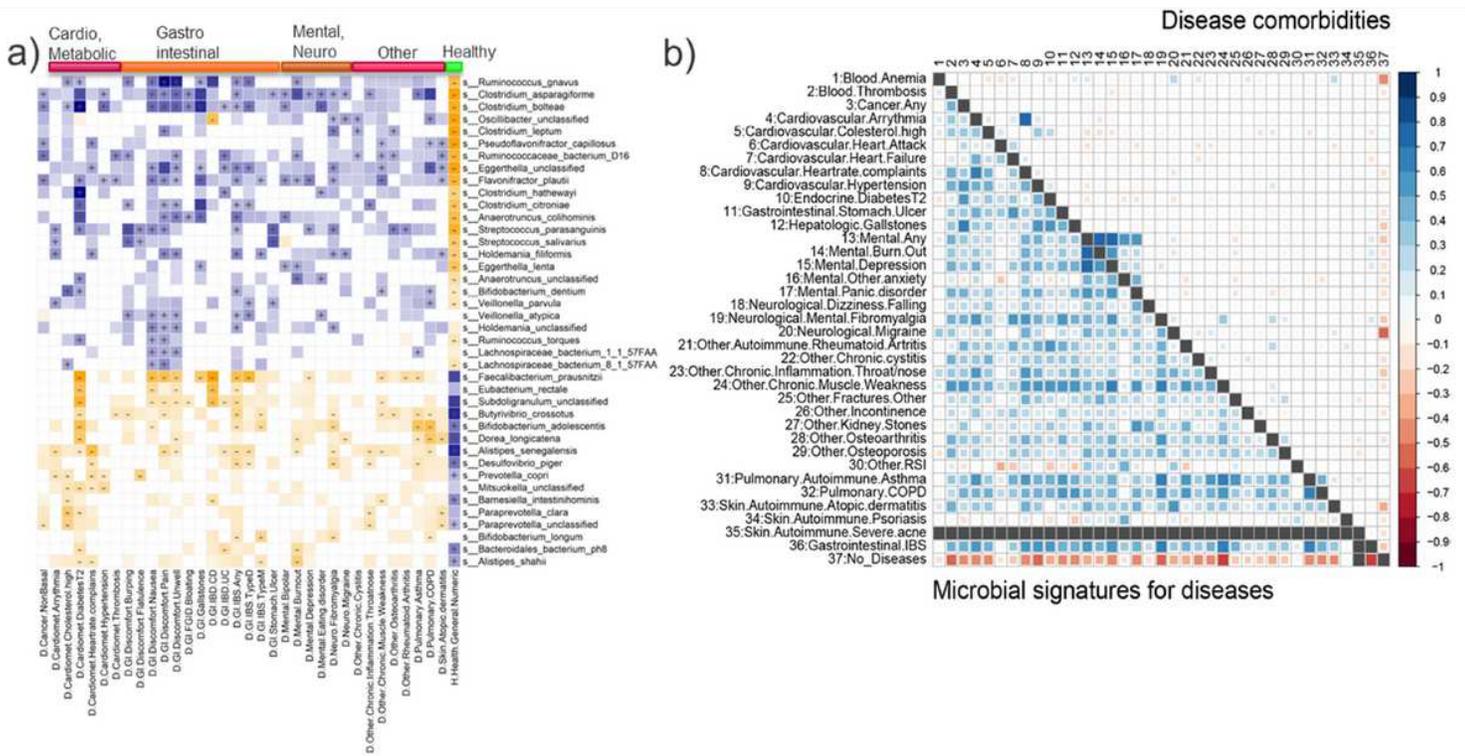


Figure 4

Microbiome signatures of health and diseases. a, Heatmap of microbial species associated with categories of diseases and health status. Diseases are sorted and labelled by disease type. Microbial species are clustered by association p-value (indicated by colour intensity) using hierarchical clustering. Associations are coloured by direction of effect (blue = positive, orange = negative), with associations significant at study-wise FDR < 0.05 marked with +/- for positive/negative correlations. Coloured associations without a label indicate nominally significant associations (p-value < 0.05, no multiple testing correction). b, Comparison of correlations between signatures predictive for diseases (lower triangle) and comorbidities of these diseases in the cohort (upper triangle).

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementarydiscussionrdy.pdf](#)
- [supplementaryfiguresrdy.pdf](#)
- [Supplementarytables.xlsx](#)