

1 **The making of the oral microbiome in Agta hunter-gatherers**

2 Begoña Dobon^{1,2+}, Federico Musciotto^{1,3+}, Alex Mira^{4,5}, Michael Greenacre^{6,7}, Abigail
3 E. Page⁸, Mark Dyble⁹, Daniel Smith¹⁰, Sylvain Viguier⁹, Rodolph Schläpfer¹,
4 Gabriela Aguilera², Leonora H. Astete¹¹, Marilyn Ngales¹¹, Vito Latora^{12,13,14}, Federico
5 Battiston^{1,15}, Lucio Vinicius^{1,9#}, Andrea B. Migliano^{1,9#*}, Jaume Bertranpetit^{2#*}

6

7 ¹Department of Anthropology, University of Zurich; Zurich, Switzerland.

8 ²Institut de Biologia Evolutiva (CSIC-Universitat Pompeu Fabra); Barcelona, Spain

9 ³Dipartimento di Fisica e Chimica, Università di Palermo; Palermo, Italy.

10 ⁴Department of Health and Genomics, Center for Advanced Research in Public Health,
11 FISABIO Foundation; Valencia, Spain.

12 ⁵CIBER Center for Epidemiology and Public Health; Madrid, Spain

13 ⁶Department of Economics and Business, Universitat Pompeu Fabra & Barcelona
14 Graduate School of Economics; Barcelona, Spain.

15 ⁷Faculty of Biosciences, Fisheries and Economics, University of Tromsø; Norway.

16 ⁸Department of Population Health, London School of Hygiene and Tropical Medicine;
17 London, UK.

18 ⁹Department of Anthropology, University College London; London, UK.

19 ¹⁰Bristol Medical School, University of Bristol; Bristol, UK.

20 ¹¹Lyceum of the Philippines University, Intramuros, Manila, Philippines.

21 ¹²School of Mathematical Sciences, Queen Mary University of London; London, UK.

22 ¹³Dipartimento di Fisica ed Astronomia, Università di Catania and INFN; Catania, Italy.

23 ¹⁴Complexity Science Hub Vienna (CSHV); Vienna, Austria.

24 ¹⁵Department of Network and Data Science, Central European University; Vienna 1100,
25 Austria.

26 ⁺These authors contributed equally to this work.

27 [#]These authors contributed equally to this work.

28 ^{*}Corresponding authors. Email: jaume.bertranpetit@upf.edu, andrea.migliano@uzh.ch

29

30 **Abstract**

31 **Ecological and genetic factors have influenced the composition of the human**
32 **microbiome during our evolutionary history. We analyzed the oral microbiota of**
33 **the Agta, a hunter-gatherer population where part of its members is adopting an**
34 **agricultural diet. We show that age is the strongest factor modulating the**
35 **microbiome, likely through immunosenescence as there is an increase of**
36 **pathogenicity with age. Biological and cultural processes generate sexual**
37 **dimorphism in the oral microbiome. A small subset of oral bacteria is influenced**
38 **by the host genome, linking host collagen genes to bacterial biofilm formation. Our**
39 **data also suggests that shifting from a fish/meat to a rice-rich diet transforms their**
40 **microbiome, mirroring the Neolithic transition. All these factors have implications**
41 **in the epidemiology of oral diseases. Thus, the human oral microbiome is**
42 **multifactorial, and shaped by various ecological and social factors that modify the**
43 **oral environment.**

44

45 **Introduction**

46 The composition and diversity of the human oral microbiota has been influenced by
47 several factors during our evolutionary history^{1,2}. Some are intrinsic biological
48 characteristics of the host, such as age, sex, and genetic composition, while others such
49 as diet, drinking water sources, oral hygiene, lifestyle and social interactions are
50 external factors²⁻⁴. These factors modulate the physiological conditions of the oral
51 cavity and affect the composition and diversity of the oral microbiota. While the oral
52 microbiota is one of the most diverse sites in the human body and shows high variability
53 between individuals, it remains stable within individuals over time⁵. Little is known

54 about how the composition of the oral microbiome is modulated in populations adapted
55 to the hunting and gathering niche, where the fully mature oral biofilm microbiome can
56 be studied without the confounding effects of tooth brushing or professional dental
57 cleaning, similar conditions to how the human oral microbiome must have evolved in
58 the past⁶.

59 To investigate the multiple ecological and genetic factors shaping the human oral
60 microbiome, we have analysed the oral microbiome of the Agta hunter-gatherers from
61 the Philippines. The Agta are predominantly hunter-gatherers (fishing, hunting, and
62 gathering)⁷, and while their main source of animal protein is obtained by riverine and
63 marine spearfishing or by hunting, other activities such as inter-tidal foraging, wild food
64 gathering, low-intensity cultivation, wage labour and trade complement their
65 economy^{8,9}. Interestingly, there is high variability among the Agta on the amount of
66 hunting, gathering and sea foraging products that is traded for rice and other items (such
67 as tobacco) with farming neighbours⁷, causing the Agta lifestyle to range from
68 completely mobile foragers with a protein-rich diet, to settled low intensity farmers,
69 with a rice-rich diet^{7,9,10}.

70 To detect fine-scale variation in the oral microbiome of Agta hunter-gatherers, we
71 collected saliva samples from 138 Agta, aged 5 to 65 years, sequencing the 16S rRNA
72 region, and identified 5430 amplicon sequence variants (ASVs)¹¹ belonging to 110
73 genera. To study the genetic host factors associated with the microbiome composition
74 we also genotyped all individuals with the Axiom Genome-wide Human Origins array.
75 We combined this information with additional individual data on household
76 composition, age, sex, and diet, measured as the proportion of meals including meat/fish
77 (any animal protein), and proportion of meals that consisted of only agricultural
78 products (rice) (Supplementary Figure 1). By using this rich dataset, we have been able

79 to discern the different contributions of age, sex, diet, and host genetics in the making of
80 the oral environment.

81

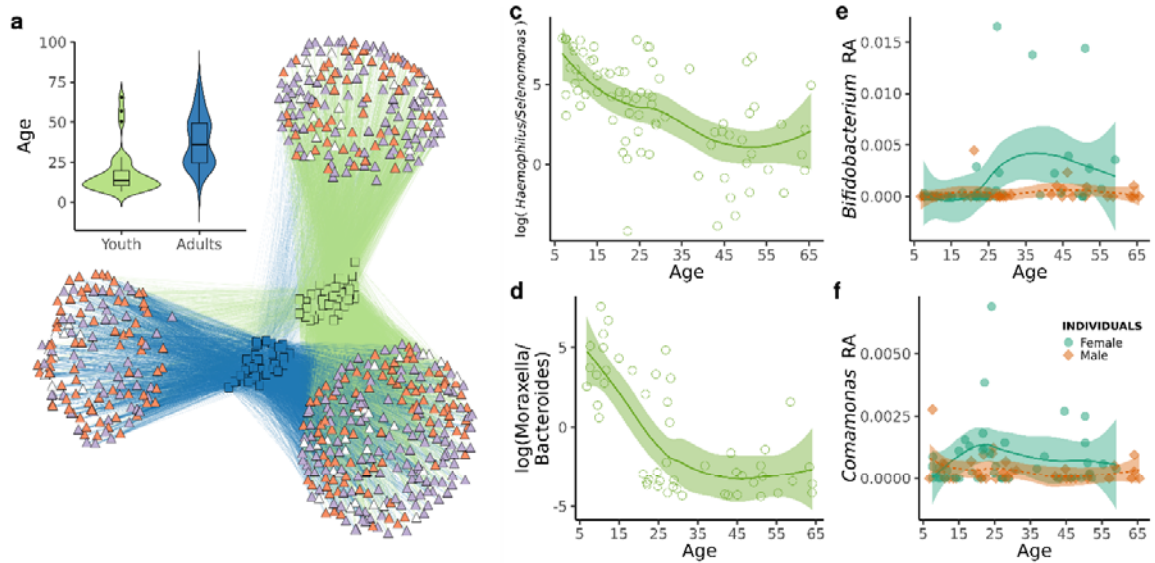
82 **Results**

83 **Factors influencing Agta oral microbiome composition.** The Agta oral microbiome is
84 mostly composed by *Firmicutes* (mean ASV prevalence = 33.1%, sd = 11.4),
85 *Proteobacteria* (27±14.6%), *Actinobacteria* (15.5±8.3%) and *Bacteroidetes*
86 (14.5±7.7%) (Supplementary Figure 2). To compare and identify the main ecological
87 and social factors contributing to microbiome variation, we performed a constrained
88 logratio analysis (LRA)¹² on the bacterial genus abundance. Marginally, age explains
89 7.2% of the total logratio variance ($P < 0.0001$, based on 9999 permutations), sex
90 explains 2.2% ($P = 0.018$), and diet 3.6% ($P = 0.015$). Altogether they explain 13.0% of
91 the total logratio variance. We also applied a bipartite stochastic block model (biSBM)
92 approach¹³ at the ASV level, where we assigned each bacterium to an individual, and
93 then clustered the individuals according to the bacteria they have in common. We
94 restricted the analysis to the Core Measurable Microbiota (CMM), that we define as
95 ASVs present in at least 10% of the Agta to reduce random errors due to low-prevalent
96 taxa (Supplementary Figure 3). The best model produced two clusters of people and
97 three clusters of bacteria (Figure 1a). While we did not find differences in diet or
98 proportions of sexes between the clusters of individuals (Supplementary Figure 4), they
99 strongly differ in their age distribution: adults (mean age = 38 years old) and youth
100 (mean age = 18 years old) (Welch t-test, $t = 5.78$, $df = 71.35$, $P < 0.0001$), with 55.48%
101 of ASVs in the CMM being more associated with one of the clusters of individuals.
102 Thus, while age, diet, and sex influence the composition of Agta microbiome, the

103 biSBM singles out age as the main modulator of the hunter-gatherer core oral

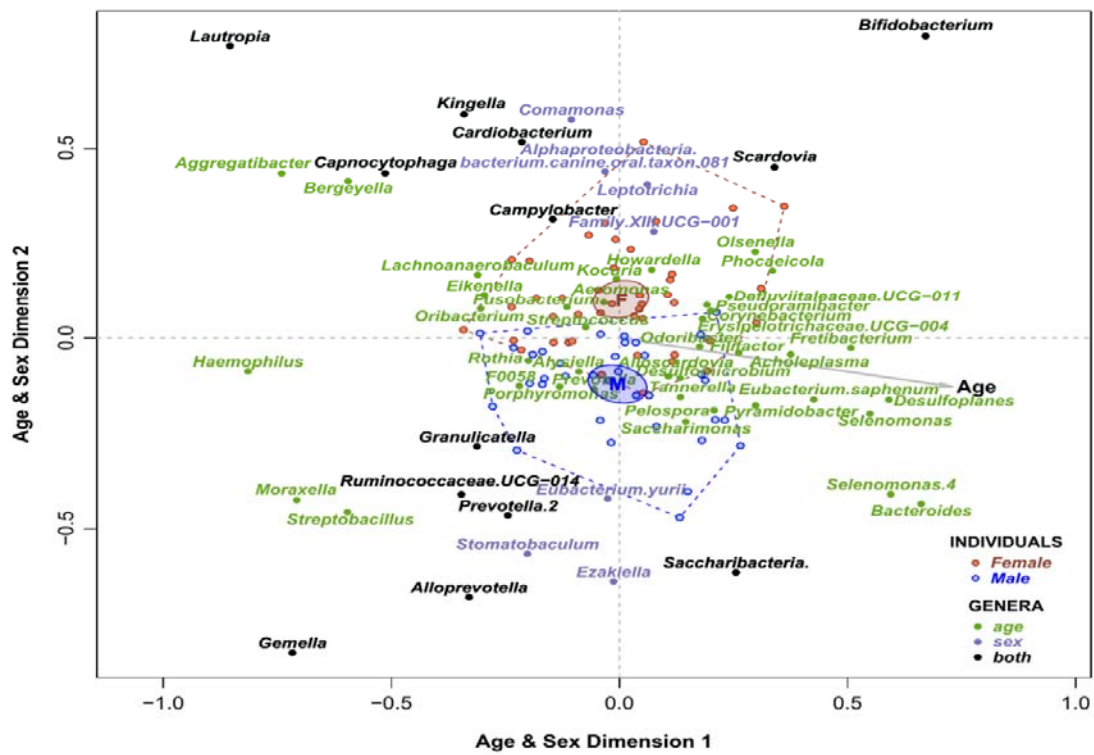
104 microbiome.

105



b

Constrained Analysis of Age & Sex



106

107

108 **Figure 1. Age and sex-related effects in the hunter-gatherer oral microbiome.** a)
109 Network representation of the hunter-gatherer CMM. ASVs (triangles) are colour-coded
110 as: putatively pathogenic (purple), non-pathogenic (orange) or unclassified (white).
111 Inset shows age distribution for the two clusters of individuals (squares). b) Logratio
112 analysis constrained to age and sex differences on the bacterial composition at genus
113 level. The effects of diet were partialled out. Only genera statistically significant in at
114 least 20 (for age) or 10 (for sex) logratios are displayed (p-value < 0.05 after Benjamini-
115 Hochberg correction). Dashed lines enclose all individuals (dots) within a category of
116 sex, with 95% confidence ellipses for their means. Taxa are colour-coded depending on
117 the associated variable: age, sex, or both. The starting point of the grey arrow indicates
118 the mean age of the population (30 years old). Log ratio of c) *Haemophilus* and
119 *Selenomonas* abundance and d) *Moraxella* and *Bacteroides* according to age. Line and
120 shaded area indicate the 95% confidence interval of the mean. Relative abundance of e)
121 *Bifidobacterium* and f) *Comamonas* according to age and sex. Lines and shaded areas
122 indicate the 95% confidence interval of the mean for each sex.

123

124 **Old age is associated with increased frequency of oral pathogens.** To investigate the
125 independent effects of ageing on the oral microbiome, we performed a LRA constrained
126 with age and sex after partialing out the effects of diet. The resulting ordination shows
127 that the effects of age and sex are mostly independent, with only few genera being
128 affected by both variables (Supplementary Figure 5): as expected, the first dimension is
129 associated with the age of the individuals, while the second dimension separates them
130 according to sex (Figure 1b).

131 There is a clear change in the composition and frequency of certain bacteria with
132 age (Figure 1c-d). At young age, we observe organisms that typically live in mucosa,
133 such as *Haemophilus* and *Moraxella*, that infect the upper and lower respiratory tract
134 but are detected in the oral cavity and saliva which are their vehicles of transmission.
135 Other genera found at younger ages include bacteria normally associated to good oral
136 health, such as *Bergeyella* and *Rothia*¹⁴. However, at older ages we observe a marked
137 decline in the abundance of those genera and an increase of important pathogens related
138 with periodontitis including the “red complex” periodontal pathogen *Tannerella*, as well

139 as other periodontitis-related bacteria (*Filifactor*, *Fretibacterium*, *Saccharimonas*,
140 *Selenomonas*, and *Phocaeicola*), consistent with a higher incidence of this disease with
141 older age¹⁵. We also found organisms associated with cavities (*Olsenella*), with dental
142 plaque and dental calculus formation (*Corynebacterium*), with pulmonary infections,
143 sepsis, or bacteremia, and with chronic diseases (*Acholeplasma*) (see Methods for in-
144 depth bacteria pathogenic classification). Another sign of ageing was the presence in the
145 oral cavity of gut bacteria (*Bacteroides*) indicating a potential age-related decline in
146 immunological function and filtering¹⁶. However, such changes are not associated with
147 a decrease in the alpha-diversity of the total oral microbiome as measured by the
148 number of bacteria observed or their phylogenetic complexity (Supplementary Figure
149 6a-b), suggesting that the overall effect of aging is a replacement of protective and
150 commensal bacteria by pathogenic ones. This is supported by an increase in the number
151 of potential pathogenic bacteria in the CCM in bacterial clusters associated with older
152 ages (Fisher exact test, $P < 0.001$) (Figure 1a).

153

154 **Sex differences shape composition but not diversity of the Agta oral microbiome.**

155 We found no differences in alpha-diversity in the Agta oral microbiome between males
156 and females (Supplementary Figure 6c-d), which may be explained by sex equality
157 within the Agta hunter-gatherer society regarding diet and social interactions^{17,18}.
158 Nevertheless, the LRA constrained to age and sex shows sex-related differences in the
159 composition of the oral microbiome (Figure 1b). For example, *Stomatobaculum* and
160 *Eubacterium yurii*, present in the oral cavity of smokers¹⁹, are associated with males,
161 consistent with Agta men chewing tobacco more frequently than women. It is also
162 interesting to mention *Comamonas* (Figure 1f), which even if it has been reported as a
163 possible contaminant in microbiome studies²⁰, its presence in females could be related

164 to its capacity of degrading the female hormone progesterone²¹. This bacterium has
165 been found in subgingival samples, where female hormones could be present either in
166 saliva or in gingival crevicular fluid.

167

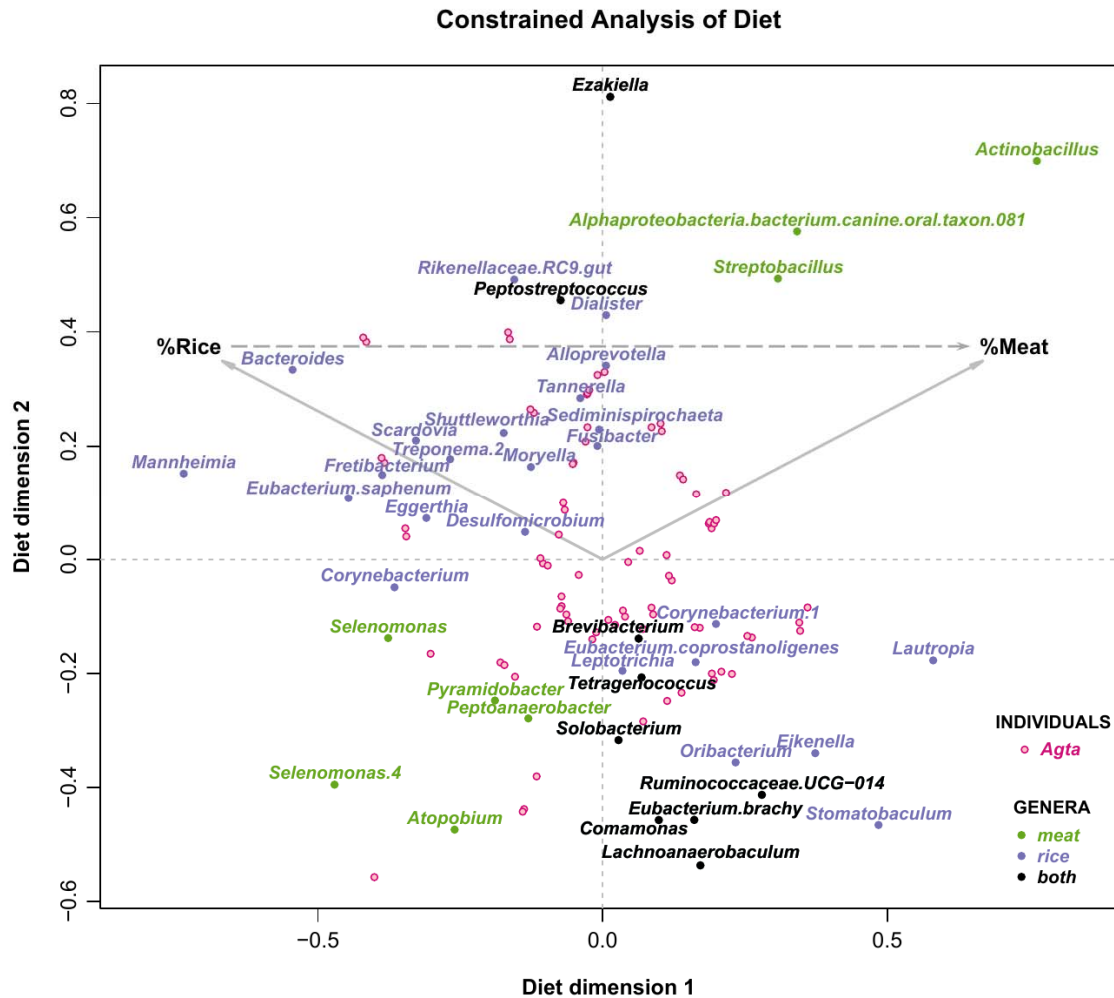
168 **Age and sex interactions in microbiome composition.** Some bacteria are significantly
169 associated to both age and sex-related differences, such as *Gemella*, which is a prevalent
170 inhabitant of the respiratory mucosa such as *Haemophilus* and *Moraxella*, supporting
171 the idea that mucosa-associated and/or respiratory-tract organisms are more frequently
172 acquired in young individuals, especially males. At older ages, the
173 *Bifidobacterium/Saccharibacteria* ratio distinguishes between sexes: while
174 *Bifidobacterium* is associated to females, the periodontal pathogen *Saccharibacteria* is
175 associated to males. Thus, the observed trend of increase of periodontal pathogens with
176 age is stronger in males, as expected by the global epidemiology of the disease^{22,23}. On
177 the other hand, we found an increase of the caries-related pathogens *Scardovia* and
178 *Bifidobacterium* associated to reproductive age females. Caries incidence increases with
179 age and is more prevalent in females²⁴, a more saccharolytic or acidic salivary
180 environment in older women, together with hormonal fluctuations and lower salivary
181 flow²⁵ could facilitate the proliferation of saccharolytic bacteria. The strong association
182 of *Bifidobacterium* with adult females could also be explained by its presence in
183 breastmilk²⁶. Also, its proliferation coincides with the start of the reproductive age and
184 the increase of childcare¹⁰ (Figure 1d).

185

186 **Influence of variation in rice consumption on the Agta oral microbiome.** While the
187 impact of diet on gut microbiome has been clearly established²⁷⁻³⁰, its role in the oral
188 microbiome is still unclear. Some studies have found little or no effect, whereas others

189 have found associations with specific nutrients^{2,25,31,32}. The variation in rice
190 consumption in the Agta allows us to assess both the relationship between the hunter-
191 gatherer diet on the oral microbiome and the effects of the recent introduction of
192 farming products. We performed a LRA on the bacterial genus abundance after
193 partialing out the effects of age and gender (Figure 2). The first dimension of the
194 ordination shows the gradient of the transition from a diet where all meals include meat
195 to where most consist of only rice. Agta following a hunter-gatherer diet, where most
196 meals contain meat, have large quantities of *Actinobacillus*, *Alphaproteobacteria* and
197 *Streptobacillus* and lower abundance of *Selenomonas*, *Atopobium*, *Peptoanaerobacter*,
198 and *Pyramidobacter*. The higher abundance of *Actinobacillus* in individuals ingesting a
199 protein-rich diet is particularly interesting, given the extraordinary proteolytic potential
200 of *A. actinomycetencomitans*, a well-known oral pathogen with destructive effects in the
201 gingival tissue and in aggressive forms of periodontitis³³. At the other extreme, in
202 individuals with a rice-rich diet, there is an increase of the highly saccharolytic dental
203 caries pathogen *Scardovia*, of *Treponema*, of gut organisms like *Butyrivibrio* and
204 *Erysipelotrichaceae*, and of *Eggerthia*, a rare organism isolated from dental abscesses.
205 We also ranked the ASV based on whether they are more present than expected in
206 individuals with high or low proportion of meals with only rice or with meat/fish. We
207 found that the scores associating each bacterial species with these two nutrients are
208 negatively correlated (Spearman's rho = -0.47, P < 0.0001). This fits with a general
209 separation of oral microorganisms in saccharolytic (caries-related, acidogenic and
210 acidophilic) and proteolytic (gum-disease and halitosis related, alkalophilic and NH₄
211 generators), as suggested in a metabolome-based study²⁵. Our results suggest that more
212 settled Agta, which consume more rice, experience a decline in oral health, confirming a

213 general pattern of health decline due to a Neolithic-like diet and a more farming-derived
 214 lifestyle^{9,34,35}.



215 **Figure 2. Effect of diet on the oral microbiome in the Agta.** Logratio analysis
 216 constrained to diet differences on the bacterial composition at genus level. The effects
 217 of age and sex were partialled out. Only genera statistically significant in more than five
 218 (for rice) or three (for meat) logratios are displayed (p-value < 0.05 after Benjamini-
 219 Hochberg correction). Taxa are colour-coded based on the variable they are associated:
 220 proportion of meals with meat (%Meat), proportion of meals with only rice (%Rice), or
 221 both. The original result was slightly rotated so that the dashed vector indicating the
 222 difference between %Meat and %Rice was horizontal, without any change in explained
 223 variance.

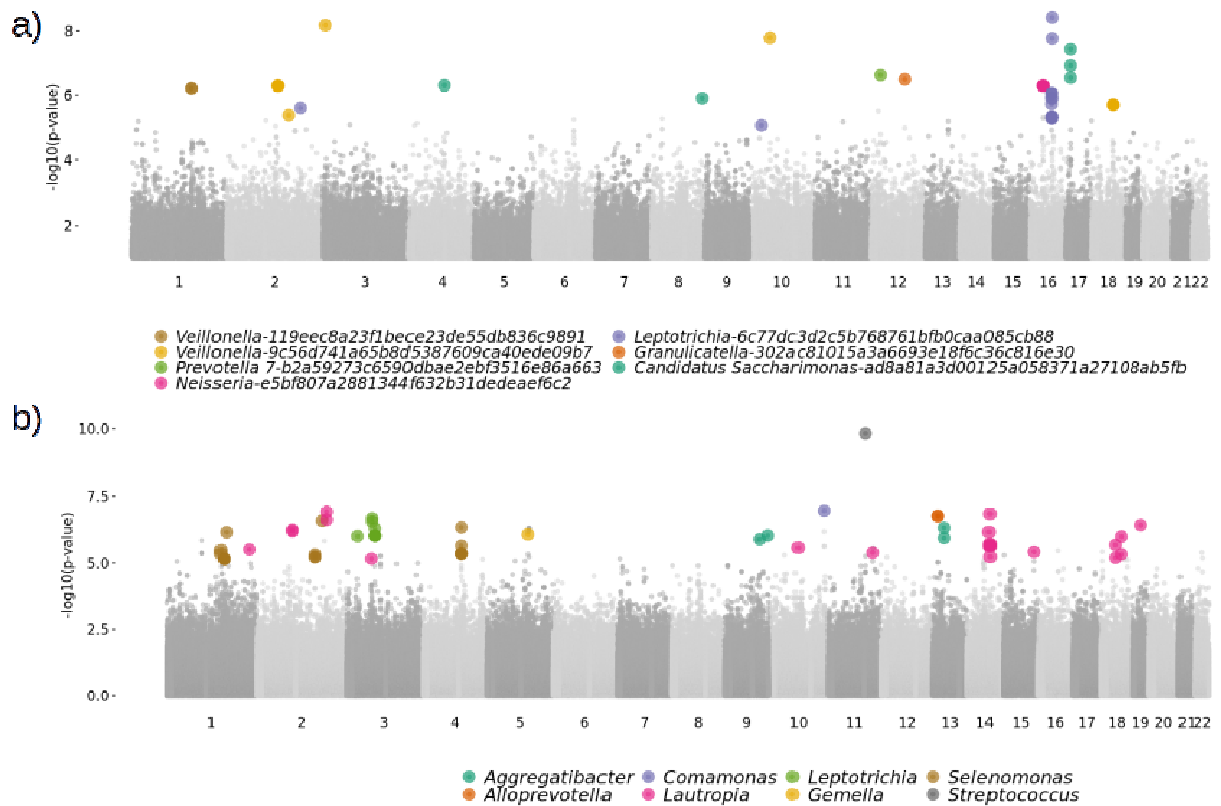
224

225 **Pathogenic oral bacteria are associated with host collagen genes.** The interaction
226 between the host genetic makeup and microbiome composition differs across body
227 sites^{36,37}, and seems especially weak in the oral cavity^{37,38}, making it difficult to assess
228 the co-evolution of our genome and the oral microbiome. To overcome this, we
229 performed a genome-wide association study (GWAS) using a mixed model approach in
230 a population that evolved in a hunter-gatherer niche and without the confounding
231 influence of antibiotics or brushing. We treated the relative abundance of each
232 bacterium as an independent trait, adding age, sex, and household as covariates and
233 kinship as a random effect. Household membership was used as proxy for the strength
234 of social interactions between individuals, as social interactions predict microbiome
235 sharing (Musciotto *et al. companion paper*). These analyses were performed using the
236 CMM, and then using 92 genera present in at least 10 Agta. All bacteria identified in the
237 Agta (Supplementary Table 1 and Supplementary Figure 7) overlap with those of other
238 oral microbiome GWAS^{3,36,37} pointing to a small subset of oral bacteria influenced by
239 the human genome (Figure 3). A pathway enrichment analysis linked this subset of the
240 oral microbiome to several biological host functions (body fat metabolism, wound
241 healing, and collagen trimmers) (Supplementary Table 2). Of relevance is an association
242 between the pathogenic bacteria *Aggregatibacter* and *Selenomonas* with genetic
243 variation in collagen genes. The ability to bind collagen is a vital feature in the oral
244 cavity, as many oral bacteria require collagen-binding proteins to attach to oral tissues³⁹
245 suggesting a genetic basis for the predisposition of biofilm formation by those bacteria.
246 We further tested if we could detect signatures of positive selection in the host genomic
247 regions associated with the oral microbiome, but we found no signals indicating recent
248 selective pressures caused by oral bacteria.

249

250

251



252

253 **Figure 3. Genome-wide association study on bacteria abundance.** Aggregated
254 Manhattan plot of the GWAS results of the a) seven ASV and b) eight genera with non-
255 zero PVE (“chip heritability”) estimates with at least one significant genetic association.
256 Each dot is a SNP and significant SNPs-bacteria associations ($q < 0.1$) are color-coded
257 according to the associated bacteria.

258

259 Discussion

260 The Agta microbiome is influenced by external factors such as social interactions
261 (Musciotto *et al. companion paper*) as well as intrinsic and ecological factors such as
262 age, sex, diet and host genetics. Among the latter, we have shown that age has the
263 strongest effect, with commensal or beneficial microbiota being replaced by potentially
264 pathogenic ones with ageing. The proliferation of oral pathogenic bacteria exhibits

265 sexual dimorphism, with caries-related (in females) and periodontitis-related (in males)
266 bacteria increasing with age, likely associated with immunosenescence⁴⁰ and with a sex-
267 specific oral environment due to biological and cultural factors. In the Agta, the increase
268 of farming-derived novel foods such as rice influences their microbiome composition
269 and health. The relatively small subset of bacteria linked to the host genome, which are
270 also found associated to other factors, suggests that the Agta oral microbiome is mainly
271 affected by environmental (diet) and intrinsic factors (age), with little influence of
272 individual host genetic variation (Supplementary Figure 5). Thus, environmental factors
273 and not host genetics are the main driving force for oral microbiota acquisition, in
274 agreement with Mukherjee *et. al*⁴¹. Based on the case study of the Agta hunter-
275 gatherers, we conclude that the human oral microbiome is multifactorial with distinct
276 subsets of bacteria shaped by specific ecological and social factors, reflecting multiple
277 adaptations in the domains of life history, sociality, and diet.

278

279

280 **Methods**

281 **Ethics approval**

282 This study was approved by UCL Ethics Committee (UCL Ethics code 3086/003) and
283 carried out with permission from local government and community members. Informed
284 consent was obtained from all participants, after group and individual explanation of
285 research objectives in the indigenous language. A small compensation (usually a
286 thermal bottle or cooking utensils) was given to each participant. The National
287 Commission for Indigenous Peoples (NCIP), advised us that the process of Free Prior
288 Informed Consent with the tribal leaders, youth and elders would be necessary to
289 validate our data collection under their supervision. It was done in 2017 with the

290 presence of all tribal leathers, elders and youth representatives at the NCIP regional
291 office, with the mediation of the regional officer and the NCIP Attorney. The validation
292 process was approved unanimously by the tribal leaders, and the NCIP, and validated
293 the full 5 years of data collection.

294

295 **Saliva sample collection**

296 Saliva samples from 155 Palanan Agta were collected over two field seasons: April-
297 June 2013 and February-October 2014. For comparative genetic studies we also used
298 saliva samples from 21 Mbendjele Bayaka, an African hunter-gather population,
299 collected in 2014, and 14 Palanan farmers collected in 2007-2009. In all cases saliva
300 was collected using the Oragene-DNA/saliva kit and participants were asked to rinse
301 their mouth with water and to spit into the vial until half full. After collection and
302 transportation, saliva samples were stored at the UCL Department of Anthropology,
303 London, UK at -20°C.

304

305 **Microbial DNA extraction and 16S rRNA gene sequencing**

306 A total of 155 Agta saliva samples were selected to study their microbiome
307 composition. DNA was extracted following the protocol for manual purification of
308 DNA for Oragene-DNA/saliva samples. The 16S rRNA gene V3-V4 region was
309 amplified by PCR with primers containing Illumina adapter overhang nucleotide
310 sequences. All PCR products were validated through an agarose gel and purified with
311 magnetic beads. Index PCR was then performed to create the final library which was
312 also validated through an agarose gel. All samples were pooled together at equimolar
313 proportions and the final pool was qPCR quantified prior to the MiSeq loading. Raw
314 Illumina pair-end sequence data were demultiplexed, quality filtered and denoised with

315 QIIME 2 2019.1⁴² and DADA2⁴³. DADA2 generates single nucleotide exact amplicon
316 sequence variants (ASVs). ASV are biological meaningful entities as they identify a
317 specific DNA sequence and allow for higher resolution than using operational
318 taxonomic units (OTUs)¹¹. Taxonomic information was assigned to ASVs using a naïve
319 Bayes taxonomy classifier against SILVA database release 132 with a 99% identity
320 sequence⁴⁴. ASVs that did not belong to the kingdom Bacteria, or that were classified as
321 mitochondrial or chloroplast sequences and samples with an extremely low number of
322 sequences (8000) were excluded from further analyses. ASVs were aligned with
323 MAFFT⁴⁵ and a rooted phylogenetic tree was constructed with FastTree2⁴⁶ using default
324 settings via QIIME 2. This resulted in a total of 5430 ASVs and 138 Agta (67 women
325 and 71 men). We generated a rarefaction curve with R package vegan (version 2.5-7)⁴⁷
326 to determine that the richness of the samples had been fully observed (Supplementary
327 Figure 8). The number of observed ASVs and Shannon Diversity index were calculated
328 with R package Phyloseq (version 1.30.0)⁴⁸. Faith's Phylogenetic Diversity index⁴⁹ was
329 calculated with R package picante (version 1.8.2)⁵⁰ using the rooted phylogenetic tree
330 generated in R⁵¹. To determine the set of microbial traits to be included in the analyses,
331 we selected ASVs with at least 10 reads in at least 2 individuals (n = 1980), then we
332 aggregated those with a taxonomic assignment at a genus level, resulting in 110 genera.
333 At the ASV level we also defined a Core Measurable Microbiota (CMM), consisting of
334 ASVs that appear in at least 10% of the Agta individuals (14 or more) resulting in 575
335 ASVs (out of 1980) that represent 90% of the composition of the Agta microbiome.

336

337 **Genotype data**

338 A total of 190 saliva samples were genotyped with the Affymetrix Axiom Genome-
339 Wide Human Origins 1 array. DNA extraction was carried out following the protocol

340 for manual purification of DNA for Oragene-DNA/saliva samples in the same
341 laboratory that sequenced the 16S rRNA data. Samples were analyzed with Axiom
342 Analysis Suite v4.0 following the Axiom genotyping best-practices workflow for saliva
343 samples. 618810 markers and 177 samples passed initial quality control. Single
344 nucleotide polymorphisms (SNPs) with less than 95% genotyping rate and samples
345 where the estimated gender from the genotypes did not match the recorded gender were
346 excluded from the analysis. Duplicated samples were identified with KING⁵² and
347 removed. This resulted in a total of 617063 markers and 174 samples: 141 Agtas, 19
348 Bayaka and 14 Palanan farmers.

349

350 **Ethnographic data collection**

351 Ethnographic data collection occurred over two field seasons from April-June 2013 and
352 February-October 2014. In the first season we censused 915 Agta individuals (54.7%
353 which were male) across 20 camps, collecting basic information on household
354 composition, sex, and estimated ages. Following relative aging protocols⁵³, accurate
355 ages were established for all individuals post data collection.

356

357 **Diet data collection**

358 Dietary recall data was collected at the household level over a period of 10 days. We
359 asked the mother and the father at the end of the day (between 17:00 – 18:00) what
360 foods they had eaten that day, including agricultural produces from trade with nearby
361 farmers. We counted the total amount of meals we had recorded for a household and
362 established what proportion of these consisted of meat, vegetables, fruits, honey, and
363 rice. Therefore, this is only a rough guide to dietary composition and does not take
364 calorific intake or absolute weighs of the different food types into account. Dietary data

365 for 80 individuals (37 males and 43 females) was annotated based on the proportion of
366 meals that consisted of only rice, and the proportion of meals that included meat
367 (primarily fish and other marine resources and game).

368

369 **Classification of oral bacteria as pathogens**

370 Bacteria were classified as potential oral pathogens if they have been reported as
371 etiological agents of periodontitis or dental caries. Assignment as periodontal pathogen
372 was performed according to the systematic review of Perez-Chaparro *et al.*⁵⁴ and
373 Socransky *et al.*⁵⁵, or if they have been previously associated with this gum disease^{56,57}.
374 Bacteria were classified as caries pathogens if they were described in transcriptomic
375 studies of human cavities, according to Simon-Soro *et al.*⁵⁸ and Simon-Soro and Mira⁵⁹,
376 previously associated with caries^{60,61}, with cavities⁶² or with dental plaque and dental
377 calculus formation^{63,64}. Bacteria reported as etiological agents of respiratory infections
378 and biofilm-mediated infections were also considered pathogens, including organisms
379 that can be present in healthy carriers. These included species described in Leung *et*
380 *al.*⁶⁵, Bellussi *et al.*⁶⁶, and Natsis and Cohen⁶⁷. Bacteria causing urinary tract infection
381 or sexually transmitted diseases which can transiently be found in the oral cavity were
382 also considered as potential pathogens and included microorganisms described in Lanao
383 *et al.*⁶⁸ and Jung *et al.*⁶⁹. Common oral commensals potentially causing endocarditis or
384 systemic infections in immunocompromised patients were not considered pathogens. If
385 a bacterium was isolated from the oral cavity of an animal, it was considered an oral
386 inhabitant for the sake of our classification. If taxonomic classification in our dataset
387 could be assigned at the genus level only, it was considered a pathogen if: i) > 90 of
388 species within the genus were pathogenic, or ii) it included a major pathogenic species
389 but the rest of species within the genus were not oral inhabitants, according to the

390 Human Oral Microbiome Database (<http://www.ehond.org/>)⁷⁰. Bacteria with taxonomic
391 assignments at higher levels than genus (family, order, class) were excluded from this
392 analysis.

393 For assignment of bacteria to pathogenic or non-pathogenic, we used species-
394 level ASVs, given that there are multiple cases where different species from the same
395 genus had a different assignment. If taxonomic classification of the ASV was only
396 possible at the genus level, it was considered a pathogen if: i) >90% of named species
397 within the genus were pathogenic, or ii) the genus included a major pathogenic species
398 but the remaining species within the genus were not classified as oral by the Human
399 Oral Microbiome Database⁷⁰. ASV with a top hit to a sequence classified as “Oral taxa”
400 in databases but without a species assignment were not considered named species and
401 were discarded from the analysis. Cases where taxonomic classification of the ASV was
402 only possible at the family level or higher were also discarded.

403

404 **Multivariate compositional data analysis on microbial composition**

405 We performed a constrained logratio analysis (LRA) using the package easyCODA¹² in
406 R⁵¹ on the Agta oral microbiome at the genus level using as constraining covariates age
407 (as continuous variable), sex (male and female), and diet (both proportion of meals with
408 meat and proportion of meals with only rice). The microbiome abundance counts of
409 each Agta individual were treated as compositional data⁷¹ and transformed to logarithms
410 of ratios (logratios). Constrained LRA is a special case of redundancy analysis¹² where
411 the total logratio variance is decomposed into parts explained by the covariates (the
412 "constrained variance") and a residual part (the "unconstrained variance", unrelated to
413 the covariates). Then, the ordination resulting from the LRA explains a maximum of the
414 constrained variance in a reduced two-dimensional solution. The statistical significance

415 of the three covariates was assessed using a multivariate permutation test (999
416 permutations) in the R package `vegan`⁴⁷. There is no correlation between these three
417 covariates, except within the diet covariate, where the two variables are negatively
418 correlated (Spearman's $\rho = -0.54$, $P < 0.0001$) (Supplementary Figure 1). To focus on
419 the genera affected only by internal factors (age and gender), we performed a
420 constrained LRA on the microbial composition after partialing out the effects of diet.
421 Similarly, to identify genera affected exclusively by the diet, we performed a
422 constrained LRA after partialing out the effects of age and gender. Taxon-covariate
423 association was ranked by counting the number of significant logratios for each of the
424 taxa, with p-value < 0.05 controlling for the false discovery rate (FDR) at level $\alpha =$
425 0.05 ^{72,73}.

426

427 **Community detection**

428 To model the relationship between the Agta and the CMM, we used a stochastic block
429 model (SBM) approach specifically suited for bipartite networks¹³. SBM infers the
430 community structure⁷⁴ that better fits the existing graph, by building a prior distribution
431 for edges that holds no information on real data and using it in the framework of
432 Bayesian inference (biSBM) to find a partition of the two types of nodes whose
433 associated entropy is maximal. In this framework, the absence of links between nodes of
434 the same type or set is not considered informative for the model, as it is expected given
435 the bipartite nature of the graph, different from the general version of SBM. We selected
436 the number of clusters in the two sets that minimize the description length⁷⁵. Robustness
437 of the clustering was assessed by calculating the average Adjusted Rand Index (ARI)
438 between iterations ($n = 100$), finding a mean ARI on Agta = 0.90 and a mean ARI on
439 ASV = 0.70. ARI measures the similarity of two partitions against a null hypothesis of

440 random assignment maintaining the size of the different clusters; the closer to 1 the
441 more robust is the classification⁷⁶. The resulting clusters were plotted with graph-tool⁷⁷.

442

443 **Ranking of bacteria associated with diet**

444 ASVs present in the Agta were ranked from -1 to 1 based on whether that ASV is more
445 present than expected in individuals from a given category: low proportion versus high
446 proportion of meals with only rice, and low proportion versus high proportion of meals
447 with meat, based on the median value of the population for each variable. Thus, a meat
448 associated score towards 1 indicates that an ASV is present more than expected in
449 individuals with high proportion of meals with meat (above the median of the
450 population), and a score towards -1 indicates that is present more in individuals with
451 low proportion of meals with meat (below the median of the population). A rice
452 associated score towards 1 indicates that an ASV is present more than expected in
453 individuals with high proportion of meals that consist of only rice, and a score towards -
454 1 indicates that is present more in individuals with low proportion of meals with only
455 rice.

456

457 **Microbiome genome-wide association studies**

458 To study the relationship between host genetics and the microbiome in the Agta, we
459 used a genome-wide association study (GWAS) approach to identify specific SNPs
460 associated with microbial abundance using GEMMA (version 0.94)⁷⁸. GWAS were
461 performed using the relative abundance of a given taxon in the Agta as a phenotype
462 trait, adding as covariates age, sex, and household (as a proxy for diet and shared
463 environment, as members of the same household share a hearth and their food on daily
464 bases). A kinship matrix calculated by KING using identical by descent segment

465 inference⁵² was included as random effects. For the GWAS analyses, we applied the
466 following quality control steps to the Agta genotypes. First, to detect ancestry outliers in
467 the dataset, we filtered the samples to keep only bi-allelic autosomal SNPs with Minor
468 Allele Frequency (MAF) > 5% and without missing data with PLINK 1.9⁷⁹. This dataset
469 was pruned for linkage disequilibrium (LD) using --indep-pairwise 50 5 0.2, and we
470 performed a principal component analysis (PCA) with EIGENSOFT (version 7.2.1)⁸⁰ to
471 identify ancestry outliers and exclude them from the analysis (Supplementary Figure 9).
472 Second, per sample heterozygosity was calculated with PLINK and samples with
473 overall increased/decreased heterozygosity rates (± 3 s.d. from the mean of the
474 population) were removed. A total of 129 Agta samples passed microbiome and
475 genotype data quality controls and were included in the microbiome GWAS analyses.
476 The analyses were done at the genus taxonomic level and on the CMM to study the
477 effect of host genetics at different taxonomic levels. Depending on the taxon analyzed,
478 as we included only samples with non-zero abundance, and SNPs with MAF < 10% and
479 with more than 5% missing data were excluded, the number of individuals tested ranged
480 from 10 to 129 and the number of SNPs tested ranged from 270569 to 313198 markers.
481 When we performed the GWAS at the genus level, we only included in the analysis 92
482 genera that are present in at least 10 Agta individuals to exclude low prevalent genera.
483 P-values were adjusted for multiple testing by FDR, and SNP-taxa associations were
484 considered significant at q-value < 0.1 on the cases where the proportion of variance in
485 the bacterial abundance explained by the genotypes (PVE or “chip heritability”) was
486 non-zero. The proportion of variance in the phenotype (bacterial abundance) explained
487 by the genotypes tested (PVE or “chip heritability”) was estimated for each taxon and
488 was considered non-zero if the standard error measurements did not intersect zero. We
489 applied genomic control to correct for cryptic relatedness and population stratification

490 and minimize false positives induced by inflated association test statistics⁸¹. To do so,
491 we estimated the genomic inflation factor as the median value of the likelihood ratio test
492 (LRT) values divided by 0.456 (median of a $\chi^2(1)$ distribution) and recalculated the p-
493 values after dividing the LRT values by the genomic inflation factor⁸². Threshold of
494 significance was set at FDR 10%, and only genomic positions having at least three
495 samples for the major homozygous genotype and for the heterozygous genotype were
496 considered. SNPs were annotated with ANNOVAR⁸³ in GRCh37 (hg19) using
497 RefSeqGene and dbSNP 147. For the enrichment analyses, we extracted the genes
498 associated to all non-intergenic SNPs and classified the genes in the background set,
499 that consisted in all genes present in the Axiom Human origin array; and the set to test,
500 that consisted of all genes that had a non-intergenic SNP significantly associated with
501 an ASV or a genus. We performed a gene ontology enrichment analysis with
502 ViSEAGO⁸⁴ and TopGo⁸⁵ R packages (Fisher exact test) and FUMA
503 GENE2FUNCTION module (Functional Mapping and Annotation of Genome-Wide
504 Association Studies)⁸⁶ to perform pathway enrichment analysis (hypergeometric test)
505 with FDR 5%.

506

507 **Selection analyses**

508 To test whether the GWAS SNPs showed any signal of recent positive selection we
509 performed a genome-wide scan of selection. We phased the Agta and Palanan farmer
510 populations independently. For each population, samples that were identified as
511 ancestry outliers by a PCA or with overall increased/decreased heterozygosity rates (± 3
512 s.d. from the mean of the population) were excluded from phasing. A total of 138 Agta
513 samples were phased using SHAPEIT2 version v2 (r900)⁸⁷ with the duoHMM method
514 to improve the phasing by integrating the known pedigree information. SNPs with

515 missing data were removed and window size was set to 5Mb for phasing. Due to the
516 small sample size, to phase the 14 unrelated Palanan farmers we used SHAPEIT2 with
517 default parameters and the 1,000 Genomes Phase 3 panel of haplotypes⁸⁸ as a reference
518 dataset. SNPs with missing data were removed. For the selection analyses, we excluded
519 one of each pair of related individuals by removing the sample in the pair with the
520 lowest call rate in the Agta phased dataset. This resulted in 38 unrelated Agta
521 individuals and 14 unrelated Palanan farmers. We ran the Integrated Haplotype Score
522 (iHS)⁸⁹ in the Agta phased dataset and the Cross-population Extended Haplotype
523 Homozygosity (XP-EHH) test⁹⁰ comparing the Agta against Palanan farmers as
524 implemented in selscan version v1.3.0⁹¹ to identify signals of positive selection in
525 GWAS SNPs. Both tests were run with default parameters and with the genetic map
526 provided by the 1,000 Genomes Phase 3⁸⁸. To identify regions under selection, for each
527 test we selected markers with scores in the 95th percentile that had at least 3 markers in
528 the 99th percentile in the surrounding area (± 10 Kb). For iHS we used absolute values,
529 while only positive scores were analyzed for XP-EHH.

530

531 **Acknowledgements**

532 A.B.M was funded by Leverhulme Trust Grant RP2011-R 045 and PLP-2017-323. J.B.
533 received grant PID2019-110933GB-I00/AEI/10.13039/501100011033 from the
534 Agencia Estatal Investigación (AEI), Spain, grant GRC 2017 SGR 702 from Secretaria
535 d'Universitats i Recerca del Departament d'Economia i Coneixement de la Generalitat
536 de Catalunya, as well as grant CEX2018-000792-M, part of the “Unidad de Excelencia
537 María de Maeztu” funded by the AEI, which granted the microbiome analysis by the
538 Servei de Genòmica at Universitat Pompeu Fabra. A.M. is funded by grant RTI2018-

539 102032-B-100 from AEI (Spain). A.E.P. is funded by grant MR/P014216/1 from the
540 Medical Research Council.

541

542 **Data availability**

543 16S amplicon data (EGAS00001005317) are deposited at the European Genome-
544 phenome Archive (EGA), which is hosted at the EBI and the CRG. Genome data
545 generated in this study has been deposited at EGA under accession number
546 EGAS00001005315. Data at the individual level on age, household composition and
547 diet that support the findings of this study are available on request from the
548 corresponding authors (JB and ABM). The individual data are not publicly available
549 due to them containing information that could compromise research participant privacy.

550

551 **Code availability**

552 Source code and data for visualization are available at

553 <https://doi.org/10.5281/zenodo.6342212>

554

555

556 **References**

- 557 1. Cornejo Ulloa, P., van der Veen, M. H. & Krom, B. P. Review: modulation of the
558 oral microbiome by the host to promote ecological balance. *Odontology* **107**,
559 437–448 (2019).
- 560 2. Weyrich, L. S. The evolutionary history of the human oral microbiota and its
561 implications for modern health. *Periodontol. 2000* **85**, 90–100 (2021).
- 562 3. Gomez, A. *et al.* Host Genetic Control of the Oral Microbiome in Health and
563 Disease. *Cell Host Microbe* **22**, 269-278.e3 (2017).
- 564 4. Willis, J. R. *et al.* Citizen science charts two major “stomatotypes” in the oral
565 microbiome of adolescents and reveals links with habits and drinking water
566 composition. *Microbiome* **6**, 218 (2018).

- 567 5. Costello, E. K. *et al.* Bacterial Community Variation in Human Body Habitats
568 Across Space and Time. *Science (80-.)*. **326**, 1694–1697 (2009).
- 569 6. Velsko, I. M. *et al.* Microbial differences between dental plaque and historic
570 dental calculus are related to oral biofilm maturation stage. *Microbiome* **7**, 102
571 (2019).
- 572 7. Page, A. E., Minter, T., Viguier, S. & Migliano, A. B. Hunter-gatherer health and
573 development policy: How the promotion of sedentism worsens the Agta’s health
574 outcomes. *Soc. Sci. Med.* **197**, 39–48 (2018).
- 575 8. Minter, T. The Agta of the Northern Sierra Madre: Livelihood Strategies and
576 Resilience Among Philippine Hunter-gatherers. (Faculty of Social Sciences,
577 Leiden University, 2010).
- 578 9. Page, A. E. *et al.* Reproductive trade-offs in extant hunter-gatherers suggest
579 adaptive mechanism for the Neolithic expansion. *Proc. Natl. Acad. Sci.* **113**,
580 4694–4699 (2016).
- 581 10. Dyble, M., Thorley, J., Page, A. E., Smith, D. & Migliano, A. B. Engagement in
582 agricultural work is associated with reduced leisure time among Agta hunter-
583 gatherers. *Nat. Hum. Behav.* **3**, 792–796 (2019).
- 584 11. Callahan, B. J., McMurdie, P. J. & Holmes, S. P. Exact sequence variants should
585 replace operational taxonomic units in marker-gene data analysis. *ISME J.* **11**,
586 2639–2643 (2017).
- 587 12. Greenacre, M. *Compositional data analysis in practice*. (Chapman and
588 Hall/CRC, 2018).
- 589 13. Larremore, D. B., Clauset, A. & Jacobs, A. Z. Efficiently inferring community
590 structure in bipartite networks. *Phys. Rev. E* **90**, 012805 (2014).
- 591 14. Rosier, B. T., Moya-Gonzalez, E. M., Corell-Escuin, P. & Mira, A. Isolation
592 and Characterization of Nitrate-Reducing Bacteria as Potential Probiotics for
593 Oral and Systemic Health. *Front. Microbiol.* **11**, 2261 (2020).
- 594 15. Kassebaum, N. J. *et al.* Global Burden of Severe Periodontitis in 1990-2010. *J.*
595 *Dent. Res.* **93**, 1045–1053 (2014).
- 596 16. De Maeyer, R. P. H. & Chambers, E. S. The impact of ageing on monocytes and
597 macrophages. *Immunol. Lett.* **230**, 1–10 (2021).
- 598 17. Dyble, M. *et al.* Networks of Food Sharing Reveal the Functional Significance of
599 Multilevel Sociality in Two Hunter-Gatherer Groups. *Curr. Biol.* **26**, 2017–2021
600 (2016).
- 601 18. Migliano, A. B. *et al.* Characterization of hunter-gatherer networks and
602 implications for cumulative culture. *Nat. Hum. Behav.* **1**, 0043 (2017).
- 603 19. Duan, X. *et al.* Smoking May Lead to Marginal Bone Loss Around Non-
604 Submerged Implants During Bone Healing by Altering Salivary Microbiome: A
605 Prospective Study. *J. Periodontol.* **88**, 1297–1308 (2017).
- 606 20. Eisenhofer, R. *et al.* Contamination in Low Microbial Biomass Microbiome
607 Studies: Issues and Recommendations. *Trends Microbiol.* **27**, 105–117 (2019).
- 608 21. Liu, S., Ying, G.-G., Liu, Y.-S., Peng, F.-Q. & He, L.-Y. Degradation of
609 Norgestrel by Bacteria from Activated Sludge: Comparison to Progesterone.

- 610 *Environ. Sci. Technol.* **47**, 130829113920003 (2013).
- 611 22. Shiau, H. J. & Reynolds, M. A. Sex Differences in Destructive Periodontal
612 Disease: A Systematic Review. *J. Periodontol.* **81**, 1379–1389 (2010).
- 613 23. Eke, P. I. *et al.* Update on Prevalence of Periodontitis in Adults in the United
614 States: NHANES 2009 to 2012. *J. Periodontol.* **86**, 611–622 (2015).
- 615 24. Ferraro, M. & Vieira, A. R. Explaining Gender Differences in Caries: A
616 Multifactorial Approach to a Multifactorial Disease. *Int. J. Dent.* **2010**, 1–5
617 (2010).
- 618 25. Zaura, E. *et al.* On the ecosystemic network of saliva in healthy young adults.
619 *ISME J.* **11**, 1218–1231 (2017).
- 620 26. Fernández, L., Pannaraj, P. S., Rautava, S. & Rodríguez, J. M. The Microbiota of
621 the Human Mammary Ecosystem. *Front. Cell. Infect. Microbiol.* **10**, 689 (2020).
- 622 27. David, L. A. *et al.* Diet rapidly and reproducibly alters the human gut
623 microbiome. *Nature* **505**, 559–563 (2014).
- 624 28. Turnbaugh, P. J. *et al.* The Effect of Diet on the Human Gut Microbiome: A
625 Metagenomic Analysis in Humanized Gnotobiotic Mice. *Sci. Transl. Med.* **1**,
626 6ra14–6ra14 (2009).
- 627 29. Schnorr, S. L. *et al.* Gut microbiome of the Hadza hunter-gatherers. *Nat.*
628 *Commun.* **5**, 3654 (2014).
- 629 30. Smits, S. A. *et al.* Seasonal cycling in the gut microbiome of the Hadza hunter-
630 gatherers of Tanzania. *Science (80-.).* **357**, 802–806 (2017).
- 631 31. Belstrøm, D. *et al.* Bacterial profiles of saliva in relation to diet, lifestyle factors,
632 and socioeconomic status. *J. Oral Microbiol.* **6**, 23609 (2014).
- 633 32. De Filippis, F. *et al.* The Same Microbiota and a Potentially Discriminant
634 Metabolome in the Saliva of Omnivore, Ovo-Lacto-Vegetarian and Vegan
635 Individuals. *PLoS One* **9**, e112373 (2014).
- 636 33. Fives-Taylor, P. M., Meyer, D. H., Mintz, K. P. & Brissette, C. Virulence factors
637 of *Actinobacillus actinomycetemcomitans*. *Periodontol. 2000* **20**, 136–67 (1999).
- 638 34. Adler, C. J. *et al.* Sequencing ancient calcified dental plaque shows changes in
639 oral microbiota with dietary shifts of the Neolithic and Industrial revolutions.
640 *Nat. Genet.* **45**, 450–455 (2013).
- 641 35. Sabbatani, S. & Fiorino, S. Dental worm disease. *Le Infez. Med.* **24**, 349–358
642 (2016).
- 643 36. Kolde, R. *et al.* Host genetic variation and its microbiome interactions within the
644 Human Microbiome Project. *Genome Med.* **10**, 6 (2018).
- 645 37. Blekhman, R. *et al.* Host genetic variation impacts microbiome composition
646 across human body sites. *Genome Biol.* **16**, 191 (2015).
- 647 38. Shaw, L. *et al.* The Human Salivary Microbiome Is Shaped by Shared
648 Environment Rather than Genetics: Evidence from a Large Family of Closely
649 Related Individuals. *MBio* **8**, e01237-17 (2017).
- 650 39. Mira, A., Artacho, A., Camelo-Castillo, A., Garcia-Esteban, S. & Simon-Soro, A.
651 Salivary Immune and Metabolic Marker Analysis (SIMMA): A Diagnostic Test

- 652 to Predict Caries Risk. *Diagnostics* **7**, 38 (2017).
- 653 40. Preshaw, P. M., Henne, K., Taylor, J. J., Valentine, R. A. & Conrads, G. Age-
654 related changes in immune function (immune senescence) in caries and
655 periodontal diseases: a systematic review. *J. Clin. Periodontol.* **44**, S153–S177
656 (2017).
- 657 41. Mukherjee, C. *et al.* Acquisition of oral microbiota is driven by environment, not
658 host genetics. *Microbiome* **9**, 54 (2021).
- 659 42. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome
660 data science using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).
- 661 43. Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina
662 amplicon data. *Nat. Methods* **13**, 581–583 (2016).
- 663 44. Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
664 processing and web-based tools. *Nucleic Acids Res.* **41**, D590–D596 (2012).
- 665 45. Katoh, K. MAFFT: a novel method for rapid multiple sequence alignment based
666 on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002).
- 667 46. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-
668 Likelihood Trees for Large Alignments. *PLoS One* **5**, e9490 (2010).
- 669 47. Oksanen, J. *et al.* Vegan: community ecology package. Ordination methods,
670 diversity analysis and other functions for community and vegetation ecologists. *R*
671 *Packag. version 2.5-7* (2020).
- 672 48. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible
673 Interactive Analysis and Graphics of Microbiome Census Data. *PLoS One* **8**,
674 e61217 (2013).
- 675 49. Faith, D. P. Conservation evaluation and phylogenetic diversity. *Biol. Conserv.*
676 **61**, 1–10 (1992).
- 677 50. Kembel, S. W. *et al.* Picante: R tools for integrating phylogenies and ecology.
678 *Bioinformatics* **26**, 1463–1464 (2010).
- 679 51. R Core Team. R: A language and environment for statistical computing. (2020).
- 680 52. Manichaikul, A. *et al.* Robust relationship inference in genome-wide association
681 studies. *Bioinformatics* **26**, 2867–2873 (2010).
- 682 53. Diekmann, Y. *et al.* Accurate age estimation in small-scale societies. *Proc. Natl.*
683 *Acad. Sci.* **114**, 8205–8210 (2017).
- 684 54. Pérez-Chaparro, P. J. *et al.* Newly Identified Pathogens Associated with
685 Periodontitis. *J. Dent. Res.* **93**, 846–858 (2014).
- 686 55. Socransky, S. S., Haffajee, A. D., Cugini, M. A., Smith, C. & Kent, R. L.
687 Microbial complexes in subgingival plaque. *J. Clin. Periodontol.* **25**, 134–144
688 (1998).
- 689 56. Camelo-Castillo, A. *et al.* Relationship between periodontitis-associated
690 subgingival microbiota and clinical inflammation by 16S pyrosequencing. *J.*
691 *Clin. Periodontol.* **42**, 1074–1082 (2015).
- 692 57. Khemwong, T. *et al.* Fretibacterium sp. human oral taxon 360 is a novel
693 biomarker for periodontitis screening in the Japanese population. *PLoS One* **14**,

- 694 e0218266 (2019).
- 695 58. Simón-Soro, A., Guillen-Navarro, M. & Mira, A. Metatranscriptomics reveals
696 overall active bacterial composition in caries lesions. *J. Oral Microbiol.* **6**, 25443
697 (2014).
- 698 59. Simón-Soro, A. & Mira, A. Solving the etiology of dental caries. *Trends*
699 *Microbiol.* **23**, 76–82 (2015).
- 700 60. Tanner, A. C. R. Anaerobic culture to detect periodontal and caries pathogens. *J.*
701 *Oral Biosci.* **57**, 18–26 (2015).
- 702 61. Kressirer, C. A. *et al.* *Scardovia wiggsiae* and its potential role as a caries
703 pathogen. *J. Oral Biosci.* **59**, 135–141 (2017).
- 704 62. Wolff, D. *et al.* Amplicon-based microbiome study highlights the loss of
705 diversity and the establishment of a set of species in patients with dentin caries.
706 *PLoS One* **14**, e0219714 (2019).
- 707 63. Mark Welch, J. L., Rossetti, B. J., Rieken, C. W., Dewhirst, F. E. & Borisy, G. G.
708 Biogeography of a human oral microbiome at the micron scale. *Proc. Natl. Acad.*
709 *Sci.* **113**, E791–E800 (2016).
- 710 64. Ferrer, M. D. & Mira, A. Oral Biofilm Architecture at the Microbial Scale.
711 *Trends Microbiol.* **24**, 246–248 (2016).
- 712 65. Leung, A. K. C., Wong, A. H. C. & Hon, K. L. Community-Acquired Pneumonia
713 in Children. *Recent Pat. Inflamm. Allergy Drug Discov.* **12**, 136–144 (2018).
- 714 66. Bellussi, L. M. *et al.* An overview on upper respiratory tract infections and
715 bacteriotherapy as innovative therapeutic strategy. *Eur. Rev. Med. Pharmacol.*
716 *Sci.* **23**, 27–38 (2019).
- 717 67. Natsis, N. E. & Cohen, P. R. Coagulase-Negative Staphylococcus Skin and Soft
718 Tissue Infections. *Am. J. Clin. Dermatol.* **19**, 671–677 (2018).
- 719 68. Lanao, A. E. & Pearson-Shaver, A. L. Mycoplasma infections. *StatPearls*
720 *[Internet]* (2020).
- 721 69. Jung, H.-S., Ehlers, M. M., Lombaard, H., Redelinguys, M. J. & Kock, M. M.
722 Etiology of bacterial vaginosis and polymicrobial biofilm formation. *Crit. Rev.*
723 *Microbiol.* **43**, 651–667 (2017).
- 724 70. Escapa, I. F. *et al.* New Insights into Human Nostril Microbiome from the
725 Expanded Human Oral Microbiome Database (eHOMD): a Resource for the
726 Microbiome of the Human Aerodigestive Tract. *mSystems* **3**, e00187-18 (2018).
- 727 71. Gloor, G. B., Macklaim, J. M., Pawlowsky-Glahn, V. & Egozcue, J. J.
728 Microbiome Datasets Are Compositional: And This Is Not Optional. *Front.*
729 *Microbiol.* **8**, 2224 (2017).
- 730 72. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical
731 and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B* **57**, 289–
732 300 (1995).
- 733 73. Storey, J. D. A direct approach to false discovery rates. *J. R. Stat. Soc. Ser. B*
734 *(Statistical Methodol.* **64**, 479–498 (2002).
- 735 74. Fortunato, S. Community detection in graphs. *Phys. Rep.* **486**, 75–174 (2010).

- 736 75. Peixoto, T. P. Nonparametric Bayesian inference of the microcanonical stochastic
737 block model. *Phys. Rev. E* **95**, 012317 (2017).
- 738 76. Hubert, L. & Arabie, P. Comparing partitions. *J. Classif.* **2**, 193–218 (1985).
- 739 77. Tiago, P. P. The graph-tool python library. *figshare* (2014).
740 doi:10.6084/m9.figshare.1164194
- 741 78. Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for
742 association studies. *Nat. Genet.* **44**, 821–824 (2012).
- 743 79. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger
744 and richer datasets. *Gigascience* **4**, 7 (2015).
- 745 80. Patterson, N., Price, A. L. & Reich, D. Population Structure and Eigenanalysis.
746 *PLoS Genet.* **2**, e190 (2006).
- 747 81. Devlin, B. & Roeder, K. Genomic Control for Association Studies. *Biometrics*
748 **55**, 997–1004 (1999).
- 749 82. Bacanu, S.-A., Devlin, B. & Roeder, K. The Power of Genomic Control. *Am. J.*
750 *Hum. Genet.* **66**, 1933–1944 (2000).
- 751 83. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of
752 genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38**,
753 e164–e164 (2010).
- 754 84. Brionne, A., Juanchich, A. & Hennequet-Antier, C. ViSEAGO: a Bioconductor
755 package for clustering biological functions using Gene Ontology and semantic
756 similarity. *BioData Min.* **12**, 16 (2019).
- 757 85. Alexa, A. & Rahnenfuhrer, J. topGO: Enrichment Analysis for Gene Ontology.
758 (2019).
- 759 86. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional
760 mapping and annotation of genetic associations with FUMA. *Nat. Commun.* **8**,
761 1826 (2017).
- 762 87. O’Connell, J. *et al.* A General Approach for Haplotype Phasing across the Full
763 Spectrum of Relatedness. *PLoS Genet.* **10**, e1004234 (2014).
- 764 88. Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68–
765 74 (2015).
- 766 89. Voight, B. F., Kudaravalli, S., Wen, X. & Pritchard, J. K. A Map of Recent
767 Positive Selection in the Human Genome. *PLoS Biol.* **4**, e72 (2006).
- 768 90. Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive
769 selection in human populations. *Nature* **449**, 913–918 (2007).
- 770 91. Szpiech, Z. A. & Hernandez, R. D. selscan: An Efficient Multithreaded Program
771 to Perform EHH-Based Scans for Positive Selection. *Mol. Biol. Evol.* **31**, 2824–
772 2827 (2014).
773