1

**Global analysis of infant gut microbiota revealed distinctive maturation dynamics across lifestyles**

Robert Bücking[1,2,3,4], Greta Pasquali[1,2,3,5], Ulrike Löber[1,2,3], Claudia Buss[6,7,8,9], Dorothee Viemann[4,10,11], Víctor Hugo Jarquín-Díaz[1,2,3]*, Sofia Kirke Forslund-Startceva[1,2,3,12,13]*

*Corresponding Author(s)
Sofia Kirke Forslund-Startceva & Víctor Hugo Jarquín-Díaz

[1] Max-Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin, Germany
[2] Charité–Universitätsmedizin Berlin, a corporate member of Freie Universität Berlin and Humboldt-Universität zu Berlin, Berlin, Germany
[3] Experimental and Clinical Research Center (ECRC), a cooperation of the Max-Delbrück Center and Charité–Universitätsmedizin, Berlin, Germany
[4] Translational Pediatrics, Department of Pediatrics, University Hospital Würzburg, Würzburg, Germany
[5] University of Padua, Department of Medicine DIMED, Padua, Italy
[6] Charité–Universitätsmedizin Berlin, Corporate Member of Freie Universität Berlin, Humboldt-Universität zu Berlin, and Berlin Institute of Health (BIH), Institute of Medical Psychology, Berlin, Germany
[7] University of California, Irvine, Development, Health and Disease Research Program, Irvine, CA, USA
[8] German Center for Child and Adolescent Health (DZKJ), partner site Berlin, Charité - Universitätsmedizin Berlin, Germany
[9] German Center for Mental Health (DZPG), partner site Berlin, Charité - Universitätsmedizin Berlin, Germany
[10] Center for Infection Research, University Würzburg, Würzburg, Germany
[11] Cluster of Excellence RESIST (EXC 2155), Hannover Medical School, Hannover, Germany
[12] German Center for Cardiovascular Research (DZHK), partner site Berlin
[13] Structural and Computational Biology Unit (SCB), EMBL Heidelberg, Germany

**Abstract**

The infant gut microbiome develops during the first years of life and influences long-term health through its interaction with immune system development. However, our understanding of early-life microbiome assembly is biased by the predominance of infants with industrialized lifestyles from North America and Europe. Here, we address this bias by assembling a globally representative dataset of infant gut microbiomes to train a microbiome maturation model that can characterize lifestyle specific patterns of microbial maturation as a function of age. Models trained exclusively on industrialized infants perform poorly when applied to non-industrialized datasets. In contrast, more diverse models including individuals from both lifestyles achieve increased correlation between microbial and chronological age. We identified differences in relevant taxa associated with the maturation in the different lifestyles. Additionally, our modeling approach detects a delay in the microbial maturation of independent cohorts of severely malnourished and preterm infants compared to healthy ones. Our results underscore the relevance of global diversity in microbiome research and provide deeper insights into context-dependent maturation dynamics of the infant gut microbiome.

**Main**

2

3

49  Early-life maturation of the gut microbiome is a key determinant of human health, influencing
50  immune system development, metabolic adjustments, and susceptibility to diseases
51  throughout life (Bisgaard et al. 2011; Willers et al. 2020; Tamburini et al. 2016). However,
52  our understanding of this process is primarily based on infants in industrialized populations,
53  which represent only a small fraction of global diversity (Abdill et al. 2022). Infants living in
54  non-industrialized or traditional communities remain largely underrepresented in microbiome
55  research, leaving a major gap in understanding how different environments and living
56  conditions shape early-life microbial community assembly.

57  Immediately after birth, the infant's gut is rapidly colonized by a variety of bacteria, and it
58  progressively develops until it resembles the adult microbiome by around three years of age
59  (Yatsunenko et al. 2012). The microbiome maturation is influenced by the simultaneous
60  development of the immune system. However, other factors including delivery mode,
61  gestational age, dietary shifts (Chu et al. 2017; Hill et al. 2017), exposure to antibiotics
62  (Bokulich et al. 2016) and lifestyle (Morandini et al. 2023) continue to shape the microbiome
63  during all stages of life. Models of microbiome maturation have shown that deviations in the
64  early microbiome development are associated with malnutrition (Subramanian et al. 2014)
65  and health outcomes later in life, such as asthma and allergies (Hoskinson et al. 2023) and
66  provide a powerful tool to detect generalizable microbial patterns across cohorts (Fahur
67  Bottino et al. 2025). Thus, the maturation of the early-life microbiome is a compelling model
68  system for investigating ecological succession and health-related microbial dynamics.

69  Despite the extensive research on the early-life microbiome, most studies remain
70  disproportionately focused on populations from North America and Europe, resulting in a
71  significant geographical bias (Abdill, Adamowicz, and Blekhman 2022). Consequently, the
72  diversity in host genetics, ethnicity, and particularly lifestyles is limited in most of the studies,
73  although those factors are known to significantly impact the adult human microbiome
74  (Clemente et al. 2015; Blekhman et al. 2015; Brooks et al. 2018; Yatsunenko et al. 2012,
75  Morandini et al. 2023). The industrialized lifestyle in those geographical regions involves
76  increased hygiene and exposure to antibiotics, reduced contact with wildlife and a dietary
77  shift toward more processed, high-caloric foods. Collectively, these factors lead to a reduced
78  microbial diversity and altered community structures in the adult human microbiome
79  (O'Keefe et al. 2015; Suez et al. 2014; Martínez et al. 2015; Almeida et al. 2019; Nayfach et
80  al. 2019; Pasolli et al. 2019). Thus, the focus on industrialized populations skews our
81  understanding of global microbiome dynamics, particularly during the critical stages of
82  microbiome maturation in infancy.

83  To address the gap driven by single cohorts and homogeneous populations in microbiome
84  maturation studies, we performed a globally representative analysis of infant gut microbiome
85  maturation across diverse populations and lifestyles. Our study integrates and analyzes
86  previously publicly available datasets using machine learning models. We specifically
87  compare the microbiome maturation trajectories between infants from industrialized and
88  various non-industrialized populations. Thereby, we contribute to a broader understanding of
89  early-life microbial colonization and its implications for human health.

90

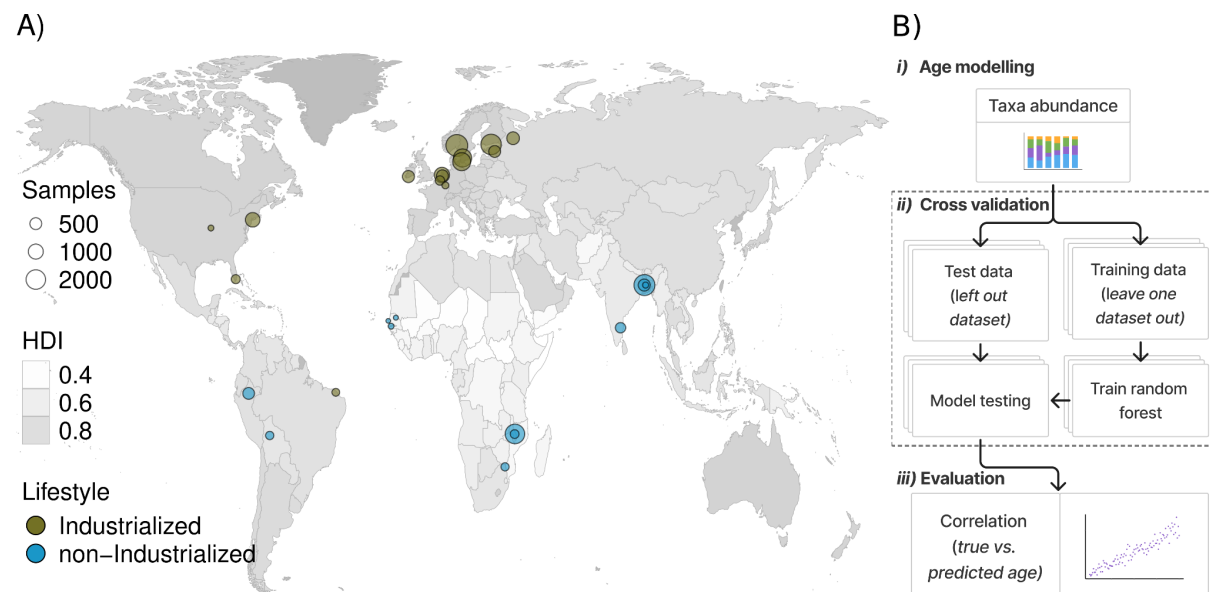91  **Materials and methods**

92  **Data collection**
93  We collected publicly available 16S rRNA gene sequencing data from the stool microbiome
94  of human infants. The inclusion of the samples in our meta-analysis was based on the
95  following characteristics: a) Samples from full-term, healthy infants under two years of
96  chronological age, not subject to any intervention b) The sequencing data was generated

4

97   using Illumina platforms c) The targeted amplicon included the V4 region of the 16S rRNA
98   gene d) Metadata with information on the age of the infant at sampling time was available.

99   Raw sequencing data was retrieved from the Sequence Read Archive (SRA), and metadata
100  was obtained either from the original publications or through direct communication with the
101  authors when necessary. The final dataset consisted of 15,077 samples from 20 studies
102  across 20 countries inAfrica, America, Asia and Europe, corresponding to 2,720 individuals
103  (Figure 1A, Supplementary Table S1). We categorized samples into two lifestyle groups:
104  Industrialized and non-industrialized, using the Human Development Index 2022 (HDI)
105  (United Nations Development Programme 2022) (Figure 1A). Samples from countries with
106  an HDI-value above the median of 0.742 were classified into the industrialized lifestyle, while
107  all other samples were classified into non-industrialized lifestyles. Samples from Peru (HDI =
108  0.762) were treated as an exception and classified as non-industrialized, as they were taken
109  from individuals living in a remote area of the Amazon rainforest with living conditions more
110  comparable to other non-industrialized individuals than to industrialized ones (Raman et al.
111  2019). The term industrialized lifestyle serves here as an umbrella term for complex,
112  multidimensional lifestyle changes compared to pre-industrial societies, as previously
113  discussed (Pasolli et al. 2019). These changes include improved hygiene and sanitized
114  environments, increased access to healthcare, and higher exposure to antibiotics and other
115  drugs, reduced contact with wildlife, and a dietary shift toward more processed, high-caloric
116  foods. These factors are known to have a huge impact on the human gut microbiome.

117  Samples classified as industrialized in this study are expected to be more affected by these
118  changes than those classified as non-industrialized. However, it is important to note that the
119  term non-industrialized here does not refer to a single lifestyle, but rather to a collection of
120  different lifestyles.
121
122



123
124  **Figure 1. Collection and processing of publicly available datasets. A)** Global distribution
125  of samples. Each circle represents the samples from one study at the specific location. Size
126  indicates the number of samples and color the lifestyle of the samples. Shading indicates
127  Human Development Index (HDI). **B)** Computational pipeline for microbial age modelling.
128  The processing has three steps: *i)* Raw data preprocessing and taxonomy annotation, *ii)*
129  Random forest regression models with leave-one-dataset-out cross-validation and *iii)*
130  Evaluation of predicted microbial age by correlation to chronological age.

**Processing of sequencing data**

The raw data was re-processed from quality check to taxonomic annotation with the same pipeline to ensure consistency across samples. The read quality was assessed with FastQC v0.11.9 ("FastQC" 2015). Where necessary, adapter contamination was removed from the sequences using Trimmomatic v0.39 (Bolger et al. 2014). As a host decontamination step, reads were mapped to the reference human genome assembly (GRCh38.p14) using bowtie2 v2.3.4.3 (Langmead and Salzberg 2012) with the --very-sensitive option. Reads mapping to the human reference were removed from further analysis. Quality checked and decontaminated reads were processed and merged to infer amplicon sequence variants (ASVs) using the pipeline from DADA2 v1.34.0 (Callahan et al. 2016) for each study independently. For study-specific details on trimming and filtering of reads, see Supplementary Table S1. Taxonomic annotation of the ASVs was done with the SILVA database v138 (Pruesse et al. 2007) using DECIPHER v3.2.0 (Wright 2016). ASVs assigned as "Mitochondria" were also discarded from further analysis. The dataset was aggregated to genus level. Based on rarefaction curves, study-specific thresholds were established to discard samples with low read counts (Supplementary Figure S1, Supplementary Table S1). The filter threshold was determined as the amount of reads where genus richness (number of observed genera) reached an asymptote. Relative abundances were calculated for the filtered datasets. Genera with a mean relative abundance below 0.005 % were removed from all analyses. The bioinformatic analysis was performed on the High Performance Computing Cluster from the Max Delbrück Centrum, Berlin (Max-Cluster).

**Alpha and beta diversity analyses**

All analyses were run in R v4.4.2. Alpha diversity indices were estimated for all samples with more than 2,000 reads after rarefication to that depth. For samples with less than 2,000 reads that passed the study-specific filter threshold, alpha diversity indices were estimated based on the raw read counts. To evaluate how the size of the dataset affects alpha diversity, each lifestyle group was randomly subsampled 10 times, increasing the number of studies by 1, individuals by 10 and samples by 100. For each subsample, the number of taxa with more than 10 rarefied reads in at least one sample was determined for each lifestyle group. Differences in the number of detected taxa between lifestyles over the number of studies, individuals and samples were tested with generalized additive models using gamlss v5.4-22(R. A. Rigby and D. M. Stasinopoulos 2005). The number of taxa was modeled as a function of lifestyle and a penalized spline smoothing term for the number of studies, individuals or samples was included. The full model was compared with a reduced model without lifestyle using a likelihood ratio test.

The effects of age and lifestyle on Shannon-diversity were tested with generalized additive models. The full model had Shannon diversity as response variable and lifestyle and age as explanatory variables, allowing for nonlinear relationships between age and diversity using penalized spline smoothing. Study was included as a random factor. The full model was compared to two reduced models using a likelihood ratio test, one without age and one without lifestyle, respectively. The proportion of variance explained by age and lifestyle was calculated as the relative decrease in deviance in the full model compared to the reduced model. The effect of individual studies on differences in Shannon diversity and the variance

179  explained by each variable was assessed by repeating the analysis and removing each
180  study from the dataset once. We compared the Shannon diversity of both groups using a 60-
181  day sliding window advancing in 7-day increments to determine the age intervals in which
182  microbial diversity differs significantly between the two lifestyles. We used Wilcoxon tests
183  with Bonferroni correction for multiple testing.
184
185  For beta diversity analysis, a principal component analysis was computed on centered-log-
186  ratio (clr) transformed raw read counts and on clr-transformed counts rarefied to 2,000 reads
187  per sample to assess the effect of differences in sequencing depth. The effects of age,
188  lifestyle and study on the composition were analyzed with a permutational analysis of
189  variance (PERMANOVA) as implemented in vegan v2.6-8. (Oksanen et al. 2024).
190  PERMANOVA was run with 999 permutations to model Euclidean distance of the clr-
191  transformed count matrix by adding the terms age, lifestyle and study sequentially.
192

**Machine learning modelling of age based on microbial composition**

194  To predict the age of the individuals based on their microbial composition, a supervised
195  machine learner was trained. Random forest regression models with 500 trees were trained
196  on a dataset including: 1) relative abundances of genera that were detected in at least five
197  samples in two studies in the training set, and 2) microbial diversity (Shannon index) and
198  richness. Random forest regression was implemented using the ranger R-package v0.17.0
199  (Wright and Ziegler 2017) with default parameters. Relative feature importance was
200  calculated for each feature as the increase in out-of-the-bag (OOB) error when the
201  respective feature was permuted and normalized to the highest importance value.
202  Significantly important features were selected using the Boruta R-package v8.0.0 (Kursa and
203  Rudnicki 2010) with permutation-based importance values and a p-value threshold of 0.01.
204  To estimate the performance of the microbial age model, a leave-one-dataset-out cross-
205  validation (LODO-CV) was implemented using the caret R-package v6.0-94. Thus, each
206  study was used once as a validation set to estimate the performance of a model trained on
207  the rest of studies. Due to the non-linear but monotonous increase in microbial age
208  alongside chronological age, microbial age was rank transformed for model performance
209  evaluation. The coefficient of determination ($R^2$) of a linear model between transformed
210  predicted microbial and the actual chronological age was used as a performance metric
211  during validation. To assess the overall effect of lifestyle on the predicted microbial age
212  obtained from the LODO-CV, a mixed model with the rank transformed predicted age as
213  response, chronological age and lifestyle as explanatory variables and study as random
214  factor was fit to the whole dataset. The full model was compared to a reduced model, without
215  lifestyle, using a likelihood ratio test (LRT). To assess the effect of individual studies on the
216  differences in predicted microbial age, the same analysis was repeated by removing each
217  study from the dataset once. To determine the time at which maturation reaches a more
218  stable state, a logistic growth model was fitted to the relationship between chronological and
219  predicted age for each lifestyle. We considered the time when the fitted model reached 90%
220  of the model's carrying capacity as the point at which maturation rate decreased and
221  reached a more stable state.
222

223  For lifestyle specific models, only subjects from the corresponding lifestyle were used in the
224  training set. The performance of lifestyle specific models was compared across studies from
225  both lifestyles using paired Wilcoxon tests with Bonferroni correction for multiple testing.

11

Lifestyle specific relevant features were defined if they were selected as relevant by Boruta in at least two models of LODO-CV for the specific lifestyle.

To assess the effect of differences in dataset size between the two lifestyles, a downsampling approach was implemented. Lifestyle specific and combined LODO-CV was performed 50 times for each study, with the training set downsampled to make the two lifestyles comparable in terms of the number of studies, individuals and samples, as well as their age distribution. First, studies and individuals in the industrialized group were randomly selected in equal numbers to those in the non-industrialized group. Then, samples above and below one year of age were randomly downsampled separately to match the size of the smaller lifestyle group.

To determine the influence of important features on the prediction of age by lifestyle, two models were trained on the complete set of samples corresponding to each lifestyle. Shapley additive explanation (SHAP) values were calculated for all samples on these two models using fastshap v0.1.1 (Greenwell 2024). Kolmogorov-Smirnof test with Bonferroni correction for multiple comparisons was used to determine significance and effect size from differences in taxon-specific SHAP-value distribution between lifestyles and between taxa within specific lifestyles. The Spearman correlation between SHAP-values and relative feature abundances for each taxon was used to assess its temporal dynamics. Unlike the direct correlation between taxon abundance and age, the correlation between SHAP-values and abundance is a more robust strategy. It is less sensitive to zero-inflated abundance distributions and nonlinear abundance-age relationships.

The longitudinal dynamics of taxa associated with lifestyle specific maturation were further investigated to detect differences between both lifestyles. Differences in trajectories of prevalence over time for those taxa were analyzed using linear models. Therefore, a full model was trained with age, lifestyle and their interaction as predictors for prevalence. The full model was tested against a reduced model without the interaction between age and lifestyle as a predictor using a likelihood ratio test as implemented in the lmtest R-package v0.9-40 (Zeileis and Hothorn 2002).

**Influence of clinical factors on maturation**
To evaluate the effect of lifestyle specific models on the prediction of microbial age under two clinical conditions, we assembled two supplementary datasets: one from infants born preterm (before 37 weeks of gestation) and another from infants diagnosed with severe acute malnutrition (SAM) (Supplementary Table S1). The raw sequences from these samples were processed as described above for the training set.

We predicted the microbial age of SAM and preterm infants employing the lifestyle-combined and lifestyle specific models. To compare differences in maturation, microbiome-for-age Z-scores (MAZ) were calculated from microbial age predictions for each model. Microbial age predictions were grouped into weekly chronological age bins. MAZ-scores were calculated by subtracting the group mean from the microbial age and dividing by the group standard deviation. Differences in MAZ-scores between healthy, preterm, and malnourished infants were tested using Wilcoxon rank-sum tests with Bonferroni correction for multiple comparisons.

To determine microbial drivers of age maturation under different clinical conditions, SHAP-values were calculated for preterm and healthy industrialized infants using the industrialized

12

274  model and SAM and healthy non-industrialized infants using the non-industrialized model.
275  Relative taxon impact was assessed by the differences in mean of SHAP-values between
276  healthy infants and infants of the same lifestyle with a clinical condition in separate age bins
277  using Wilcoxon tests. P-values were adjusted for multiple testing across all taxa and age
278  bins within each lifestyle group using Bonferroni correction.

279  **Results**

280  **Age and lifestyle drive the maturation of the gut microbiome**
281  To investigate the development of the human infant gut microbiome across different
282  lifestyles, we analyzed 11,255 samples from infants with an industrialized lifestyle and 3,873
283  samples from infants with non-industrialized lifestyles. We detected 890 genera, of which
284  only 61 were present in all studies regardless of the lifestyle. Among the remaining genera,
285  506 genera were detected in microbiomes from both lifestyles, while 365 genera were
286  exclusive to industrialized microbiomes and 19 were found only in non-industrialized
287  microbiomes (Supplementary Figure S2A). Of the lifestyle specific taxa, 97% had a
288  prevalence below 1% within the respective lifestyle. The relatively small number of taxa
289  specific to non-industrialized samples may result from differences in sample size between
290  the two lifestyles. Therefore, we created 1,000 downsampled versions of the dataset in
291  which the number of studies, individuals and samples above and below one year of age
292  were equal for both lifestyle groups. In the downsampled datasets, we still detected more
293  lifestyle specific genera in industrialized microbiomes compared to non-industrialized ones,
294  although the difference was less pronounced (mean = 80.35 and 71.02, respectively, paired
295  Wilcoxon test, p-value < 0.001) (Supplementary Figure S3A).
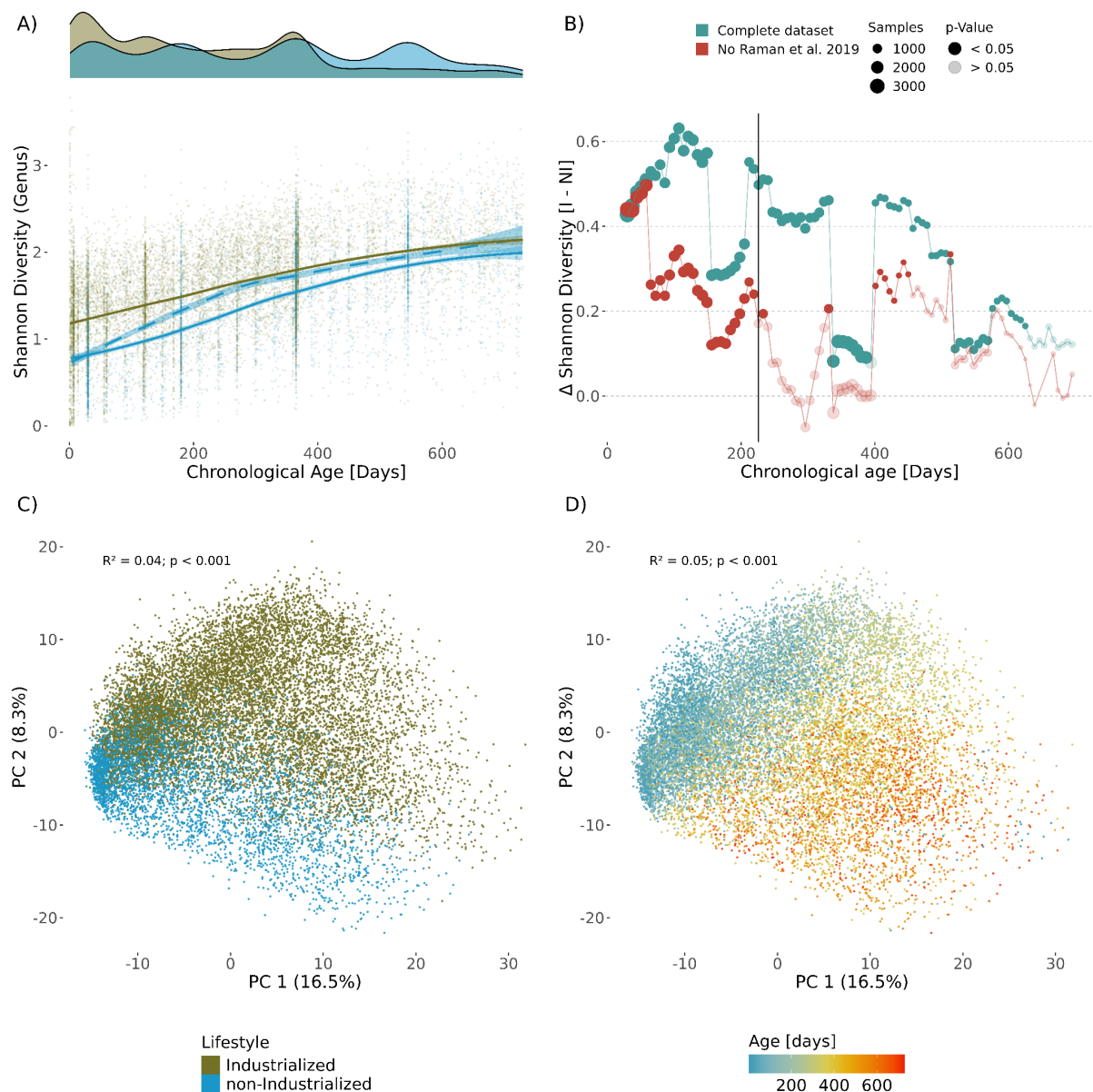296  We observed that industrialized samples had on average a two times higher sequencing
297  depth compared to non-industrialized samples (Supplementary Figure S2B). Therefore, we
298  rarefied to 2,000 reads per sample for alpha diversity analyses. We observed a positive
299  relationship between the genus richness and the number of studies, individuals, and
300  samples within each lifestyle, independently of the sequencing depth (Supplementary
301  Figures S3B, C and D). Overall, in any of the lifestyles, the number of detected genera
302  reached the saturation of the community, suggesting an unexplored microbial diversity in
303  infants from both industrialized and non-industrialized lifestyles. Genus richness was
304  significantly lower in non-industrialized than in industrialized samples, regardless of the
305  number of studies, individuals and samples included in both groups (Likelihood ratio-test, $\chi^2$
306  (DF) = 1.72, 2.32, 4.45, $p$-value<0.01, explained variance: 8.7%, 16.0%, 15.6% for the
307  number of studies, individuals and samples, respectively).
308
309  Microbial diversity (Shannon index, Figure 2A) and richness (Supplementary Figure S2C)
310  increased along the chronological age (Likelihood ratio-test, $\chi^2$ (DF) = 0.85, $p$-value<0.01,
311  explained variance: 18.8%), consistent with previous studies (Bokulich et al. 2016; Stewart
312  et al. 2018; Yatsunenko et al. 2012). While the increase in alpha diversity metrics was
313  independent of lifestyle, non-industrialized microbiomes showed reduced overall diversity
314  compared to industrialized ones (Likelihood ratio-test, $\chi^2$ (DF) = 0.80, $p$-value < 0.01,
315  explained variance: 0.1%). The non-industrialized samples had a significantly lower fraction
316  of ASVs unassigned at genus level compared to industrialized samples (Likelihood ratio-test,
317  $\chi^2$ (DF) = 0.49, $p$-value < 0.01, explained variance: 0.2%) (Supplementary Figure S2D),
318  indicating that bias in the annotation database is not the cause of a reduced alpha diversity
319  in non-industrialized samples. A study by Raman et al. from 2019 with 1,891 samples mainly

15

320  from Bangladesh drove this significant effect on alpha diversity between lifestyles. The
321  variance explained by the remaining non-industrialized studies was lower when this study
322  was not included (Likelihood ratio-test, $\chi^2$ (DF) = 0.88, $p$-value = 0.78, explained variance <
323  0.001%). However, alpha diversity was consistently reduced in non-industrialized individuals,
324  independently of the study, until around 200 days of chronological age (Figure 2B). While
325  this demonstrates the influence of a single study, it also highlights the importance of
326  gathering multiple studies to assess differences in alpha diversity between lifestyles that are
327  not observable in single cohorts with restricted age spans.
328
329  A total of 4% of the microbial community structure was determined by the lifestyle
330  (PERMANOVA, $p$<0.001), with a clear separation between non-industrialized and
331  industrialized samples along PC2 (Figure 2C, Supplementary Figure S2E). Among the most
332  influential taxa driving this separation are *Veilonella*, *Erysipelatoclostridium* and *Bacteroides*
333  in the industrialized direction and *Faecalibacterium* in the non-industrialized direction
334  (Supplementary Figure S2E). An additional 5% was explained by the age (PERMANOVA,
335  $p$<0.001), with the distinction between lifestyles becoming clearer at older chronological
336  ages, following parallel trajectories along PC1 over time (Figure 2D, Supplementary Figure
337  S2E). This trajectory was mainly driven by *Faecalibacterium*, *Bacteroides*,
338  *Lachnoclostridium, Intestinibacter* and *Ruminococcus* (Supplementary Figure S2E).
339  Although the study specific effects accounted for 7% of the residual variation
340  (Supplementary Figure S2F), there was a consistent influence of age and lifestyle. While the
341  average sequencing depth differed between both lifestyle groups, a PCA on rarefied counts
342  showed the same pattern in sample distribution. The first two principal components were
343  largely driven by the age and lifestyle, respectively, indicating a biological rather than a
344  technical effect driven by sequencing depth difference (Supplementary Figure S2G, H).
345

16

17



**Figure 2. Gut microbial diversity and community structure follows lifestyle specific maturation trajectories. A)** Shannon diversity calculated from rarefied counts aggregated to genus level, plotted over chronological age and stratified by lifestyle. Locally estimated scatterplot smoothing (LOESS) curves are fit to visualize temporal trends. The solid lines represent the independent increase in alpha diversity by lifestyle. The dashed line shows the increase in diversity for non-industrialized microbiomes, excluding one study with a large sample size (Raman et al., 2019). When all non-industrialized studies are included, diversity is reduced overall compared to industrialized studies throughout the entire chronological age span. However, when all samples from Raman et al., 2019 are removed, the difference in diversity is significant only during the first 200 days of life. **B)** Difference in Shannon diversity between industrialized (W) and non-industrialized (NW) populations in a 60 days sliding window. Red represents comparisons on the complete industrialized vs non-industrialized datasets and blue comparisons between complete industrialized and a reduced non-industrialized dataset without the study of Raman et al., 2019. The number of samples in each comparison is represented by size. Transparent points indicate no significant differences between lifestyles. **C)** Principal Component Analysis (PCA) on centered log ratio transformed counts colored by lifestyle and D) age in days of life. $R^2$ values indicate the

18

19

364    proportion of variance explained in a Permutational multivariate analysis of variance
365    (PERMANOVA) using a sequential model with the terms age, lifestyle and study.

366

367    **A model for microbiome maturation in human infants**
368    We developed a microbiome age index by training machine learning models that predict an
369    individual's age based on their microbial composition and diversity at the time of sampling to
370    identify microbial features that characterize the maturation process in different lifestyles. The
371    model's evaluation using LODO-CV demonstrated a strong correlation between predicted
372    microbial age and chronological age in both industrialized and non-industrialized samples
373    ($R^2$: 0.85 and 0.76, respectively).

374

375    Microbial age initially increased with chronological age, reaching a stable phase at 455 and
376    524 days for industrialized and non-industrialized lifestyles, respectively (90% of carrying
377    capacity, logistic growth model), following the pattern observed for microbial diversity. This
378    reflects a phase of rapid initial microbial colonization directly after birth, followed by a more
379    stable phase where the community reaches stability. Although the overall maturation
380    trajectory remained consistent across lifestyles, non-industrialized infants exhibited slightly
381    delayed maturation compared to industrialized infants (Likelihood ratio-test, $\chi^2$ (DF) = 1, *p*-
382    value < 0.01, explained variance = 0.01 %). Similarly, as for alpha diversity, the 1,891
383    samples from Raman et al., 2019 were highly influential on the delayed maturation effect,
384    further confirming the high variability among non-industrialized lifestyles and difficulties of
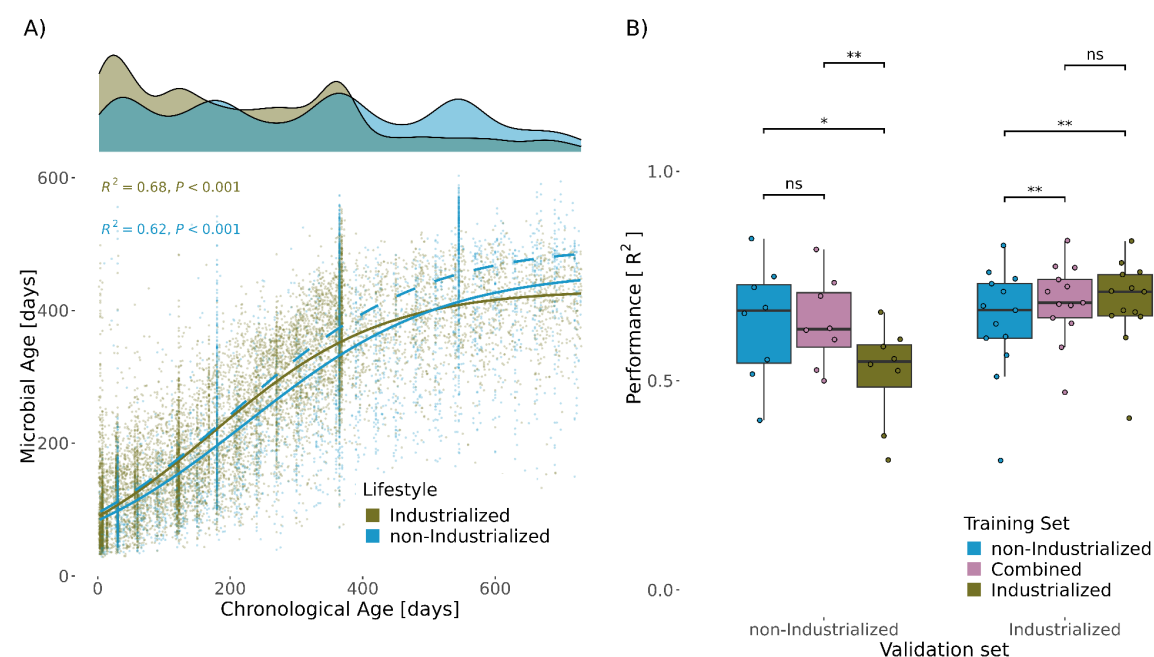385    generalization among studies with non-industrialized lifestyles.

386

387    To determine whether microbial age prediction differs between industrialized and non-
388    industrialized infants, we trained specific models on three types of datasets: 1) a combined
389    dataset including samples from both lifestyles, 2) a dataset with samples from the
390    industrialized lifestyle and 3) a dataset from a non-industrialized lifestyle and evaluated their
391    performance on lifestyle specific datasets as validation (Figure 3B). Models trained on a
392    combined dataset did not significantly outperform lifestyle specific models, when the lifestyle
393    of the validation and training set matched. However, models trained on a single lifestyle and
394    evaluated on a different lifestyle validation dataset showed a significant decrease in
395    performance. The performance in microbial age prediction of non-industrialized samples
396    using an industrialized model significantly declined (paired Wilcoxon test, Bonferroni
397    adjusted p-value < 0.1) compared to models that included non-industrialized samples in the
398    training set. For industrialized datasets, the performance in predicting the microbial age
399    showed a significant reduction when using models trained only on non-industrialized
400    samples (paired Wilcoxon test, Bonferroni adjusted p-value < 0.05).

401

402    To determine whether differences in sequencing depth had an effect on the observed
403    differences between the two lifestyles, we rarefied the samples to the same sequencing
404    depth and trained a new model on the rarefied dataset (Supplementary Figure S3E). The
405    significant differences in model performance between different lifestyles were reproducible.
406    While the combined model did not perform significantly better than the industrialized model
407    for non-industrialized data, the direction of the effect remained consistent. Studies including
408    non-industrialized individuals were underrepresented in our study. To assess the effect of
409    lower sample size on the performance of microbial age prediction of models trained on non-
410    industrialized samples, we repeated the analysis for each validation set on 50 downsampled
411    versions of the training set where both lifestyles had a similar number of studies, individuals,

20

21

412  and samples above and below one year of age. We found a reproducible difference in model
413  performance between lifestyles when the downsampled industrialized training sets were
414  used (Supplementary Figure S3G).

415

416  Non-industrialized samples represent a heterogeneous lifestyle group, defined by the
417  absence of industrialized lifestyle characteristics rather than sharing a similar lifestyle. Living
418  conditions include inhabitants of a slum in Bangladesh, as well as infants living in the
419  Amazon rainforest. Despite the expected variability among non-industrialized samples, our
420  results indicate that training models on a more diverse dataset improved their generalizability
421  for all non-industrialized datasets.

422



423
424  **Figure 3. Lifestyle determines microbial age modelling but not development stages. A)**
425  Microbial age development over chronological age. Microbial age was predicted using
426  random forest regression based on relative abundances of bacterial genera against
427  chronological age using a leave-one-study-out cross-validation (LODO-CV) approach. A
428  logistic growth curve was fit to visualize the nonlinear trajectory of the maturation dynamic
429  over chronological age. The solid lines represent the microbial age by lifestyle with the
430  complete datasets. The dashed line shows the microbial age in non-industrialized
431  individuals, excluding the study by Raman et al., 2019. **B)** Performance of microbial age
432  modelling with lifestyle specific models. Separate models were trained either on
433  microbiomes from industrialized samples, non-industrialized samples, or a combined
434  dataset. Performance was assessed using LODO-CV, where each study served as a
435  validation set once for each model. Coefficient of determination ($R^2$) from a linear model of
436  rank transformed microbial age versus chronological age was used as performance metrics
437  and was calculated for each combination of model and validation set separately. Differences
438  in model performances were tested with a paired Wilcoxon test and adjusted for multiple
439  comparisons using Bonferroni correction (*$p<0.1$, **$p<0.05$, ***$p<0.01$).
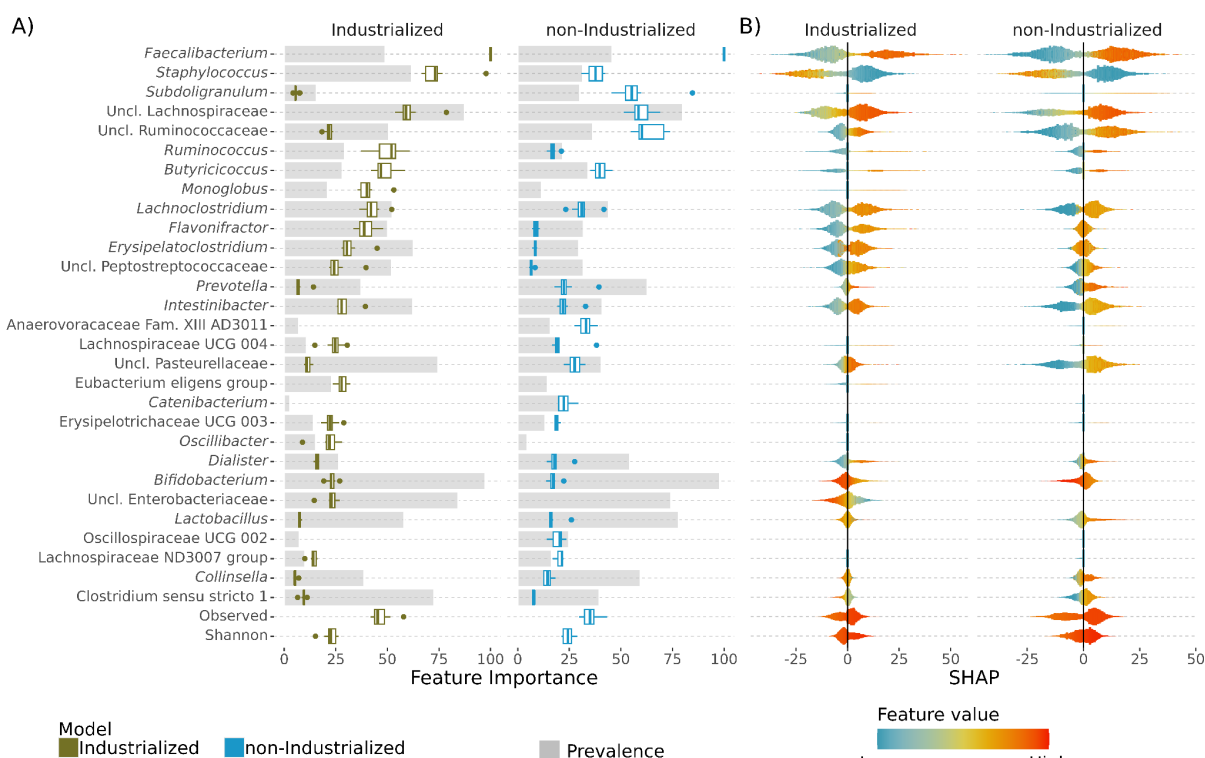
440

**Specific bacteria drive the microbiome maturation in different lifestyles**
442  Having identified consistent differences in microbial maturation between lifestyles, we aimed
443  to investigate which microbial taxa characterize these differences. Thus, we calculated the

22

444 feature importance scores for the random forest models (Figure 4A, Supplementary Figure
445 S4A, Supplementary Table S2). Overall, we identified 105 important taxa that account for a
446 mean combined relative abundance of 98 % across all samples. However, only 40 of those
447 taxa were important in both lifestyles and the majority were lifestyle specific. Among the
448 generalist taxa shared across the lifestyles, *Faecalibacterium* was the most important taxon
449 in both lifestyles. In contrast, *Ruminococcus* and *Staphylococcus* were more important in
450 industrialized microbiota than in non-industrialized, and *Subdoligranulum* and *Prevotella*
451 showed greater importance in non-industrialized models, suggesting them as lifestyle
452 specialists. Additionally, 65 taxa were important only for one lifestyle. *Monoglobus* and
453 *Flavonifractor* were examples of taxa highly relevant only in industrialized microbiota, while
454 *Oscillospiraceae UCG 002* was an exclusive specialist in non-industrialized individuals. All
455 taxa important in only one lifestyle specific model were still present in the other lifestyle,
456 although less prevalent (Supplementary Figure S5A). Feature importance was independent
457 of prevalence for the non-industrialized dataset. For the industrialized dataset, however,
458 feature importance was only correlated with prevalence for taxa with a prevalence below
459 0.15 (Supplementary Figure S5B).
460 More taxa were classified as important in industrialized than in non-industrialized models (95
461 and 50, respectively). Our selection of important taxa could have been affected by the
462 number of studies and thus number of models in the LODO-CV for each lifestyle. Therefore,
463 we compared the number of taxa identified as significantly important by Boruta in each
464 lifestyle specific model trained on the downsampled data sets described above.
465 Industrialized models still had a significantly higher number of important taxa on average
466 than non-industrialized models (103 and 89 respectively; Wilcoxon test: p < 0.001)
467 (Supplementary Figure S3F).
468



**Figure 4. Specific bacterial taxa distinguish the microbial maturation between lifestyles. A)** Relative feature importance of specific bacteria is different depending on the lifestyle and its relevance is independent from its prevalence. Relative feature importance in

25

473 random forest models trained separately on samples from individuals with an industrialized
474 or a non-industrialized lifestyle. Grey bars indicate the mean prevalence of a specific taxon
475 across studies in both lifestyles. Importances are displayed only for features significantly
476 important for the respective lifestyle. **B)** Shapley additive explanation (SHAP) values for
477 specific bacteria show consistent predictive impact of a given feature between the two
478 lifestyles. However, the abundance of *Lactobacillus* or *Collinsella* has the opposite effect on
479 the prediction of microbial age between lifestyles: A higher abundance increases the
480 predicted age in non-industrialized individuals, and decreases it in industrialized individuals.
481 Violin plots represent the distribution of SHAP-values per feature. Positive SHAP-values
482 indicate an increase in predicted microbial age driven by a given feature, whereas negative
483 SHAP-values indicate a decrease in the predictive effect. Color indicates relative taxon
484 abundance or alpha diversity, scaled per feature. Only values for features significantly
485 important for the respective lifestyle are shown.
486

487 **Lifestyle specific taxonomic drivers differ in their colonization patterns**
488 Once we identified specialist and generalist features for both lifestyles, we aimed to further
489 characterize the temporal dynamics of these features and define the changes in maturation
490 between the two lifestyles. Therefore, we performed SHAP analysis to assess the
491 contribution of each feature to the predicted microbial age (Shapley 1953). Positive SHAP-
492 values indicate that a feature in a given sample increases the predicted microbial age,
493 whereas negative SHAP-values indicate a decreasing effect. The three taxa with higher
494 relative feature importance values, *Staphylococcus*, *Faecalibacterium*, and one unclassified
495 Lachnospiraceae taxon, also showed the highest range of SHAP-values, underlining their
496 importance in our model (Figure 4B, Supplementary Figure 4B). In general, feature values
497 and their corresponding SHAP-values exhibited a positive monotonic relationship across
498 taxa, meaning that the increase in feature value corresponds to an increase in SHAP-values,
499 and vice versa. However, taxa like *Bifidobacterium* and *Staphylococcus* showed negative
500 relationships. While the relationship between feature values and SHAP-values was
501 consistent across lifestyles for most taxa, *Lactobacillus* and *Collinsella* showed opposite
502 relationships in each lifestyle group, indicating lifestyle specific ecological roles of those taxa.
503 All predictive taxa had a significant difference in SHAP-value distribution between lifestyles
504 ($p < 0.001$, Kolmogorov-Smirnov test, Bonferroni correction, Kolmogorov-Smirnov D = 0.26
505 between models, 0.35 within the industrialized model and 0.39 within the non-industrialized
506 model) (Supplementary Figure S5C), which further confirmed the difference in predictive
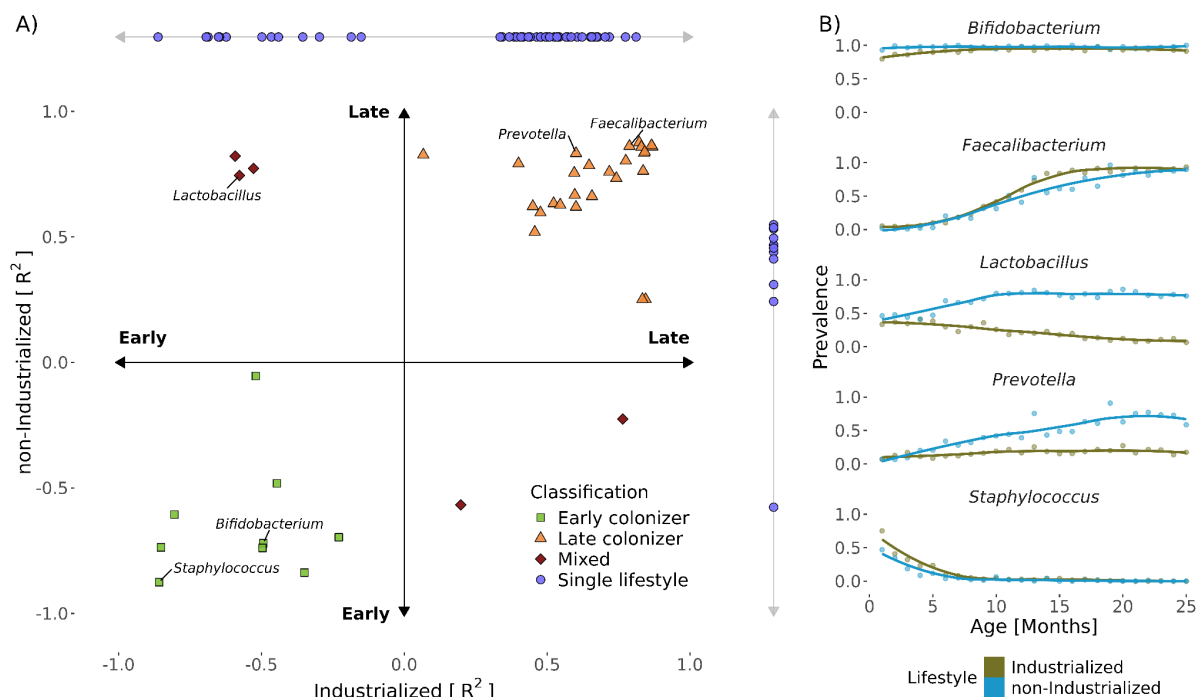507 importance for each taxa between lifestyles and within lifestyles.
508

509 We further evaluated the relationship between the predictive influence of each feature and
510 the predicted microbial age and calculated Spearman correlations between feature
511 abundance values and their corresponding SHAP-values. A positive correlation indicates
512 that higher abundance values for a feature are associated with an increased predicted
513 microbial age and vice versa. Based on the association between abundance and predictive
514 influence of each taxon, we grouped taxa into four colonization patterns. Early colonizer taxa
515 were more abundant in early-life and declined over time, whereas late colonizers were
516 largely absent in the first months of life and increased in abundance later in life. Mixed
517 colonizers include taxa whose colonization patterns differed between lifestyles and single
518 lifestyle taxa were only important in models from one lifestyle.
519

26

520  The majority (89%) of taxa important in both lifestyles had a consistent colonization pattern
521  across lifestyles (Figure 5A, Supplementary Table S2). Late colonizers accounted for 66% of
522  these taxa. Similarly, the majority (79%) of the single lifestyle taxa exhibited a late colonizer
523  pattern. These features combined drive the increase in alpha diversity described above
524  (Figure 2A) and the maturation of microbial communities with age.

525

526  *Faecalibacterium*, as an example of a late colonizer in both lifestyles, was nearly absent in
527  the first months of life and increased in prevalence and abundance over time with a similar
528  rate of increase in both lifestyles (Figure 5B, Supplementary Figure S6). *Prevotella* also
529  displayed a late colonizer pattern in both lifestyles. However, in non-industrialized
530  individuals, the rate of increase in prevalence was higher than for industrialized individuals
531  (Likelihood ratio-test, $\chi^2$ (DF) = 1.00, p-value < 0.001), indicating a different colonization
532  dynamic between lifestyles. *Staphylococcus*, as an early colonizer, was predominantly
533  detected during the first seven months of life and declined in prevalence thereafter, while
534  *Bifidobacterium* remained present in most individuals but decreased in relative abundance
535  over time. *Lactobacillus*, classified as a mixed colonizer, initially had a similar prevalence in
536  both lifestyles, however, while it declined in industrialized infants, it became increasingly
537  prevalent in non-industrialized infants (Likelihood ratio-test, $\chi^2$ (DF) = 1.00, p-value < 0.001).

538



539
540  **Figure 5. Bacterial taxa with significant global importance in lifestyle specific models**
541  **differ in their colonization dynamics. A)** The majority of taxa with high model importance
542  in both lifestyles correspond to late (eg. *Faecalibacterium* and *Prevotella*) or early (eg.
543  *Staphylococcus*) colonizers for both lifestyles. A set of taxa shows a mixed dynamic,
544  meaning that colonization patterns differ between lifestyles. Spearman correlation between
545  SHAP-values and relative taxon abundance in models trained on industrialized and non-
546  industrialized samples for all significantly important taxa. Taxa consistently showing a
547  negative correlation across both models are classified as early colonizers, while those with a
548  positive correlation are classified as late colonizers. Taxa with discordant colonization
549  patterns between lifestyles are classified as mixed and those significantly important in only

550 one lifestyle as single lifestyle. Correlations for taxa not significantly important in the
551 respective lifestyle or with Bonferroni-adjusted *p*-values > 0.05 are displayed on the top and
552 right side of the plot. **B)** Prevalence of taxa representative of the different colonization
553 dynamics in samples from both lifestyles plotted over chronological age, binned in months.
554
555

**Health conditions affect microbiome maturation**
557 We applied the combined and lifestyle specific models to two independent datasets from
558 clinically relevant conditions: one comprising infants diagnosed with SAM from two cohorts in
559 Bangladesh, and another consisting of preterm infants from three industrialized cohorts.
560 Both clinical conditions represent microbial communities clearly distinct from healthy
561 industrialized and non-industrialized infants (Supplementary Figure S7A, B).
562

563 SAM infants exhibited delayed microbial maturation compared to healthy infants from both
564 lifestyle groups across models (Figure 6A, B, Supplementary Figure S8A, B), consistent with
565 previous findings (Subramanian et al. 2014; Gehrig et al. 2019). The difference in microbial
566 age prediction for SAM infants and healthy non-industrialized infants was stronger when
567 predictions were based on a model including non-industrialized samples (Δ mean MAZ =
568 1.44 and 1.40 for the non-industrialized and combined model, respectively) compared to the
569 industrialized (1.15) model. The industrialized model underestimated microbial age for SAM
570 and healthy non-industrialized infants, supporting the need for lifestyle specific models to
571 assess microbial age not only for healthy infants but also for those that deviate from healthy
572 conditions.
573

574 To further investigate the main drivers of this delay, we calculated SHAP-values for healthy
575 and SAM infants within the non-industrialized dataset based on the non-industrialized model.
576 We identified 69 microbial features significantly associated with delayed maturation. The
577 most influential were the abundance of *Faecalibacterium* and microbial richness (Figure 6E,
578 Supplementary Figure S7C). When compared to healthy non-industrialized infants, SAM
579 infants had very low levels of the important late colonizer *Faecalibacterium* and did not show
580 any increase in microbial richness over time (Supplementary Figure S8C).
581

582 Preterm infants showed delayed microbiome maturation compared to industrialized full-term
583 infants, which persisted at least up to three months of age (Figure 6C). The delay in
584 maturation was larger in models including industrialized full-term infants (Δ mean MAZ $_{combined}$
585 $_{model}$ = 0.93 and Δ mean MAZ $_{Industrialized\ model}$ = 0.88) than in the non-industrialized model (0.76).
586 The non-industrialized model predicts increased microbial age in all preterms and, in
587 general, industrialized infants. This resulted in lower MAZ in non-industrialized full-terms
588 compared to industrialized preterms (Figure 6C, D). Although the non-industrialized model
589 predicted similar maturation trajectories for industrialized preterms and non-industrialized
590 full-terms, both groups showed a distinct compositional difference, with preterm infants
591 clustering independently from lifestyle (Supplementary Figure S7B). In contrast, both the
592 combined and industrialized models distinguished the maturation trajectories of preterm-born
593 industrialized infants and non-industrialized infants.
594

595 To investigate the drivers of delayed microbial maturation in preterm infants, we calculated
596 SHAP-values for preterm and full-term industrialized samples based on the industrialized
597 model. We identified 14 microbial features associated with delayed maturation in preterm

598 infants. Microbial richness was the most influential feature, which was reduced in preterms at
599 all ages (Figure 6F, Supplementary Figures S7D and S8D). Unlike in SAM infants, we
600 identified 17 microbial features associated with accelerated maturation in preterms. While for
601 SAM infants the drivers had a consistent effect across chronological age, most of the drivers
602 in preterm infants were more influential at specific time points during maturation. We
603 detected a persistent increase in abundance of *Staphylococcus* in preterm infants that led to
604 a decreased microbial age only at older chronological ages. One caveat of analysing
605 preterm infants is the reduced microbial diversity and rapid temporal shifts at early time
606 points in life.

607



608
609 **Figure 6. Health conditions affect microbiome maturation differently in lifestyle**
610 **specific models. A)** Microbial age development over chronological age in combined and
611 lifestyle specific models for healthy industrialized and non-industrialized infants and non-
612 industrialized infants with acute severe malnutrition (SAM). LOESS curves are fit to indicate
613 temporal trends. **B)** Differences in microbiome-for-age Z-scores (MAZ) between SAM and
614 healthy infants from both lifestyles predicted with a combined and the lifestyle specific
615 models. Diamonds represent medians, thick lines 50th percentiles, thin lines 95th
616 percentiles. Differences in MAZ were tested with a Wilcoxon test and adjusted for multiple
617 comparisons using Bonferroni correction (*$p<0.1$, **$p<0.05$, ***$p<0.01$, ****$p<0.001$). **C)**
618 Microbial age development over chronological age in combined and lifestyle specific models
619 for full-term healthy industrialized and non-industrialized infants and preterm born
620 industrialized infants. **D)** Differences in microbiome-for-age Z-scores (MAZ) between preterm
621 and healthy full-term infants from both lifestyles predicted with a combined and the lifestyle
622 specific models. **E)** Differences in mean SHAP-values between non-industrialized infants
623 with SAM and healthy non-industrialized infants for selected taxa binned by month. SHAP-

624 values are based on the non-industrialized model. Differences in SHAP-values were tested
625 for each age bin with a Wilcoxon test and adjusted for multiple comparisons using Bonferroni
626 correction (*$p<0.1$, **$p<0.05$, ***$p<0.01$, ****$p<0.001$). **F)** Differences in mean SHAP-values
627 between preterm and full-term industrialized infants for selected taxa binned by week.
628 SHAP-values are based on the industrialized model.
629
630
631

## Discussion

633 In this study, we analyzed the maturation dynamics of the infant gut microbiome between
634 industrialized and non-industrialized populations using machine learning models and a
635 globally diverse dataset of more than 15,000 early-life microbiomes from 20 different
636 countries. This comprehensive approach allowed us to detect significant lifestyle specific
637 maturation signatures. Here, we highlight the importance of diverse datasets in microbiome
638 research, demonstrating that more homogenous datasets predominantly based on
639 industrialized populations significantly decrease model performance for non-industrialized
640 populations. Our analysis included a particularly diverse non-industrialized dataset,
641 consisting of infants living in vastly different conditions, including a slum in Bangladesh and
642 the Amazon rainforest. Despite the wide diversity of living conditions within these groups, a
643 greater lifestyle diversity significantly improved model performance for all non-industrialized
644 lifestyles. This enabled us to identify both distinct microbiome maturation dynamics
645 compared to those of industrialized populations, and characteristics generalizable between
646 all lifestyles, as well as deviations from a healthy microbiome maturation in clinically relevant
647 conditions such as severe acute malnutrition (SAM) or preterm birth.
648
649 At the intra-individual level, age was the strongest driver of microbial diversity. It rapidly
650 increased during the first months of life, followed by a stabilization phase around 18 months
651 consistent across lifestyles and in line with previous studies (Bokulich et al. 2016;
652 Yatsunenko et al. 2012; Roswall et al. 2021; Kuang et al. 2016). While previous studies
653 observed a higher alpha diversity in non-industrialized infants after two years (De Filippo et
654 al. 2010; Yatsunenko et al. 2012), in our study, alpha diversity was reduced in non-
655 industrialized individuals compared to industrialized ones early in life. This might reflect the
656 higher prevalence of formula feeding in high income countries (Zong et al. 2021), which is
657 associated with an increased microbial diversity before the introduction of solid food
658 (Dwijayanti et al. 2025; Ho et al. 2018). Our findings highlight the need for caution when
659 generalizing diversity trends in early-life microbiomes from understudied non-industrialized
660 populations.
661
662 In contrast to the minimal impact on intra-individual diversity, lifestyle had a strong and age-
663 independent effect on the inter-individual microbiome variation and taxonomic composition.
664 These compositional differences underscore the importance of lifestyle as a primordial
665 ecological factor shaping the early-life microbiota rather than geography (Olm et al. 2022).
666 Our findings emphasise that while the microbial maturation process follows a conserved
667 longitudinal trajectory, its community diversity and composition at an early age are distinctly

668 modulated by lifestyle. This provides the foundation for exploring the specific colonization
669 patterns that differentiate these microbial trajectories, as well as for predicting microbial age.
670

671 Microbial age prediction models revealed distinct colonization dynamics shaped by lifestyle.
672 In non-industrialized infants, *Prevotella* appeared as a relevant late colonizer with high
673 predictive importance, aligning with its known underrepresentation in industrialized adult gut
674 microbiome (Tett et al. 2019). Our findings suggest that the diminished prevalence of
675 *Prevotella* in industrialized adults may result from early-life colonization patterns.
676 Conversely, *Staphylococcus* was identified as a highly relevant early colonizer in
677 industrialized populations, possibly reflecting its higher prevalence in infants born via c-
678 section, which is more common in countries with higher HDI (Ye et al. 2016; Dominguez-
679 Bello et al. 2010; Reyman et al. 2019). Interestingly, while *Prevotella* and *Staphylococcus*
680 showed consistent trajectories across lifestyles, *Lactobacillus* had divergent colonization
681 patterns, increasing in prevalence in non-industrialized infants but declining over time in
682 industrialized ones. The observed contrasting patterns may reflect both ecological
683 succession events during the maturation (Pasolli et al. 2020) and differences in the
684 prevalence of c-sections (Tamburini et al. 2016; Ye et al. 2016).
685

686 Contrary to our hypothesis, not all the taxa known to play important roles in the early-life
687 microbiota were informative for age prediction. Despite its well-documented biological role in
688 the metabolism of breast milk oligosaccharides (Ennis et al. 2024; Ojima et al. 2022),
689 *Bifidobacterium* had little contribution to the performance of our microbial age prediction
690 model, likely due to its ubiquity in infants. This highlights a key challenge in using amplicon
691 based data for age modeling and emphasizes the need for higher taxonomic resolution to
692 disentangle the functionality of distinct species or strains at different stages of an infant's life
693 (Vatanen et al. 2022; Ennis et al. 2024).
694

695 Our modeling framework can detect clinically relevant deviations in microbiome maturation
696 by applying it to infants diagnosed with SAM and preterm born infants. Delayed microbiome
697 maturation in SAM infants and microbial signatures distinguishing preterm from full-term
698 infants have been reported previously (Subramanian et al. 2014; Van Rossum et al. 2024).
699 Here, we identify specific microbial taxa associated with delayed maturation in both clinical
700 conditions, providing a deeper insight into the drivers of these conditions. It has also been
701 reported that malnutrition during pregnancy is associated with preterm birth (Bloomfield
702 2011) and preterm born infants show increased susceptibility to malnutrition (Harding et al.
703 2017), particularly in low income countries (Sania et al. 2014). Our results indicate a
704 potential developmental and nutritional link between both conditions. Although the preterm
705 and SAM infants differed in lifestyles and age span, further studies should examine whether
706 delayed microbiome maturation represents a shared signature of both conditions. Notably,
707 the observed differences in maturation delay between lifestyle specific models highlight the
708 importance of using reference populations that reflect the diversity of the target population.
709

710 Despite the limited taxonomic resolution of amplicon-based data, our study advances our
711 understanding of global microbiome maturation and highlights the value of integrating
712 underrepresented populations into microbiome research. By compiling one of the most
713 diverse infant microbiome datasets to date, we improved microbial age prediction models
714 that can be applied across lifestyles and help address persistent geographic biases in the
715 study of early-life microbiomes. Previous work has shown that microbiomes from non-

37

716 industrialized populations also differ markedly between those populations(Abdill et al. 2025).
717 However, limited data from non-industrialized infant cohorts prevented deeper comparisons
718 across lifestyles or larger geographic regions. Expanding studies of infant microbiome
719 maturation in globally diverse populations will be essential to refine these models and
720 broaden their applicability. Achieving more equitable and globally represented early-life
721 microbial studies will also require going beyond sample inclusion to foster close
722 collaborations and promote capacity building with local researchers and ensure reciprocal
723 benefits for the communities (Armenteras 2021). While these goals go beyond the scope of
724 this study, our work illustrates the potential and responsibility of microbial research to draw
725 more globally representative conclusions.
726

727 **Reference**

728 Abdill, Richard J., Elizabeth M. Adamowicz, and Ran Blekhman. 2022. "Public
729 Human Microbiome Data Are Dominated by Highly Developed Countries." *PLOS*
730 *Biology* 20 (2): e3001536. https://doi.org/10.1371/journal.pbio.3001536.

731 Abdill, Richard J., Samantha P. Graham, Vincent Rubinetti, et al. 2025. "Integration
732 of 168,000 Samples Reveals Global Patterns of the Human Gut Microbiome." *Cell*
733 188 (4): 1100-1118.e17. https://doi.org/10.1016/j.cell.2024.12.017.

734 Almeida, Alexandre, Alex L. Mitchell, Miguel Boland, et al. 2019. "A New Genomic
735 Blueprint of the Human Gut Microbiota." *Nature* 568 (7753): 499–504.
736 https://doi.org/10.1038/s41586-019-0965-1.

737 Armenteras, Dolors. 2021. "Guidelines for Healthy Global Scientific Collaborations."
738 *Nature Ecology & Evolution* 5 (9): 1193–94. https://doi.org/10.1038/s41559-021-
739 01496-y.

740 Bisgaard, Hans, Nan Li, Klaus Bonnelykke, et al. 2011. "Reduced Diversity of the
741 Intestinal Microbiota during Infancy Is Associated with Increased Risk of Allergic
742 Disease at School Age." *Journal of Allergy and Clinical Immunology* 128 (3): 646-
743 652.e5. https://doi.org/10.1016/j.jaci.2011.04.060.

744 Bloomfield, Frank H. 2011. "How Is Maternal Nutrition Related to Preterm Birth?"
745 *Annual Review of Nutrition* 31 (August): 235–61. https://doi.org/10.1146/annurev-
746 nutr-072610-145141.

747 Bokulich, Nicholas A., Jennifer Chung, Thomas Battaglia, et al. 2016. "Antibiotics,
748 Birth Mode, and Diet Shape Microbiome Maturation during Early Life." *Science*
749 *Translational Medicine* 8 (343): 343ra82-343ra82.
750 https://doi.org/10.1126/scitranslmed.aad7121.

751 Bolger, Anthony M., Marc Lohse, and Bjoern Usadel. 2014. "Trimmomatic: A Flexible
752 Trimmer for Illumina Sequence Data." *Bioinformatics* 30 (15): 2114–20.
753 https://doi.org/10.1093/bioinformatics/btu170.

754 Callahan, Benjamin J., Paul J. McMurdie, Michael J. Rosen, Andrew W. Han, Amy Jo
755 A. Johnson, and Susan P. Holmes. 2016. "DADA2: High-Resolution Sample
756 Inference from Illumina Amplicon Data." *Nature Methods* 13 (7): 7.
757 https://doi.org/10.1038/nmeth.3869.

38

39

Chu, Derrick M., Jun Ma, Amanda L. Prince, Kathleen M. Antony, Maxim D. Seferovic, and Kjersti M. Aagaard. 2017. "Maturation of the Infant Microbiome Community Structure and Function Across Multiple Body Sites and in Relation to Mode of Delivery." *Nature Medicine* 23 (3): 314–26. https://doi.org/10.1038/nm.4272.

De Filippo, Carlotta, Duccio Cavalieri, Monica Di Paola, et al. 2010. "Impact of Diet in Shaping Gut Microbiota Revealed by a Comparative Study in Children from Europe and Rural Africa." *Proceedings of the National Academy of Sciences* 107 (33): 14691–96. https://doi.org/10.1073/pnas.1005963107.

Dominguez-Bello, Maria G., Elizabeth K. Costello, Monica Contreras, et al. 2010. "Delivery Mode Shapes the Acquisition and Structure of the Initial Microbiota across Multiple Body Habitats in Newborns." *Proceedings of the National Academy of Sciences* 107 (26): 11971–75. https://doi.org/10.1073/pnas.1002601107.

Dwijayanti, Ira, Farah Nuriannisa, Laura Navika Yamani, Catur Wulandari, Fasty Arum Utami, and Trias Mahmudiono. 2025. "Changes in Gut Microbiota Diversity and Composition during Feeding Transitions in Infants: A Scoping Review." *Nutrition* 138 (October): 112814. https://doi.org/10.1016/j.nut.2025.112814.

Ennis, Dena, Shimrit Shmorak, Evelyn Jantscher-Krenn, and Moran Yassour. 2024. "Longitudinal Quantification of Bifidobacterium Longum Subsp. Infantis Reveals Late Colonization in the Infant Gut Independent of Maternal Milk HMO Composition." *Nature Communications* 15 (1): 894. https://doi.org/10.1038/s41467-024-45209-y.

Fahur Bottino, Guilherme, Kevin S. Bonham, Fadheela Patel, et al. 2025. "Early Life Microbial Succession in the Gut Follows Common Patterns in Humans across the Globe." *Nature Communications* 16 (1): 660. https://doi.org/10.1038/s41467-025-56072-w.

"FastQC." 2015. June. https://qubeshub.org/resources/fastqc.

Gehrig, Jeanette L., Siddarth Venkatesh, Hao-Wei Chang, et al. 2019. "Effects of Microbiota-Directed Foods in Gnotobiotic Animals and Undernourished Children." *Science* 365 (6449): eaau4732. https://doi.org/10.1126/science.aau4732.

Greenwell, Brandon. 2024. *Fastshap: Fast Approximate Shapley Values*. https://CRAN.R-project.org/package=fastshap.

Harding, Jane E, Barbara E Cormack, Tanith Alexander, Jane M Alsweiler, and Frank H Bloomfield. 2017. "Advances in Nutrition of the Newborn Infant." *The Lancet* 389 (10079): 1660–68. https://doi.org/10.1016/S0140-6736(17)30552-4.

Hill, Cian J., Denise B. Lynch, Kiera Murphy, et al. 2017. "Evolution of Gut Microbiota Composition from Birth to 24 Weeks in the INFANTMET Cohort." *Microbiome* 5 (1): 4. https://doi.org/10.1186/s40168-016-0213-y.

Ho, Nhan T., Fan Li, Kathleen A. Lee-Sarwar, et al. 2018. "Meta-Analysis of Effects of Exclusive Breastfeeding on Infant Gut Microbiota across Populations." *Nature Communications* 9 (1): 4169. https://doi.org/10.1038/s41467-018-06473-x.

Kuang, Ya-Shu, Sheng-Hui Li, Yong Guo, et al. 2016. "Composition of Gut Microbiota in Infants in China and Global Comparison." *Scientific Reports* 6 (1): 36666. https://doi.org/10.1038/srep36666.

40

41

Kursa, Miron B., and Witold R. Rudnicki. 2010. "Feature Selection with the Boruta Package." *Journal of Statistical Software* 36 (September): 1–13. https://doi.org/10.18637/jss.v036.i11.

Langmead, Ben, and Steven L. Salzberg. 2012. "Fast Gapped-Read Alignment with Bowtie 2." *Nature Methods* 9 (4): 357–59. https://doi.org/10.1038/nmeth.1923.

Martínez, Inés, James C. Stegen, Maria X. Maldonado-Gómez, et al. 2015. "The Gut Microbiota of Rural Papua New Guineans: Composition, Diversity Patterns, and Ecological Processes." *Cell Reports* 11 (4): 527–38. https://doi.org/10.1016/j.celrep.2015.03.049.

Nayfach, Stephen, Zhou Jason Shi, Rekha Seshadri, Katherine S. Pollard, and Nikos C. Kyrpides. 2019. "New Insights from Uncultivated Genomes of the Global Human Gut Microbiome." *Nature* 568 (7753): 505–10. https://doi.org/10.1038/s41586-019-1058-x.

Ojima, Miriam N, Lin Jiang, Aleksandr A Arzamasov, et al. 2022. "Priority Effects Shape the Structure of Infant-Type Bifidobacterium Communities on Human Milk Oligosaccharides." *The ISME Journal* 16 (9): 2265–79. https://doi.org/10.1038/s41396-022-01270-3.

O'Keefe, Stephen J. D., Jia V. Li, Leo Lahti, et al. 2015. "Fat, Fibre and Cancer Risk in African Americans and Rural Africans." *Nature Communications* 6 (1): 6342. https://doi.org/10.1038/ncomms7342.

Oksanen, Jari, Gavin L. Simpson, F. Guillaume Blanchet, et al. 2024. *Vegan: Community Ecology Package*. https://CRAN.R-project.org/package=vegan.

Olm, Matthew R., Dylan Dahan, Matthew M. Carter, et al. 2022. "Robust Variation in Infant Gut Microbiome Assembly across a Spectrum of Lifestyles." *Science* 376 (6598): 1220–23. https://doi.org/10.1126/science.abj2972.

Pasolli, Edoardo, Francesco Asnicar, Serena Manara, et al. 2019. "Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle." *Cell* 176 (3): 649-662.e20. https://doi.org/10.1016/j.cell.2019.01.001.

Pasolli, Edoardo, Francesca De Filippis, Italia E. Mauriello, et al. 2020. "Large-Scale Genome-Wide Analysis Links Lactic Acid Bacteria from Food with the Gut Microbiome." *Nature Communications* 11 (1): 2610. https://doi.org/10.1038/s41467-020-16438-8.

Pruesse, Elmar, Christian Quast, Katrin Knittel, et al. 2007. "SILVA: A Comprehensive Online Resource for Quality Checked and Aligned Ribosomal RNA Sequence Data Compatible with ARB." *Nucleic Acids Research* 35 (21): 7188–96. https://doi.org/10.1093/nar/gkm864.

R. A. Rigby and D. M. Stasinopoulos. 2005. "Generalized Additive Models for Location, Scale and Shape,(with Discussion)." *Applied Statistics* 54: 507–54.

Raman, Arjun S., Jeanette L. Gehrig, Siddarth Venkatesh, et al. 2019. "A Sparse Covarying Unit That Describes Healthy and Impaired Human Gut Microbiota Development." *Science* 365 (6449): eaau4735. https://doi.org/10.1126/science.aau4735.

42

43

843  Reyman, Marta, Marlies A. van Houten, Debbie van Baarle, et al. 2019. "Impact of
844  Delivery Mode-Associated Gut Microbiota Dynamics on Health in the First Year of
845  Life." *Nature Communications* 10 (1): 1. https://doi.org/10.1038/s41467-019-13014-7.

846  Roswall, Josefine, Lisa M. Olsson, Petia Kovatcheva-Datchary, et al. 2021.
847  "Developmental Trajectory of the Healthy Human Gut Microbiota during the First 5
848  Years of Life." *Cell Host & Microbe* 29 (5): 765-776.e3.
849  https://doi.org/10.1016/j.chom.2021.02.021.

850  Sánchez-Quinto, Andrés, Daniel Cerqueda-García, Luisa I. Falcón, et al. 2020. "Gut
851  Microbiome in Children from Indigenous and Urban Communities in México: Different
852  Subsistence Models, Different Microbiomes." *Microorganisms* 8 (10): 10.
853  https://doi.org/10.3390/microorganisms8101592.

854  Sania, Ayesha, Donna Spiegelman, Janet Rich-Edwards, et al. 2014. "The
855  Contribution of Preterm Birth and Intrauterine Growth Restriction to Childhood
856  Undernutrition in Tanzania." *Maternal & Child Nutrition* 11 (4): 618–30.
857  https://doi.org/10.1111/mcn.12123.

858  Shapley, L. S. 1953. "17. A Value for n-Person Games." In *Contributions to the
859  Theory of Games, Volume II*, edited by Harold William Kuhn and Albert William
860  Tucker. Princeton University Press. https://doi.org/doi:10.1515/9781400881970-018.

861  Stewart, Christopher J., Nadim J. Ajami, Jacqueline L. O'Brien, et al. 2018.
862  "Temporal Development of the Gut Microbiome in Early Childhood from the TEDDY
863  Study." *Nature* 562 (7728): 7728. https://doi.org/10.1038/s41586-018-0617-x.

864  Subramanian, Sathish, Sayeeda Huq, Tanya Yatsunenko, et al. 2014. "Persistent
865  Gut Microbiota Immaturity in Malnourished Bangladeshi Children." *Nature* 510
866  (7505): 7505. https://doi.org/10.1038/nature13421.

867  Suez, Jotham, Tal Korem, David Zeevi, et al. 2014. "Artificial Sweeteners Induce
868  Glucose Intolerance by Altering the Gut Microbiota." *Nature* 514 (7521): 181–86.
869  https://doi.org/10.1038/nature13793.

870  Tamburini, Sabrina, Nan Shen, Han Chih Wu, and Jose C. Clemente. 2016. "The
871  Microbiome in Early Life: Implications for Health Outcomes." *Nature Medicine* 22 (7):
872  713–22. https://doi.org/10.1038/nm.4142.

873  Tett, Adrian, Kun D. Huang, Francesco Asnicar, et al. 2019. "The Prevotella Copri
874  Complex Comprises Four Distinct Clades Underrepresented in Westernized
875  Populations." *Cell Host & Microbe* 26 (5): 666-679.e7.
876  https://doi.org/10.1016/j.chom.2019.08.018.

877  United Nations Development Programme. 2022. *Human Development Index: 2022
878  Statistical Update*. (New York). http://hdr.undp.org/en/data/.

879  Vatanen, Tommi, Qi Yan Ang, Léa Siegwald, et al. 2022. "A Distinct Clade of
880  Bifidobacterium Longum in the Gut of Bangladeshi Children Thrives during Weaning."
881  *Cell* 0 (0). https://doi.org/10.1016/j.cell.2022.10.011.

882  Willers, Maike, Thomas Ulas, Lena Völlger, et al. 2020. "S100A8 and S100A9 Are
883  Important for Postnatal Development of Gut Microbiota and Immune System in Mice
884  and Infants." *Gastroenterology* 159 (6): 2130-2145.e5.
885  https://doi.org/10.1053/j.gastro.2020.08.019.

44

45

886     Wright, Erik S. 2016. "Using DECIPHER v2.0 to Analyze Big Biological Sequence
887         Data in R." *The R Journal* 8 (1): 352–59. https://doi.org/10.32614/RJ-2016-025.

888     Wright, Marvin N., and Andreas Ziegler. 2017. "Ranger: A Fast Implementation of
889         Random Forests for High Dimensional Data in C++ and R." *Journal of Statistical*
890         *Software* 77 (1): 1–17. https://doi.org/10.18637/jss.v077.i01.

891     Yatsunenko, Tanya, Federico E. Rey, Mark J. Manary, et al. 2012. "Human Gut
892         Microbiome Viewed across Age and Geography." *Nature* 486 (7402): 7402.
893         https://doi.org/10.1038/nature11053.

894     Ye, J, J Zhang, R Mikolajczyk, Mr Torloni, Am Gülmezoglu, and Ap Betran. 2016.
895         "Association between Rates of Caesarean Section and Maternal and Neonatal
896         Mortality in the 21st Century: A Worldwide Population-Based Ecological Study with
897         Longitudinal Data." *BJOG: An International Journal of Obstetrics & Gynaecology* 123
898         (5): 745–53. https://doi.org/10.1111/1471-0528.13592.

899     Zeileis, Achim, and Torsten Hothorn. 2002. "Diagnostic Checking in Regression
900         Relationships." *R News* 2 (3): 7–10.

901     Zong, Xin'nan, Han Wu, Min Zhao, Costan G. Magnussen, and Bo Xi. 2021. "Global
902         Prevalence of WHO Infant Feeding Practices in 57 LMICs in 2010–2018 and Time
903         Trends since 2000 for 44 LMICs." *eClinicalMedicine* 37 (July).
904         https://doi.org/10.1016/j.eclinm.2021.100971.

905     **Acknowledgements**

915

916     **Author contributions**
917     Conceptualization: SKFS, CB. Data curation, Formal Analysis and Investigation: RB.
918     Supervision: CB, DV, VHJD, SKFS. Writing – Original Draft Preparation: RB, GP, VHJD.
919     Writing – Review & Editing: RB, GP, UL, CB, DV, VHJD, SKFS. Funding Acquisition: SKFS,
920     CB, DV.


921     **Conflict of Interest**
922     The authors declare that there is no conflict of interests
923


924     **Supplementary Data**

925     Table S1. List of studies with number of samples and individuals, assigned lifestyle,
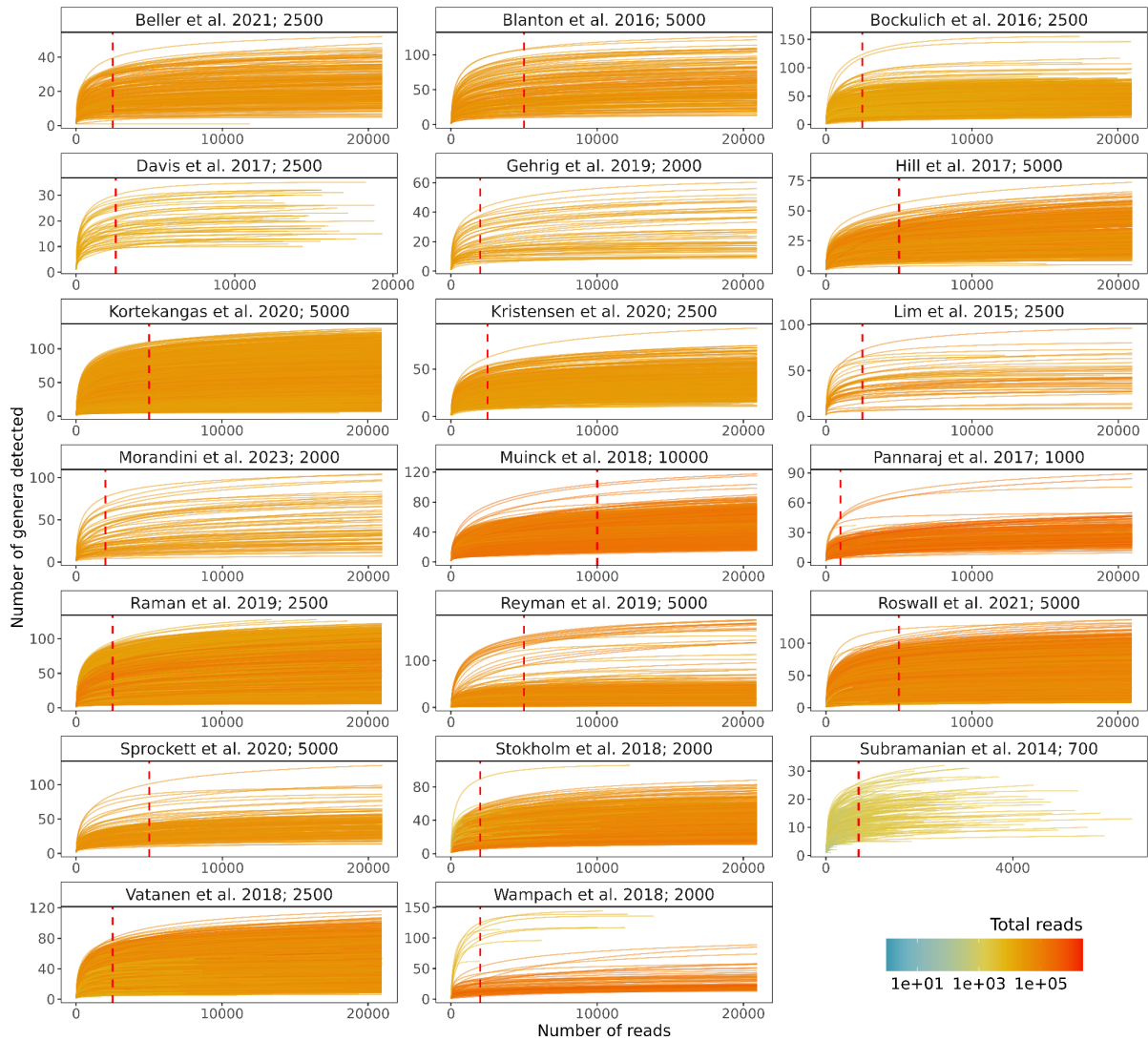926     countries of origin and filtering thresholds


46

47

927
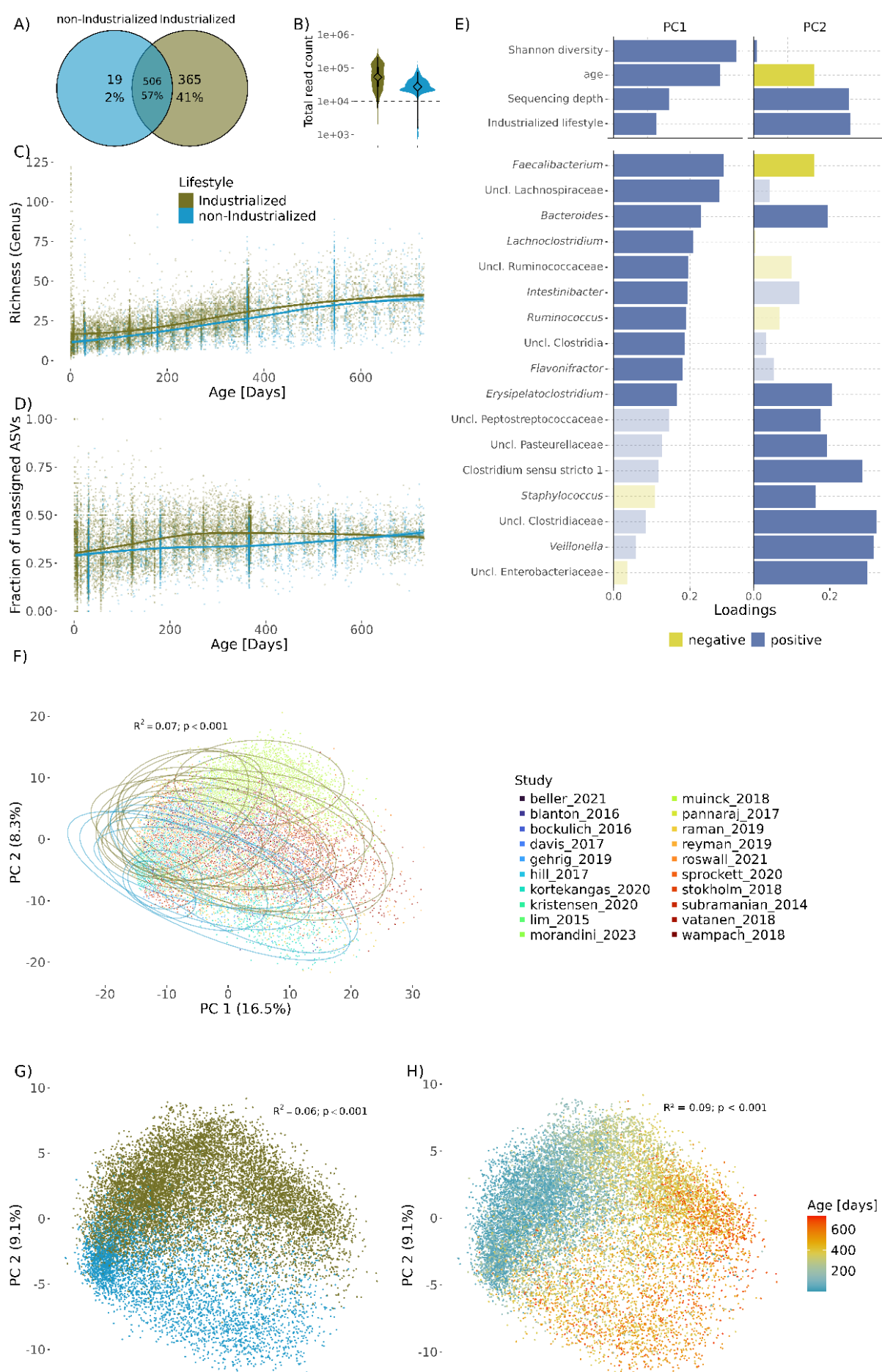928     Table S2. Classification of taxa into early and late colonizers for each lifestyle



929
930 **Figure S1.** Study specific rarefaction curves to determine study-specific filtering threshold.
931 Filter thresholds are indicated by red dashed vertical lines and specified after the study
932 name in the pane title. Color indicates the total number of taxonomically assigned reads for
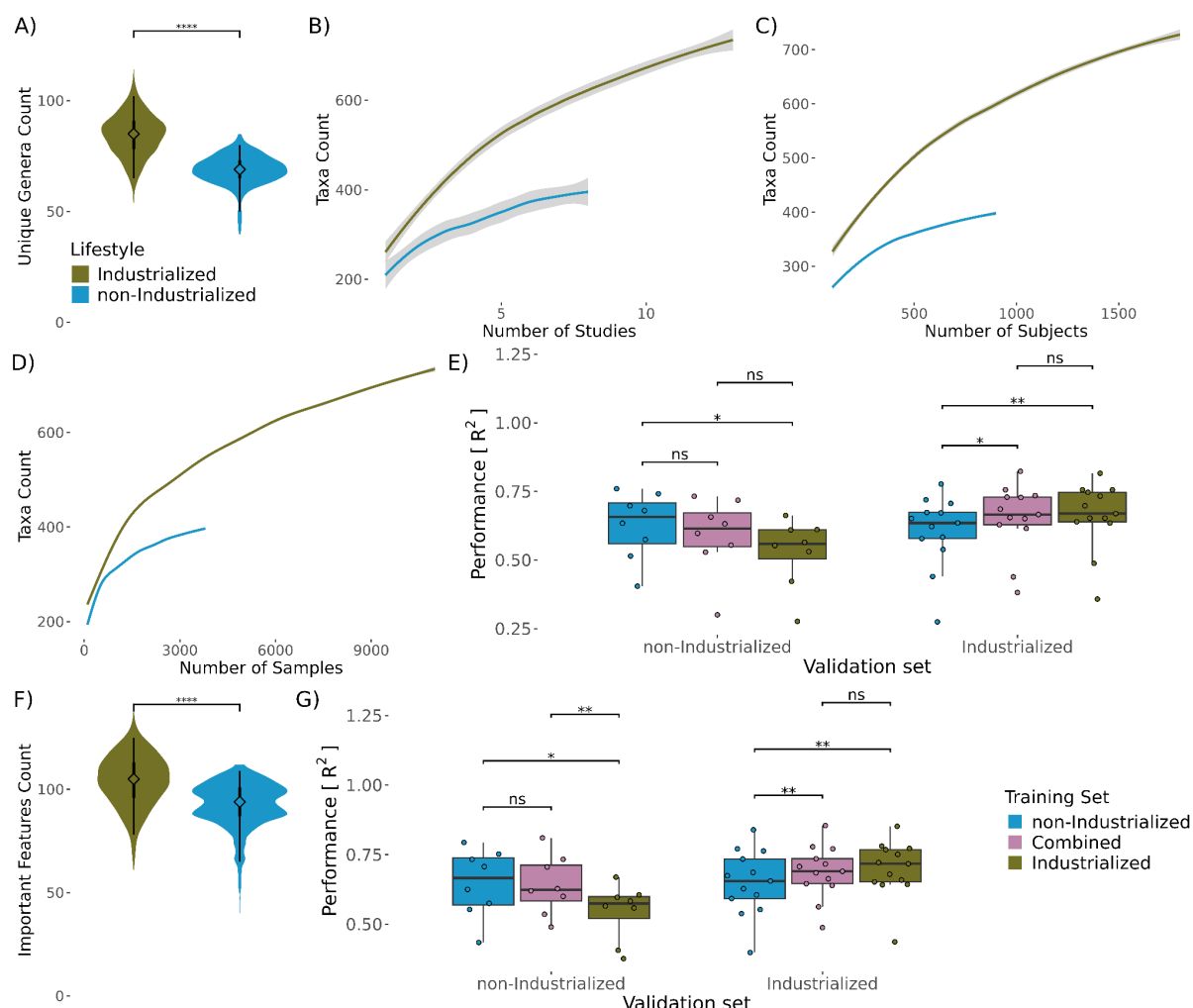933 each sample.

934
935

48

49



50

51

**Figure S2. Gut microbial diversity across lifestyles and datasets A)** Venn diagram for the amount of genera detected in each lifestyle. **B)** Difference in sequencing depth between lifestyles. Violin plot with the distributions of total read counts per sample. Dashed line marks 10,000 reads. **C)** Microbial richness over time calculated from rarefied counts aggregated to genus level, plotted over chronological age and stratified by lifestyle. Locally estimated scatterplot smoothing (LOESS) curves are fit to visualize temporal trends. **D)** Fraction of ASVs unassigned at genus level per samples plotted over chronological age and stratified by lifestyle. Locally estimated scatterplot smoothing (LOESS) curves are fit to visualize temporal trends. **E)** Drivers of beta diversity in a PCA on clr-transformed raw counts. Upper panel: Spearman correlations of sample variables with the first two principal components (PCs). Lover panel: Top 10 loadings of taxa on the first two PCs. Loadings not in the top 10 for one PC are transparent. **F)** PCA on clr-transformed raw counts colored by study. $R^2$ value indicates the proportion of variance explained in a PERMANOVA using a sequential model with the terms age, lifestyle and study. **G)** PCA on clr-transformed rarefied counts colored by lifestyle and **H)** age. Read counts were rarefied to 2,000 per sample.



**Figure S3.** Results for downsampling both lifestyle groups to the same amount of studies, subjects and samples. **A)** Difference in the number of genera detected only in industrialized or non-industrialized samples. Each value represents one downsampling iteration. Diamonds represent medians, thick lines 50th percentiles, thin lines 95th percentiles. Significance was
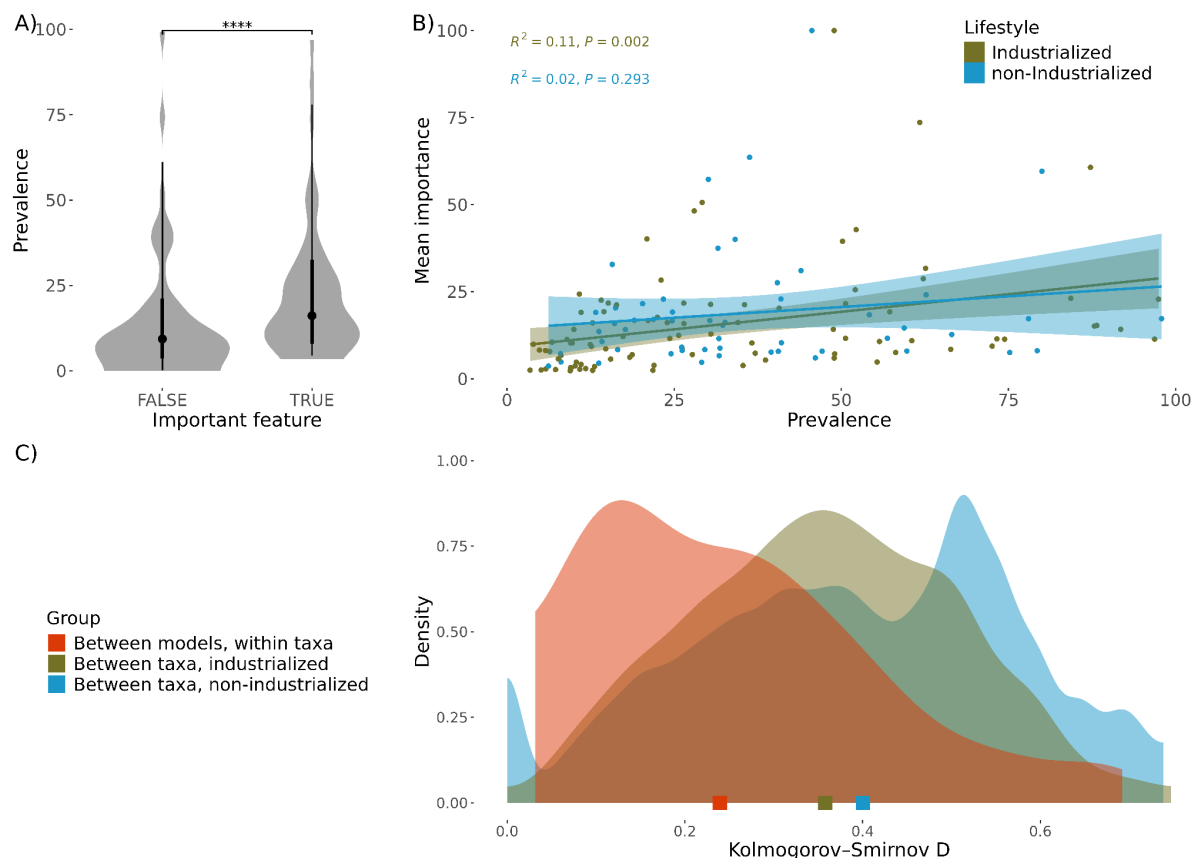
52

959 tested with a paired Wilcoxon test (*$p$<0.1, **$p$<0.05, ***$p$<0.01, ****p<0.001). **B, C, D)**
960 Number of taxa detected in rarefied count data by lifestyle in randomly subsampled datasets
961 to specific numbers of studies, subjects and samples. **E)** Performance of microbial age
962 modelling with lifestyle specific models on relative abundances of read counts that were
963 rarefied to 2,000 reads per sample. Performance was assessed using LODO-CV, where
964 each study served as a validation set for models trained on a downsampled dataset.
965 Coefficient of determination ($R^2$) from a linear model of rank transformed microbial age
966 versus chronological age is used as performance metrics and was calculated for each
967 combination of model and validation set separately. Differences in model performances were
968 tested with a paired Wilcoxon test and adjusted for multiple comparisons using Bonferroni
969 correction. **F)** Difference in the number of important features between industrialized and non-
970 industrialized models trained on downsampled data. Each value represents one
971 downsampling iteration. Significance was tested with a paired Wilcoxon test. **G)**
972 Performance of microbial age modelling with lifestyle specific models. Separate models were
973 trained either on microbiomes from industrialized samples, non-industrialized samples, or a
974 combined dataset. Each point represents the mean out of 50 downsampling iterations.
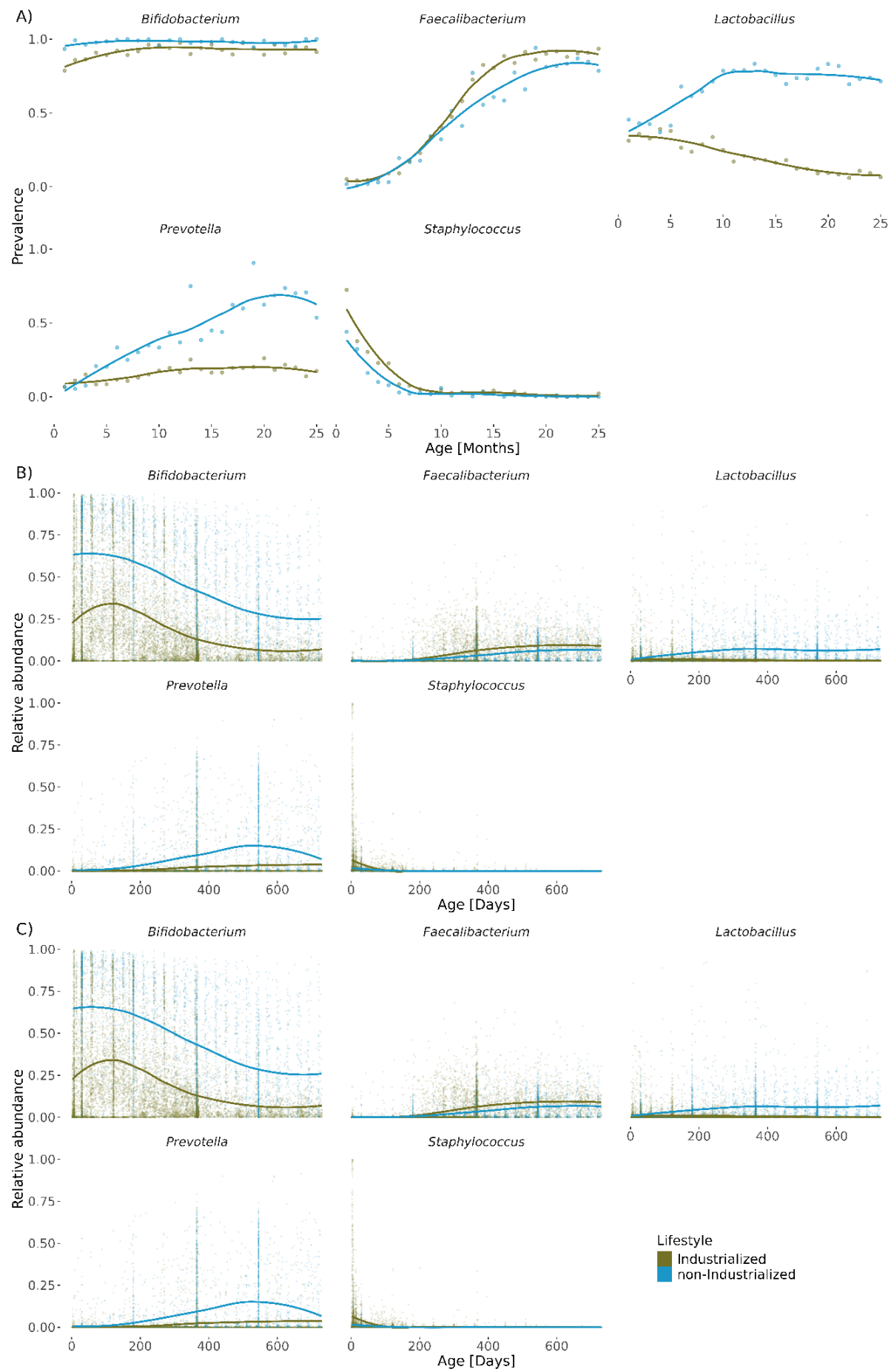975 (*$p$<0.1, **$p$<0.05, ***$p$<0.01, ****p<0.001).

57

**Figure S4.** Relevance of important taxa in lifestyle specific models. All taxa identified as significantly important by Boruta in at least two models of LODO-CV for the specific lifestyle are displayed. **A)** Relative feature importance in random forest models trained separately on samples from individuals with an industrialized or a non-industrialized lifestyle. Grey bars indicate the prevalence of a specific taxon in each of the lifestyles. Importances are displayed only for features significantly important for the respective lifestyle. **B)** Violin plot for Shapley additive explanation (SHAP) values. Color indicates relative taxon abundance or alpha diversity, scaled per feature. Only values for features significantly important for the respective lifestyle are shown.



**Figure S5: A)** Density plots for the distributions of effect sizes derived from Kolmogorov-Smirnov tests comparing SHAP-value distributions. Comparisons were made either between models for the same taxa or between different taxa within individual models. Diamonds represent medians, thick lines 50th percentiles, thin lines 95th percentiles. Significance was tested with a paired Wilcoxon test (*$p<0.1$, **$p<0.05$, ***$p<0.01$, ****$p<0.001$). **B)** Correlation of mean relative feature importance and mean prevalence over all studies for all taxa from Figure S4. **C)** Difference in prevalences of taxa important only in one lifestyle group. Prevalences are compared between the lifestyle where the feature is important and the lifestyle where it is not important. Squares represent medians.
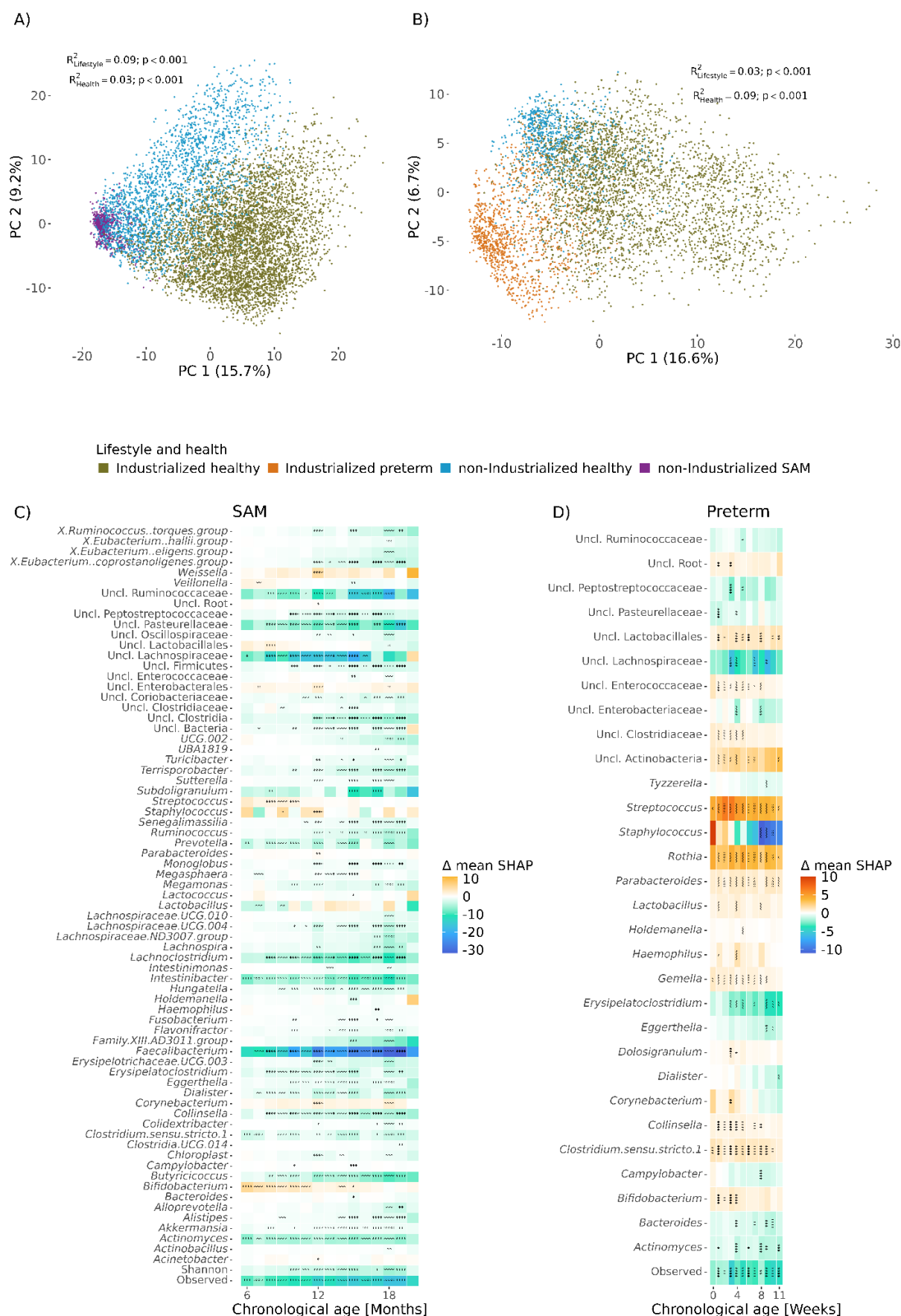
58

59



1003

60

61

1004 **Figure S6.** Longitudinal dynamics of taxa representative of the different colonization
1005 dynamics in samples from both lifestyles. **A)** Prevalence over chronological age binned in
1006 months in rarefied data. **B)** Relative abundance over chronological age in un-rarefied and **C)**
1007 rarefied data. LOESS curves are fit to visualize temporal trends.
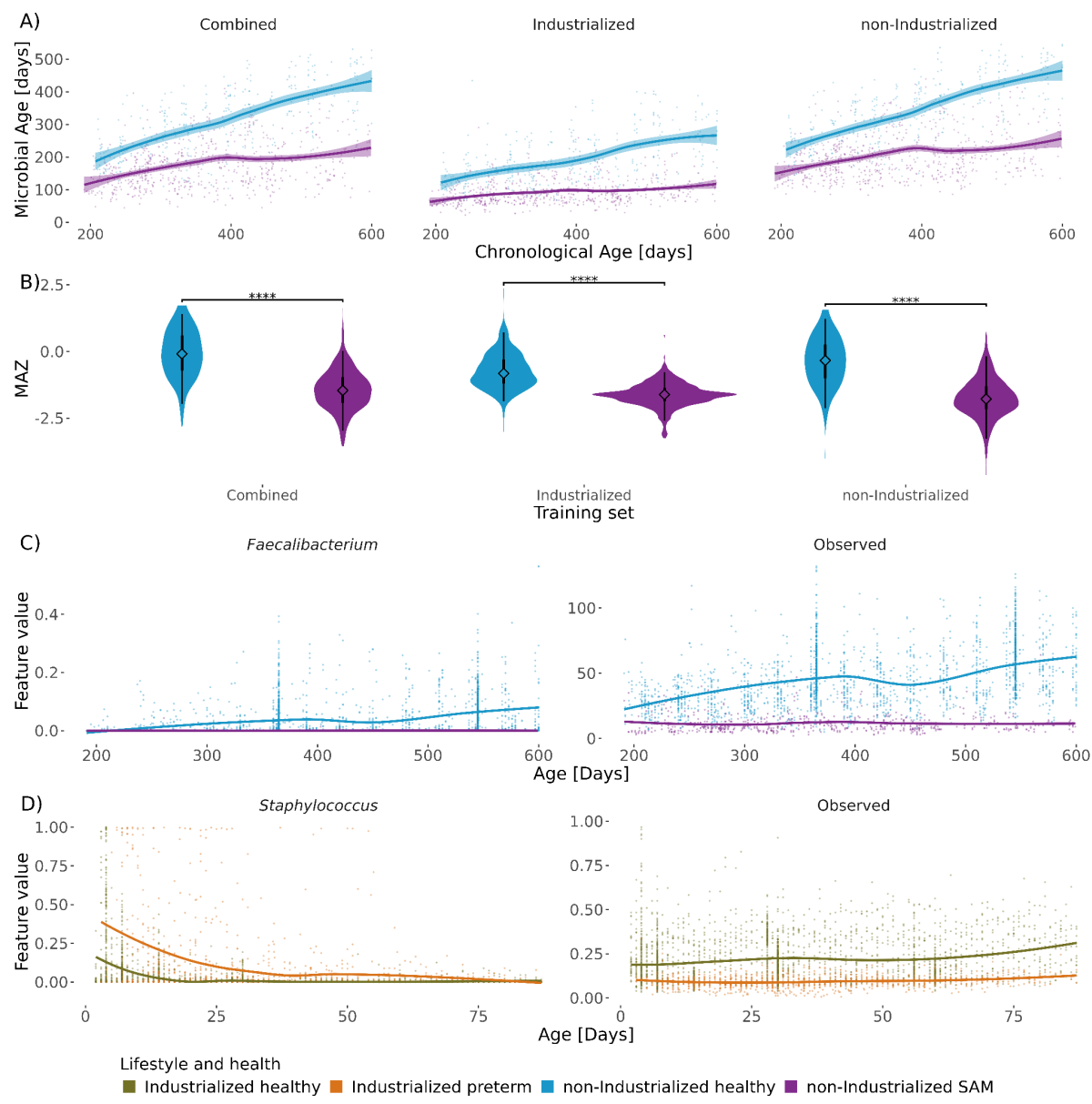1008

62

63



**Figure S7:** Application of the microbial age model to health conditions. **A, B)** PCA on clr-
transformed raw counts colored by lifestyle and health condition for healthy infants and non-

64

65

1012   industrialized infants with SAM (A) and preterm industrialized infants (B). Healthy samples
1013   outside the age range of non-healthy samples were excluded. $R^2$ values indicate the
1014   proportion of variance explained in a Permutational multivariate analysis of variance
1015   (PERMANOVA) using a sequential model with the terms age, lifestyle, health and study. **C)**
1016   Differences in mean SHAP-values between non-industrialized infants with SAM and healthy
1017   non-industrialized infants by month. Depicted are all taxa with a significant difference in at
1018   least one month. SHAP-values are based on the non-industrialized model. Differences in
1019   SHAP-values were tested for each age bin with a Wilcoxon test and adjusted for multiple
1020   comparisons using Bonferroni correction (*$p<0.1$, **$p<0.05$, ***$p<0.01$, ****$p<0.001$). **D)**
1021   Differences in mean SHAP-values between preterm and full-term industrialized infants
1022   binned by week. Depicted are all taxa with a significant difference in at least one week.
1023   SHAP-values are based on the industrialized model.

1024
1025
1026
1027
1028



1029

66

67

1030 **Figure S8: A)** Predicted age for healthy and SAM infants in studies containing both healthy
1031 and SAM infants. **B)** Differences in MAZ between healthy and malnourished infants in
1032 studies containing both healthy and SAM infants. Diamonds represent medians, thick lines
1033 50th percentiles, thin lines 95th percentiles. Significance was tested with a paired Wilcoxon
1034 test (*$p$<0.1, **$p$<0.05, ***$p$<0.01, ****p<0.001). **C)** Relative abundance over time for selected
1035 taxa and observed richness driving a delay in maturation in SAM and **D)** preterm infants
1036 plotted over chronological age. LOESS curves are fit to visualize temporal trends.
1037
1038
1039
1040
1041
1042
1043
1044
1045

68