

1 Supplementary material:

2 Sample preparation for 16S amplicon sequencing. Sample processing was performed as recently
3 described [1]. In brief, the DNA from all samples was extracted using the Qiagen DNA Minikit (Qiagen,
4 Hilden, Germany), following the spin protocol for DNA purification from body fluids. Subsequently,
5 the V4 region of the 16S rRNA gene was amplified using forward (5'-GTGCCAGCMGCCGCGGTAA-
6 3') and reverse (5'-GGACTACHVGGGTWTCTAAT-3') primers [2] and modified with an Illumina
7 adaptor sequence at the 5' end. PCR cycling conditions were 95 °C for 6 minutes and 40 cycles at 95
8 °C for 30 seconds, at 59 °C for 30 seconds and at 72 °C for 1.5 minutes (with a final elongation step at
9 72 °C for 5 minutes). PCR products were purified by QIAquick PCR purification kit (Qiagen, Hilden,
10 Germany). The samples were semi-quantified using the DNA 7500 kit with an Agilent 2100 Bioanalyzer
11 (Agilent Technologies, Palo Alto, CA). As suggested previously, samples with less than 1 ng/μl after
12 PCR and purification should be excluded from further analyses [3]. As part of our quality control, four
13 swabs were mock-exposed for several seconds during each sampling procedure and processed together
14 with the samples from this study. All negative control samples were below 1 ng/μl after PCR and
15 purification and, therefore, mock samples were not sent for sequencing. Samples were submitted for
16 indexing and pair-end 2x250 bp sequencing (Reagent Kit v2) on the Illumina MiSeq platform (San
17 Diego, USA).

18 Inferring sequence variants (SVs) of 16S amplicon sequencing data. Reads were analysed using the
19 *dada2* package version 1.5.0 and *workflow* in R version 3.1.2. as previously described [4] In brief,
20 forward and reverse reads were trimmed at 200 bp and at 150 bp to remove low quality regions,
21 respectively. The amplicon errors were corrected using the *dada2* algorithm with default parameters.
22 SVs shorter than 245 or longer than 257 base pairs were removed as were chimeras. Taxonomy was
23 assigned using the *assignTaxonomy* function [5].

24 Alpha and beta diversity analyses of SVs (for the non-infection analysis). Alpha diversity was assessed
25 with the functions *estimate* (richness) and *diversity* (Shannon diversity indices [SDIs]). In brief, the
26 richness was calculated using the command *specnumber* in R, and the SDI via the *diversity* function. As
27 for analysing the SVs, the distance matrices (DM; square matrix containing the distances, taken pairwise
28 between two different samples) for Beta diversity analyses were calculated by unweighted Jaccard index
29 (presence-absence based) and weighted Ružička index (abundance-based) of dissimilarity, and the
30 distances were calculated using the *vegdist* function from the *vegan* package. The multivariate dispersion
31 of each sample group was determined by calculating the average distance (based on Jaccard and Ružička
32 indices) to the sample type's centroid using the *betadisper* function. Outputs were visualised in principal
33 component analyses (PCoA) (*plot* function) and beta dispersion boxplots (*boxplot* function).

34 For the purpose of additional differential analysis, amplicon SVs were inferred with *dada2* 1.16.0 and
35 contaminating sequences across the entire data set were identified with *decontam* 1.8.0 using
36 comparison to sequence prevalences in a set of negative controls that were collected at the same time

37 and treated the same way as the samples. Batch effect correction between cages was carried out
38 separately for each timepoint by batch mean centering (BMC) on centered log-ratio transformed (CLR)
39 abundance data. The corrected data sets underwent differential analysis for each time point separately
40 with the R package DESeq 1.28.1. (model: ~ cage + group, test: Wald).

41 A particular strength from our study is the parallel inclusion of culture data. Species from culture were
42 identified and semi-quantitatively counted. The resulting numbers were transferred to absolute and
43 relative values for comparison to sequencing data. Furthermore, we calculated the DM for absolute
44 (CFUs) and relative values by the weighted Ružička index (abundance-based analyses) using the *vegdist*
45 function. We then computed the PCoA plots as described for the SVs above. Linear correlation of paired
46 dissimilarity value of weighted (Ružička) and culture data (absolute or relative values) DM were
47 calculated using *aes* and *geom_smooth* (method = "lm") function of *ggplot2* in R.

48 Alpha and beta diversity analyses of SVs (for the infection analyses). Alpha diversity values of the
49 animals which were inoculated by any of the 4 different pneumococcal strains were pooled. Relative
50 abundances of bacterial families were also pooled and visualised as average relative abundances per
51 time point.

52 As for Beta diversity analyses average relative abundances of each SV for each of the four groups was
53 first calculated. Subsequently, we received the DM by weighted Ružička index (abundance-based)
54 calculation (*vegdist* function) and clustered the results to receive a dendrogram (*hclust* function). Results
55 were then visualised by a dendrogram (function as. Dendrogram) and a heat map of the most abundant
56 SVs (log10 values) using the *heatmap.2* function and *ggplot2* package in R.

57

58

- 59 1. **Kraemer JG, Ramette A, Aebi S, Oppliger A, Hilty M.** Influence of pig farming on the
60 human's nasal microbiota: The key role of the airborne microbial communities. *Appl Environ*
61 *Microbiol* 2018.
- 62 2. **Caporaso JG, Lauber CL, Walters WA, Berg-Lyons D, Lozupone CA et al.** Global
63 patterns of 16S rRNA diversity at a depth of millions of sequences per sample. *Proc Natl Acad Sci U S*
64 *A* 2011;108 Suppl 1:4516-4522.
- 65 3. **Biesbroek G, Sanders EA, Roeselers G, Wang X, Caspers MP et al.** Deep sequencing
66 analyses of low density microbial communities: working at the boundary of accurate microbiota
67 detection. *PLoS One* 2012;7(3):e32942.
- 68 4. **Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ et al.** DADA2: High-
69 resolution sample inference from Illumina amplicon data. *Nat Methods* 2016;13(7):581-583.
- 70 5. **Wang Q, Garrity GM, Tiedje JM, Cole JR.** Naive Bayesian classifier for rapid assignment
71 of rRNA sequences into the new bacterial taxonomy. *Appl Environ Microbiol* 2007;73(16):5261-5267.

72

73

74 **Supplementary figure legends**

75

76 **Supplementary Figure 1: Bacterial density in the oropharynx of mice.** PCR concentrations were
77 measured by the Bioanalyzer and are shown for CTRL and SMK exposed mice at indicated time
78 points. Box-plots indicate median and interquartile range with mean indicated as + and outliers are
79 shown. Data are from 8-10 mice per group sampled longitudinally

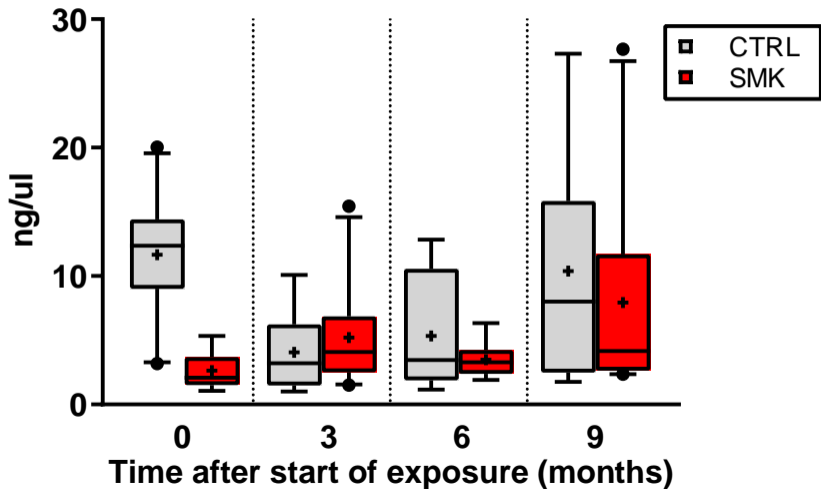
80 **Supplementary Figure 2: Differential abundance of SVs in SMK and CTRL groups prior to**
81 **smoke exposure.** DESeq analysis indicates the fold-change of SVs in bacterial communities of mice
82 prior to the start of smoke or air exposure ($p < 0.001$). Mice were assigned in systematic random
83 manner into cages and acclimatized for 4 weeks prior to the first oropharyngeal swab sampling (time
84 point 0) taken before the initiation of exposure to cigarette smoke or room air. The differential
85 abundances of SVs between SMK and CTRL groups with $p < 0.001$ only are presented. The taxa of the
86 SVs are labeled at genus level and the family belongings are indicated in the identical colors as in
87 Figure 2C and 2D. The data shown here includes bacterial families that were regrouped as 'others' and
88 shown in grey in in Figure 2C and 2D.

89 **Supplementary Figure 3: Differential abundance of SVs in SMK and CTRL groups after three**
90 **months.** DESeq analysis indicates the fold-change of SVs in bacterial communities of mice after three
91 months of smoke or air exposure ($p < 0.001$). The differential abundances of SVs between SMK and
92 CTRL groups with $p < 0.001$ only are presented. The taxa of the SVs are labeled at genus level and the
93 family belongings are indicated in the identical colors as in Figure 2C and 2D. The data shown here
94 includes bacterial families that were regrouped as 'others' and shown in grey in in Figure 2C and 2D.

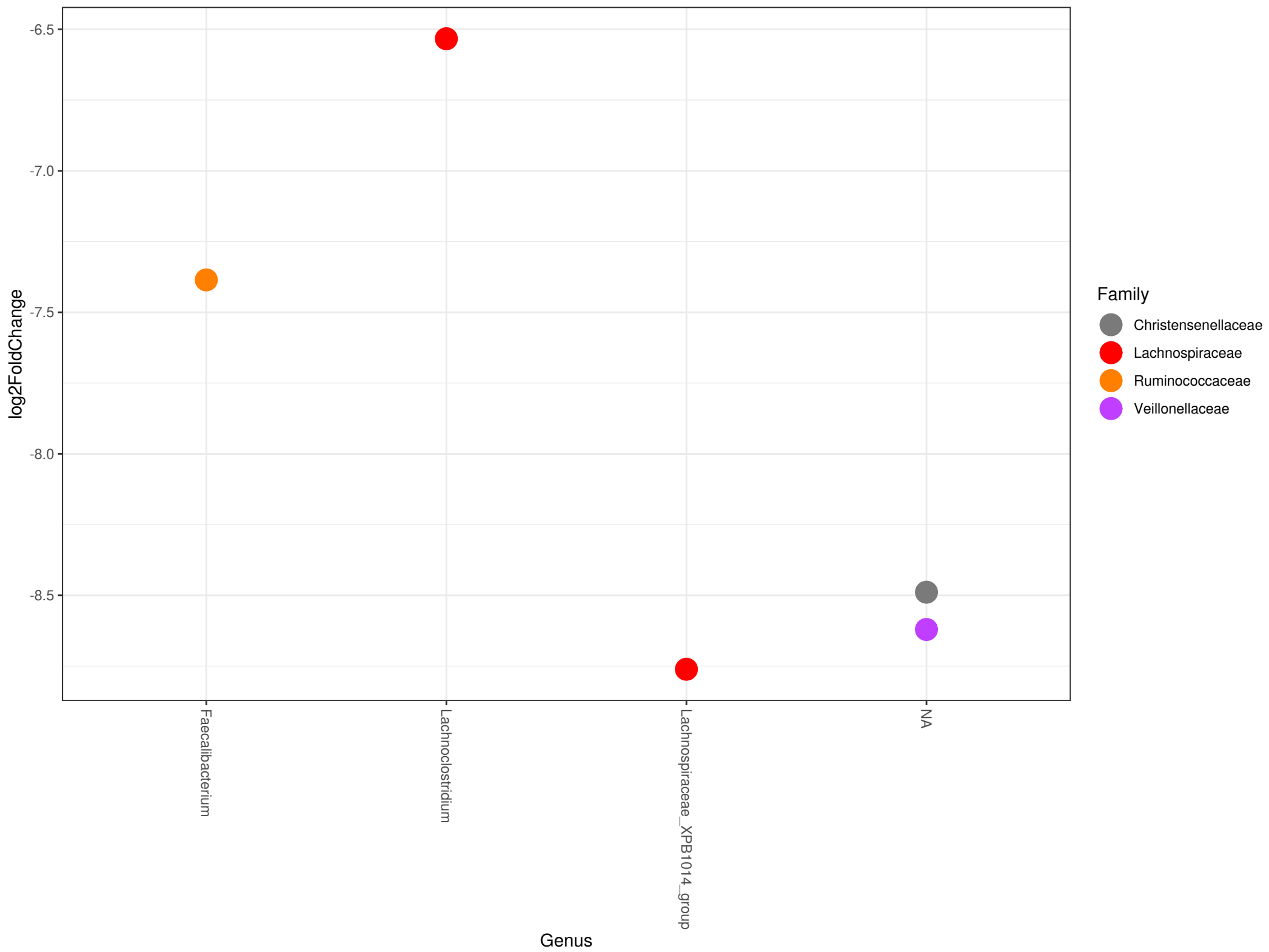
95 **Supplementary Figure 4: Differential abundance of SVs in SMK and CTRL groups after six**
96 **months.** DESeq analysis indicates the fold-change of SVs in bacterial communities of mice after six
97 months of smoke or air exposure ($p < 0.001$). The differential abundances of SVs between SMK and
98 CTRL groups with $p < 0.001$ only are presented. The taxa of the SVs are labeled at genus level and the
99 family belongings are indicated in the identical colors as in Figure 2C and 2D. The data shown here
100 includes bacterial families that were regrouped as 'others' and shown in grey in in Figure 2C and 2D.

101 **Supplementary Figure 5: Differential abundance of SVs in SMK and CTRL groups after nine**
102 **months.** DESeq analysis indicates the fold-change of SVs in bacterial communities of mice after nine
103 months ($p < 0.001$). One group of mice were smoke exposed for 6 months followed by a period of 3
104 months of smoke cessation (9 months in total) (SMK) while the control mice were air exposed for the
105 whole time for 9 months (CTRL). The differential abundances of SVs between SMK and CTRL
106 groups with $p < 0.001$ only are presented. The taxa of the SVs are labeled at genus level and the family
107 belongings are indicated in the identical colors as in Figure 2C and 2D. The data shown here includes
108 bacterial families that were regrouped as 'others' and shown in grey in in Figure 2C and 2D.

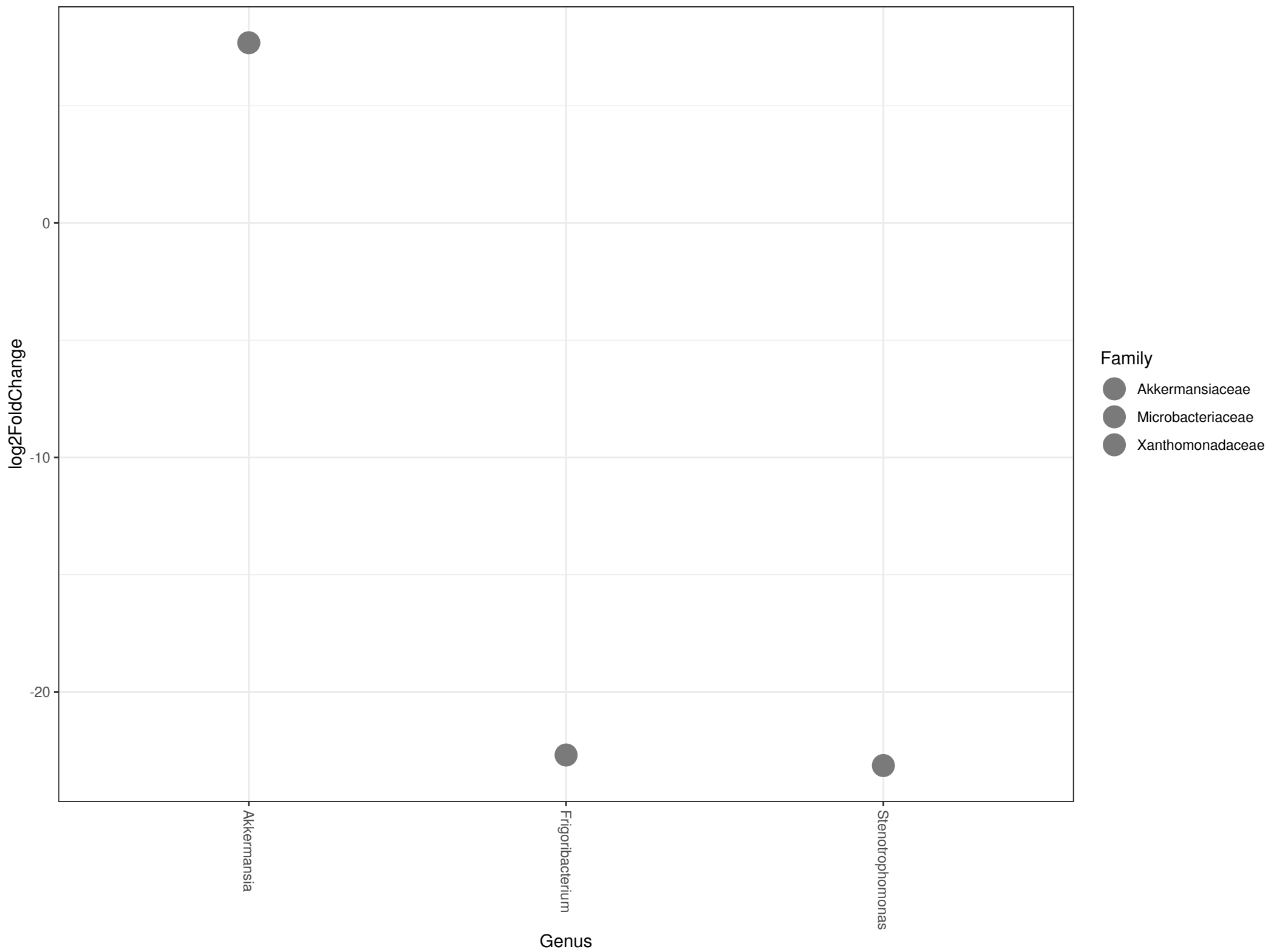
Supplementary figure 1



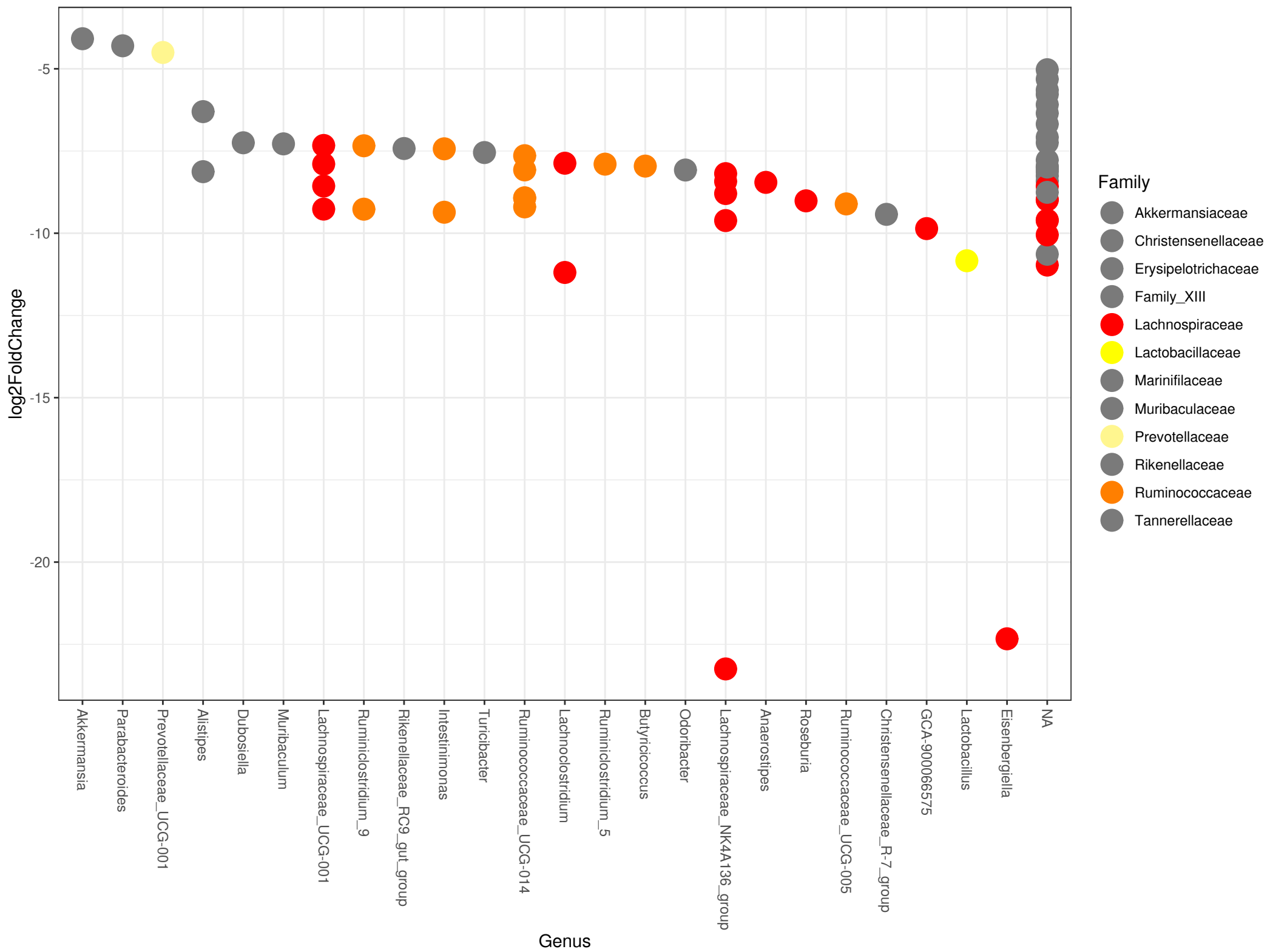
Supplementary figure 2: Differential abundance at time point 0 (adj. p-value = 0.001)



Supplementary figure 3: Differential abundance at time point 3 months (adj. p-value = 0.001)



Supplementary figure 4: Differential abundance at time point 6 months (adj. p-value = 0.001)



Supplementary figure 5 Differential abundance at time point 9 months (adj. p-value = 0.001)

