

# Microbiota long-term dynamics and prediction of acute graft-versus-host-disease in pediatric allogeneic stem cell transplantation

**Anna Cécilia Ingham**

Technical University of Denmark

**Katrine Kielsen**

Rigshospitalet

**Hanne Mordhorst**

Technical University of Denmark

**Marianne Ifversen**

Rigshospitalet

**Frank M Aarestrup**

Technical University of Denmark

**Klaus G Müller**

Rigshospitalet

**Sünje Johanna Pamp** (✉ [sjpa@food.dtu.dk](mailto:sjpa@food.dtu.dk))

Technical University of Denmark <https://orcid.org/0000-0002-6236-1763>

---

## Research

**Keywords:** Holobiont, Gut, oral, and nasal microbiota, HSCT, acute GvHD, immune reconstitution, microbiome, antibiotics, amplicon sequence variants, machine learning, prediction

**Posted Date:** March 3rd, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-258775/v1>

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

[Read Full License](#)

---

1           **Microbiota long-term dynamics and prediction of acute graft-versus-host-**  
2           **disease in pediatric allogeneic stem cell transplantation**

3  
4  
5   **Anna Cäcilia Ingham<sup>1#</sup>, Katrine Kielsen<sup>2, 3</sup>, Hanne Mordhorst<sup>1</sup>, Marianne Ifversen<sup>3</sup>, Frank M.**  
6   **Aarestrup<sup>1</sup>, Klaus Gottlob Müller<sup>2, 3, 4</sup>, Sünje Johanna Pamp<sup>1\*</sup>**

7  
8  
9  
10  
11  
12   <sup>1</sup> Research Group for Genomic Epidemiology, Technical University of Denmark, Kongens Lyngby,  
13   Denmark.

14  
15   <sup>2</sup> Institute for Inflammation Research, Department of Rheumatology and Spine Disease, Copenhagen  
16   University Hospital, Rigshospitalet, Copenhagen, Denmark.

17  
18   <sup>3</sup> Department of Pediatrics and Adolescent Medicine, Copenhagen University Hospital Rigshospitalet,  
19   Copenhagen, Denmark.

20  
21   <sup>4</sup> Institute of Clinical Medicine, University of Copenhagen, Copenhagen, Denmark.

22  
23  
24  
25  
26   # Present address: Department of Bacteria, Parasites and Fungi, Statens Serum Institut, Copenhagen,  
27   Denmark.

28  
29   \*Correspondence: sjpa@dtu.dk  
30   Technical University of Denmark, Research Group for Genomic Epidemiology, 2800 Kongens Lyngby,  
31   Denmark.

32  
33  
34   **Keywords**

35   Holobiont; Gut, oral, and nasal microbiota; HSCT; acute GvHD; immune reconstitution; microbiome;  
36   antibiotics; amplicon sequence variants; machine learning; prediction.

41 **Abstract**

42

43 **Background**

44 Patients undergoing allogeneic hematopoietic stem cell transplantation (HSCT) exhibit changes in their  
45 gut microbiota and are experiencing a range of complications, including acute graft-versus-host disease  
46 (aGvHD). It is unknown if, when, and under which conditions a re-establishment of microbial and  
47 immunological homeostasis occurs. It is also unclear whether microbiota long-term dynamics occur at  
48 other body sites than the gut such as the mouth or nose. Moreover, it is not known whether the patients'  
49 microbiota prior to HSCT holds clues to whether the patient would suffer from severe complications  
50 subsequent to HSCT. Here, we take a holobiont perspective and performed an integrated host-microbiota  
51 analysis of the gut, oral, and nasal microbiotas in 29 children undergoing allo-HSCT.

52

53 **Results**

54 The bacterial diversity decreased in the gut, nose, and mouth during the first month and reconstituted  
55 again 1-3 months after allo-HSCT. The microbial community composition traversed three phases over one  
56 year. Distinct taxa discriminated the microbiota temporally at all three body sides, including *Enterococcus*  
57 spp., *Lactobacillus* spp., and *Blautia* spp. in the gut. Of note, certain microbial taxa appeared already  
58 changed in the patients prior to allo-HSCT as compared to healthy children. Acute GvHD occurring after  
59 allo-HSCT could be predicted from the microbiota composition at all three body sites prior to HSCT. The  
60 reconstitution of CD4+ T cells, T<sub>H</sub>17 and B cells was associated with distinct taxa of the gut, oral, and nasal  
61 microbiota.

62

63 **Conclusions**

64 This study reveals for the first time bacteria in the mouth and nose that may predict aGvHD. Monitoring  
65 of the microbiota at different body sites in HSCT patients, and particularly through involvement of samples  
66 prior to transplantation, may be of prognostic value and could assist in guiding personalized treatment  
67 strategies. The identification of distinct bacteria that have a potential to predict post-transplant aGvHD  
68 might provide opportunities for an improved preventive clinical management, including a modulation of  
69 microbiomes. The host-microbiota associations shared between several body sites might also support an  
70 implementation of more feasible oral and nasal swab sampling-based analyses. Altogether, the findings  
71 suggest that the microbiota and host factors together could provide actionable information to guiding  
72 precision medicine.

## 73 **Background**

74 In allogeneic hematopoietic stem cell transplantation (allo-HSCT), the infusion of donor derived stem cells  
75 is employed as a curative treatment for various types of hematologic and non-hematologic disorders [1].  
76 In allo-HSCT patients, the human gut microbiota changes subsequent to transplantation, which may in  
77 part be attributable to antimicrobial treatment and conditioning regimens [2–4]. Butyrate-producing  
78 bacteria affiliated with the order *Clostridiales* are depleted in the gut early after transplantation, while  
79 *Proteobacteria*, and *Lactobacillales* such as *Enterococcus* spp. expand, possibly due to both increased  
80 oxygen levels in the intestinal lumen in the absence of butyrate, and antimicrobial resistance [2–5].  
81 However, microbiota dynamics in HSCT patients have so far mainly been monitored in detail during the  
82 first month post HSCT and not over longer periods of time. Hence, it is unclear whether and when the  
83 microbiota re-establishes to similar microbial community structures as prior to HSCT.

84 Conditioning-induced intestinal epithelial permeability might promote bacterial translocation and  
85 bacteremia [6]. This is recognized as the initial step in the pathogenesis of acute graft-versus-host disease  
86 (aGvHD) [7]. Acute GvHD is a common side effect of allo-HSCT in which alloreactive donor T cells exhibit  
87 cytotoxic activity against healthy tissue in the host, including the gut epithelium [7]. Acute GvHD severity  
88 can be distinguished in four grades dependent on the extent of organs affected: Grade 0-I presents as no  
89 or mild, and grade II-IV as moderate to severe aGvHD. Recently, studies have suggested that a lower gut  
90 microbiota diversity is associated with aGvHD and aGvHD-related mortality and that certain bacterial taxa  
91 dominating post HSCT may be involved in promoting aGvHD [3,8–12]. However, it has not been examined  
92 whether microbiota composition prior to HSCT has a predictive value in forecasting possible aGvHD  
93 severity, and which is addressed in the present study.

94 The microbiota exerts immunomodulatory function on the host's adaptive immune system, for example  
95 on T cells [13]. For instance, human commensal gut strains affiliated with *Bacteroides* and *Clostridia* can  
96 induce T regulatory ( $T_{reg}$ ) cells in germ-free mice [14]. Recent findings suggest that functionally different  
97 T cell subsets, such as T helper 17 ( $T_H17$ ) and  $T_{reg}$  cells are involved in the pathogenesis of aGVHD [15–  
98 17]. The microbiota at body sites other than the gut, such as the oral and nasal cavities, have also been  
99 suggested to be involved in immunomodulation [18]. We have previously proposed that the gut  
100 microbiota is associated with immune cell reconstitution after allo-HSCT [4]. However, it is unknown if the  
101 microbiotas at other mucosal sites are affected by allo-HSCT, whether they are associated with aGvHD,  
102 and whether they are associated with recovery of the patients' immune system.

103 Here, we monitored the microbiota dynamics in the gut, oral, and nasal cavities in pediatric allogeneic  
104 HSCT patients over a period of one year. At all three body sites, we identify distinct temporal bacterial  
105 abundance trajectories. In a machine learning approach, we predict aGvHD severity from pre-transplant  
106 microbiotas in the gut, oral, and nasal cavities which may be useful for early preventive managements in  
107 the clinical setting. By relating the microbiota composition to immune cell counts, inflammation and  
108 infection markers, antibiotic treatment, clinical outcomes, and patients' baseline parameters, we uncover  
109 similarities in host-microbial associations at different body sites.

110

## 111 **Results**

112 We characterized long-term microbiota dynamics in pediatric allo-HSCT at three body sites: the gut, and  
113 oral and nasal cavities (Figure 1). Fecal samples, buccal swabs, and anterior naris swabs were collected

114 from 29 children at 10 time points over a one-year period: Twice prior to HSCT, on the day of HSCT, weekly  
115 during the first month after HSCT, and at three follow-up time points up to twelve months post HSCT  
116 (Figure 1). Microbial community dynamics in these samples were determined by 16S rRNA gene profiling.  
117 A total of 709 patient samples (212 fecal samples, 248 oral swabs, and 249 nasal swabs from 10 time  
118 points) were characterized. Upon sequence filtering (see Methods), we retained 2465 ASVs for the fecal,  
119 377 ASVs for the oral, and 197 ASVs for the nasal core microbiota sets. We predicted the development of  
120 aGvHD severity from pre-transplant gut, oral, and nasal microbial abundances using machine learning. In  
121 addition, we assessed multivariate associations between the microbiota at the different body sites and  
122 immune reconstitution, immune markers, and clinical outcomes. Immune reconstitution was determined  
123 through quantitative measurements of T, B, and NK cells, and other leukocyte subpopulations in  
124 peripheral blood (Figure 1). We assessed systemic inflammation through levels of C-reactive protein (CRP),  
125 and measured procalcitonin as an approximation of infection (Figure 1, see Methods).

126

### 127 **Patient cohort and outcomes**

128 The 29 children had a median age of 8.2 years (range: 2.5-16.4) at the time of HSCT. Nine patients (31%)  
129 had no or mild aGvHD (grade 0 or I) and 20 patients (69%) developed moderate to severe aGvHD (grade  
130 II-IV) at median +14 days following HSCT (range: day +9 to day +61) (Supplementary Table S1, Additional  
131 file 1; and <https://doi.org/10.6084/m9.figshare.13567502>). The main organs involved in aGvHD included  
132 the skin (all), intestinal tract (n=3), and the liver (n=2). During the follow-up period of 21.4 months on  
133 average (range: 10.1 – 32.7 months), two patients (7%) relapsed and one patient underwent a donor  
134 lymphocyte infusion. Three patients (10%) died (one relapse-related death at day +91 and two treatment-  
135 related deaths at days +111 and +241, respectively). Due to their low incidence, we did not focus our  
136 analysis on relapse and mortality. For 25 patients (86.2%)  $\geq 1$  bacterial infection indicated by positive  
137 microbial culture was reported throughout the monitored period. All patients were treated  
138 prophylactically with trimethoprim and sulfamethoxazole prior to HSCT. In cases of fever or clinical signs  
139 of infections, antibiotic treatment with meropenem (28 patients), vancomycin (24 patients), ciprofloxacin  
140 (20 patients), phenoxymethylpenicillin (14 patients), or other antibiotics was commenced according to  
141 culture-based results or clinical presentation.

142

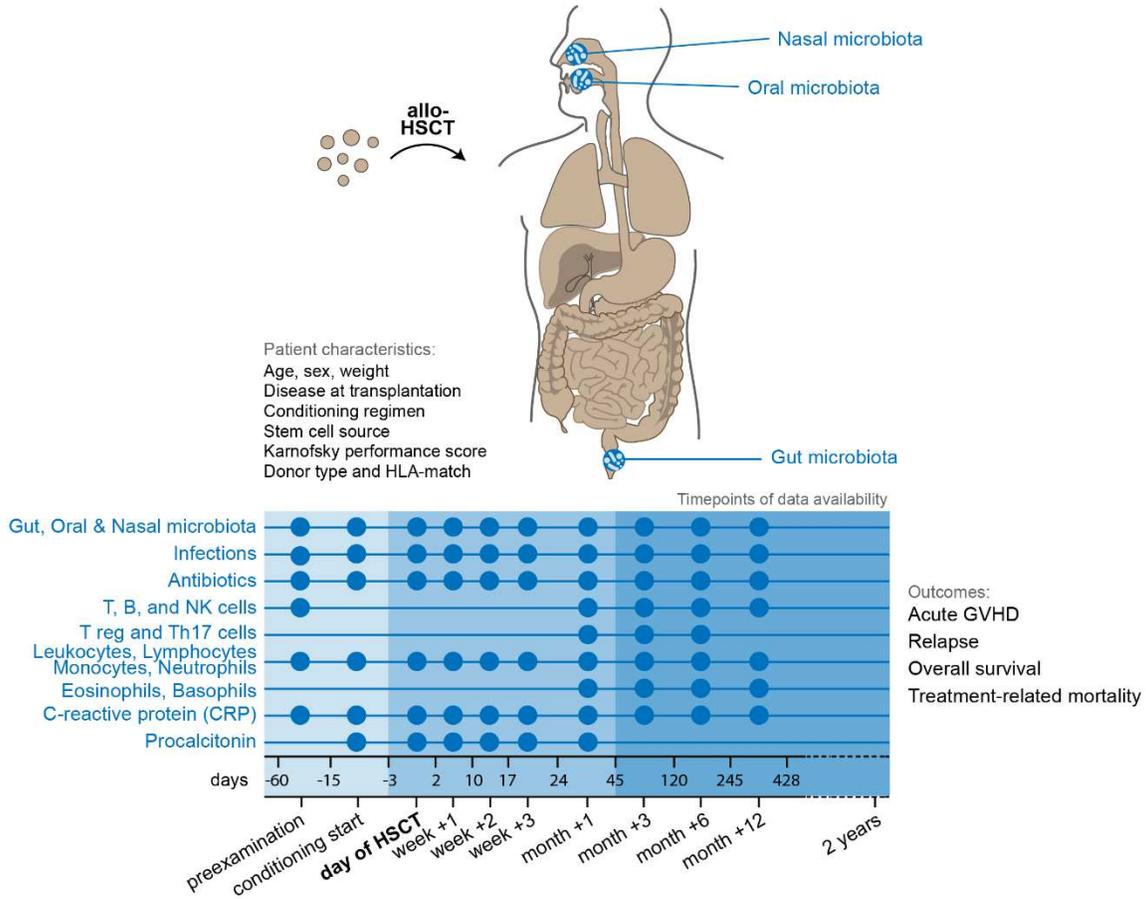
### 143 **Bacterial alpha diversity decreases in relation to allo-HSCT at all three body sites**

144 Alpha diversity (Inverse Simpson) in the gut was overall the highest, followed by the oral cavity, and the  
145 nose (Figure 1B). The lowest alpha diversity was observed within the first month post HSCT for all three  
146 body sites. However, the exact time points were somewhat different for each body site: the day of HSCT  
147 to week +3 for the gut, week +3 for the oral cavity, and week +1 for the nasal cavity. The decrease in  
148 microbial diversity was significant for the nasal cavity, where the median alpha diversity decreased from  
149 4.43 at the start of conditioning to 2.65 in week +1 ( $P = 0.02$ ) (Figure 1B). Alpha diversity increased again  
150 at all body sites thereafter. However, alpha diversity was lower again at month +12 in the nasal cavity.

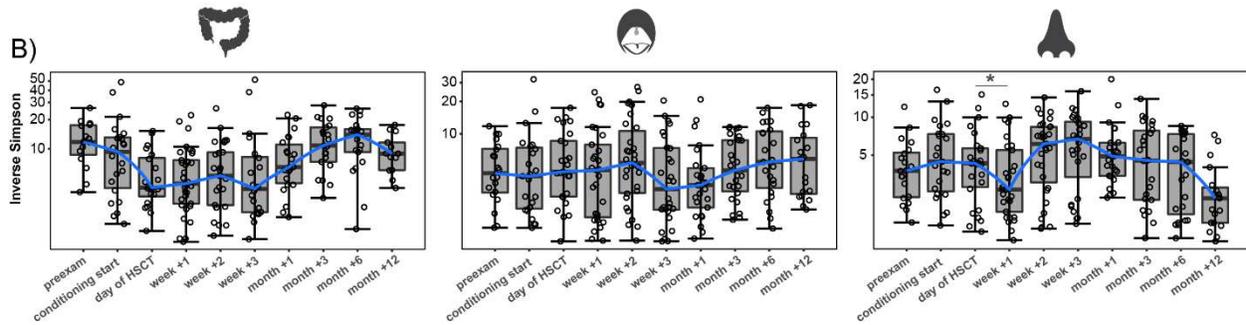
151

152

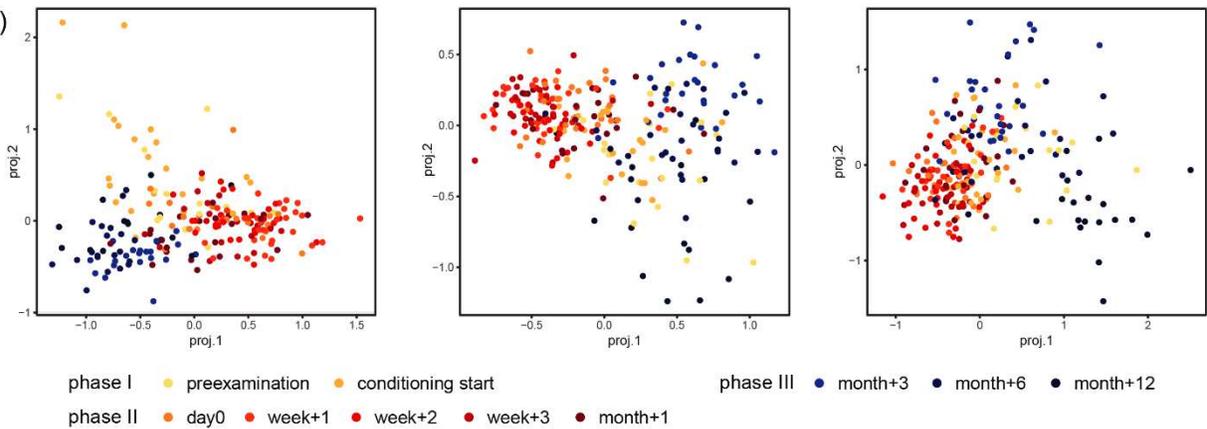
A)



B)



C)



154 **Figure 1. Monitoring gut, oral, and nasal microbiota and the host immune system in allogeneic hematopoietic**  
155 **stem cell transplantation (HSCT).** A) Twenty-nine children were monitored before, at the time of, and immediately  
156 post allogeneic HSCT, as well as at late follow-up time points. Patients' baseline characteristics, clinical outcomes, as  
157 well as immune cell counts, and inflammation and infection markers over time were monitored. Patient  
158 characteristics are described in detail in Table S1 (Additional File 1). Host immune system parameters were related  
159 to longitudinal dynamics of the gut, oral, and nasal microbiota that was assessed at the denoted time points. B)  
160 Bacterial alpha diversity before, at the time of, and after HSCT at each body site, displayed on a log<sub>10</sub> transformed  
161 y-axis for visualization purposes. Asterisks indicate significant differences in median inverse Simpson index between  
162 time points \*  $P < 0.05$ . C) Tree-based sparse linear discriminant (LDA) analyses by time point in relation to HSCT. For  
163 fecal samples, the positive LDA scores were observed for samples collected immediately post HSCT. For both oral  
164 and nasal samples, the positive LDA scores were observed for samples from before HSCT and from late follow up-  
165 time points.  
166

### 167 **Microbial community composition in patients prior to HSCT differs from healthy controls**

168 We hypothesized that the bacterial alpha diversity at the first sampling time point (preexamination) might  
169 already be lower in these patients as compared to age-matched healthy children due to the treatment  
170 given prior to the referral to allo-HSCT and enrolment in this study. To assess this, we compared the gut  
171 microbiota at preexamination to that of healthy children (median age 6.8 years) [19]. As expected, the  
172 alpha diversity was 2.4-fold lower in the patients at preexamination (median InvSimpson 11.7) as  
173 compared to the healthy children (median InvSimpson 28.2) (Supplementary Figure S1A, Additional File  
174 2). Bacterial composition differed between the two groups (anosim,  $p=0.001$ ,  $R=0.44$ , Figure S1B). This  
175 difference was to a certain extent due to a larger variation within the HSCT group (betadisper,  $p<0.001$ )  
176 (Supplementary Figure S1 B, Additional File 2). Through linear discriminant analysis (LEfSe) and differential  
177 abundance analysis (DeSeq2), we found taxa that were significantly more abundant in the patients already  
178 at preexamination as compared to the healthy controls: these included *Bacilli* (e.g. *Lactobacillus*,  
179 *Enterococcus*), *Erysipelotrichaceae*, and *Enterobacteriaceae* (e.g. *Klebsiella*). In contrast, certain taxa were  
180 more abundant in the healthy children, such as *Prevotella*, *Ruminococcaceae* (e.g. *Ruminococcus*), and  
181 *Akkermansia*, as compared to the patients at preexamination (Supplementary Figure S1 C and D,  
182 Additional File 2; and <https://doi.org/10.6084/m9.figshare.13614230>).  
183

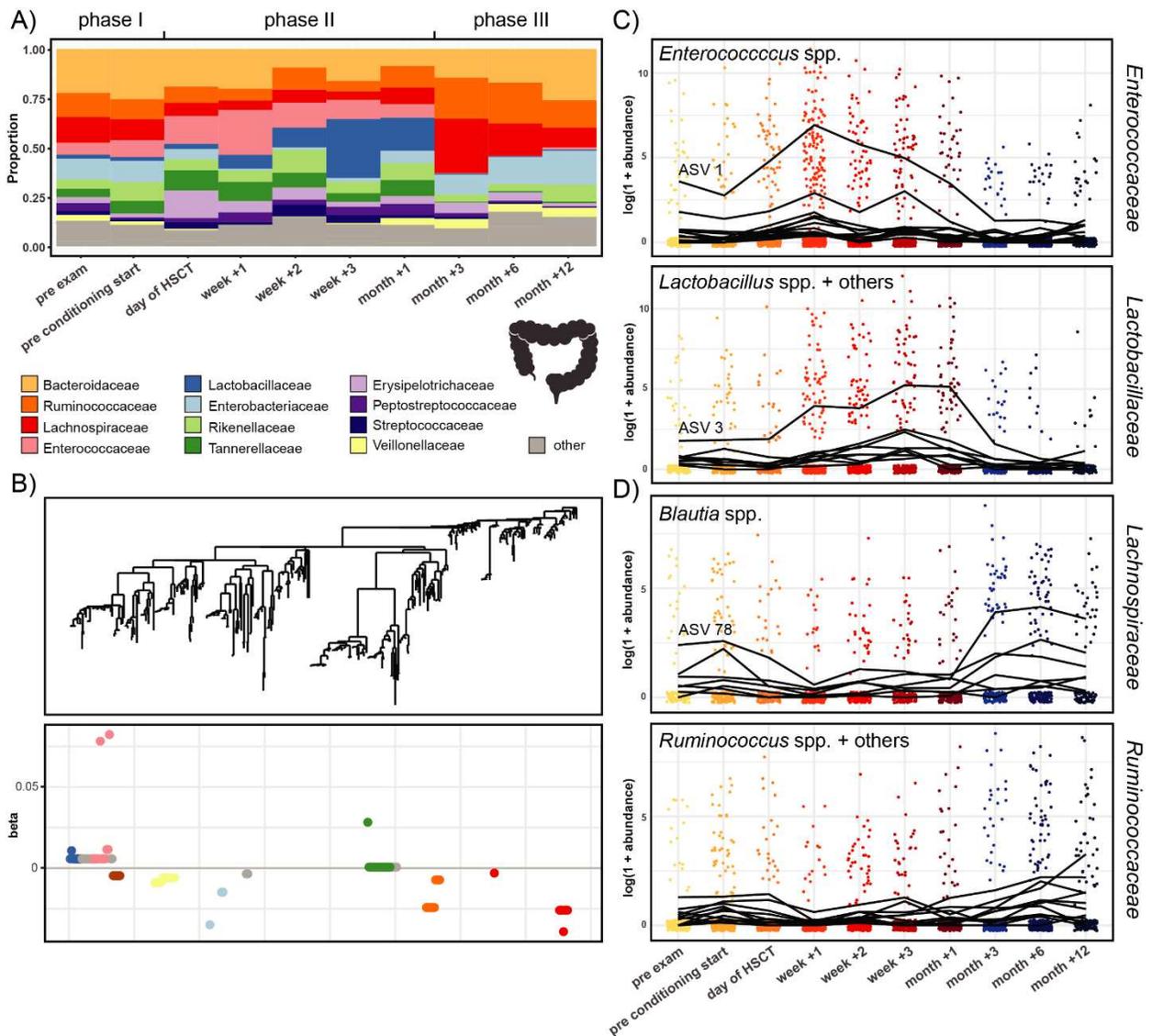
### 184 **Temporal microbial community dynamics appear in three interlaced phases over one year**

185 For a more detailed assessment of gut, oral, and nasal ASVs that best characterized samples from different  
186 time points, we performed tree-based sparse linear discriminant analyses (LDA). We observed at all three  
187 body sites that samples divided into three partly interlaced phases; phase I: samples at pre-examination  
188 and conditioning start, phase II: day of HSCT to month +1, and phase III: month +3 to month +12 (Figure  
189 1C). Interestingly, samples from phase I and III overlapped for the oral and nasal cavities, suggesting a  
190 possible return of microbial communities from later time points to a state similar to before HSCT. Of note,  
191 the nasal community composition at month +12, that exhibited low alpha diversity, was different from  
192 samples of week +1 (phase II) that also exhibited low alpha diversity (Figure 1B and C).

193 To get a more detailed view of the microbial abundance dynamics, we examined the 12 most abundant  
194 families at each body site, respectively (Figures 2A, 3A, Additional File 2: Figures S2 and S3A). In the gut,  
195 we observed a reduction in *Lachnospiraceae* in phase II, immediately after HSCT, from 13% at pre-  
196 examination to 4.7% in week +1, followed by a recovery to 27.5% in month +3 at the start of phase III  
197 (Figure 2A). Concurrently, an expansion of *Enterococcaceae* in phase II (pre-examination: 6.1%; week +1:

218 22.8%) and *Lactobacillaceae* in phase II (pre-examination: 2%; week +1: 7%) occurred, followed by a  
 219 reduction in phase III from month +3 onwards to 0.2% and 0.6%, respectively (Figure 2A).  
 220 In the oral cavity, we observed a reduced relative abundance of *Actinomycetaceae* for several time points  
 221 in phase II as compared to the time points in phase I (prior to HSCT) and at later follow-up time points.  
 222 For example, *Actinomycetaceae* abundances were 9.7% at pre-examination and 2.9% in week +3 (Figure  
 223 3A). Furthermore, *Streptococcaceae* abundances were lower from the day of HSCT until week +2  
 224 compared with before HSCT and late follow-up time points (pre-examination: 44.6%; week +1: 23.3%;  
 225 month +3: 51.3%, Figure 3A).  
 226 In the nasal cavity, we observed a reduced relative abundance of *Corynebacteriaceae* and *Moraxellaceae*  
 227 at most time points in phase II, as compared to samples from phase I and III (Additional File 2: Figure S3).  
 228 For example, *Corynebacteriaceae* abundances were 28.7% at pre-examination and 0.7% in week +1  
 229 (Additional File 2: Figure S3).

210  
 211



212

213 **Figure 2. Temporal microbial community dynamics in the gut.** A) Relative abundances over time of the 12 most  
214 abundant families in the gut. B) Tree-based sparse linear discriminant analysis (LDA). Coefficients of discriminating  
215 clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the phylogenetic tree. C)  
216 Trajectories of ASVs affiliated with the families *Enterococcaceae* and *Lactobacillaceae*, with increasing abundances  
217 after HSCT. The most abundant discriminating ASV for each family is indicated. D) Trajectories of ASVs affiliated with  
218 the families *Lachnospiraceae* and *Ruminococcaceae*, with decreasing abundances after HSCT and recovery at late  
219 follow-up time points. The most abundant discriminating ASV for *Blautia* spp. is indicated. Detailed taxonomic  
220 information and LDA-coefficients of the displayed ASVs are listed in Additional File 1: Table S2.

221

### 222 **Distinct *Enterococcus*, *Lactobacillus*, and *Blautia* lineages discriminate the gut microbiota temporally**

223 In order to determine which specific taxa in the gut were driving the differences between samples in the  
224 LDA (Figure 1C), we examined the individual discriminating ASVs. In general, in tree-based sparse LDA,  
225 ASVs with positive LDA coefficients are overrepresented in samples with positive LDA scores, while ASVs  
226 with negative LDA coefficients likewise are associated with samples with negative LDA scores (Figures 1C,  
227 2B, 2C, and 2D). The LDA revealed 19 clades (total 102 ASVs) in the gut that best separated samples by  
228 time point (Figure 2B). The two most discriminating clades with positive LDA-coefficients comprised ASVs  
229 of the family *Enterococcaceae* and *Lactobacillaceae* (Figure 2B). The ASVs of these two clades increased  
230 in abundance from the day of HSCT (*Enterococcaceae*) and week +1 (*Lactobacillaceae*), respectively, in  
231 support of the family abundances and in line with the positive LDA scores of phase II samples (Figures 2A,  
232 2C, and 1C). Of note, the order *Lactobacillales* and genus *Lactobacillus* (family *Lactobacillaceae*) appeared  
233 already to be higher at pre-examination as compared to healthy children (Supplementary Figure S1D,  
234 Additional File 2). From month +3 onwards, their abundances decreased again to levels comparable to the  
235 time of pre-examination (i.e. pre-treatment) (Figure 2C). All members of the *Enterococcaceae* clade, with  
236 the exception of one ASV, were *Enterococcus* spp. (Additional File 1: Table S2). The most abundant and  
237 most frequently observed *Enterococcus* was ASV 1 (Figure 2C and Additional File 1: Table S2). More  
238 detailed sequence analysis of the partial 16S rRNA gene sequence using SINA and BLAST alignments  
239 revealed that it belonged to the *Enterococcus faecium* group. The most abundant and most frequently  
240 observed *Lactococcus* was ASV 3 (Figure 2C and Additional File 1: Table S2), and its partial 16S rRNA gene  
241 sequence exhibited a high sequence similarity to *Lactobacillus rhamnosus*.

242 The two most discriminative clades with negative LDA-coefficients included two individual ASVs and one  
243 clade of the *Lachnospiraceae* family, and two *Ruminococcaceae* clades (Figure 2B, Additional File 1: Table  
244 S2). The abundances of these ASVs decreased in week +1 and recovered from month +3 onwards,  
245 returning to abundances comparable with before HSCT or higher (Figure 2D), in agreement with the  
246 abundance patterns for those families (Figure 2A). Of note, the family *Ruminococcaceae* appears already  
247 to be lower at pre-examination as compared to healthy children (Supplementary Figure S1 C and D,  
248 Additional File 2). All ASVs within the *Lachnospiraceae* group belonged to the genus *Blautia* (Additional  
249 File 1: Table S2). The most abundant and most frequently observed *Blautia* was ASV 78 (Figure 2D and  
250 Additional File 1: Table S2), and its partial 16S rRNA gene sequence exhibited a high sequence similarity  
251 to *Blautia wexlerae*.

252

### 253 **Distinct *Actinomyces* and *Streptococcus* lineages discriminate the oral microbiota temporally**

254 The tree-based sparse LDA identified 10 clades of in total 71 ASVs in the oral cavity that best separated  
255 samples by time points along the first axis (Figure 3B). The two largest discriminating groups of ASVs were

256 affiliated with *Actinomycetaceae* and *Streptococcaceae* (Figure 3B, Additional File 1: Table S2). The most  
257 abundant and among the most frequently observed ASVs were *Actinomyces* ASV 18 and *Streptococcus*  
258 ASV 28 (Figure 3C and Additional File 1: Table S2), and their partial 16S rRNA gene sequence exhibited a  
259 high sequence similarity to the *Actinomyces viscosus* and *Streptococcus mitis* groups, respectively.  
260 Additional discriminating ASVs were affiliated with *Prevotellaceae*, and *Bacillales* Family XI (*Gemella* spp.),  
261 respectively. The most abundant and frequently observed ASVs were affiliated with *Prevotella*  
262 *melaninogenica* (ASV 42) and *Gemella sanguis* (ASV 208). In agreement with the relative family abundance  
263 dynamics, these clades shared a pattern of depletion from the day of HSCT or week +1 onwards (phase  
264 II), until their abundances recovered from month +3 onwards (phase III) (Figures 3A and 3C) to an  
265 abundance similar to before HSCT, as observed for *Ruminococcaceae* and *Lachnospiraceae* in the gut.

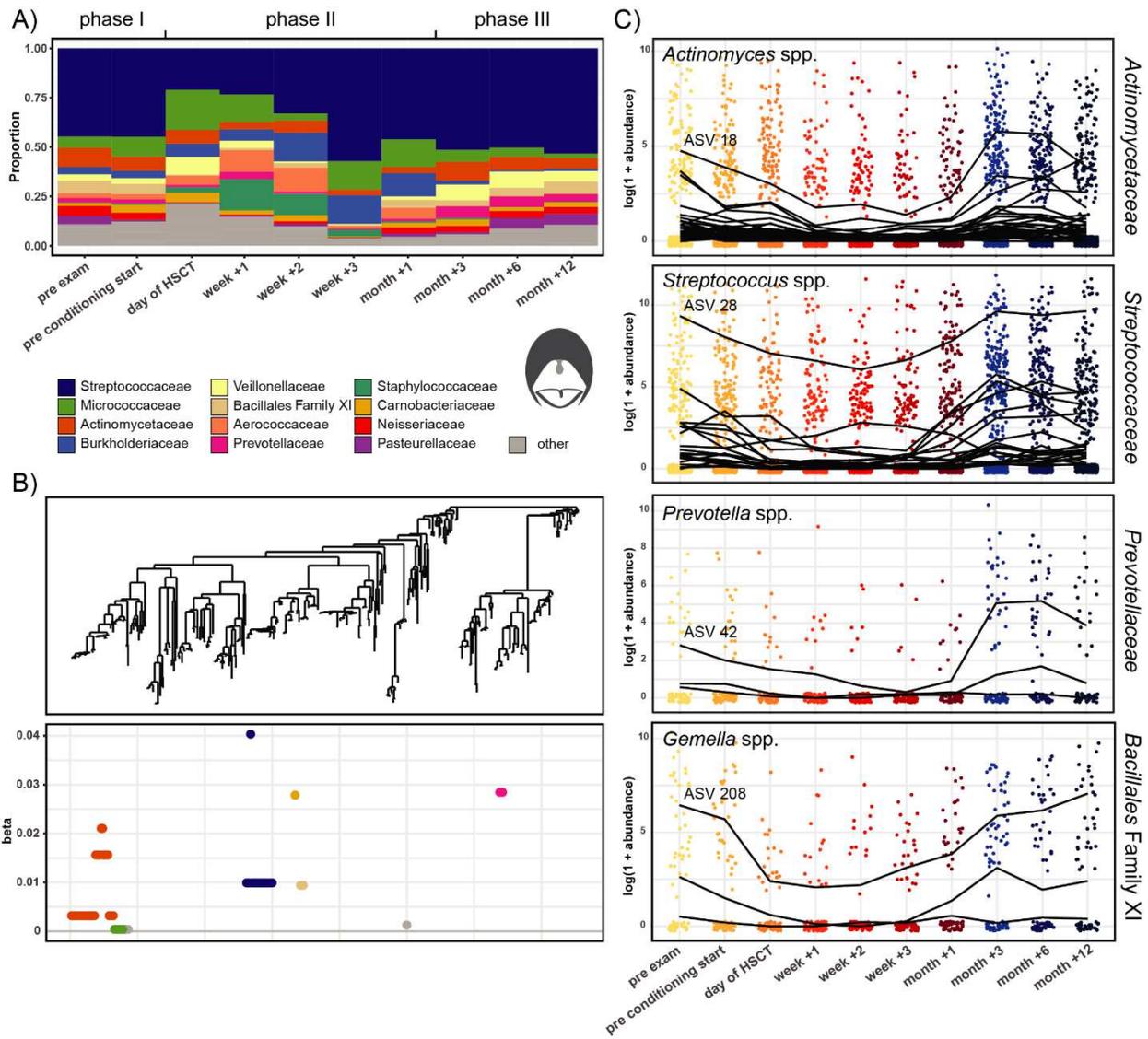
266

### 267 **Distinct *Corynebacteriaceae* and *Streptococcaceae* lineages discriminate the nasal microbiota** 268 **temporally**

269 The LDA revealed 30 discriminating nasal clades on axis 1 (comprising in total 36 ASVs), many of which  
270 consisted of individual ASVs (Additional File 2: Figure S3B). ASVs affiliated with the same family did not  
271 always covary in abundance. The *Corynebacteriaceae*, *Streptococcaceae*, and *Moraxellaceae* ASVs all had  
272 positive LDA-coefficients, i.e. their abundances decreased after HSCT and increased again from month+3  
273 onwards (Additional File 2: Figures S3B and S3C). The most abundant and most frequently observed  
274 *Corynebacteriaceae* was ASV 14 (Additional File 2: Figure S3C and Additional File 1: Table S2), and its  
275 partial 16S rRNA gene sequence exhibited a high sequence similarity to *Corynebacterium propinquum*.

276

277



278  
 279 **Figure 3. Temporal microbial community dynamics in the oral cavity.** A) Relative abundances over time of the 12  
 280 most abundant families in the oral cavity. B) Tree-based sparse linear discriminant analysis (LDA). Coefficients of  
 281 discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the phylogenetic  
 282 tree. C) Trajectories of ASVs affiliated with the families *Actinomycetaceae*, *Streptococcaceae*, *Prevotellaceae*, and  
 283 Family XI (Class *Bacillales*), with decreasing abundances after HSCT and recovery at late follow-up time points. The  
 284 most abundant discriminating ASV for each family is indicated. Detailed taxonomic information and LDA-coefficients  
 285 of the displayed ASVs are listed in Additional File 1: Table S2.

286  
 287

288 **Acute GvHD severity can be predicted from gut microbiota composition prior to HSCT**

289 To reveal potential associations between the gut microbiota and the severity of acute GvHD, we examined  
 290 the 12 most abundant families at each body site in patients with no or mild (grade 0-I) and moderate to  
 291 severe (grade II-IV) aGvHD. In the gut, *Tannerellaceae* were less abundant at time points before HSCT in  
 292 patients with grade 0-I compared to grade II-IV, especially at pre-examination and at start of conditioning  
 293 (Figure 4A). In order to predict aGvHD (grade 0-I versus grade II-IV) from microbial abundances at time  
 294 points up until the time of stem cell infusion, we implemented machine-learning models (see Methods –

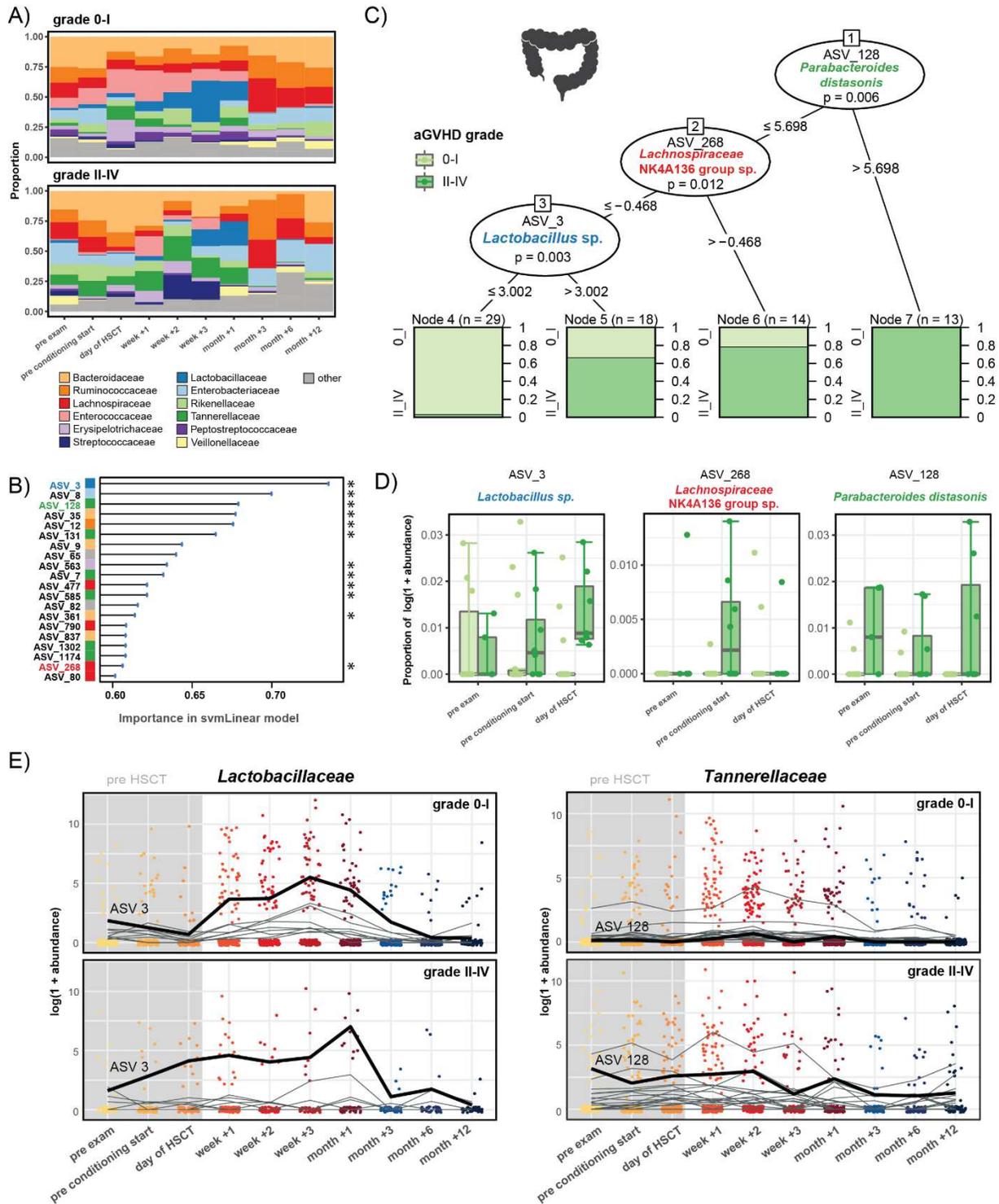
295 Statistical analysis). This analysis revealed 3 significant predictive ASVs in the gut: ASV 128 (*Parabacteroides*  
296 *distasonis*, Tannerellaceae,  $P < 0.01$ ), ASV 268 (*Lachnospiraceae* NK4A136 group sp., Lachnospiraceae,  $P$   
297 = 0.01) and ASV 3 (*Lactobacillus* sp., Lactobacillaceae,  $P < 0.01$ ) (Figures 4B and 4C, and Additional File 1:  
298 Table S3). This means, high abundances of these ASVs before HSCT were associated with the subsequent  
299 development of aGvHD grade II-IV post HSCT (Figure 4C). For instance, all pre-transplant samples with a  
300 variance stabilized abundance  $>5.7$  of ASV 128 (*Parabacteroides distasonis*) and 67% with a variance  
301 stabilized abundance  $>3$  of ASV 3 (*Lactobacillus* sp.) originated from patients who later developed aGvHD  
302 grade II-IV (Figure 4C). In agreement, log transformed relative abundances of these ASVs were mostly  
303 higher at pre-examination, conditioning start, and the day of HSCT in patients who later developed aGvHD  
304 grade II-IV compared with those exhibiting grade 0-I (Figure 4D). For instance, the average abundance of  
305 ASV 128 (*Parabacteroides distasonis*) was 5.5 times higher at pre-examination in grade II-IV versus in  
306 grade 0-I patients (Figure 4D). The temporal trajectory of ASV 3 (*Lactobacillus* sp.) also revealed a higher  
307 abundance at time points up to the transplantation in patients with grade II-IV aGvHD compared to those  
308 with grade 0-I (Figure 4E). Within the *Lactobacillaceae* identified by the LDA, this pattern seemed to be  
309 restricted to ASV3 (Figure 4E). ASV 128 (*Parabacteroides distasonis*) was part of the discriminating group  
310 of *Tannerellaceae* identified in the LDA (Figure 4E, and Additional File 1: Table S3). Its trajectory faceted  
311 by aGvHD severity confirmed the observation of increased pre-HSCT abundances in patients with  
312 subsequent development of aGvHD grade II-IV (Figure 4E).

313

#### 314 **Acute GvHD severity can be predicted from oral microbiota composition prior to HSCT**

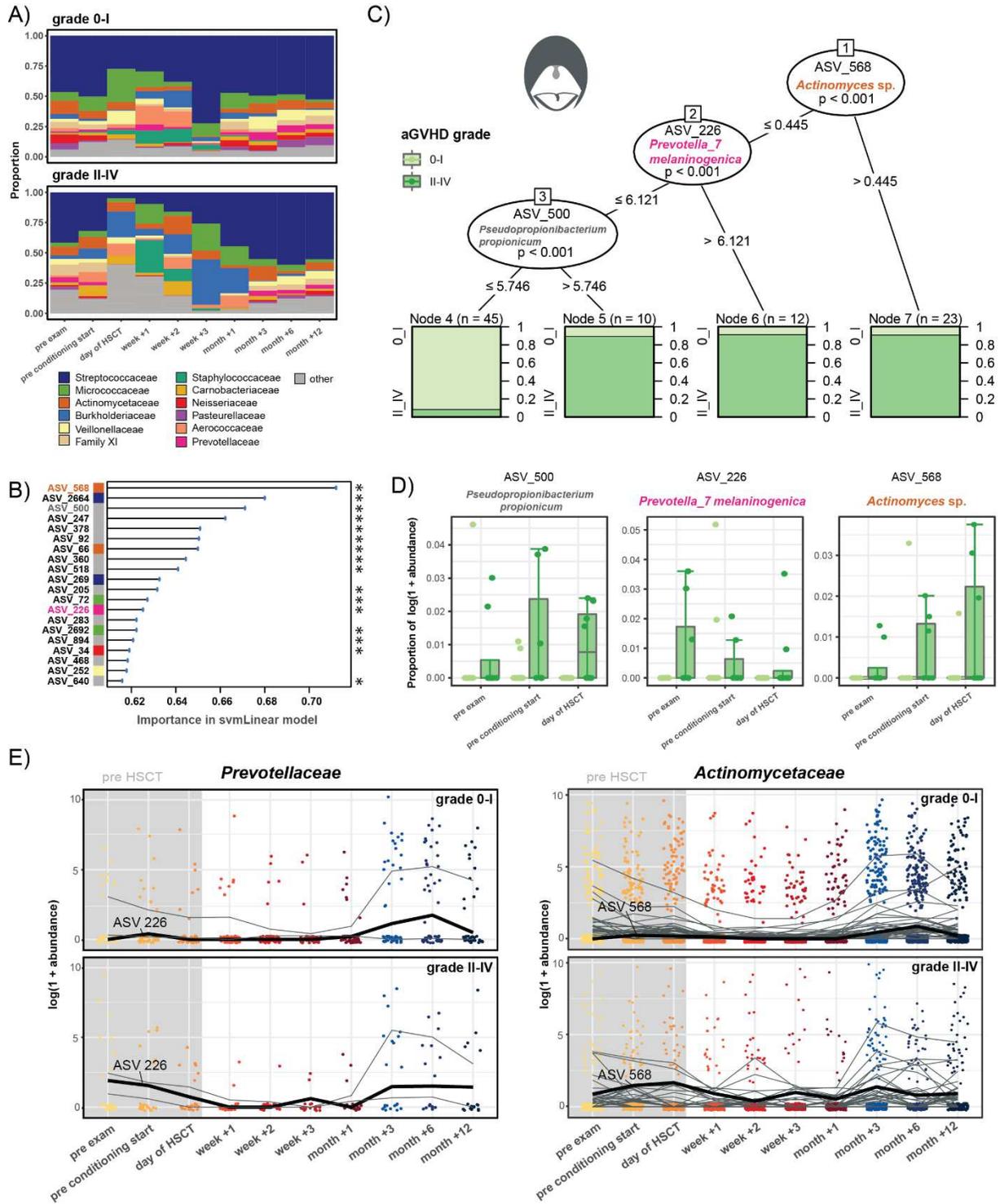
315 In the oral cavity, the bacterial community before HSCT in patients with grade II-IV aGvHD was  
316 characterized by a lower relative abundance of *Neisseriaceae*, and higher relative abundances of  
317 *Aerococcaceae* and *Prevotellaceae*, compared with grade 0-I aGvHD, especially at pre-examination and  
318 conditioning start (Figure 5A). Our machine learning approach predicted aGvHD severity (grade 0-I versus  
319 II-IV) from the abundances of 3 significant oral ASVs pre-HSCT: ASV 568 (*Actinomyces* sp.,  
320 Actinomycetaceae,  $P < 0.001$ ), ASV 226 (*Prevotella melaninogenica*, Prevotellaceae,  $P < 0.001$ ) and ASV  
321 500 (*Pseudopropionibacterium propionicum*, Propionibacteriaceae,  $P < 0.001$ ) (Figures 5B and 5C, and  
322 Additional File 1: Table S3). High abundances of these ASVs before transplantation predicted the  
323 development of aGvHD grade II-IV after HSCT (Figure 5C). For instance, 91% of samples with a variance  
324 stabilized abundance  $>0.4$  of ASV 568 (*Actinomyces* sp.) and 92% of samples with a variance stabilized  
325 abundance  $>6.1$  of ASV 226 (*Prevotella melaninogenica*) originated from patients with subsequent  
326 development of aGvHD grade II-IV (Figure 5C). In support, pre-HSCT log transformed relative abundances  
327 of these ASVs were higher in those patients. For example, the median relative abundance of ASV 500  
328 (*Pseudopropionibacterium propionicum*) on the day of HSCT was 10 times higher in grade II-IV versus in  
329 grade 0-I patients (Figure 5D). Temporal trajectories of oral *Actinomycetaceae* and *Prevotellaceae*,  
330 identified also in the LDA, showed that the abundances of ASV 226 (*Prevotella melaninogenica*) and ASV  
331 568 (*Actinomyces* sp.) were higher at time points up to the transplantation in patients with grade II-IV  
332 versus those with grade 0-I (Figure 5E).

333



334  
 335 **Figure 4. Machine learning-based prediction of aGVHD severity from the pre-HSCT gut microbiota composition.** A)  
 336 Relative abundances of the 12 most abundant families over time in the gut in patients with aGVHD grade 0-I versus  
 337 II-IV. B) Importance plot of top 20 predictive gut ASVs identified by the svmLinear model with importance scores  
 338 indicating the mean decrease in prediction accuracy in case the respective ASV would be excluded from the model.  
 339 The final cross-validated svmLinear model predicted aGVHD (0-I versus II-IV) from the abundances of gut ASVs pre-  
 340 HSCT with 86% accuracy (95% CI: 65% to 97%). The ASVs that were also confirmed by Boruta feature selection are  
 341 indicated with asterisk. C) Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by

342 nonparametric regression for prediction of aGvHD. Numbers along the branches indicate split values of variance  
 343 stabilized bacterial abundances. The terminal nodes show the proportion of samples originating from patients (n =  
 344 number of samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depicting the log transformed relative abundances of  
 345 the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared with grade II-IV patients.  
 346 E) Trajectories of *Lactobacillaceae* and *Tannerellaceae* ASVs that were identified by tree-based sparse LDA, including  
 347 ASV 3 and ASV 128 that were predictive for aGvHD (bold lines), in patients with aGvHD grade 0-I vs II-IV.  
 348



349

350 **Figure 5. Machine learning-based prediction of aGvHD severity from the pre-HSCT oral microbiota composition.**  
351 A) Relative abundances the 12 most abundant families over time in the oral cavity in patients with aGvHD grade 0-I  
352 versus II-IV. B) Importance plot of top 20 predictive oral ASVs identified by the svmLinear model with importance  
353 scores indicating the mean decrease in prediction accuracy in case the respective ASV would be excluded from the  
354 model. The final cross-validated svmLinear model predicted aGvHD (0-I versus II-IV) from the abundances of oral  
355 ASVs pre-HSCT with 92% accuracy (95% CI: 73% to 99%). The ASVs that were also confirmed by Boruta feature  
356 selection are indicated with asterisk. C) Conditional inference tree (CTREE) displaying ASVs identified as significant  
357 split nodes by nonparametric regression for prediction of aGvHD. Numbers along the branches indicate split values  
358 of variance stabilized bacterial abundances. The terminal nodes show the proportion of samples originating from  
359 patients (n = number of represented samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depict the log transformed  
360 relative abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared  
361 with grade II-IV patients. E) Trajectories of *Prevotellaceae* and *Actinomycetaceae* ASVs that were identified by tree-  
362 based sparse LDA, including ASV 226 and ASV 568 that were predictive for aGvHD (bold lines), in patients with aGvHD  
363 grade 0-I vs II-IV.  
364  
365

### 366 **Acute GvHD severity can be predicted from nasal microbiota composition prior to HSCT**

367 The proportion of nasal *Neisseriaceae* prior to HSCT was higher in patients with aGvHD grade 0-I as  
368 compared to grade II-IV (Additional File 2: Figure S4A). In contrast, *Actinomycetaceae* and  
369 *Corynebacteriaceae* exhibited a higher abundance in aGvHD grade II-IV patients prior to HSCT compared  
370 to those with grade 0-I (Additional File 2: Figure S4A). We found two ASVs significantly predicting aGvHD  
371 grade with opposite effects, ASV 66 and ASV 47. A high pre-HSCT abundance of ASV 66 (*Actinomyces* sp.,  
372 *Actinomycetaceae*,  $P = 0.03$ ) predicted development of aGvHD grade II-IV. The partial 16S rRNA gene  
373 sequence of ASV 66 exhibited a high sequence similarity to *Actinomyces viscosus*. A total of 94% of  
374 samples with a variance stabilized abundance >6.4 of ASV 66 originated from patients with subsequent  
375 development of aGvHD grade II-IV (Additional File 2: Figures S4B and S4C). In support, pre-HSCT log  
376 transformed relative abundances of ASV 66 (*Actinomyces* sp.) were 2.3 times higher in patients with  
377 aGvHD grade II-IV compared to those with grade 0-I (Additional File 2: Figure S4C). In contrast, high pre-  
378 HSCT abundance of ASV 47 (*Rothia* sp.,  $P = 0.03$ ) predicted that patients would be spared from aGvHD.  
379 The partial 16S rRNA gene sequence of ASV 47 exhibited a high sequence similarity to *Rothia aeria*. All  
380 nasal samples with a variance stabilized pre-HSCT abundance >-3.05 of ASV 47 (*Rothia* sp.) originated from  
381 patients who subsequently developed no or mild aGvHD (grade 0-I) (Additional File 2: Figure S4B and S3C).  
382

### 383 **Reconstitution of CD4+ T cells and the T<sub>H</sub>17 subpopulation is associated with gut, oral, and nasal 384 microbiota**

385 In order to characterize associations between the microbiota and immune cell counts, immune markers,  
386 and clinical outcomes in HSCT that potentially might impact our predictions of aGvHD, we implemented  
387 two multivariate multi-table approaches, namely sparse partial least squares (sPLS) regression and  
388 canonical correspondence analyses (CCpNA). Using sPLS regression, we identified three clusters of ASVs  
389 for each body site, respectively (Figures 6A, and Additional File 2: S5A and S6A), which was supported by  
390 the CCpNA (Figure 6B, and Additional File 2: S5B and S6B). Several cell populations of the adaptive immune  
391 response were associated with one cluster each at all three body sites according to the sPLS analysis.  
392 These included T cell counts at late follow-up time points, particularly CD4+ T cells in months +3 and +6,  
393 and the subpopulation of T<sub>H</sub>17 cells in months +1 and +3. In the gut, high numbers of these adaptive  
394 immune cell populations were associated with high abundances of mainly *Lachnospiraceae*,

395 *Ruminococcaceae*, and *Lactobacillaceae* ASVs (gut cluster 1, Figure 6A). Of note, two of the *Lactobacillus*  
396 spp. ASVs in gut cluster 1 (ASV 31 and ASV 586) were also observed as members of the group of  
397 *Lactobacillaceae* that discriminated samples from different time points in the LDA (Figure 2C). In the oral  
398 cavity, the same lymphocyte subsets were positively correlated with specific *Flavobacteriaceae*,  
399 *Prevotellaceae*, *Veillonellaceae*, and *Neisseriaceae* ASVs (oral cluster 3, Additional File 2: Figure S5A). The  
400 nasal cluster 1 that was affiliated with high T cell counts comprised predominantly *Veillonellaceae*  
401 (Additional File 2: Figure S5A). The nasal cluster 3 was characterized by high T cell counts at pre-  
402 examination and exhibited a high abundance of ASV 47 (*Rothia* sp.) and other taxa that were associated  
403 with no to mild aGvHD (grade 0-I) (Additional File 2: Figure S4).

404 In the CCpNA, we observed that samples in gut cluster 1 (mainly from months +3 and +6) belonged to  
405 patients with benign primary diseases, who received conditioning regimens involving fludarabine (Figure  
406 6B). Moreover, these patients had a high number of bacterial and viral infections and were treated often  
407 with phenoxymethylpenicillin compared to the overall patient population. In the oral cavity, samples  
408 associated with CD4+ T cell reconstitution similarly stemmed from late follow-up time points and from  
409 pre-examination. Patients in oral cluster 3 were generally treated with few antibiotics. The CCpNA of the  
410 nasal data set indicated that patients with high CD4+ T cell and T<sub>H</sub>17 cell counts at late follow-up time  
411 points exhibited moderate to severe aGvHD (grade II-IV). Furthermore, these patients were treated with  
412 meropenem, ciprofloxacin, and vancomycin more often compared with the remaining patient population  
413 (Figure S6B). Most samples in the nasal cluster 1 were collected in weeks +2 and +3.

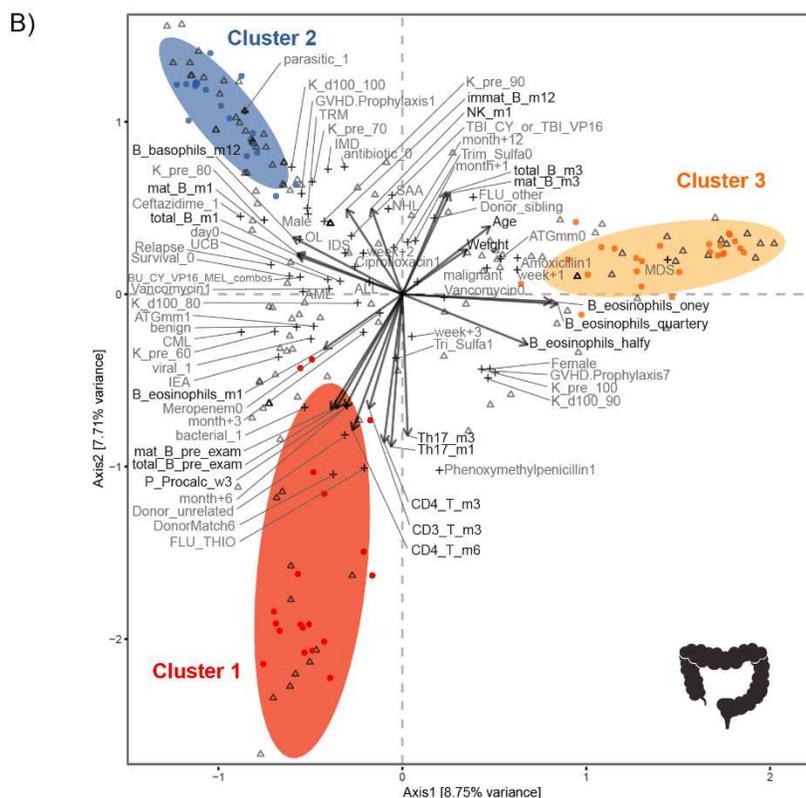
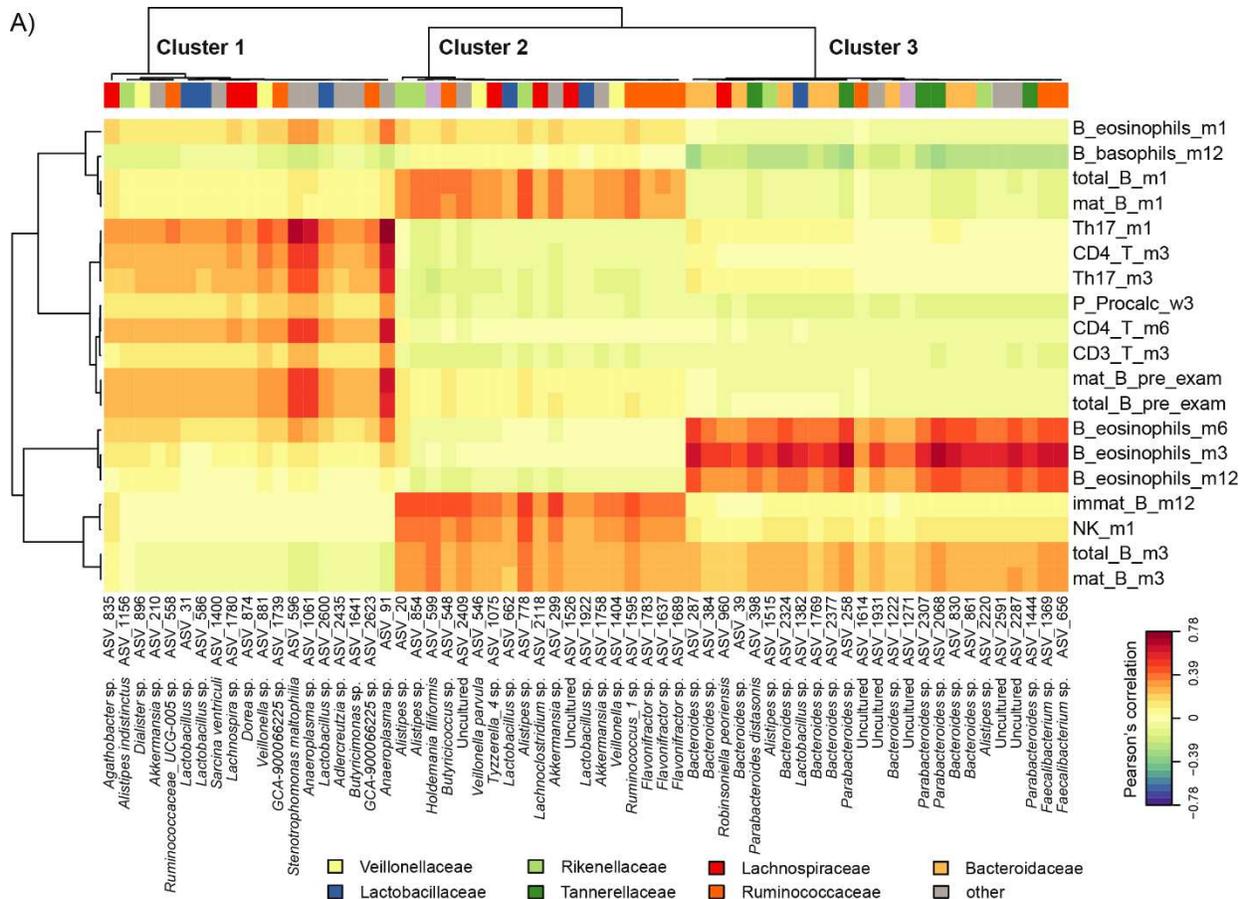
414

#### 415 **Reconstitution of B cells is associated with gut, oral, and nasal microbiota**

416 At all three body sites, B cell counts at several late follow-up time points exhibited associations with  
417 microbial abundances. High B cell counts were positively correlated with high abundances of  
418 *Ruminococcaceae*, *Lachnospiraceae*, and *Rikenellaceae*, as well as few *Veillonellaceae* and  
419 *Lactobacillaceae* in the gut (cluster 2, Figure 6A). In addition, the gut cluster 2 was associated with high  
420 NK cell counts in month +1. In the oral cavity, ASVs within the small cluster 1, particularly ASV 422  
421 (*Actinomyces odontolyticus*) and ASV 546 (*Veillonella parvula*), were positively correlated with these cell  
422 counts, whereas ASVs affiliated with *Staphylococcaceae* and *Lactobacillaceae* (oral cluster 2) exhibited  
423 negative correlations (Additional File 2: Figure S5A). ASV 422 (*Actinomyces odontolyticus*) was also  
424 observed within the group of *Actinomycetaceae* ASVs in the LDA of the oral microbiota. In the nasal cavity,  
425 abundances of *Streptococcaceae*, *Moraxellaceae*, and *Corynebacteriaceae* within nasal cluster 3 were  
426 positively correlated with high B cell counts, particularly in month +3 (Additional File 2: Figure S6A). The  
427 CCpNA indicated that samples in gut cluster 2 were taken predominantly in week +2, whereas samples in  
428 oral cluster 1 were mainly collected in months +3 and +6 (Figures 6B and Additional File 2: Figure S5B).

429 Both the gut and oral CCpNA indicated that the associations between B cell counts and microbial  
430 abundances predominantly occurred in patients who underwent a conditioning regimen without TBI and  
431 without fludarabine (in contrast to conditioning regimens involving TBI or fludarabine). Furthermore,  
432 these patients were treated with ceftazidime, vancomycin, and ciprofloxacin, but sparsely with other  
433 antimicrobial agents (Figure 6B and Additional File 2: Figure S5B). The CCpNA on the gut data set revealed  
434 that samples in this cluster (gut cluster 2) originated from both patients diagnosed with malignant diseases  
435 and benign diseases (Figure 6B).

436



438 **Figure 6. Multivariate associations of the gut microbiota with immune and clinical parameters in HSCT.** A)  
439 Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis (dimensions 1, 2, and 3)  
440 displaying pairwise correlations  $>0.3$ / $<-0.3$  between ASVs (bottom) and continuous immune and clinical parameters  
441 (right). Red indicates a positive correlation, and blue indicates a negative correlation, respectively. Based on the sPLS  
442 regression model, hierarchical clustering (clustering method: complete linkage, distance method: Pearson's  
443 correlation) was performed resulting in the three depicted clusters. B) Canonical correspondence analysis (CCpNA)  
444 relating gut microbial abundances (circles) to continuous (arrows) and categorical (+) immune and clinical  
445 parameters. ASVs and variables with at least one correlation  $>0.3$ / $<-0.3$  in the sPLS analysis were included in the  
446 CCpNA. The triplot shows variables and ASVs with a score  $>0.3$ / $<-0.3$  on at least one of the first three CCpNA axes,  
447 displayed on axis 1 versus 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence  
448 interval) correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. Abbreviations not  
449 mentioned in text: ATGmm, anti-thymocyte globulin; B\_, blood; BU, busulfan; CY, Cyclophosphamide; DonorMatch6,  
450 matched unrelated donor; FLU\_other, fludarabine combinations without thiotepa; GvHD.Prophylaxis1, treatment  
451 with cyclosporine; GvHD.Prophylaxis7, treatment with cyclosporine and methotrexate; immat\_B, immature B cells;  
452 K\_d100, Karnofsky score on day +100; K\_pre, Karnofsky score before HSCT; m1, month+1; m3, month+3; m6,  
453 month+6; m12, month+12; mat\_B, mature B cells; MEL, melphalan; total\_B, total B cells; P\_, plasma; parasitic,  
454 parasitic infection; pre\_cond, before conditioning start; pre\_exam, pre-examination; THIO, thiotepa; viral, viral  
455 infection; VP16, Etoposide.

456

### 457 **Body site-specific immune-microbial associations**

458 In addition to immune-microbial associations shared between two or three of the examined body sites,  
459 we observed a few patterns that were exclusive to individual sites. In cluster 3 in the gut, we observed  
460 ASVs primarily affiliated with *Bacteroidaceae* and *Tannerellaceae* whose abundances showed positive  
461 correlations with eosinophil counts in months +3, +6, and +12. In the oral cavity, the sPLS analysis revealed  
462 a sub-cluster of oral cluster 3 comprising ASVs affiliated with various families, e.g. ASV 1172 (*Actinomyces*  
463 *sp.*), which was also identified as one of the discriminating *Actinomycetaceae* ASVs in the LDA. In the sPLS  
464 analysis, this sub-cluster was associated with high counts of T<sub>reg</sub> and T<sub>H</sub>17 cells at late follow-up time points  
465 (Additional File 2: Figure S6A).

466

### 467 **Discussion**

468 Both the microbiota and the immune system are subject to major changes during allogeneic HSCT. Failure  
469 to re-establish host-microbial homeostasis might have adverse consequences for the patients, such as  
470 prolonged immune deficiency. Long-term surveillance of microbial dynamics is required to understand i)  
471 the shifts in the microbial community structure induced by HSCT and its accompanying treatments, and  
472 ii) at which time points and under which conditions re-establishment of immunological and microbial  
473 homeostasis occurs. Such knowledge may be of great prognostic value and may assist in guiding  
474 personalized treatment strategies. Here, we present a comprehensive assessment of temporal microbial  
475 abundance trajectories from before, at the time of, and after HSCT, to late follow-up time points up to  
476 one year.

477

478 We have identified a group of *Ruminococcaceae*, and a clade of *Blautia* spp. (*Lachnospiraceae*), temporally  
479 discriminating microbial community structure in the gut in relation to HSCT. We show a clear pattern of  
480 depletion of fecal *Blautia* spp. immediately post HSCT, as well as their recovery from month +3 post HSCT  
481 onwards. One could describe the trajectories of these potentially beneficial taxa as a “smile”-shape.  
482 Previous studies have associated the taxonomic families of *Ruminococcaceae* and *Lachnospiraceae* (both

483 class *Clostridia*), and especially the genus *Blautia* (family *Lachnospiraceae*), with lower mortality, lower  
484 GvHD, and higher bacterial diversity in adult allo-HSCT recipients [4,9,20–22]. In turn, a loss of those taxa  
485 after HSCT was associated with subsequent adverse outcomes. Our findings extend the potential of  
486 *Blautia* spp. abundances as an indicator of favorable clinical outcomes, as we characterize abundance  
487 dynamics in children and provide important insight into the time point for the expected return to  
488 abundances comparable to pre-HSCT time points (i.e. between month +1 and +3).

489  
490 Adverse effects, like bacteremia and GvHD, have been found to accompany an expansion of the genus  
491 *Enterococcus* post transplantation [2,3,6,23]. We have found a characteristic expansion of this genus, as  
492 well as of certain *Lactobacillaceae* after HSCT, in agreement with other recent studies [4,6,11]. In addition,  
493 we were able to show a decrease of *Enterococcus* spp. and *Lactobacillaceae* from month +3 to abundances  
494 comparable to pre-HSCT levels. The abundance of these taxa over the course of one year might be  
495 described as a “frown”-shaped trajectory. As for the “smile”-trajectory of potentially beneficial taxa, the  
496 “frown”-trajectories of these taxa could be the first step towards a novel basis to evaluate the re-  
497 establishment of patients’ microbial homeostasis and associated convalescence. Importantly,  
498 *Enterococcus* was already higher abundant in the patient cohort at preexamination prior to HSCT as  
499 compared to the healthy age-matched cohort, most likely due to prior chemotherapy and antibiotic  
500 treatment given before referral to HSCT. Knowledge about the abundance level of *Enterococcus* before  
501 HSCT could therefore provide valuable information about potential high-risk individuals already prior to  
502 transplantation. It should be noted, however, that despite the observed different abundance levels in  
503 patients and healthy controls, and the further expansion of *Enterococcus* post HSCT being in line with  
504 previous studies, our multivariate analyses did not reveal direct detrimental host-microbial associations  
505 of *Enterococcus* in the present cohort.

506  
507 We have to our knowledge for the first time determined long-term dynamics of the oral and nasal  
508 microbiota in allogeneic HSCT patients. Interestingly, we identified abundance trajectories of  
509 phylogenetically closely related groups of *Actinomycetaceae*, *Streptococcaceae*, *Prevotellaceae*, and  
510 Family XI (*Gemella* spp., Class *Bacillales*) in the oral cavity, resembling the “smile”-shaped trajectories  
511 observed in gut. These taxa are part of the normal oral microbiota. Our findings are in agreement with  
512 previous studies reporting the detection of fewer *Prevotella* spp. and *Streptococcus* spp. in the oral cavity  
513 during the first month post HSCT [24,25]. In addition, our current study provides insight into the time of  
514 recovery of these taxa in month +3 after HSCT.

515  
516 For the oral cavity, a post-transplant expansion of *Enterococcus* spp. and *Staphylococcus* spp. has been  
517 reported previously [25,26]. Consistently, we observed an increased relative abundance of  
518 *Staphylococcaceae* during the first month post HSCT, but we did not identify *Enterococcus* spp. or  
519 *Staphylococcus* spp. as significant drivers of temporal dynamics in the oral cavity. Previously, increased  
520 *Enterococcus* abundances post HSCT were found predominantly in patients who developed oral mucositis,  
521 which was not directly assessed in our study [25,27]. Therefore, our findings suggest that further  
522 investigation of taxa that exhibit “smile”-like abundance trajectories could be relevant in direct relation  
523 to oral mucositis. Especially *Actinomycetaceae*, *Streptococcaceae*, and *Prevotellaceae*, when low-

524 abundant, might be candidates for bacterial predictors of oral mucositis, and furthermore might be  
525 employed to facilitate preventive management.

526  
527 In the nasal cavity, the microbiota did not exhibit temporal patterns as distinct as the “smile“- and  
528 “frown“- shaped trajectories in the gut and the oral cavity. One could speculate that nasal bacterial  
529 abundance patterns might be more individualized, which might in turn conceal pronounced patterns  
530 when looking at the patient population as a whole. However, certain host-microbial associations observed  
531 in the gut were reflected in the nasal cavity. For instance, reconstitution of CD4+ T cells and the T<sub>H</sub>17  
532 subset were associated with distinct groups of ASVs at all three body sites.

533  
534 Together, these findings suggest that the oral and potentially also the nasal cavity might constitute easily  
535 accessible microbial niches suitable for investigating host-microbial associations in the context of HSCT,  
536 similar to current strategies for the gut. While mucous membranes that are in close association with  
537 distinct microbial communities characterize all three niches, it is more feasible to collect buccal and  
538 anterior naris swabs during clinical routine as compared to collecting fecal samples. Fecal sample  
539 collection is dependent on bowel movements, which often are impaired in this patient group. Therefore,  
540 our study provides valuable knowledge for possible future applications that could include the monitoring  
541 of oral microbial dynamics in clinical routine, which might be easier to implement than routine fecal  
542 sampling.

543  
544 We identified AVSs with the potential to predict post-transplant aGvHD, which might open opportunities  
545 to improved preventive clinical management, for example by intensified prophylactic immunosuppression  
546 for patients at increased risk. Some ASVs were significant for both, discriminating the microbiota in long-  
547 term dynamics as well as in the prediction of aGvHD severity from the microbiotas prior to HSCT, such as  
548 ASV 3 (*Lactobacillus* sp.) in the gut, as well as ASV 568 (*Actinomyces* sp.) and ASV 226 (*Prevotella*  
549 *melaninogenica*) in the oral cavity. While we do not yet understand the biological mechanisms underlying  
550 this observation, these taxa could be of particular interest for a long-term monitoring in pediatric HSCT  
551 patients, starting prior to HSCT. Like the gut microbiota, the oral and nasal commensal residents might be  
552 of systemic relevance, and a more holistic picture of microbial influences might be drawn by examining  
553 various niches with bacterial communities potentially interacting across body sites. In light of intimate  
554 host-microbiota interactions, the microbial community patterns might also be a marker for underlying  
555 changes occurring in the immune system.

556  
557 High abundances at late follow-up time points of two fecal *Lactobacillus* spp. that expanded after HSCT  
558 showed positive correlations with T cell reconstitution. This is in line with previous studies suggesting that  
559 the expansion of *Lactobacillus*, a genus commonly associated with probiotic properties, might promote  
560 immune homeostasis and thereby exert a protective effect to limit *Enterococcus* expansion [4,23,28]. A  
561 potential explanation indicated by our results might be that high *Lactobacillus* abundances outlasting  
562 enterococcal dominance promotes T cell reconstitution. However, the associated cell populations include  
563 T<sub>H</sub>17 cells which can facilitate inflammation, and therefore it is difficult to determine whether the  
564 observed *Lactobacillus* expansion is exclusively beneficial [13]. However, Th17 cells could perhaps add to

565 the host defense in these patients and therefore be beneficial for local homeostasis, although with the  
566 unusual cost of harmful inflammation.

567  
568 Furthermore, we found associations of high *Lachnospiraceae* and *Ruminococcaceae* in the gut with rapid  
569 B and NK cell reconstitution, which is in support of our previous study [4]. These two *Clostridiales* families  
570 play an important role in providing the host with short-chain fatty acids (SCFAs), such as butyrate [5,29].  
571 A study demonstrated that SCFAs can facilitate the differentiation of human naïve B cells to plasma cells  
572 in culture [30]. Whether SCFAs also directly influence B cell proliferation is yet unknown.

573  
574 We have made several observations in which infections and/or antibiotic treatments were associated with  
575 the abundance of specific bacterial clusters at certain body sites, immune cell counts, and aGvHD. For  
576 example, patients whose samples were represented by gut microbiota cluster 1 experienced a high  
577 number of infections and were treated often with phenoxymethylpenicillin compared to the overall  
578 patient population. In contrast, patients affiliated with gut microbiota cluster 2 experienced treatment  
579 with ceftazidime, vancomycin, and ciprofloxacin, but sparsely with other antimicrobial agents.  
580 Furthermore, patients affiliated with oral microbiota cluster 3 were generally treated with few antibiotics,  
581 and, patients whose sample were represented by the nasal microbiota cluster 1 were treated often with  
582 meropenem, ciprofloxacin, and vancomycin compared with the remaining patient population. However,  
583 it is challenging to interpret these observations, as these patient samples were also associated with other  
584 features, such as an increased or decreased abundance of certain immune cells (see Additional file 3 for  
585 further discussion), or the patients were exposed to other treatments as well, such as TBI or fludarabine.  
586 Overall, however, our observations are consistent with previous reports that antimicrobial treatment is  
587 associated with changes in microbiota composition in patients undergoing allo-HSCT and might impact  
588 clinical outcomes [4,11,31–33]. It will be important to gain a more mechanistic understanding of the  
589 possible effects of antimicrobial treatment to disentangle the effect of antibiotics from that of other  
590 medications and host responses. Such insight could for example allow selecting more suitable  
591 antimicrobials for treatment in HSCT patients that spare the elimination of beneficial taxa, whose decline  
592 might be associated with more severe clinical outcomes. The choice of antibiotic treatment might also be  
593 important to take into consideration in patients that might potentially be referred to HSCT eventually,  
594 given that we already observed certain changes in the microbiota in the patients at referral compared to  
595 healthy controls. The microbiota at referral already exhibited some features that were associated with  
596 more severe side effects.

597  
598 Associations between aGvHD severity and the microbiota have to date merely been based on logistic  
599 regression and correlation analyses [8,34–36]. In addition, microbial abundances at the time of neutrophil  
600 recovery or engraftment were assessed, i.e. at time points shortly before, concurrent to, or potentially  
601 after aGvHD onset [8,16,36]. Here, we have implemented machine learning techniques to take the  
602 assessment of microbiota-aGvHD relations from correlative to predictive modeling: We presented  
603 evidence that aGvHD severity may be predicted from pre-HSCT microbial abundances in the gut, as well  
604 as in the oral and nasal cavities. This could open up opportunities for the future where microbial markers  
605 guide early interventions to prevent aGvHD. This could include a modulation of the microbiota of patients  
606 predicted to be at high risk with synthetic microbiotas containing beneficial bacteria, including probiotics.

607 Notably, we have to our knowledge for the first time revealed microbial taxa in the oral and nasal cavity  
608 that may predict aGvHD. A further discussion on possible connections between specific microbial taxa of  
609 the gut, oral, and nasal cavity, immune responses, and aGvHD can be found in Additional file 3.

610

## 611 **Conclusions**

612

613 With the present study we bring forward a comprehensive framework of host-microbial associations in  
614 allogeneic HSCT. We focused on long-term microbial dynamics, demonstrating distinct microbial  
615 abundance patterns of disturbance and recovery, as well as making predictions about aGvHD from the  
616 pre-transplant microbiota. We discovered that the microbial community composition in patients prior to  
617 HSCT already differs somewhat from healthy controls in regard to key microbial taxa, opening up  
618 opportunities for potential preventive measure in the future. Moreover, we confirmed the depletion of  
619 *Blautia* spp. and expansion of *Enterococcus* spp. in the gut after HSCT and expand this knowledge by  
620 precisely defining which phylogenetically closely related sequence variants of these genera are  
621 characteristic for those patterns, and when they return to pre-HSCT levels. We identified similar patterns  
622 for members of the oral and nasal microbiota and propose month +3 post-transplant as a possible  
623 universally crucial time point for microbiota reconstitution after HSCT. We demonstrate that high  
624 abundances of for example an intestinal *P. distasonis* ASV, and an oral *P. melaninogenica* ASV pre-HSCT  
625 predict the development of moderate to severe aGvHD post-transplant. When relating microbial  
626 abundances with immune cell counts, we found rapid B and NK cell reconstitution to be associated with  
627 high abundances of *Lachnospiraceae* and *Ruminococcaceae*, which also depended on antibiotics treatment.  
628 Distinct ASVs at all three body sites were associated with T<sub>H</sub>17 cell counts, suggesting future research on  
629 a potential immunomodulatory involvement of the microbiota in inflammation regulation, which might  
630 play a role for aGvHD development. We have discovered host-microbial associations shared between two  
631 or more of the examined body sites. This may open up opportunities for implementing a more feasible  
632 oral and nasal swab sampling into research and clinical diagnostic activities to design more precise patient  
633 treatment strategies to reduce serious side effects and improve immune and microbiota reconstitution.

634

635

## 636 **Materials and Methods**

637

### 638 **Patient recruitment and sample collection**

639 We recruited 29 children (age range: 2.5 - 16.4 years) who underwent their first myeloablative allogeneic  
640 hematopoietic stem cell transplantation at Copenhagen University Hospital Rigshospitalet (Denmark)  
641 between November 2015 and October 2017. We provide detailed information about the patients' clinical  
642 characteristics in Table S1 (Additional File 1). Every patient underwent a myeloablative conditioning  
643 regimen starting on day -10 for patients receiving a graft from a haploidentical donor, and on day -7 for  
644 patients with sibling or matched unrelated donors (Additional File 1: Table S1). One patient had a donor  
645 lymphocyte infusion on day +223 after the first transplantation. Immune cell count date of this patient  
646 was excluded from our analysis from the time of donor lymphocyte infusion. We grouped the patients  
647 into four categories of conditioning regimens: 1. TBI\_CY\_or\_TBI\_VP16 (n=6; TBI + cyclophosphamide or

648 TBI + etoposide), 2. BU\_CY\_VP16\_MEL\_combos (n=6; Combinations of busulfan, cyclophosphamide,  
649 etoposide and melphalan), 3. FLU\_THIO (n=12; subgroups: fludarabine + busulfan + thiotepa (n=6);  
650 fludarabine + treosulfan + thiotepa (n=4); fludarabine + thiotepa (n=1); fludarabine + cyclophosphamide  
651 + thiotepa (n=1)), and 4. FLU\_other (n=5; subgroups: fludarabine + busulfan (n=2); fludarabine +  
652 cyclophosphamide (n=2); fludarabine + treosulfan (n=1)) (Additional File 1: Table S1). The following  
653 sampling time points were defined: pre-examination (between day -57 and day -15), around the start of  
654 conditioning (between day -14 and day -3 and latest 2 days after conditioning start), at time of HSCT  
655 (between day -2 and day +2), and weekly during the first 3 weeks after transplantation (week +1: day +3  
656 to day +10, week +2: day +11 to day +17, week +3: day +18 to day +24) (Figure 1A). Broader intervals  
657 applied to follow-up time points: Month +1 (between days +25 and +45), month +3 (between days +46  
658 and +120), month +6 (between days +121 and +245), and month +12 (between day +246 and +428). Acute  
659 GvHD was graded by daily clinical assessment of skin, liver and gastro-intestinal manifestations according  
660 to the Glucksberg criteria [37]. We group aGvHD severity into grade 0-I and grade II-IV, reflecting clinical  
661 practice where grade I represents limited alloreactivity with no (or very limited) impact on the overall  
662 clinical outcome of HSCT, and therefore no need for medical treatment of these patients, such as the use  
663 of glucocorticoids, which is first-line treatment for grade II-IV aGvHD.  
664 To address certain specific questions, we also analyzed the microbiota (from time point 0) of a cohort of  
665 18 healthy children that were part of a previous study [19]. The median age of these children was 6.8  
666 years (interquartile range 4.6 to 9.6). A total of 30 fecal samples were obtained (11 children provided two  
667 samples each within an interval of six months). The children did not receive any antibiotics within the  
668 month prior to sample collection. The samples were processed in the same way as the fecal samples of  
669 the patients of this study (described below).

670

#### 671 **Infections and antibiotics**

672 Records of bacterial, fungal, viral, and parasitic infections and antibiotic treatment from before HSCT  
673 (from day -30 or at the collection time of the first microbiota sample in case this was earlier) until month  
674 +12 (day +428) were taken into consideration (or as long as data was available for the most recent  
675 patients; data accessed in July 2018). This corresponds to the length of the sampling period of fecal and  
676 swab samples.

677

#### 678 **Analysis of immune cell subpopulations**

679 Leukocyte counts were recorded daily during hospitalization starting prior to HSCT, and later weekly in  
680 the outpatient clinic by flow cytometry (Sysmex XN) or microscopy (CellaVision DM96 microscope) in case  
681 of very low counts. Monitored subpopulations included lymphocytes, monocytes, neutrophils, basophils,  
682 and eosinophils.

683

#### 684 **Analysis of T, B and NK cells in peripheral blood**

685 T, B, and NK cell counts in  $\times 10^9/L$  were determined at pre-examination, and in month +1, +3, +6, and +12.  
686 Trucount Tubes (Becton Dickinson, Albertslund, Denmark) were used to quantify these cell types in  
687 peripheral blood on a FC500 flow cytometer (Beckman Coulter, Copenhagen, Denmark). For  
688 immunofluorescence staining, the following conjugated monoclonal antibodies were used for CD3+ T  
689 cells, CD3+CD4+ T cells and CD3+CD8+ T cell quantification: CD3-PerCP, CD3-FITC, CD4-FITC, CD8-PE

690 (Becton Dickinson). CD45-PerCP, CD16/56-PE antibodies were used to determine NK cells based on their  
691 CD45+CD16+CD56+ phenotype. For B cells, total B cells (CD45+CD19+), mature B cells  
692 (CD45+CD19+CD20+) and immature B cells (CD45+CD19+CD20-) were differentiated by using CD20-FITC  
693 and CD19-PE antibodies.

694

#### 695 **Subtyping of T cells**

696 Peripheral blood samples were collected in month +1, +3 and +6 for isolation of peripheral blood  
697 mononuclear cells (PBMCs) by gradient centrifugation of heparinized blood with Lymphoprep™ (Axis-  
698 Shield, Oslo, Norway). PBMCs were washed in PBS (Life Technologies, Invitrogen, Paisley, U.K.) three times  
699 and then resuspended in RPMI 1640 buffer containing HEPES (Biological Industries Israel Beit-Haemek Ltd,  
700 Kibbutz Beit-Haemek, Israel), L-glutamine (GIBCO, Invitrogen, Carlsbad, CA) and Gentamycin (GIBCO), 30%  
701 fetal bovine serum (Biological Industries) and 10% Dimethyl Sulfoxide (VWR, Herlev, Denmark) for cryo-  
702 preservation in liquid nitrogen.

703 T cell subsets, i.e. T<sub>H</sub>17 cells and T<sub>reg</sub> cells, were quantified from frozen PBMCs by flow cytometry on a  
704 FACS Fortessa III flow cytometer (Becton Dickinson, Albertslund, Denmark). PBMCs were thawed and  
705 washed before incubation with Fixable viability stain 620 (Becton Dickinson) and a set of conjugated  
706 monoclonal antibodies for 30 minutes on ice: CD3-APC-A750 (Beckmann Coulter), CD4-PE-Cy7 (Beckmann  
707 Coulter), CD8-A700 (Becton Dickinson), CD25-PE (Becton Dickinson), CD39-PerCP-Cy5.5 (Beckmann  
708 Coulter), CD196-BV510 (Biolegend, San Diego, USA), CD127-BV711 (Biolegend), CD161-BV650 (Becton  
709 Dickinson) and CD45RA-BV786 (Becton Dickinson). Next, PBMCs were washed and incubated with  
710 transcription factor buffer set (BD) for 45 min on ice. Afterwards, PBMCs were washed twice and  
711 intracellular monoclonal antibodies were added and incubated for 45 minutes on ice: RORγT-A488  
712 (Becton Dickinson), FOXP3-A647 (Becton Dickinson) and Helios-PB (Beckmann Coulter). TH17 cells were  
713 determined by the CD4+RORγT+ phenotype, and T<sub>reg</sub> cells by the CD4+CD25<sup>high</sup>FOXP3+ phenotype.  
714 Absolute cell counts in x10<sup>9</sup>/L were obtained by multiplying the frequency of T<sub>H</sub>17 and T<sub>reg</sub> cells with the  
715 CD4+ T cell count from the same time point.

716

#### 717 **Quantification of inflammation and infection markers**

718 Markers were measured at the Department of Clinical Biochemistry, Copenhagen University Hospital  
719 Rigshospitalet, Denmark. As a marker of infection, plasma procalcitonin was determined by sandwich  
720 electrochemiluminescence immunoassays (ECLIA). As a marker of systemic inflammation, CRP was  
721 measured by latex immunoturbidimetric assays (LIA).

722

#### 723 **DNA isolation from fecal, oral, and nasal samples and 16S rRNA gene sequencing**

724 A total of 212 fecal samples for analysis of the intestinal microbiota were collected from 29 patients at  
725 the 10 time points described above. The gut microbiota was characterized at ≤6 time points in 9 patients  
726 (31%), at 7-8 time points in 13 patients (45%) and at 9-10 time points in 7 patients (24%) (Additional File  
727 1: Table S1). DNA from fecal samples, one blank control per extraction round (thereof sequenced: 14),  
728 one mock community sample (Biodefense and Emerging Infectious Research (BEI) Resources of the  
729 American Type Culture Collection (ATCC) (Manassas, VA, USA), Catalog No. HM-276D) per sequencing run

730 and two collection tube controls was isolated using the QIAamp Fast DNA Stool Mini kit (Qiagen, Venlo,  
731 Netherlands), following the manufacturer's instructions with modifications according to [38].  
732 We collected 248 buccal swabs (3x at ≤6 time points (10%), 11x at 7-8 time points (38%), 15x at 9-10 time  
733 points (52%)) and 249 anterior naris swabs (3x at ≤6 time points (10%), 9x at 7-8 time points (31%), 17x at  
734 9-10 time points (59%)). DNA from swab samples, one blank control per extraction round (thereof  
735 sequenced: 28), one mock community sample per run, two collection tube controls, and two sampling  
736 swab controls was isolated using the QIAamp UCP Pathogen Mini kit (Qiagen, Venlo, Netherlands), with  
737 the 'Protocol: Pretreatment of Microbial DNA from Eye, Nasal, Pharyngeal, or other Swabs (Protocol  
738 without Pre-lysis)' and subsequently the 'Protocol: Sample Prep (Spin Protocol)', following the  
739 manufacturer's instructions with the following modifications: 550µl instead of 500µl Buffer ATL was used  
740 during pretreatment; DNA was eluted twice with 20µl Buffer AVE into 1.5 ml DNA LoBind tubes  
741 (Eppendorf, Hamburg, Germany) instead of the tubes provided with the kits.  
742 Library construction and sequencing on an Illumina MiSeq instrument (Illumina Inc., San Diego, CA, USA)  
743 was performed at the Multi Assay Core facility (DMAC), Technical University of Denmark. DNA  
744 concentration of each sample was measured using a NanoDrop spectrophotometer (Thermo Scientific,  
745 Waltham, MA, USA). Library construction was performed according to the *16S Metagenomic Sequencing*  
746 *Library Preparation* protocol by Illumina [39]: The V3-V4 region of the 16S ribosomal RNA gene were  
747 amplified in a PCR in each sample and in the controls, using the following previously evaluated primers,  
748 preceded by Illumina adapters [40]: 341F (5'-  
749 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGCCTACGGGNGGCWGCAG-3') and 805R (5'-  
750 GTCTCGTGGGCTCGGAGATGTGTATAAGAGACAGGACTACHVGGGTATCTAATCC-3'). Amplicons were then  
751 analyzed for quantity and quality in an Agilent 2100 Bioanalyzer with the use of an Agilent RNA 1000 Nano  
752 Kit (Agilent Technology, Santa Clara, CA, USA). Subsequently, the amplicons were purified on AMPure XP  
753 Beads (Beckman Coulter, Copenhagen, Denmark) according to the manufacturer's instructions. Illumina  
754 adapters and dual-index barcodes were then added to the amplicon target in a PCR according to Illumina  
755 [39] using the 96 sample Nextera XT Index Kit (Illumina, FC-131–1002). A final clean-up of the libraries was  
756 performed in another PCR step, using AMPure XP Beads (Beckman Coulter, Copenhagen, Denmark)  
757 according to the manufacturer's instructions, followed by a confirmation of the target size in an Agilent  
758 2100 Bioanalyzer (Agilent Technologies). Before sequencing, DNA concentration was determined with a  
759 Qubit (Life Technologies, Carlsbad, CA, USA) and libraries were pooled. In preparation for sequencing, the  
760 pooled libraries were denatured with NaOH, diluted with hybridization buffer, and heat denatured. 5%  
761 PhiX was included as an internal control for low-diversity libraries. Paired-end sequencing with 2 × 300bp  
762 reads was performed with a MiSeq v3 reagent kit on an Illumina MiSeq instrument (Illumina Inc., San  
763 Diego, CA, USA).

764

#### 765 **16S rRNA gene sequence pre-processing**

766 Raw sequence reads were demultiplexed based on sample-specific barcodes and 'read 1' and 'read 2'  
767 FASTQ files for each sample were generated on the Illumina MiSeq instrument by the MiSeq reporter  
768 software. Primers were removed by using cutadapt (version 1.16) [41] at a tolerated maximum error rate  
769 of 15% for matching the primer sequence anchored in the beginning of each read. In the case that at least  
770 one read of a pair did not contain the primer, the pair was discarded. Only pairs in which the forward read

771 contained the forward primer (341F) and the reverse read contained the reverse primer (805R) were  
772 retained.

773 The resulting reads were further processed using the R package DADA2 (version 1.8) to infer high-  
774 resolution amplicon sequence variants (ASV) [42]. Forward and reverse reads were truncated at 280 bp  
775 and 200 bp respectively. This way, the majority of reads retained a quality score >25 according to MultiQC  
776 analysis [43]. These truncation thresholds also ensured an overlap of 480 bp (expected amplicon length  
777 of 460 bp + 20 bp), allowing to merge forward and reverse reads. Samples were pooled for the sample  
778 inference step (*dada()* function) to increase the power for detecting rare variants. Default values were  
779 used for all other quality filtration parameters in DADA2. DNA from samples with a read count <10,000  
780 after preliminary chimera and contaminant removal were re-sequenced. DNA from feces samples with a  
781 read count <5,000 were re-extracted. Eventually, chimeras were identified by sample and removed from  
782 the whole data set (over all sequencing runs) based on a consensus decision (*removeBimeraDenovo()*  
783 function, method “consensus”). Taxonomic assignment on ASVs was done by using the Silva reference  
784 data base (version 132), formatted for DADA2 [44]. Additional species assignment by exact reference  
785 strain matching was performed using the Silva species-assignment training data base, formatted for  
786 DADA2 [44].

787 The resulting ASV and taxonomy tables were integrated with the R package phyloseq and its dependencies  
788 (version 1.24.0) [45]. The data was split into two data sets, one containing feces sample data and one  
789 containing nasal and oral swab data. Subsequently, contaminant removal was performed with the R  
790 package decontam [46]. Potential technical batch effects by sequencing run, 96-well plate, extraction kit,  
791 extraction round, experimenter, and extraction date were assessed by ordination (Principal Coordinates  
792 Analysis (PCoA)).

793 For both, the fecal sample data set and the swab data set, contaminants were identified by sequencing  
794 run as a batch effect and a subsequent calculation of a consensus probability. For the feces sample data  
795 set, contaminants were identified by both, increased prevalence in 14 blank extraction controls and by  
796 relating ASV frequency to post-PCR sample DNA concentration, assuming inverse correlation (method  
797 “both”, frequency threshold: 0.2, prevalence threshold: 0.075) [46]. After manual evaluation of edge  
798 cases, 89 ASVs were removed from the fecal sample data set as contaminants. In an additional step, we  
799 identified 7 contaminants from 2 sampling tube controls (method and thresholds as stated above). In  
800 total, 96 ASVs were removed as contaminants from the fecal sample data set.

801 For the swab sample data set, contaminants were identified by both, increased prevalence in 28 blank  
802 extraction controls and by relating ASV frequency to post-PCR sample DNA concentration (method “both”,  
803 frequency threshold: 0.1, prevalence threshold: 0.6) [46]. A more stringent threshold for prevalence  
804 compared to frequency was chosen here, given the low biomass of the swab samples, accompanied by  
805 post-PCR DNA concentrations similar to those in blank controls. After manual evaluation of edge cases,  
806 1137 ASVs were removed from the swab sample data set as contaminants. In an additional step, we  
807 identified 16 contaminants from 2 sampling tube controls and 2 swab controls (method “both”, frequency  
808 threshold: 0.075, prevalence threshold: 0.5). In total, 1153 ASVs were removed as contaminants from the  
809 oral and nasal swab sample data set.

810 For each subset, we created a phylogenetic tree by de novo alignment of the inferred ASVs, following a  
811 previously described workflow [47]. First, we performed multiple alignment with the package DECIPHER  
812 [48]. Subsequently, we built a neighbor-joining tree using the package phangorn [49], based on which we

813 fitted a GTR+G+I (Generalized time-reversible with Gamma rate variation) maximum likelihood tree. The  
814 phylogenetic tree for each data set (fecal, oral, and nasal) was then integrated with the respective  
815 phyloseq object.

816 Next, we took core subsets of the ASVs remaining after contaminant removal using the function *kOverA()*  
817 from R package *genefilter* [50]. In the fecal set, 2465 ASVs with  $\geq 5$  reads in  $\geq 2$  samples were retained.  
818 With  $\geq 5$  reads in  $\geq 10$  samples, 509 ASVs were retained from the oral sample set, and 602 ASVs from the  
819 nasal sample set. Additional manual contaminant filtering was applied to the oral and nasal core sets.  
820 ASVs affiliated with taxonomic families commonly found in both the oral or nasal cavity and the gut were  
821 only retained in the oral sample set in case they had  $\geq 10$  reads in  $\geq 10$  samples. ASVs of families only  
822 expected in the gut were removed from the oral and nasal sample sets after manually assessing their  
823 abundances. Subsequently, we retained 377 ASVs in the oral sample set, and 197 ASVs in the nasal sample  
824 set.

825 For the comparison of the fecal microbiota in preexamination samples ( $n=15$ ) of HSCT patients and  
826 healthy children ( $n=18$ ), these data were combined in a phyloseq object. The same set of putative  
827 contaminants was removed from the healthy data set as were identified within the full fecal data set of  
828 HSCT patients. Subsequently, a core subset was taken as described above (retaining ASVs with  $\geq 5$  reads  
829 in  $\geq 2$  samples).

830

### 831 **Statistical analysis**

832 Statistical analyses and generation of graphs was performed in R (version 3.5.1, R Foundation for  
833 Statistical Computing, Vienna, Austria) [51]. The R scripts documenting the major steps of our statistical  
834 analyses are available from figshare (<https://doi.org/10.6084/m9.figshare.12280001>). Sequencing data,  
835 and experimental and clinical data (<https://doi.org/10.6084/m9.figshare.12280028>) were integrated for  
836 analysis by using the R package *phyloseq* and its dependencies [45]. We also provide the resulting  
837 *phyloseq* objects through figshare (<https://doi.org/10.6084/m9.figshare.12280004>). Plots were  
838 generated with the packages *ggplot2* [52], *mixOmics* [53], *treeDA* [54], *caret* [55], and *partykit* [56,57].  
839 From the core sets of ASV counts for each body site, bacterial alpha diversity (denoted by the inverse  
840 Simpson index) was calculated and compared between time points by using a Friedman test with  
841 Benjamini-Hochberg correction for multiple testing, and a post-hoc Conover test. To gain insight into  
842 changes of microbial abundances over time in relation to HSCT, we agglomerated ASV counts on family  
843 levels with the function *tax\_glom()* in *phyloseq* [45]. Thereafter, we displayed the relative abundances of  
844 the 12 most abundant families at each body site for each time point. We also depicted relative abundances  
845 over time on family level in patients with aGvHD grade 0-I versus grade II-IV.

846 In order to determine which particular ASVs are relevant in temporal microbial abundance dynamics at  
847 each body site, we implemented tree-based sparse linear discriminant analysis (LDA) with the package  
848 *treeDA* [54]. This supervised method implements prior information about phylogenetic relationships  
849 between ASVs to perform supervised discrimination of classes, here time points, and induces sparsity  
850 constraints to increase interpretability [58]. Leaves and nodes of the phylogenetic tree, representing  $\log+1$   
851 transformed ASV abundances and the sums thereof respectively, were used as predictive features. The  
852 core oral and nasal sets were used as input as described above, while the fecal set was further reduced to  
853 389 ASVs with  $> 5$  reads in  $> 10$  samples for this analysis. Leave-one-out cross validation (LOOCV) was  
854 performed to choose the optimal minimum number of predictive features ensuring sparse, interpretable

855 models. The resulting LDA models had 9 components. By default, this number corresponds to the number  
856 of predicted classes (here 10 time points) less one. To identify relevant components, we plotted sample  
857 scores colored by time points along each component and plotted the components pairwise against each  
858 other (Figure 1C). Thereby, we revealed that the first LDA-component for each body site showed the  
859 highest sample scores and best separated the samples by time point. Therefore, we proceeded with  
860 displaying temporal trajectories of clades of predictive features (ASVs) on the first component. For  
861 selected groups of predictive ASVs we displayed trajectories for patients with aGvHD grade 0-I versus with  
862 grade II-IV.

863 Next, we implemented machine learning models to predict aGvHD grade post-transplant from preceding  
864 ASV abundances. The strategy and R code for the machine learning approach was partially adapted from  
865 a previous approach [59,60]. As a preparative step for this analysis, we variance-stabilized the ASV count  
866 data. To do so, we first performed size factor estimation for zero-inflated data on the core data sets for  
867 each body site with the package GMPR [61]. Subsequently, we transformed the data by using the  
868 function *varianceStabilizingTransformation()* in the package DESeq2 [62]. The function implements a  
869 Gamma-Poisson mixture model to account for both library size differences and biological variability [63].  
870 For the prediction of aGvHD grade, we compared the performances of four different classifiers (random  
871 forest (rf), boosted logistic regression (LogitBoost), support vector machines with linear kernel  
872 (svmLinear), and support vector machines with radial basis function kernel (svmRadial)) using the package  
873 caret [55]. We took subsets of the phyloseq objects comprising only the time points preceding aGvHD  
874 onset: pre-examination, conditioning start, and at the time of HSCT. Prior to fitting the models, we  
875 excluded ASVs with near zero variance, i.e. those that were not differentially abundant between any  
876 samples, by using the function *nearZeroVar()* in package caret [55]. Thereby we obtained sets of 238, 186,  
877 and 100 ASVs for the fecal, oral, and nasal data set, respectively, which were then assessed as potential  
878 predictors of subsequent aGvHD. All classifiers were trained on a randomly chosen subset of 70% of the  
879 data to build a predictive model evaluated on a test set (30% of the data). Splitting was performed in a  
880 way that samples from the same patient at different time points were kept together in either the testing  
881 or training set to ensure that the outcome of a patient can only appear in either the testing set or the  
882 training set, but not both. Thirty iterations of 10-fold cross-validation were performed for each classifier,  
883 both with and without up-sampling. Up-sampling refers to the process of replacement-based sampling of  
884 the class with fewer samples (here aGvHD grade II-IV) to the same size as the class with more samples  
885 (here aGvHD grade 0-I) to achieve a balanced design. SvmLinear on up-sampled data was chosen as the  
886 best performing predictive model for all three data sets (gut, oral, and nasal). Subsequently, we performed  
887 Boruta feature selection using the package Boruta [64]. The Boruta algorithm is a Random Forest  
888 classification based wrapper that compares the importance of real features to that of so called 'shadow  
889 attributes' with randomly shuffled values. Features that are less important than the 'shadow attributes'  
890 are iteratively removed. Here, we retained those ASVs in each data set that were both, among the 50  
891 most important predictors in the svmLinear model and confirmed by the Boruta algorithm (Additional File  
892 1: Table S3). Subsequently, we fitted a CTREE on each set of selected predictors (17 gut, 26 oral, and 12  
893 nasal ASVs) by using the package partykit (Additional File 1: Table S3) [56,57]. In the CTREE analysis the  
894 effect of the predictive ASVs on aGvHD grade is evaluated in a nonparametric regression framework. Using  
895 CTREE, we found 3 significant ASVs each in the gut and in the oral data set, and two significant ASVs in the  
896 nasal data set. CTREE iteratively tests if the abundance of any ASV has a significant effect on aGvHD grade.

897 In the case that a significant relation is found, the ASVs with the largest effect is picked as a node for the  
898 tree. The procedure is then recursively repeated until no further significant effect of any ASV on aGvHD is  
899 found. We plotted the result as a tree featuring the significant split nodes, represented by the ASVs and  
900 the Bonferroni-corrected p-values indication significant influence of their abundance on aGvHD grade.  
901 The terminal nodes of the tree show the proportion of samples stemming from patients with aGvHD grade  
902 0-I versus II-IV, under the condition of the abundance split criterion described on each branch. Since we  
903 used variance stabilized bacterial abundances as input for the machine learning analyses, abundances can  
904 be presented as negative values in some cases and are therefore not easy to interpret intuitively.  
905 Therefore, we additionally displayed the log-transformed relative abundances of all ASVs significantly  
906 predicting aGvHD in boxplots at the three investigated time points (pre-examination, conditioning start,  
907 and at the time of HSCT).

908 Subsequently, we were interested in associations between the fecal, oral, and nasal microbiota and  
909 immune cell counts, and clinical outcomes in HSCT. Records of immune markers, and immune cell counts  
910 contained left- and right-censored measurements, i.e. observations below or above the detection (or  
911 recording) limit, respectively. In order to use these data in analyses that do not tolerate censored records,  
912 we needed to impute the censored data. Therefore, we first fitted the non-parametric maximum  
913 likelihood estimator (NPMLE, also called Turnbull estimator) for univariate interval censored data on each  
914 variable that contained censored records, using the function *ic\_np()* in the R package *icenReg* [65].  
915 Subsequently, censored records were imputed, informed by the model that was fitted on the entirety of  
916 observed and censored data of each variable, using the *imputeCens()* function [65]. Next, we took the  
917 median of measurements for the time points defined above for those immune markers, and immune cell  
918 counts that have been measured more frequently than that. This way, we obtained comparable data sets.  
919 Continuous immune marker and cell count data that was systematically missing for certain sampling time  
920 points was split by time points and unavailable time points were excluded. Missing values in continuous  
921 immune marker and cell count data were imputed for variables with  $\leq 50\%$  missingness. Simultaneous  
922 multivariate non-parametric imputation was performed using the R package *missForest* [66]. Variables  
923 with more than 50% missing values were excluded from the analysis.

924 Next, we implemented two multivariate multi-table approaches to gain a detailed understanding of how  
925 the fecal, oral, and nasal microbiota might be associated with immune cell counts, immune markers, and  
926 clinical outcomes in HSCT. Evaluated clinical outcomes comprised acute GvHD (grade 0-I versus II-IV),  
927 relapse, overall survival, and treatment-related mortality. Furthermore, we included bacterial alpha  
928 diversity (inversed Simpson index), antibiotic treatment, infections, Karnofsky scores before conditioning  
929 and at day +100, and patients' baseline parameters (age, weight, sex, primary disease, malignant versus  
930 benign primary disease, conditioning regimen (including ATG treatment), chemotherapeutic agents'  
931 dosages, TBI treatment and dosage, stem cell source, GvHD prophylactic regimen, donor type  
932 (sibling/matched unrelated/haploidentical), donor HLA-match, and donor sex).

933 For each body site, we performed sparse partial least squares (sPLS) regression by using the function *sp/s()*  
934 in the package *mixOmics* [53]. In sPLS regression, two matrices are being integrated and both their  
935 structures are being modelled. Here, we used variance stabilized ASV abundances as explanatory variables  
936 and all continuous clinical and immune parameters as response variables. The method allows multiple  
937 response variables. Collinear, and noisy data can be handled by this method as well [67]. We did not limit  
938 the number of response variables to be kept for each component (*keepY*) prior to model calculation. The

939 number of explanatory variables (ASVs) to be kept on each component (keepX) was set to 25 after running  
940 the sPLS regression models for each body site with a range of values between 20 and 40 for keepX,  
941 showing results robust to keepX. The *perf()* function was used to inform the choice of 3 relevant  
942 components. Based on the sPLS regression models for each body site, we then performed hierarchical  
943 clustering with the *cim()* function, using the clustering method “complete linkage” and the distance  
944 method “Pearson’s correlation”. Thereby, we generated matrices of coefficients indicating correlations  
945 between ASV abundances and continuous clinical and immune parameters.  
946 Subsequently, we carried out canonical (i.e. bidirectional) correspondence analysis (CCpNA), which is a  
947 multivariate constrained ordination method. This method allow us to assess associations of both  
948 categorical and continuous clinical and immune parameters to ASV abundances. We included ASVs and  
949 variables with a correlation of  $>0.2$ / $<-0.2$  (oral and nasal data set) or  $>0.3$ / $<-0.3$  (fecal data set) in the sPLS  
950 analysis into the CCpNA, and additionally included categorical variables that could not be included in the  
951 sPLS. The method was implemented with the *cca()* function in package *vegan* [68]. It implements a Chi-  
952 square transformation of the log+1 transformed ASV count matrix and subsequent weighted linear  
953 regression, followed by singular value decomposition. We depicted the CCpNA results as a triplot with plot  
954 dimensions corresponding in length to the percentage of variance explained by each axis. At each body  
955 site, we identified three clusters of ASVs through hierarchical clustering based on the first three latent  
956 dimensions of each sPLS analysis (Figure 6A, and Additional File 2: Figures S5A and S6A). The CCpNA  
957 analyses reinforced the cluster separations and additionally provided insight into associations with  
958 categorical variables, including patient baseline parameters, the occurrence of infections, antibiotics  
959 treatment, and clinical outcomes (Figure 6B, and Additional File 2: Figures S5B and S6B).  
960 We compared bacterial alpha diversity and community composition in the gut of HSCT patients at  
961 preexamination with that of healthy children. Alpha diversity (inverse Simpson index) between the two  
962 groups was compared by a Kruskal-Wallis test. Community composition was visualized in a principal  
963 coordinates analysis (PCoA), and analysis of similarities (ANOSIM, package *vegan*) was used to assess  
964 significant differences in the means of rank dissimilarities between the two groups. DESeq2 was employed  
965 for identification of differentially abundant genera among the top 100 most abundant genera with  $>10$   
966 total reads [62]. Differences in relative abundance of genera identified as differentially abundant were  
967 visualized in a heat tree (package *metacoder*) [69]. Higher taxonomic level differential abundance was  
968 assessed by linear discriminant analysis effect size (LEfSe) on centered-log ratio (CLR) transformed data  
969 with an LDA cutoff of 4 (package *microbiomeMarker*) [70]. LefSe accounts for the hierarchical structure  
970 of bacterial phylogeny, thereby allowing identification of differentially abundant taxa on several  
971 taxonomic levels (here kingdom to genus). For additional information see  
972 <https://doi.org/10.6084/m9.figshare.13614230>.

973

974

## 975 **List of abbreviations**

976 AML: Acute myeloid leukemia

977 ASV: Amplicon sequence variant

978 ATG: Anti-thymocyte globulin

979 CCpNA: Canonical correspondence analysis

980 CML: Chronic myeloid leukemia

981 CRP: C-reactive protein  
982 CTREE: Conditional inference tree  
983 ECLIA: Electrochemiluminescence immunoassays  
984 GvL effect: Graft-versus-leukemia effect  
985 (a)GvHD: (Acute) graft-versus-host disease  
986 HSCT: Hematopoietic stem cell transplantation  
987 IDS: Immunodeficiency syndromes  
988 IEA: Inherited abnormalities of erythrocyte differentiation or function  
989 IMD: Inherited disorders of metabolism  
990 LIA: Latex immunoturbidimetric assay  
991 LDA: Linear discriminant analysis  
992 LogitBoost: Boosted logistic regression  
993 LOOCV: Leave-one-out cross validation  
994 MDS: Myelodysplastic or myeloproliferative disorders  
995 MM: Multiple myeloma  
996 NHL: Non-Hodgkin lymphomas  
997 NPMLE: Non-parametric maximum likelihood estimator  
998 OL: Other leukemia  
999 OTU: Operational taxonomic unit  
1000 PBMC: Peripheral blood mononuclear cell  
1001 PCoA: Principal Coordinates Analysis  
1002 Rf: Random forest  
1003 SAA: Severe aplastic anemia  
1004 SCFA: Short-chain fatty acid  
1005 sPLS: Sparse partial least squares analysis  
1006 svmLinear: Support vector machines with linear kernel  
1007 svmRadial: Support vector machines with radial basis function kernel  
1008 TBI: Total body irradiation  
1009  $T_H17$  cell: T helper 17 cell  
1010  $T_{reg}$  cell: T regulatory cell  
1011 UCB: Umbilical cord blood  
1012

## 1013 **Declarations**

1014

## 1015 **Acknowledgements**

1016 We thank the patients and their families for their participation in this study. We also thank Marlene  
1017 Danner Dalgaard and Neslihan Bicen (Multi Assay Core facility (DMAC), Technical University of Denmark)  
1018 for library construction and sequencing. Sequence pre-processing described in this paper was performed  
1019 using the DeiC National Life Science Supercomputer at DTU. Furthermore, we would like to thank Patrick  
1020 Murigu Kamau Njage (Technical University of Denmark) for helpful discussions related to machine  
1021 learning models.  
1022

## 1023 **Author's contributions**

1024 A.C.I., K.K., K.G.M., and S.J.P. designed the research; A.C.I., K.K., H.M, M.I., and S.J.P. performed the  
1025 research; A.C.I., and S.J.P. contributed analytic tools; A.C.I., and S.J.P. analyzed the data; A.C.I. and S.J.P.  
1026 wrote the manuscript; and K.K., M.I., F.M.A., and K.G.M. edited the manuscript.

1027

1028 **Funding**

1029 This work was supported by the European Union’s Framework program for Research and Innovation,  
1030 Horizon2020 (643476), and by the National Food Institute, Technical University of Denmark.

1031

1032 **Availability of data and materials**

1033 The 16S rRNA gene sequences are available through the European Nucleotide Archive (ENA) at the  
1034 European Bioinformatics Institute (EBI) under accession number PRJEB30894. The datasets generated  
1035 and/or analysed in this study as well as the R code used to analyze the data are available from the figshare  
1036 repository [https://figshare.com/projects/Microbiota\\_long-](https://figshare.com/projects/Microbiota_long-term_dynamics_and_prediction_of_acute_graft-versus-host-disease_in_allogeneic_stem_cell_transplantation/80366)

1037 [term\\_dynamics\\_and\\_prediction\\_of\\_acute\\_graft-versus-host-](https://figshare.com/projects/Microbiota_long-term_dynamics_and_prediction_of_acute_graft-versus-host-disease_in_allogeneic_stem_cell_transplantation/80366)

1038 [disease\\_in\\_allogeneic\\_stem\\_cell\\_transplantation/80366](https://figshare.com/projects/Microbiota_long-term_dynamics_and_prediction_of_acute_graft-versus-host-disease_in_allogeneic_stem_cell_transplantation/80366) (see also individual links in the Methods  
1039 section).

1040

1041 **Ethics approval and consent to participate**

1042 Written informed consent was obtained from the patients and/or their legal guardians. The study protocol  
1043 was approved by the local ethics committee (H-7-2014-016) and the Danish Data Protection Agency.

1044

1045 **Consent for publication**

1046 Not applicable.

1047

1048 **Competing interests**

1049 The authors declare that they have no competing interests.

1050

1051 **References**

- 1052 1. Chabannon C, Kuball J, Bondanza A, Dazzi F, Pedrazzoli P, Toubert A, et al. Hematopoietic stem cell  
1053 transplantation in its 60s: A platform for cellular therapies. *Sci Transl Med* [Internet]. American  
1054 Association for the Advancement of Science; 2018 [cited 2018 Aug 2];10:eaap9630. Available from:  
1055 <http://www.ncbi.nlm.nih.gov/pubmed/29643233>
- 1056 2. Shono Y, van den Brink MRM. Gut microbiota injury in allogeneic haematopoietic stem cell  
1057 transplantation. *Nat Rev Cancer* [Internet]. Nature Publishing Group; 2018 [cited 2018 Feb 21]; Available  
1058 from: <http://www.nature.com/doi/10.1038/nrc.2018.10>
- 1059 3. Holler E, Butzhammer P, Schmid K, Hundsrucker C, Koestler J, Peter K, et al. Metagenomic Analysis of  
1060 the Stool Microbiome in Patients Receiving Allogeneic Stem Cell Transplantation: Loss of Diversity Is  
1061 Associated with Use of Systemic Antibiotics and More Pronounced in Gastrointestinal Graft-versus-Host  
1062 Disease. *Biol Blood Marrow Transplant* [Internet]. 2014 [cited 2015 Oct 22];20:640–5. Available from:  
1063 <http://www.sciencedirect.com/science/article/pii/S1083879114000755>
- 1064 4. Ingham AC, Kielsen K, Cilieborg MS, Lund O, Holmes S, Aarestrup FM, et al. Specific gut microbiome  
1065 members are associated with distinct immune markers in pediatric allogeneic hematopoietic stem cell  
1066 transplantation. *Microbiome* [Internet]. BioMed Central; 2019 [cited 2019 Sep 18];7:131. Available from:  
1067 <https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-019-0745-z>
- 1068 5. Rivera-Chávez F, Lopez CA, Bäumlner AJ. Oxygen as a driver of gut dysbiosis. *Free Radic Biol Med*  
1069 [Internet]. Pergamon; 2017 [cited 2018 Feb 18];105:93–101. Available from:  
1070 <https://www.sciencedirect.com/science/article/pii/S0891584916304361?via%3Dihub>
- 1071 6. Taur Y, Xavier JB, Lipuma L, Ubeda C, Goldberg J, Gobourne a., et al. Intestinal Domination and the  
1072 Risk of Bacteremia in Patients Undergoing Allogeneic Hematopoietic Stem Cell Transplantation. *Clin*  
1073 *Infect Dis* [Internet]. 2012;55:905–14. Available from:  
1074 <http://cid.oxfordjournals.org/lookup/doi/10.1093/cid/cis580>
- 1075 7. Ghimire S, Weber D, Mavin E, Wang X nong, Dickinson AM, Holler E. Pathophysiology of GvHD and  
1076 Other HSCT-Related Major Complications. *Front Immunol* [Internet]. Frontiers; 2017 [cited 2018 Oct  
1077 19];8:79. Available from: <http://journal.frontiersin.org/article/10.3389/fimmu.2017.00079/full>
- 1078 8. Golob JL, Pergam SA, Srinivasan S, Fiedler TL, Liu C, Garcia K, et al. Stool Microbiota at Neutrophil  
1079 Recovery Is Predictive for Severe Acute Graft vs Host Disease After Hematopoietic Cell Transplantation.  
1080 *Clin Infect Dis* [Internet]. Oxford University Press; 2017 [cited 2018 Nov 23];65:1984–91. Available from:  
1081 <https://academic.oup.com/cid/article/65/12/1984/4085173>
- 1082 9. Jenq RR, Taur Y, Devlin SM, Ponce DM, Goldberg JD, Ahr KF, et al. Intestinal *Blautia* Is Associated with  
1083 Reduced Death from Graft-versus-Host Disease. *Biol Blood Marrow Transplant* [Internet]. 2015 [cited  
1084 2016 May 9];21:1373–83. Available from:  
1085 <http://www.sciencedirect.com/science/article/pii/S1083879115002931>
- 1086 10. Han L, Zhao K, Li Y, Han H, Zhou L, Ma P, et al. A gut microbiota score predicting acute graft-versus-  
1087 host disease following myeloablative allogeneic hematopoietic stem cell transplantation. *Am J*  
1088 *Transplant*. 2020;20:1014–27.

- 1089 11. Peled JU, Gomes ALC, Devlin SM, Littmann ER, Taur Y, Sung AD, et al. Microbiota as predictor of  
1090 mortality in allogeneic hematopoietic-cell transplantation. *N Engl J Med.* 2020;382:822–34.
- 1091 12. Stein-Thoeringer CK, Nichols KB, Lazrak A, Docampo MD, Slingerland AE, Slingerland JB, et al. Lactose  
1092 drives *Enterococcus* expansion to promote graft-versus-host disease. *Science* (80- ). 2019;366:1143–9.
- 1093 13. Honda K, Littman DR. The microbiota in adaptive immune homeostasis and disease. *Nature*  
1094 [Internet]. Nature Publishing Group; 2016 [cited 2018 Aug 21];535:75–84. Available from:  
1095 <http://www.nature.com/articles/nature18848>
- 1096 14. Atarashi K, Tanoue T, Oshima K, Suda W, Nagano Y, Nishikawa H, et al. Treg induction by a rationally  
1097 selected mixture of *Clostridia* strains from the human microbiota. *Nature* [Internet]. Nature Publishing  
1098 Group; 2013 [cited 2018 Sep 6];500:232–6. Available from:  
1099 <http://www.nature.com/articles/nature12331>
- 1100 15. Kielsen K, Ryder LP, Lennox-Hvenekilde D, Gad M, Nielsen CH, Heilmann C, et al. Reconstitution of  
1101 Th17, Tc17 and Treg cells after paediatric haematopoietic stem cell transplantation: Impact of  
1102 interleukin-7. *Immunobiology* [Internet]. 2018 [cited 2018 Feb 7];223:220–6. Available from:  
1103 <http://www.ncbi.nlm.nih.gov/pubmed/29033080>
- 1104 16. Han L, Jin H, Zhou L, Zhang X, Fan Z, Dai M, et al. Intestinal Microbiota at Engraftment Influence  
1105 Acute Graft-Versus-Host Disease via the Treg/Th17 Balance in Allo-HSCT Recipients. *Front Immunol*  
1106 [Internet]. 2018 [cited 2018 May 17];9:669. Available from:  
1107 <http://www.ncbi.nlm.nih.gov/pubmed/29740427>
- 1108 17. Ratajczak P, Janin A, Peffault de Latour R, Leboeuf C, Desveaux A, Keyvanfar K, et al. Th17/Treg ratio  
1109 in human graft-versus-host disease. *Blood* [Internet]. American Society of Hematology; 2010 [cited 2018  
1110 Nov 19];116:1165–71. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/20484086>
- 1111 18. Larsen JM. The immune response to *Prevotella* bacteria in chronic inflammatory disease.  
1112 *Immunology* [Internet]. Wiley/Blackwell (10.1111); 2017 [cited 2018 Nov 25];151:363–74. Available  
1113 from: <http://doi.wiley.com/10.1111/imm.12760>
- 1114 19. De Pietri S, Ingham AC, Frandsen TL, Rathe M, Krych L, Castro-Mejía JL, et al. Gastrointestinal toxicity  
1115 during induction treatment for childhood acute lymphoblastic leukemia: The impact of the gut  
1116 microbiota. *Int J Cancer.* Wiley-Liss Inc.; 2020;147:1953–62.
- 1117 20. Weber D, Oefner PJ, Hiergeist A, Koestler J, Gessner A, Weber M, et al. Low urinary indoxyl sulfate  
1118 levels early after transplantation reflect a disrupted microbiome and are associated with poor outcome.  
1119 *Blood* [Internet]. 2015 [cited 2015 Oct 22];126:1723–8. Available from:  
1120 <http://www.bloodjournal.org/content/126/14/1723>
- 1121 21. Taur Y, Jenq RR, Perales M, Littmann ER, Morjaria S, Ling L, et al. The effects of intestinal tract  
1122 bacterial diversity on mortality following allogeneic hematopoietic stem cell transplantation.  
1123 *Transplantation* [Internet]. 2014;124:1174–82. Available from:  
1124 <http://www.bloodjournal.org/content/bloodjournal/124/7/1174.full.pdf?sso-checked=true>
- 1125 22. Andermann TM, Peled JU, Ho C, Reddy P, Riches M, Storb R, et al. The Microbiome and  
1126 Hematopoietic Cell Transplantation: Past, Present, and Future. *Biol Blood Marrow Transplant* [Internet].

1127 Elsevier; 2018 [cited 2018 May 29]; Available from:  
1128 <https://www.sciencedirect.com/science/article/pii/S1083879118300879?via%3Dihub>

1129 23. Jenq RR, Ubeda C, Taur Y, Menezes CC, Khanin R, Dudakov J a., et al. Regulation of intestinal  
1130 inflammation by microbiota following allogeneic bone marrow transplantation. *J Exp Med* [Internet].  
1131 2012;209:903–11. Available from: <http://www.jem.org/cgi/doi/10.1084/jem.20112408>

1132 24. Verma D, Garg PK, Dubey AK. Insights into the human oral microbiome. *Arch Microbiol* [Internet].  
1133 Springer Berlin Heidelberg; 2018;200:525–40. Available from: <http://dx.doi.org/10.1007/s00203-018-1505-3>

1135 25. Osakabe L, Utsumi A, Saito B, Okamatsu Y, Kinouchi H, Nakamaki T, et al. Influence of Oral Anaerobic  
1136 Bacteria on Hematopoietic Stem Cell Transplantation Patients: Oral Mucositis and General Condition.  
1137 *Transplant Proc* [Internet]. Elsevier; 2017 [cited 2018 Mar 12];49:2176–82. Available from:  
1138 <http://www.ncbi.nlm.nih.gov/pubmed/29149979>

1139 26. Soga Y, Maeda Y, Ishimaru F, Tanimoto M, Maeda H, Nishimura F, et al. Bacterial substitution of  
1140 coagulase-negative staphylococci for streptococci on the oral mucosa after hematopoietic cell  
1141 transplantation. *Support Care Cancer*. 2011;19:995–1000.

1142 27. Olczak-Kowalczyk D, Daszkiewicz M, Krasuska-Sławińska, Dembowska-Bagińska B, Gozdowski D,  
1143 Daszkiewicz P, et al. Bacteria and Candida yeasts in inflammations of the oral mucosa in children with  
1144 secondary immunodeficiency. *J Oral Pathol Med*. 2012;41:568–76.

1145 28. Chung H, Pamp SJ, Hill JA, Surana NK, Edelman SM, Troy EB, et al. Gut immune maturation depends  
1146 on colonization with a host-specific microbiota. *Cell* [Internet]. Elsevier; 2012 [cited 2018 Aug  
1147 16];149:1578–93. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/22726443>

1148 29. Mathewson ND, Jenq R, Mathew A V, Koenigsknecht M, Hanash A, Toubai T, et al. Gut microbiome-  
1149 derived metabolites modulate intestinal epithelial cell damage and mitigate graft-versus-host disease.  
1150 *Nat Immunol* [Internet]. NIH Public Access; 2016 [cited 2018 May 15];17:505–13. Available from:  
1151 <http://www.ncbi.nlm.nih.gov/pubmed/26998764>

1152 30. Kim M, Qie Y, Park J, Kim CH. Gut Microbial Metabolites Fuel Host Antibody Responses. *Cell Host*  
1153 *Microbe* [Internet]. Elsevier Inc.; 2016;20:202–14. Available from:  
1154 <http://dx.doi.org/10.1016/j.chom.2016.07.001>

1155 31. Shono Y, Docampo MD, Peled JU, Perobelli SM, Velardi E, Tsai JJ, et al. Increased GVHD-related  
1156 mortality with broad-spectrum antibiotic use after allogeneic hematopoietic stem cell transplantation in  
1157 human patients and mice. *Sci Transl Med* [Internet]. 2016 [cited 2016 May 23];8:339ra71-339ra71.  
1158 Available from: <http://stm.sciencemag.org/content/8/339/339ra71>

1159 32. Weber D, Jenq RR, Peled JU, Taur Y, Hiergeist A, Koestler J, et al. Microbiota Disruption Induced by  
1160 Early Use of Broad Spectrum Antibiotics is an Independent Risk Factor of Outcome after Allogeneic Stem  
1161 Cell Transplantation [Internet]. *Biol. Blood Marrow Transplant*. Elsevier Inc.; 2017. Available from:  
1162 <http://linkinghub.elsevier.com/retrieve/pii/S1083879117302756>

1163 33. Weber D, Hiergeist A, Weber M, Dettmer K, Wolff D, Hahn J, et al. Detrimental effect of broad-  
1164 spectrum antibiotics on intestinal microbiome diversity in patients after allogeneic stem cell

- 1165 transplantation: Lack of commensal sparing antibiotics. *Clin Infect Dis* [Internet]. 2018 [cited 2018 Sep  
1166 20]; Available from: <http://www.ncbi.nlm.nih.gov/pubmed/30124813>
- 1167 34. Liu C, Frank DN, Horch M, Chau S, Ir D, Horch EA, et al. Associations between acute gastrointestinal  
1168 GvHD and the baseline gut microbiota of allogeneic hematopoietic stem cell transplant recipients and  
1169 donors. *Bone Marrow Transplant Adv online Publ* [Internet]. 2017;doi:1–8. Available from:  
1170 <https://www.nature.com/bmt/journal/vaop/ncurrent/pdf/bmt2017200a.pdf>
- 1171 35. Biagi E, Zama D, Nastasi C, Consolandi C, Fiori J, Rampelli S, et al. Gut microbiota trajectory in  
1172 pediatric patients undergoing hematopoietic SCT. *Bone Marrow Transplant* [Internet]. Nature Publishing  
1173 Group; 2015 [cited 2018 Jul 2];50:992–8. Available from: <http://www.nature.com/articles/bmt201516>
- 1174 36. Mancini N, Greco R, Pasciuta R, Barbanti MC, Pini G, Morrow OB, et al. Enteric Microbiome Markers  
1175 as Early Predictors of Clinical Outcome in Allogeneic Hematopoietic Stem Cell Transplant: Results of a  
1176 Prospective Study in Adult Patients. *Open Forum Infect Dis* [Internet]. Oxford University Press; 2017  
1177 [cited 2018 Dec 8];4. Available from:  
1178 <http://academic.oup.com/ofid/article/doi/10.1093/ofid/ofx215/4367678>
- 1179 37. Glucksberg H, Storb R, Fefer A, Buckner CD, Neiman PE, Clift RA, et al. Clinical manifestations of  
1180 graft-versus-host disease in human recipients of marrow from HL-A-matched sibling donors.  
1181 *Transplantation* [Internet]. 1974;18:295–304. Available from:  
1182 <http://www.ncbi.nlm.nih.gov/pubmed/4153799>
- 1183 38. Knudsen BE, Bergmark L, Munk P, Lukjancenko O, Prieme A, Aarestrup FM, et al. Impact of Sample  
1184 Type and DNA Isolation Procedure on Genomic Inference of Microbiome Composition. *bioRxiv*  
1185 [Internet]. 2016;1:064394. Available from: <http://biorxiv.org/lookup/doi/10.1101/064394>
- 1186 39. 16S Metagenomic Sequencing Library Preparation. [Internet]. [cited 2018 Apr 17]. Available from:  
1187 [https://support.illumina.com/content/dam/illumina-](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf)  
1188 [support/documents/documentation/chemistry\\_documentation/16s/16s-metagenomic-library-prep-](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf)  
1189 [guide-15044223-b.pdf](https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/16s/16s-metagenomic-library-prep-guide-15044223-b.pdf)
- 1190 40. Klindworth A, Pruesse E, Schweer T, Peplies J, Quast C, Horn M, et al. Evaluation of general 16S  
1191 ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies.  
1192 *Nucleic Acids Res* [Internet]. Oxford University Press; 2013 [cited 2018 Apr 17];41:e1–e1. Available from:  
1193 [http://academic.oup.com/nar/article/41/1/e1/1164457/Evaluation-of-general-16S-ribosomal-RNA-](http://academic.oup.com/nar/article/41/1/e1/1164457/Evaluation-of-general-16S-ribosomal-RNA-gene-PCR)  
1194 [gene-PCR](http://academic.oup.com/nar/article/41/1/e1/1164457/Evaluation-of-general-16S-ribosomal-RNA-gene-PCR)
- 1195 41. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.  
1196 *EMBnet.journal* [Internet]. 2011 [cited 2018 Jun 26];17:10. Available from:  
1197 <http://journal.embnet.org/index.php/embnetjournal/article/view/200>
- 1198 42. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution  
1199 sample inference from Illumina amplicon data. *Nat Methods* [Internet]. 2016 [cited 2016 Jul 28];13:581–  
1200 3. Available from: <http://www.nature.com/nmeth/journal/v13/n7/full/nmeth.3869.html>
- 1201 43. Ewels P, Magnusson M, Lundin S, Källner M. MultiQC: summarize analysis results for multiple tools  
1202 and samples in a single report. *Bioinformatics* [Internet]. Oxford University Press; 2016 [cited 2018 Sep  
1203 11];32:3047–8. Available from: <https://academic.oup.com/bioinformatics/article->

- 1204 lookup/doi/10.1093/bioinformatics/btw354
- 1205 44. Callahan B. Silva taxonomic training data formatted for DADA2 (Silva version 132). 2018 [cited 2018  
1206 Jun 26]; Available from: <https://doi.org/10.5281/zenodo.1172783#.WzJRh15uQOA.mendeley>
- 1207 45. McMurdie PJ, Holmes S. phyloseq: An R Package for Reproducible Interactive Analysis and Graphics  
1208 of Microbiome Census Data. Watson M, editor. PLoS One [Internet]. Public Library of Science; 2013  
1209 [cited 2018 Jan 24];8:e61217. Available from: <http://dx.plos.org/10.1371/journal.pone.0061217>
- 1210 46. Davis NM, Proctor DM, Holmes SP, Relman DA, Callahan BJ. Simple statistical identification and  
1211 removal of contaminant sequences in marker-gene and metagenomics data. Microbiome [Internet].  
1212 BioMed Central; 2018 [cited 2018 Dec 27];6:226. Available from:  
1213 <https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-018-0605-2>
- 1214 47. Callahan BJ, Sankaran K, Fukuyama JA, McMurdie PJ, Holmes SP. Bioconductor workflow for  
1215 microbiome data analysis: from raw reads to community analyses. F1000Research [Internet].  
1216 2016;5:1492. Available from:  
1217 <http://www.ncbi.nlm.nih.gov/pubmed/27508062%5Cnhttp://www.pubmedcentral.nih.gov/articlerende>  
1218 [r.fcgi?artid=PMC4955027](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4955027)
- 1219 48. Wright ES. Using DECIPHER v2.0 to Analyze Big Biological Sequence Data in R. R J. 2016;8:352–9.
- 1220 49. Schliep KP. phangorn: phylogenetic analysis in R. Bioinformatics [Internet]. Oxford University Press;  
1221 2011 [cited 2018 Nov 27];27:592–3. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/21169378>
- 1222 50. Gentleman R, Carey V, Huber W, Hahne F. genefilter: methods for filtering genes from microarray  
1223 experiments. R package version 1.58.1. 2017.
- 1224 51. R Core Team. R: A Language and Environment for Statistical Computing [Internet]. R Foundation for  
1225 Statistical Computing, Vienna, Austria; 2018. Available from: <https://www.r-project.org/>
- 1226 52. Wickham H. ggplot2: Elegant Graphics for Data Analysis. Springer Verlag New York; 2016.
- 1227 53. Rohart F, Gautier B, Singh A, Lê Cao K-A. mixOmics: An R package for ‘omics feature selection and  
1228 multiple data integration. Schneidman D, editor. PLOS Comput Biol [Internet]. Public Library of Science;  
1229 2017 [cited 2017 Dec 11];13:e1005752. Available from: <http://dx.plos.org/10.1371/journal.pcbi.1005752>
- 1230 54. Fukuyama J. treeDA: Tree-Based Discriminant Analysis. [Internet]. 2017. Available from:  
1231 <https://github.com/jfukuyama/treedda>
- 1232 55. Kuhn M, Wing J, Weston S, Williams A, Keefer C, Engelhardt A, et al. caret: Classification and  
1233 Regression Training. R package version 6.0-80. [Internet]. 2018. Available from: [https://cran.r-](https://cran.r-project.org/package=caret)  
1234 [project.org/package=caret](https://cran.r-project.org/package=caret)
- 1235 56. Hothorn T, Zeileis A, Cheng E, Ong S. partykit: A modular toolkit for recursive partytioning in R. J  
1236 Mach Learn Res. 2015;16:3905–9.
- 1237 57. Hothorn T, Hornik K, Zeileis A. Unbiased Recursive Partitioning: A Conditional Inference Framework. J  
1238 Comput Graph Stat [Internet]. Taylor & Francis; 2006 [cited 2018 Nov 24];15:651–74. Available from:

- 1239 <http://www.tandfonline.com/doi/abs/10.1198/106186006X133933>
- 1240 58. Fukuyama J, Rumker L, Sankaran K, Jegannathan P, Dethlefsen L, Relman DA, et al. Multidomain  
1241 analyses of a longitudinal human microbiome intestinal cleanout perturbation experiment. *PLoS Comput*  
1242 *Biol* [Internet]. Public Library of Science; 2017 [cited 2018 Mar 2];13:e1005706. Available from:  
1243 <http://www.ncbi.nlm.nih.gov/pubmed/28821012>
- 1244 59. Njage PMK, Henri C, Leekitcharoenphon P, Mistou M, Hendriksen RS, Hald T. Machine Learning  
1245 Methods as a Tool for Predicting Risk of Illness Applying Next-Generation Sequencing Data. *Risk Anal*  
1246 [Internet]. John Wiley & Sons, Ltd (10.1111); 2018 [cited 2018 Dec 24];risa.13239. Available from:  
1247 <https://onlinelibrary.wiley.com/doi/abs/10.1111/risa.13239>
- 1248 60. Njage PMK, Leekitcharoenphon P, Hald T. Improving hazard characterization in microbial risk  
1249 assessment using next generation sequencing data and machine learning: Predicting clinical outcomes in  
1250 shigatoxigenic *Escherichia coli*. *Int J Food Microbiol* [Internet]. Elsevier; 2019 [cited 2018 Dec  
1251 24];292:72–82. Available from:  
1252 <https://www.sciencedirect.com/science/article/pii/S0168160518308936#f0005>
- 1253 61. Chen J, Zhang L. GMPR: Geometric Mean of Pairwise Ratios. R package version 0.1.3. 2017.
- 1254 62. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data  
1255 with DESeq2. *Genome Biol* [Internet]. BioMed Central; 2014 [cited 2018 Jan 24];15:550. Available from:  
1256 <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-014-0550-8>
- 1257 63. McMurdie PJ, Holmes S. Waste Not, Want Not: Why Rarefying Microbiome Data Is Inadmissible.  
1258 McHardy AC, editor. *PLoS Comput Biol* [Internet]. 2014 [cited 2016 Apr 13];10:e1003531. Available from:  
1259 <http://dx.plos.org/10.1371/journal.pcbi.1003531>
- 1260 64. Kurska MB, Rudnicki WR. Feature Selection with the Boruta Package. *J. Stat. Softw.* 2010. p. 1–13.
- 1261 65. Anderson-Bergman C. **icenReg** : Regression Models for Interval Censored Data in *R*. *J Stat Softw*  
1262 [Internet]. 2017;81. Available from: <http://www.jstatsoft.org/v81/i12/>
- 1263 66. Stekhoven DJ, Buhlmann P. MissForest--non-parametric missing value imputation for mixed-type  
1264 data. *Bioinformatics* [Internet]. Oxford University Press; 2012 [cited 2018 Sep 18];28:112–8. Available  
1265 from: <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btr597>
- 1266 67. Lee D, Lee W, Lee Y, Pawitan Y. Sparse partial least-squares regression and its applications to high-  
1267 throughput data analysis. *Chemom Intell Lab Syst* [Internet]. Elsevier; 2011 [cited 2018 Jan 24];109:1–8.  
1268 Available from: <https://www.sciencedirect.com/science/article/pii/S016974391100150X>
- 1269 68. Oksanen J, Blanchet FG, Friendly M, Kindt R, Legendre P, McGlinn D, et al. *vegan*: Community  
1270 Ecology Package. R package version 2.5-2. [Internet]. 2018. Available from: [https://cran.r-](https://cran.r-project.org/package=vegan)  
1271 [project.org/package=vegan](https://cran.r-project.org/package=vegan)
- 1272 69. Foster ZSL, Shapton TJ, Grünwald NJ. Metacoder: An R package for visualization and manipulation of  
1273 community taxonomic diversity data. *PLoS Comput Biol*. 2017;13:1–15.
- 1274 70. Cao Y. *microbiomeMarker*: microbiome biomarker analysis. R package version 0.0.1.9000.

1275 <https://github.com/yiluheihei/microbiomeMarker>. 2021.

1276

1277

1278 **Figure and Table Legends**

1279

1280 **Figure 1. Monitoring gut, oral, and nasal microbiota and the host immune system in allogeneic**  
1281 **hematopoietic stem cell transplantation (HSCT).** A) Twenty-nine children were monitored before, at the  
1282 time of, and immediately post allogeneic HSCT, as well as at late follow-up time points. Patients' baseline  
1283 characteristics, clinical outcomes, as well as immune cell counts, and inflammation and infection markers  
1284 over time were monitored. Patient characteristics are described in detail in Table S1 (Additional File 1).  
1285 Host immune system parameters were related to longitudinal dynamics of the gut, oral, and nasal  
1286 microbiota that was assessed at the denoted time points. B) Bacterial alpha diversity before, at the time  
1287 of, and after HSCT at each body site, displayed on a log<sub>10</sub> transformed y-axis for visualization purposes.  
1288 Asterisks indicate significant differences in median inverse Simpson index between time points \*  $P < 0.05$ .  
1289 C) Tree-based sparse linear discriminant (LDA) analyses by time point in relation to HSCT. For fecal  
1290 samples, positive LDA scores were observed for samples collected immediately post HSCT. For both oral  
1291 and nasal samples, positive LDA scores were observed for samples from before HSCT and from late follow  
1292 up-time points.

1293

1294 **Figure 2. Temporal microbial community dynamics in the gut.** A) Relative abundances over time of the  
1295 12 most abundant families in the gut. B) Tree-based sparse linear discriminant analysis (LDA). Coefficients  
1296 of discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the  
1297 phylogenetic tree. C) Trajectories of ASVs affiliated with the families *Enterococcaceae* and  
1298 *Lactobacillaceae*, with increasing abundances after HSCT. The most abundant discriminating ASV for each  
1299 family is indicated. D) Trajectories of ASVs affiliated with the families *Lachnospiraceae* and  
1300 *Ruminococcaceae*, with decreasing abundances after HSCT and recovery at late follow-up time points. The  
1301 most abundant discriminating ASV for *Blautia* spp. is indicated. Detailed taxonomic information and LDA-  
1302 coefficients of the displayed ASVs are listed in Additional File 1: Table S2.

1303

1304 **Figure 3. Temporal microbial community dynamics in the oral cavity.** A) Relative abundances over time  
1305 of the 12 most abundant families in the oral cavity. B) Tree-based sparse linear discriminant analysis (LDA).  
1306 Coefficients of discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted  
1307 along the phylogenetic tree. C) Trajectories of ASVs affiliated with the families *Actinomycetaceae*,  
1308 *Streptococcaceae*, *Prevotellaceae*, and Family XI (Class *Bacillales*), with decreasing abundances after HSCT  
1309 and recovery at late follow-up time points. The most abundant discriminating ASV for each family is  
1310 indicated. Detailed taxonomic information and LDA-coefficients of the displayed ASVs are listed in  
1311 Additional File 1: Table S2.

1312

1313 **Figure 4. Machine learning-based prediction of aGvHD severity from the pre-HSCT gut microbiota**  
1314 **composition.** A) Relative abundances of the 12 most abundant families over time in the gut in patients  
1315 with aGvHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive gut ASVs identified by the  
1316 svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the  
1317 respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted  
1318 aGvHD (0-I versus II-IV) from the abundances of gut ASVs pre-HSCT with 86% accuracy (95% CI: 65% to

1319 97%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C)  
1320 Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric  
1321 regression for prediction of aGvHD. Numbers along the branches indicate split values of variance stabilized  
1322 bacterial abundances. The terminal nodes show the proportion of samples originating from patients (n =  
1323 number of samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depicting the log transformed relative  
1324 abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared  
1325 with grade II-IV patients. E) Trajectories of *Lactobacillaceae* and *Tannerellaceae* ASVs that were identified  
1326 by tree-based sparse LDA, including ASV 3 and ASV 128 that were predictive for aGvHD (bold lines), in  
1327 patients with aGvHD grade 0-I vs II-IV.

1328

1329 **Figure 5. Machine learning-based prediction of aGvHD severity from the pre-HSCT oral microbiota**

1330 **composition.** A) Relative abundances the 12 most abundant families over time in the oral cavity in patients

1331 with aGvHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive oral ASVs identified by the

1332 svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the

1333 respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted

1334 aGvHD (0-I versus II-IV) from the abundances of oral ASVs pre-HSCT with 92% accuracy (95% CI: 73% to

1335 99%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C)

1336 Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric

1337 regression for prediction of aGvHD. Numbers along the branches indicate split values of variance stabilized

1338 bacterial abundances. The terminal nodes show the proportion of samples originating from patients (n =

1339 number of represented samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depict the log transformed

1340 relative abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I

1341 compared with grade II-IV patients. E) Trajectories of *Prevotellaceae* and *Actinomycetaceae* ASVs that

1342 were identified by tree-based sparse LDA, including ASV 226 and ASV 568 that were predictive for aGvHD

1343 (bold lines), in patients with aGvHD grade 0-I vs II-IV.

1344

1345 **Figure 6. Multivariate associations of the gut microbiota with immune and clinical parameters in HSCT.**

1346 A) Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis (dimensions

1347 1, 2, and 3) displaying pairwise correlations  $>0.3$ / $<-0.3$  between ASVs (bottom) and continuous immune

1348 and clinical parameters (right). Red indicates a positive correlation, and blue indicates a negative

1349 correlation, respectively. Based on the sPLS regression model, hierarchical clustering (clustering method:

1350 complete linkage, distance method: Pearson's correlation) was performed resulting in the three depicted

1351 clusters. B) Canonical correspondence analysis (CCpNA) relating gut microbial abundances (circles) to

1352 continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables with at least

1353 one correlation  $>0.3$ / $<-0.3$  in the sPLS analysis were included in the CCpNA. The triplot shows variables

1354 and ASVs with a score  $>0.3$ / $<-0.3$  on at least one of the first three CCpNA axes, displayed on axis 1 versus

1355 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval)

1356 correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. Abbreviations not

1357 mentioned in text: ATGmm, anti-thymocyte globulin; B\_, blood; BU, busulfan; CY, Cyclophosphamide;

1358 DonorMatch6, matched unrelated donor; FLU\_other, fludarabine combinations without thiotepa;

1359 GvHD.Prophylaxis1, treatment with cyclosporine; GvHD.Prophylaxis7, treatment with cyclosporine and

1360 methotrexate; immat\_B, immature B cells; K\_d100, Karnofsky score on day +100; K\_pre, Karnofsky score

1361 before HSCT; m1, month+1; m3, month+3; m6, month+6; m12, month+12; mat\_B, mature B cells; MEL,  
1362 melphalan; total\_B, total B cells; P\_, plasma; parasitic, parasitic infection; pre\_cond, before conditioning  
1363 start; pre\_exam, pre-examination; THIO, thiotepa; viral, viral infection; VP16, Etoposide.

1364

1365 **Supplementary Table S1. Patient characteristics.** Abbreviations: HLA, human leukocyte antigen; TBI, total  
1366 body irradiation; CY, Cyclophosphamide; VP16, Etoposide; BU, Busulfan; MEL, Melphalan; GvHD, graft-  
1367 versus-host disease.

1368

1369 **Supplementary Table S2. Taxonomy of a subset of LDA clade members and corresponding LDA-**  
1370 **coefficients in the gut, oral cavity, and nasal cavity.**

1371

1372 **Supplementary Table S3. Taxonomy of aGvHD predictors within the fecal, oral, and nasal microbiota.**

1373 ASVs that were significantly predicting aGvHD severity according to the conditional inference tree  
1374 regression model are highlighted in bold. Of the 50 most important gut ASVs identified by the svmLinear  
1375 model, 17 were confirmed by Boruta feature selection and are listed here. In the oral and nasal cavities,  
1376 26 and 12 ASVs were confirmed by Boruta selection, respectively. Listed in bold are those ASVs with a  
1377 significant predictive effect on aGvHD severity, tested in a regression framework with CTREE (see  
1378 Methods).

1379

1380 **Figure S1. The gut microbiota in the HSCT patients at pre-exam differs from the gut microbiota of age-**  
1381 **matched healthy children.** A) Fecal bacterial alpha diversity (inverse Simpson index) was 2.4-fold higher  
1382 in healthy children (n=18) compared to children at pre-examination before HSCT (n=15). B) Fecal bacterial  
1383 composition was significantly different between the two groups (anosim, p=0.001, R=0.44), and within-  
1384 group variance was significantly greater in the HSCT group (betadisper, p<0.001). C) The taxa which best  
1385 explain differences in community structure between HSCT patients at preexamination and healthy  
1386 children were identified by analysis of LEfSe (Linear discriminant analysis Effect Size). LEfSe accounts for  
1387 the hierarchical structure of bacterial phylogeny, thereby allowing identification of differentially abundant  
1388 taxa on several taxonomic levels (here: kingdom to genus). Count data was centered-log ratio (CLR)  
1389 transformed within the LEfSe analysis. The higher the LDA score (log10), the higher the effect size of the  
1390 respective taxon in explaining group difference. Here, we show taxa with an LDA score >4. D) Differentially  
1391 abundant genera between the two groups were additionally identified by DESeq2. Of the top 100 most  
1392 abundant genera (of the whole gut microbiota data set), eighteen genera were significantly more  
1393 abundant in healthy children (yellow), and 15 genera were significantly more abundant in the patients at  
1394 preexam (purple). Differences in median proportions of these genera (and their supertaxa) are displayed  
1395 in a heat tree. See also additional information at <https://doi.org/10.6084/m9.figshare.13614230>.

1396

1397

1398 **Figure S2. Most abundant taxonomic families in the gut, oral cavity, and nasal cavity in allo-HSCT**  
1399 **patients.** Rank abundance curves displaying the proportions of the 12 most abundant taxonomic families  
1400 at each body site (gut, oral cavity, and nasal cavity).

1401

1402 **Figure S3. Tree-based sparse linear discriminant analysis revealing nasal ASVs that distinguish time**  
1403 **points from each other in relation to HSCT.** A) Relative abundances over time of the 12 most abundant  
1404 families in the nasal cavity. B) Coefficients of discriminating clades of ASVs on the first LDA axis, colored  
1405 by taxonomic family, and plotted along the phylogenetic tree. C) Trajectories of ASVs in one discriminating  
1406 group, affiliated with the family *Corynebacteriaceae*, with decreasing abundances after HSCT and recovery  
1407 at late follow-up time points. The most abundant discriminating ASV is indicated. Detailed taxonomic  
1408 information and LDA-coefficients of the displayed ASVs are listed in Table S2.

1409  
1410 **Figure S4. Machine learning-based prediction of aGvHD severity from nasal microbial abundances pre-**  
1411 **HSCT.** A) Relative abundances of the 12 most abundant families over time in the nasal cavity in patients  
1412 with aGvHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive nasal ASVs identified by the  
1413 svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the  
1414 respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted  
1415 aGvHD (0-I versus II-IV) from the abundances of nasal ASVs pre-HSCT with 76% accuracy (95% CI: 56% to  
1416 90%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C)  
1417 Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric  
1418 regression for prediction of aGvHD. Numbers along the branches indicate split values of variance stabilized  
1419 bacterial abundances. The terminal nodes show the proportion of samples originating from patients with  
1420 aGvHD grade 0-I vs II-IV (n = number of samples). D) Boxplots depict the log transformed relative  
1421 abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared  
1422 with grade II-IV patients.

1423  
1424 **Figure S5. Multivariate associations of the oral microbiota with immune and clinical parameters in HSCT.**  
1425 A) Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis dimensions  
1426 1, 2, and 3, displaying pairwise correlations  $>0.2/ <-0.2$  between oral ASVs (bottom), and continuous  
1427 immune and clinical parameters (right). Red indicated positive correlation, and blue indicates negative  
1428 correlation, respectively. Based on the sPLS regression model, hierarchical clustering (clustering method:  
1429 complete linkage, distance method: Pearson's correlation) was performed resulting in the three depicted  
1430 clusters. B) Canonical correspondence analysis (CCpNA) relating oral microbial abundances (circles) to  
1431 continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables with at least  
1432 one correlation  $>0.2/ <-0.2$  in the sPLS analysis were included in the CCpNA. The triplot shows variables  
1433 and ASVs with a score  $>0.3/ <-0.3$  on at least one the first three CCpNA axes, displayed on axis 1 versus 2  
1434 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval)  
1435 correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. For visualization  
1436 purposes, a focused section of the CCpNA triplot is shown. Abbreviations are described in Figure 6.  
1437 Additional abbreviations: fungal, fungal infection; haploid, haploidentical donor; hemo, hemoglobin;  
1438 leuko, leukocytes; lympho, lymphocytes; w1, week+1; w2, week+2; w3, week+3.

1439  
1440 **Figure S6. Multivariate associations of the nasal microbiota with immune and clinical parameters in**  
1441 **HSCT.** A) Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis  
1442 dimensions 1, 2, and 3, displaying pairwise correlations  $>0.2/ <-0.2$  between nasal ASVs (bottom), and  
1443 continuous immune and clinical parameters (right). Red indicated positive correlation, and blue indicates

1444 negative correlation, respectively. Based on the sPLS regression model, hierarchical clustering (clustering  
1445 method: complete linkage, distance method: Pearson's correlation) was performed resulting in the three  
1446 depicted clusters. B) Canonical correspondence analysis (CCpNA) relating nasal microbial abundances  
1447 (circles) to continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables  
1448 with at least one correlation  $>0.2$ / $<-0.2$  in the sPLS analysis were included in the CCpNA. The triplot shows  
1449 variables and ASVs with a score  $>0.3$ / $<-0.3$  on at least one the first three CCpNA axes, displayed on axis 1  
1450 versus 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval)  
1451 correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. For visualization  
1452 purposes, a focused section of the CCpNA triplot is shown. Abbreviations are described in Figures 6 and  
1453 S5. Additional abbreviations: DonorMatch8, unrelated donor with 1 HLA mismatch; PB, peripheral blood.

## Additional File 1 – Supplementary Tables

**Supplementary Table S1. Patient characteristics.** Abbreviations: HLA, human leukocyte antigen; TBI, total body irradiation; CY, Cyclophosphamide; VP16, Etoposide; BU, Busulfan; MEL, Melphalan; GvHD, graft-versus-host disease.

Characteristics	Number of patients	Percentage of all patients (%)	
Transplant recipients	29	100	
Average recipient age in years	9.5 (Range: 2.5-16.4)	NA	
Gut microbiota characterized	At ≤ 6 timepoints	9	31
	At 7-8 timepoints	13	44.8
	At 9-10 timepoints	7	24.1
Oral microbiota characterized	At ≤ 6 timepoints	3	10.3
	At 7-8 timepoints	11	37.9
	At 9-10 timepoints	15	51.7
Nasal microbiota characterized	At ≤ 6 timepoints	3	10.3
	At 7-8 timepoints	9	31
	At 9-10 timepoints	17	58.6
Patient sex	Female	13	44.8
	Male	16	55.2
Disease at transplantation	Malignant hematologic diseases	20	69
	Severe aplastic anemia	1	3.4
	Other benign disorders including immunodeficiencies	8	27.6
Donor type	HLA-matched sibling	14	48.3
	HLA-matched unrelated donor (9/10 or 10/10 match)	12	41.4
	Haplo-identical related donor (5/10 match)	3	10.3
Stem cell source	Bone marrow	23	79.3
	Umbilical cord blood	2	6.9
	Peripheral blood	4	13.8
Conditioning regimen	TBI + CY / TBI + VP16	6	20.7
	Combinations of BU, CY, VP16 and MEL	6	20.7
	Fludarabine-thiothepa combinations	12	41.4

	<b>Other combinations with fludarabine</b>	5	17.2
<b>Anti-thymocyte globulin treatment</b>		16	55.2
<b>Antibiotics pre- and post-HSCT</b>		29	100
<b>Sex mismatch (female donor to male recipient)</b>		7	24.1
<b>Acute GvHD</b>	<b>Grade 0-I</b>	20	69
	<b>Grade II-IV</b>	9	31
<b>Infections</b>	<b>At least one registered bacterial infection</b>	25	86.2
	<b>At least one registered viral infection</b>	29	100
	<b>At least one registered fungal infection</b>	14	48.3
<b>Overall Survival*</b>	<b>Alive</b>	26	89.7
	<b>Dead</b>	3	10.3
<b>Relapse of primary disease*</b>		2	6.9
<b>Non-relapse mortality</b>		2	6.9
<b>Re-transplantation</b>		1	3.4

\*Mean follow-up time after HSCT: 21.4 months (range: 10.1 – 32.7 months)

**Supplementary Table S2. Taxonomy of a subset of LDA clade members and corresponding LDA-coefficients in the gut, oral cavity, and nasal cavity.**

ASV number	Body site	Phylum	Family	Genus/Species/Description	LDA coefficient (beta)	Number of patients with presence of this ASV at $\geq 1$ time point
	<b>Gut</b>					
ASV_1	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.082	29
ASV_30	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	26
ASV_158	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	13
ASV_178	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	12
ASV_344	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	8
ASV_395	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.012	14
ASV_424	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.012	9
ASV_543	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	7
ASV_552	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus avium</i>	0.006	7
ASV_720	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.078	15
ASV_730	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	12
ASV_784	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	8
ASV_951	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Enterococcus</i> sp.	0.006	12
ASV_1186	gut	<i>Firmicutes</i>	<i>Enterococcaceae</i>	<i>Melissococcus</i> sp.	0.006	9
ASV_3	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus</i> sp.	0.011	26
ASV_31	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus</i> sp.	0.006	11
ASV_67	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus sakei</i>	0.006	15
ASV_74	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus</i> sp.	0.006	16
ASV_144	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Pediococcus</i> sp.	0.006	10
ASV_151	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Pediococcus pentosaceus</i>	0.006	12
ASV_189	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus</i> sp.	0.006	10
ASV_222	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus rhamnosus</i>	0.006	8
ASV_586	gut	<i>Firmicutes</i>	<i>Lactobacillaceae</i>	<i>Lactobacillus</i> sp.	0.006	9
ASV_78	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia</i> sp.	-0.039	26
ASV_237	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia faecis</i>	-0.026	23
ASV_265	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia</i> sp.	-0.026	24
ASV_271	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia obeum</i>	-0.026	12
ASV_375	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia</i> sp.	-0.026	15
ASV_554	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia</i> sp.	-0.026	16
ASV_1291	gut	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Blautia</i> sp.	-0.026	9
ASV_132	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcus bromii</i>	-0.007	16
ASV_192	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcus bromii</i>	-0.007	24
ASV_206	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	DTU089	-0.024	23
ASV_306	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcus</i> sp.	-0.007	13

ASV_348	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcus</i> sp.	-0.007	11
ASV_556	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Caproiciproducens</i> sp.	-0.024	9
ASV_566	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminoclostridium</i> sp.	-0.024	13
ASV_583	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminoclostridium</i> sp.	-0.024	23
ASV_682	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Caproiciproducens</i> sp.	-0.024	9
ASV_750	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Caproiciproducens</i> sp.	-0.024	13
ASV_792	gut	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminoclostridium</i> sp.	-0.024	15
	<b>Oral</b>					
ASV_18	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.021	27
ASV_66	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	26
ASV_117	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces naeslundii</i>	0.016	27
ASV_126	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces naeslundii</i>	0.016	19
ASV_138	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	22
ASV_155	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces odontolyticus</i>	0.003	20
ASV_227	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.016	7
ASV_235	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	F0332	0.003	18
ASV_262	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	18
ASV_345	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.016	9
ASV_389	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces odontolyticus</i>	0.003	7
ASV_403	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	15
ASV_407	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.016	7
ASV_2693	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.021	5
ASV_422	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces odontolyticus</i>	0.003	15
ASV_431	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.016	12
ASV_2697	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces naeslundii</i>	0.016	9
ASV_461	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	11
ASV_475	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces gerencseriae</i>	0.003	6
ASV_484	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	4
ASV_501	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces odontolyticus</i>	0.003	13
ASV_516	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces odontolyticus</i>	0.003	12
ASV_2700	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces gerencseriae</i>	0.003	14
ASV_568	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	12
ASV_600	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	7
ASV_642	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	7
ASV_798	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	F0332	0.003	10
ASV_871	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces massiliensis</i>	0.003	13
ASV_2729	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces graevenitzii</i>	0.003	9
ASV_1055	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	6

ASV_1172	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.003	7
ASV_2751	oral	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> sp.	0.016	8
ASV_2664	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	29
ASV_10	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	29
ASV_16	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	29
ASV_27	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	29
ASV_28	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.040	28
ASV_37	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	19
ASV_48	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	27
ASV_173	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus salivarius</i>	0.010	15
ASV_183	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	19
ASV_188	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	14
ASV_2674	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	14
ASV_230	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	8
ASV_269	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus cristatus</i>	0.010	29
ASV_282	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus parasanguinis</i>	0.010	13
ASV_2683	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus mitis</i>	0.010	11
ASV_480	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	18
ASV_481	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus peroris</i>	0.010	12
ASV_802	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	24
ASV_1531	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	11
ASV_1599	oral	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	0.010	13
ASV_42	oral	<i>Bacteroidetes</i>	<i>Prevotellaceae</i>	<i>Prevotella melaninogenica</i>	0.028	27
ASV_226	oral	<i>Bacteroidetes</i>	<i>Prevotellaceae</i>	<i>Prevotella melaninogenica</i>	0.028	15
ASV_800	oral	<i>Bacteroidetes</i>	<i>Prevotellaceae</i>	<i>Prevotella</i> sp.	0.028	6
ASV_2665	oral	<i>Firmicutes</i>	Family XI	<i>Gemella</i> sp.	0.009	29
ASV_208	oral	<i>Firmicutes</i>	Family XI	<i>Gemella sanguinis</i>	0.009	28
ASV_2701	oral	<i>Firmicutes</i>	Family XI	<i>Gemella</i> sp.	0.009	5
	<b>Nasal</b>					
ASV_14	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium</i> sp.	0.117	24
ASV_354	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium durum</i>	0.054	15
ASV_360	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium</i> sp.	0.054	17
ASV_2704	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium</i> sp.	0.022	12
ASV_2707	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium durum</i>	0.054	14
ASV_1206	nasal	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Lawsonella</i> sp.	0.007	16

**Supplementary Table S3. Taxonomy of aGvHD predictors within the fecal, oral, and nasal microbiota.**

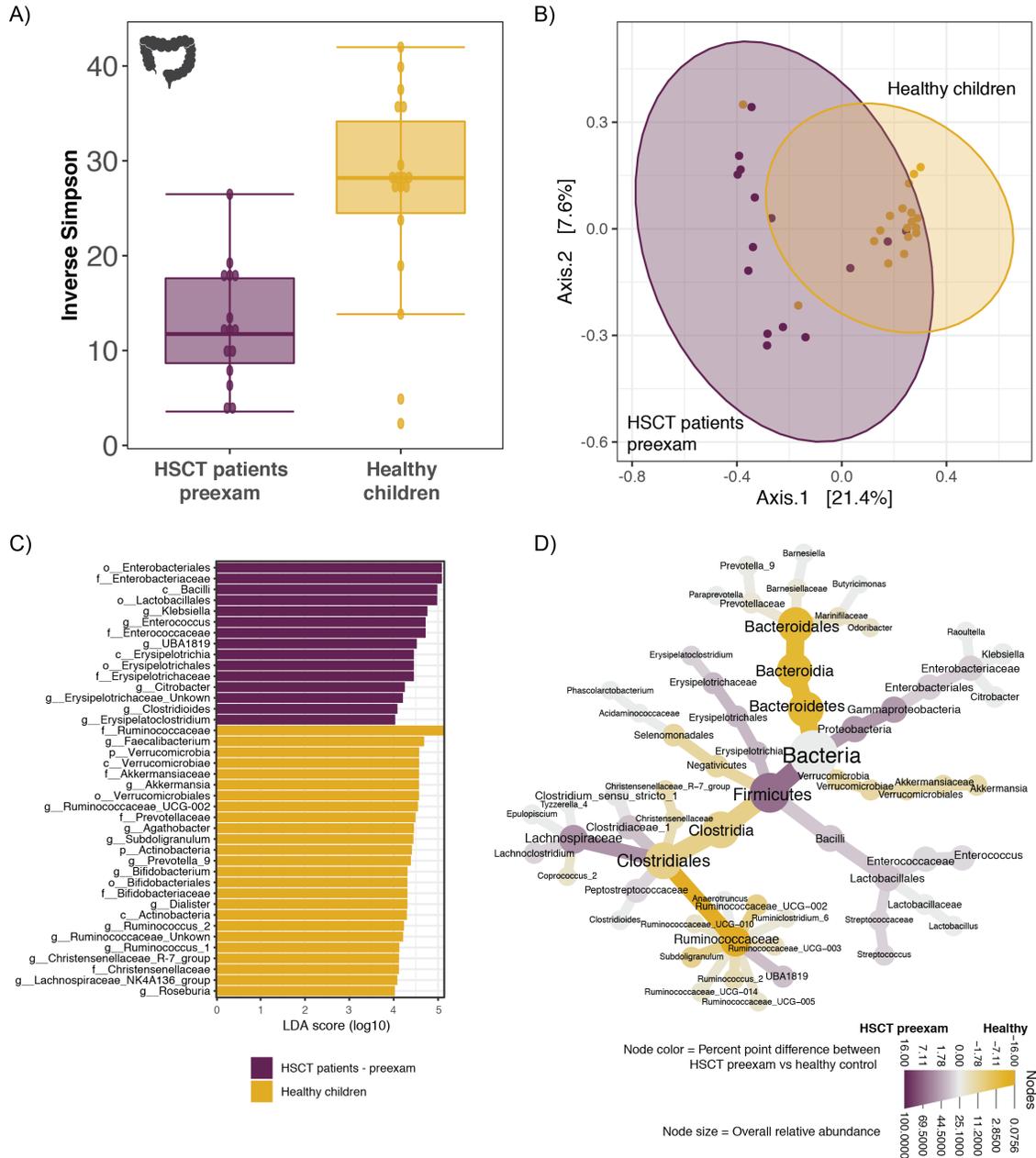
ASVs that were significantly predicting aGvHD severity according to the conditional inference tree regression model are highlighted in bold. Of the 50 most important gut ASVs identified by the svmLinear model, 17 were confirmed by Boruta feature selection and are listed here. In the oral and nasal cavities, 26 and 12 ASVs were confirmed by Boruta selection, respectively. Listed in bold are those ASVs with a significant predictive effect on aGvHD severity, tested in a regression framework with CTREE (see Methods).

ASV number	Phylum	Family	Genus/Species/Description	Body Site	Importance in svmLinear model	Mean Importance in Boruta
	<b>Gut</b>					
<b>ASV_3</b>	<b>Firmicutes</b>	<b>Lactobacillaceae</b>	<b>Lactobacillus sp.</b>	gut	<b>0.74</b>	<b>8.85</b>
ASV_7	Bacteroidetes	Tannerellaceae	Parabacteroides merdae	gut	0.63	4.62
ASV_8	Proteobacteria	Enterobacteriaceae	Escherichia/Shigella sp.	gut	0.70	4.16
ASV_12	Firmicutes	Ruminococcaceae	UBA1819 sp.	gut	0.68	5.52
ASV_35	Bacteroidetes	Bacteroidaceae	Bacteroides uniformis	gut	0.68	4.66
ASV_50	Firmicutes	Ruminococcaceae	Subdoligranulum sp.	gut	0.60	3.78
<b>ASV_128</b>	<b>Bacteroidetes</b>	<b>Tannerellaceae</b>	<b>Parabacteroides distasonis</b>	gut	<b>0.68</b>	<b>7.42</b>
ASV_131	Bacteroidetes	Tannerellaceae	Parabacteroides distasonis	gut	0.66	6.49
ASV_166	Fusobacteria	Fusobacteriaceae	Fusobacterium sp.	gut	0.60	4.08
<b>ASV_268</b>	<b>Firmicutes</b>	<b>Lachnospiraceae</b>	<b>Lachnospiraceae_NK4A136_group sp.</b>	gut	<b>0.61</b>	<b>5.02</b>
ASV_361	Bacteroidetes	Tannerellaceae	Parabacteroides sp.	gut	0.61	4.40
ASV_477	Firmicutes	Lachnospiraceae	Eisenbergiella tayi	gut	0.62	3.64
ASV_507	Firmicutes	Lachnospiraceae	Lachnospira pectinoschiza	gut	0.60	3.91
ASV_563	Firmicutes	Erysipelotrichaceae	Coprobacillus cateniformis	gut	0.63	5.38
ASV_567	Firmicutes	Ruminococcaceae	Flavonifractor sp.	gut	0.60	2.81
ASV_585	Bacteroidetes	Tannerellaceae	Parabacteroides sp.	gut	0.62	4.53
ASV_687	Firmicutes	Lachnospiraceae	NA	gut	0.59	3.83
	<b>Oral</b>					
ASV_15	Firmicutes	Streptococcaceae	Lactococcus sp.	oral	0.61	7.31
ASV_34	Proteobacteria	Neisseriaceae	Neisseria sp.	oral	0.62	5.08
ASV_66	Actinobacteria	Actinomycetaceae	Actinomyces sp.	oral	0.65	4.84
ASV_72	Actinobacteria	Micrococcaceae	Rothia sp.	oral	0.63	4.49
ASV_92	Bacteroidetes	Bacteroidaceae	Bacteroides sp.	oral	0.65	5.36
ASV_98	Actinobacteria	Micrococcaceae	Rothia mucilaginoso	oral	0.58	5.87
ASV_135	Firmicutes	Streptococcaceae	Streptococcus mutans	oral	0.60	4.24

ASV_205	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcacea</i> <i>e_UCG-002</i> sp.	oral	0.63	4.38
<b>ASV_226</b>	<b><i>Bacteroidetes</i></b>	<b><i>Prevotellaceae</i></b>	<b><i>Prevotella_7</i></b> <b><i>melaninogenica</i></b>	<b>oral</b>	<b>0.63</b>	<b>8.19</b>
ASV_247	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcacea</i> <i>e_UCG-013</i> sp.	oral	0.66	4.80
ASV_270	<i>Firmicutes</i>	<i>Veillonellaceae</i>	<i>Veillonella</i> sp.	oral	0.59	3.47
ASV_338	<i>Bacteroidetes</i>	<i>Prevotellaceae</i>	<i>Alloprevotella</i> sp.	oral	0.61	4.26
ASV_360	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium</i> sp.	oral	0.64	6.50
ASV_374	<i>Fusobacteria</i>	<i>Leptotrichiaceae</i>	<i>Leptotrichia</i> <i>wadei</i>	oral	0.60	4.39
ASV_378	<i>Firmicutes</i>	<i>Lachnospiraceae</i>	<i>Oribacterium</i> <i>sinus</i>	oral	0.65	7.17
<b>ASV_500</b>	<b><i>Actinobacteria</i></b>	<b><i>Propionibacteriaceae</i></b>	<b><i>Pseudopropionib</i></b> <b><i>acterium</i></b> <b><i>propionicum</i></b>	<b>oral</b>	<b>0.67</b>	<b>8.51</b>
ASV_518	<i>Bacteroidetes</i>	<i>Bacteroidaceae</i>	<i>Bacteroides</i> sp.	oral	0.64	4.67
<b>ASV_568</b>	<b><i>Actinobacteria</i></b>	<b><i>Actinomycetaceae</i></b>	<b><i>Actinomyces</i> sp.</b>	<b>oral</b>	<b>0.71</b>	<b>8.79</b>
ASV_640	<i>Bacteroidetes</i>	<i>Tannerellaceae</i>	<i>Parabacteroides</i> sp.	oral	0.62	2.81
ASV_871	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> <i>massiliensis</i>	oral	0.61	4.65
ASV_894	<i>Firmicutes</i>	<i>Ruminococcaceae</i>	<i>Ruminococcacea</i> <i>e_UCG-014</i> sp.	oral	0.62	2.74
ASV_140 3	<i>Firmicutes</i>	<i>Veillonellaceae</i>	<i>Selenomonas_3</i> sp.	oral	0.59	2.82
ASV_266 4	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	oral	0.68	6.11
ASV_266 6	<i>Actinobacteria</i>	<i>Micrococcaceae</i>	<i>Rothia</i> <i>dentocariosa</i>	oral	0.61	4.47
ASV_267 6	<i>Firmicutes</i>	<i>Aerococcaceae</i>	<i>Abiotrophia</i> <i>defectiva</i>	oral	0.61	5.74
ASV_269 2	<i>Actinobacteria</i>	<i>Micrococcaceae</i>	<i>Rothia</i> sp.	oral	0.62	6.30
	<b>Nasal</b>					
ASV_266 4	<i>Firmicutes</i>	<i>Streptococcaceae</i>	<i>Streptococcus</i> sp.	nasal	0.61	11.49
ASV_266 6	<i>Actinobacteria</i>	<i>Micrococcaceae</i>	<i>Rothia</i> <i>dentocariosa</i>	nasal	0.56	5.80
<b>ASV_47</b>	<b><i>Actinobacteria</i></b>	<b><i>Micrococcaceae</i></b>	<b><i>Rothia</i> sp.</b>	<b>nasal</b>	<b>0.59</b>	<b>6.41</b>
ASV_51	<i>Proteobacteria</i>	<i>Pasteurellaceae</i>	<i>Haemophilus</i> sp.	nasal	0.64	5.67
ASV_52	<i>Actinobacteria</i>	<i>Micrococcaceae</i>	<i>Rothia</i> <i>mucilaginoso</i>	nasal	0.61	8.58
<b>ASV_66</b>	<b><i>Actinobacteria</i></b>	<b><i>Actinomycetaceae</i></b>	<b><i>Actinomyces</i> sp.</b>	<b>nasal</b>	<b>0.67</b>	<b>9.05</b>
ASV_267 0	<i>Proteobacteria</i>	<i>Neisseriaceae</i>	NA	nasal	0.61	6.69
ASV_117	<i>Actinobacteria</i>	<i>Actinomycetaceae</i>	<i>Actinomyces</i> <i>naeslundii</i>	nasal	0.58	4.73
ASV_125	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Corynebacterium</i> <i>_1</i> <i>accolens</i>	nasal	0.57	5.43
ASV_252	<i>Firmicutes</i>	<i>Veillonellaceae</i>	<i>Veillonella</i> <i>rogosae</i>	nasal	0.58	5.54
ASV_270	<i>Firmicutes</i>	<i>Veillonellaceae</i>	<i>Veillonella</i> sp.	nasal	0.57	5.41

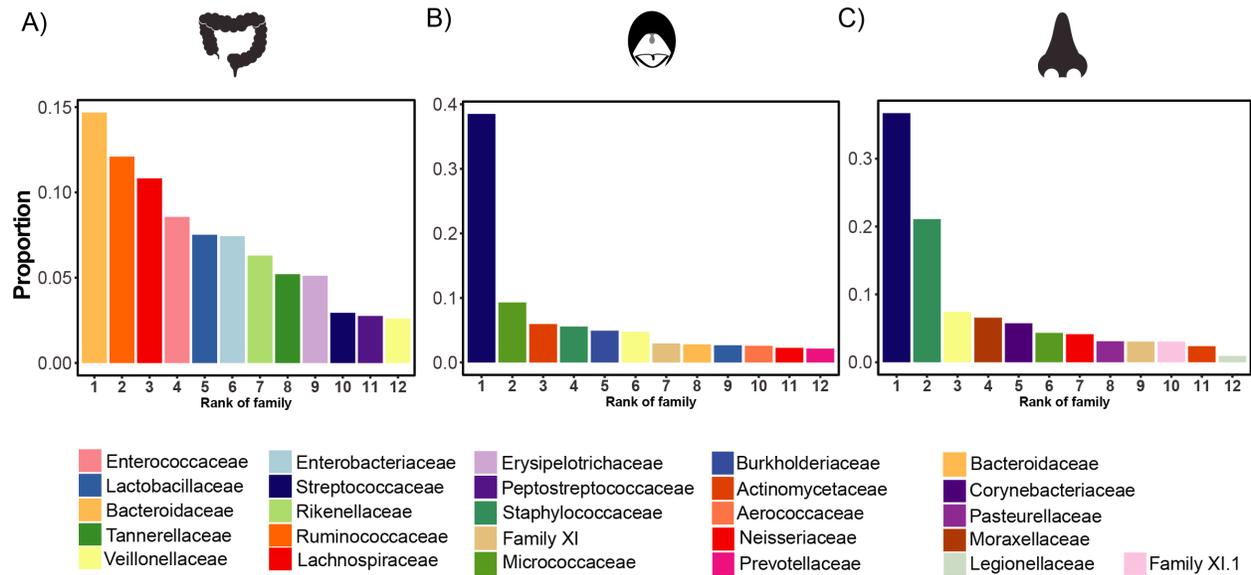
ASV_269 4	<i>Actinobacteria</i>	<i>Corynebacteriaceae</i>	<i>Lawsonella</i> sp.	nasal	0.65	8.42
--------------	-----------------------	---------------------------	-----------------------	-------	------	------

## Additional File 2 – Supplementary Figures

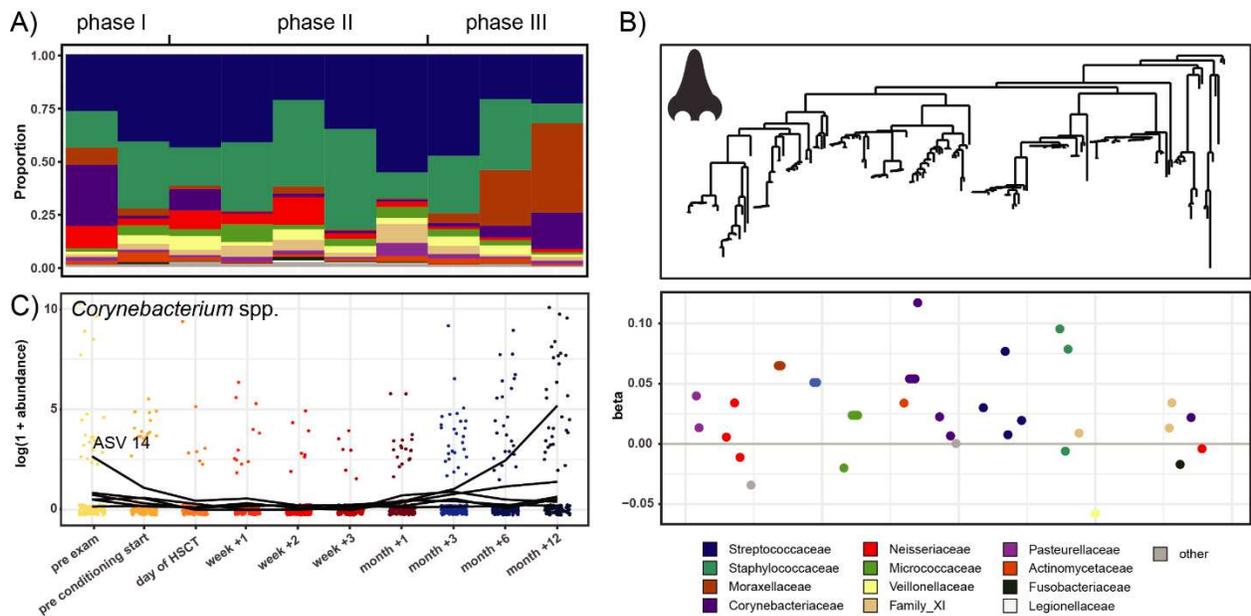


**Figure S1. The gut microbiota in the HSCT patients at pre-exam differs from the gut microbiota of age-matched healthy children.** A) Fecal bacterial alpha diversity (inverse Simpson index) was 2.4-fold higher in healthy children (n=18) compared to children at pre-examination before HSCT (n=15). B) Fecal bacterial composition was significantly different between the two groups (anosim,  $p=0.001$ ,  $R=0.44$ ), and within-group variance was significantly greater in the HSCT group (betadisper,  $p<0.001$ ). C) The taxa which best explain differences in community structure between HSCT patients at preexamination and healthy children were identified by analysis of LefSe (Linear discriminant analysis Effect Size). LefSe accounts for the hierarchical structure of bacterial phylogeny, thereby allowing identification of differentially abundant taxa on several taxonomic levels (here: kingdom to genus). Count data was centered-log ratio (CLR) transformed within the LefSe analysis. The higher the LDA score ( $\log_{10}$ ), the higher the effect

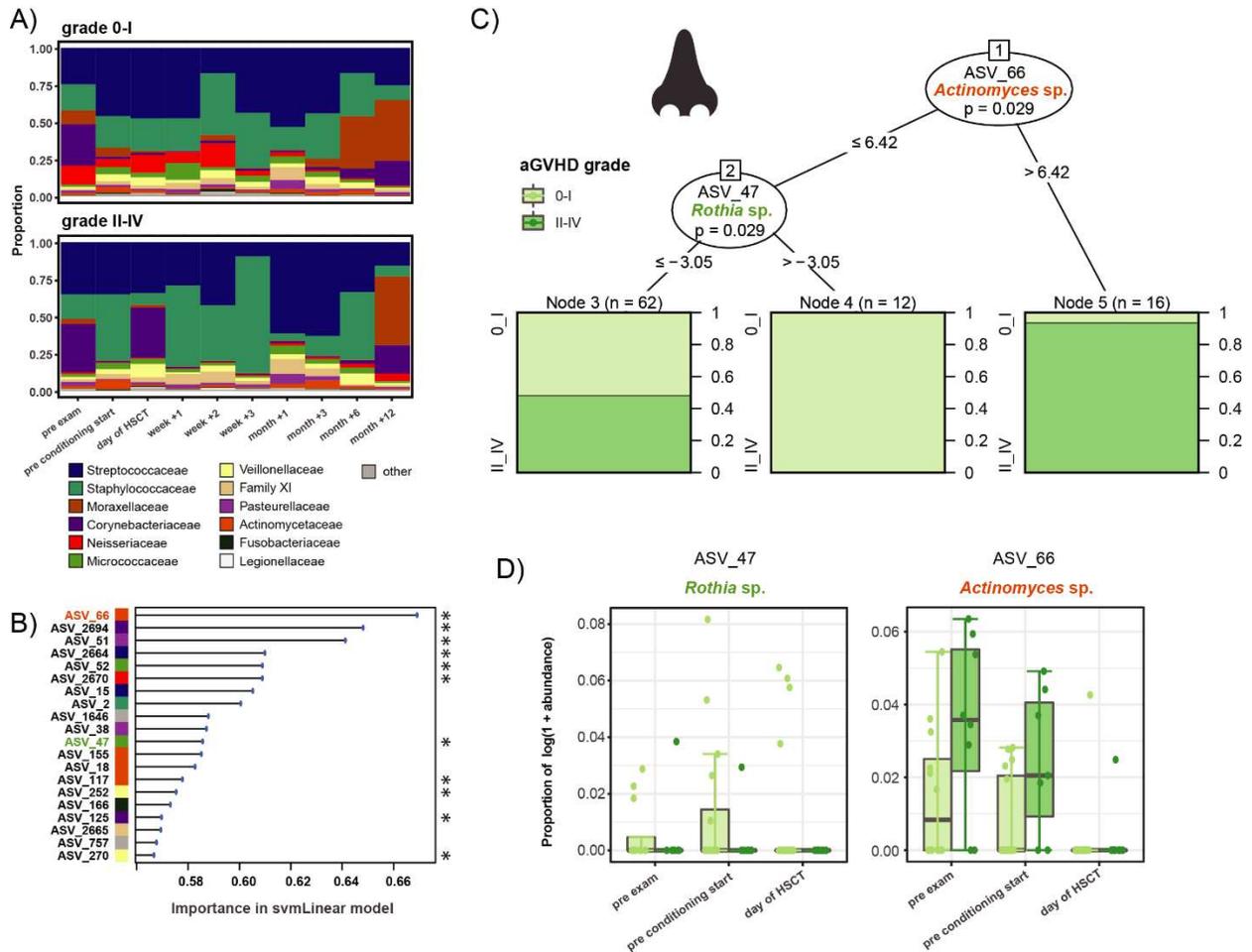
size of the respective taxon in explaining group difference. Here, we show taxa with an LDA score >4. D) Differentially abundant genera between the two groups were additionally identified by DESeq2. Of the top 100 most abundant genera (of the whole gut microbiota data set), eighteen genera were significantly more abundant in healthy children (yellow), and 15 genera were significantly more abundant in the patients at preexam (purple). Differences in median proportions of these genera (and their supertaxa) are displayed in a heat tree. See also additional information at <https://doi.org/10.6084/m9.figshare.13614230>.



**Figure S2. Most abundant taxonomic families in the gut, oral cavity, and nasal cavity in allo-HSCT patients.** Rank abundance curves displaying the proportions of the 12 most abundant taxonomic families at each body site (gut, oral cavity, and nasal cavity).

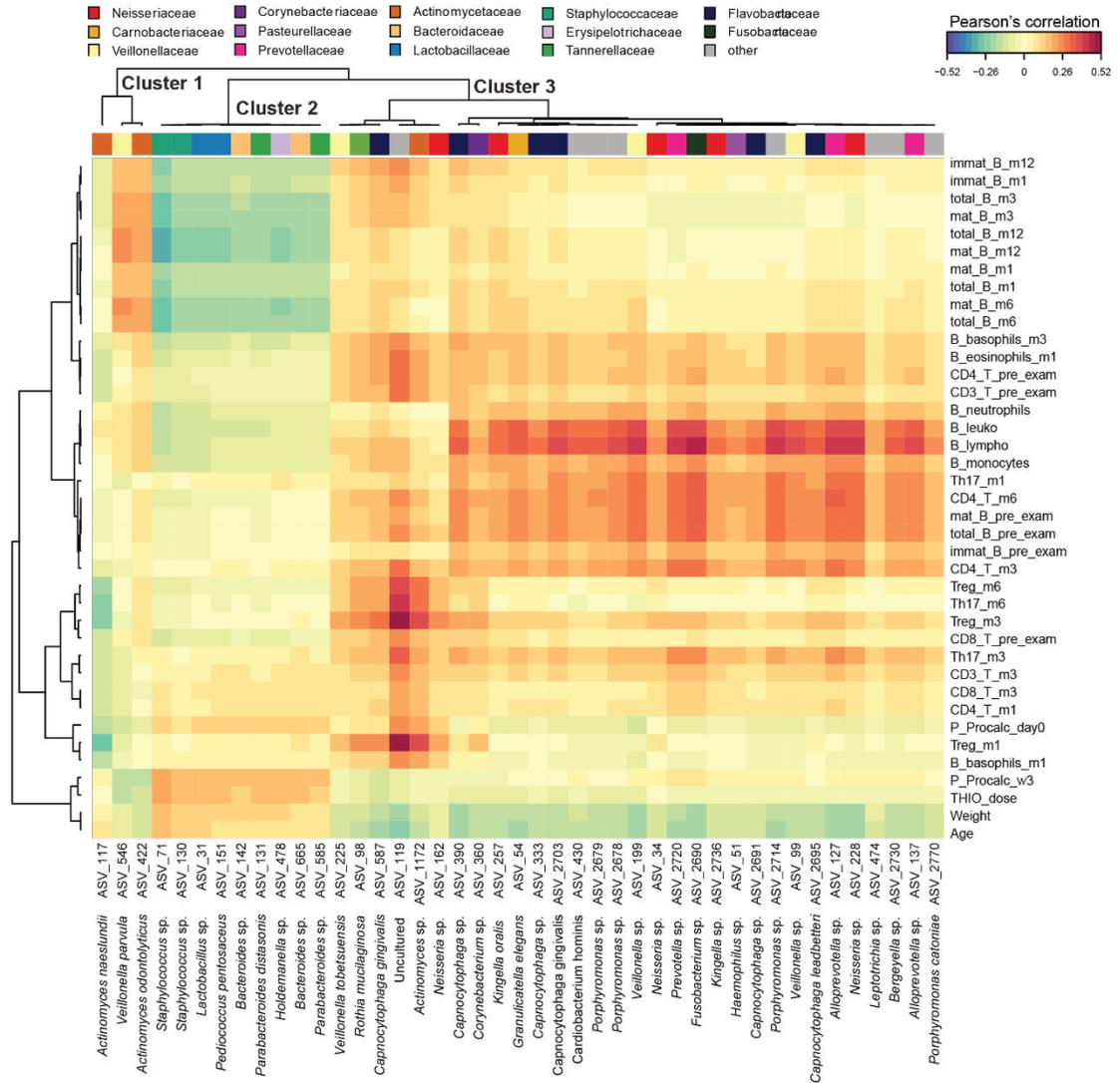


**Figure S3. Tree-based sparse linear discriminant analysis revealing nasal ASVs that distinguish time points from each other in relation to HSCT.** A) Relative abundances over time of the 12 most abundant families in the nasal cavity. B) Coefficients of discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the phylogenetic tree. C) Trajectories of ASVs in one discriminating group, affiliated with the family *Corynebacteriaceae*, with decreasing abundances after HSCT and recovery at late follow-up time points. The most abundant discriminating ASV is indicated Detailed taxonomic information and LDA-coefficients of the displayed ASVs are listed in Table S2.

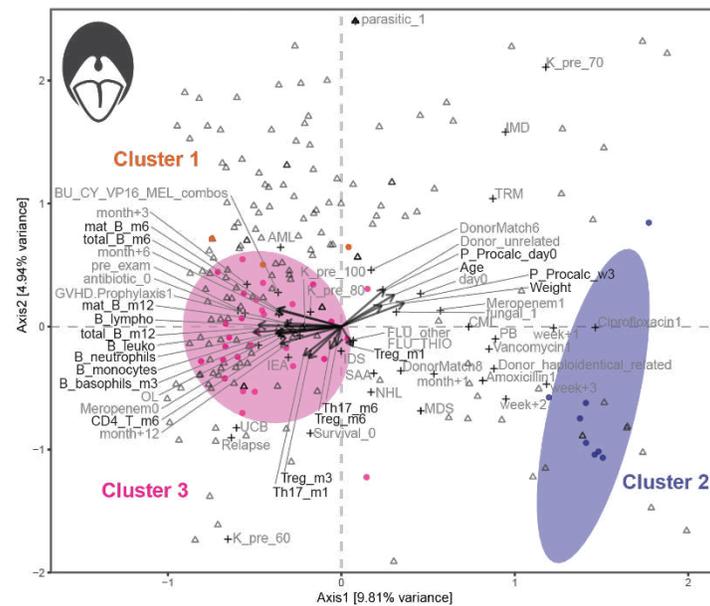


**Figure S4. Machine learning-based prediction of aGVHD severity from nasal microbial abundances pre-HSCT.** A) Relative abundances of the 12 most abundant families over time in the nasal cavity in patients with aGVHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive nasal ASVs identified by the svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted aGVHD (0-I versus II-IV) from the abundances of nasal ASVs pre-HSCT with 76% accuracy (95% CI: 56% to 90%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C) Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric regression for prediction of aGVHD. Numbers along the branches indicate split values of variance stabilized bacterial abundances. The terminal nodes show the proportion of samples originating from patients with aGVHD grade 0-I vs II-IV (n = number of samples). D) Boxplots depict the log transformed relative abundances of the predictive ASVs at time points up to the transplantation in aGVHD grade 0-I compared with grade II-IV patients.

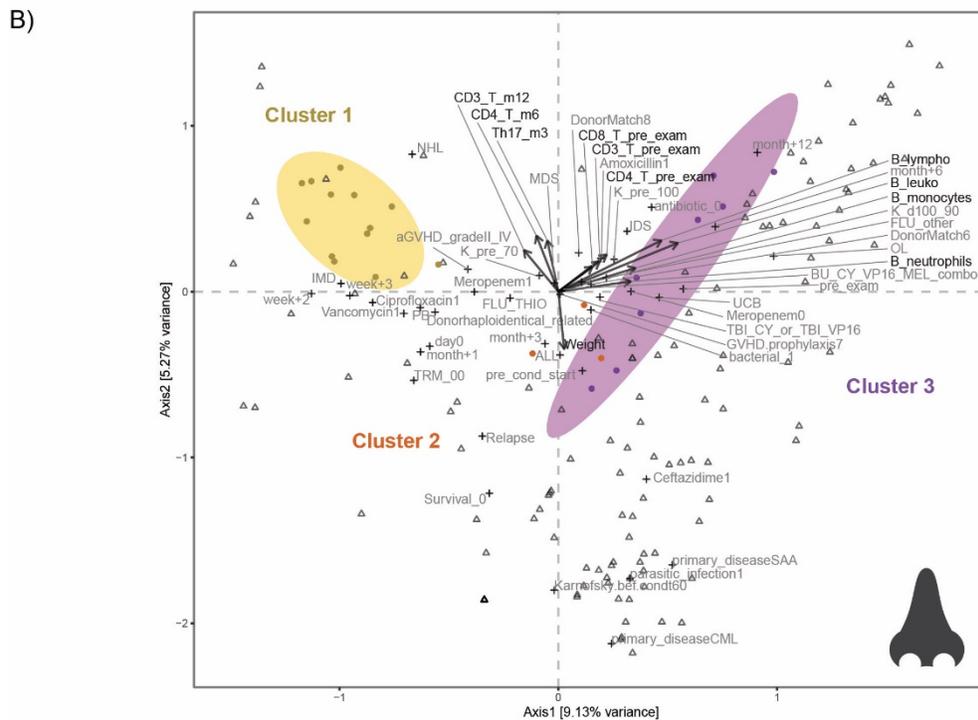
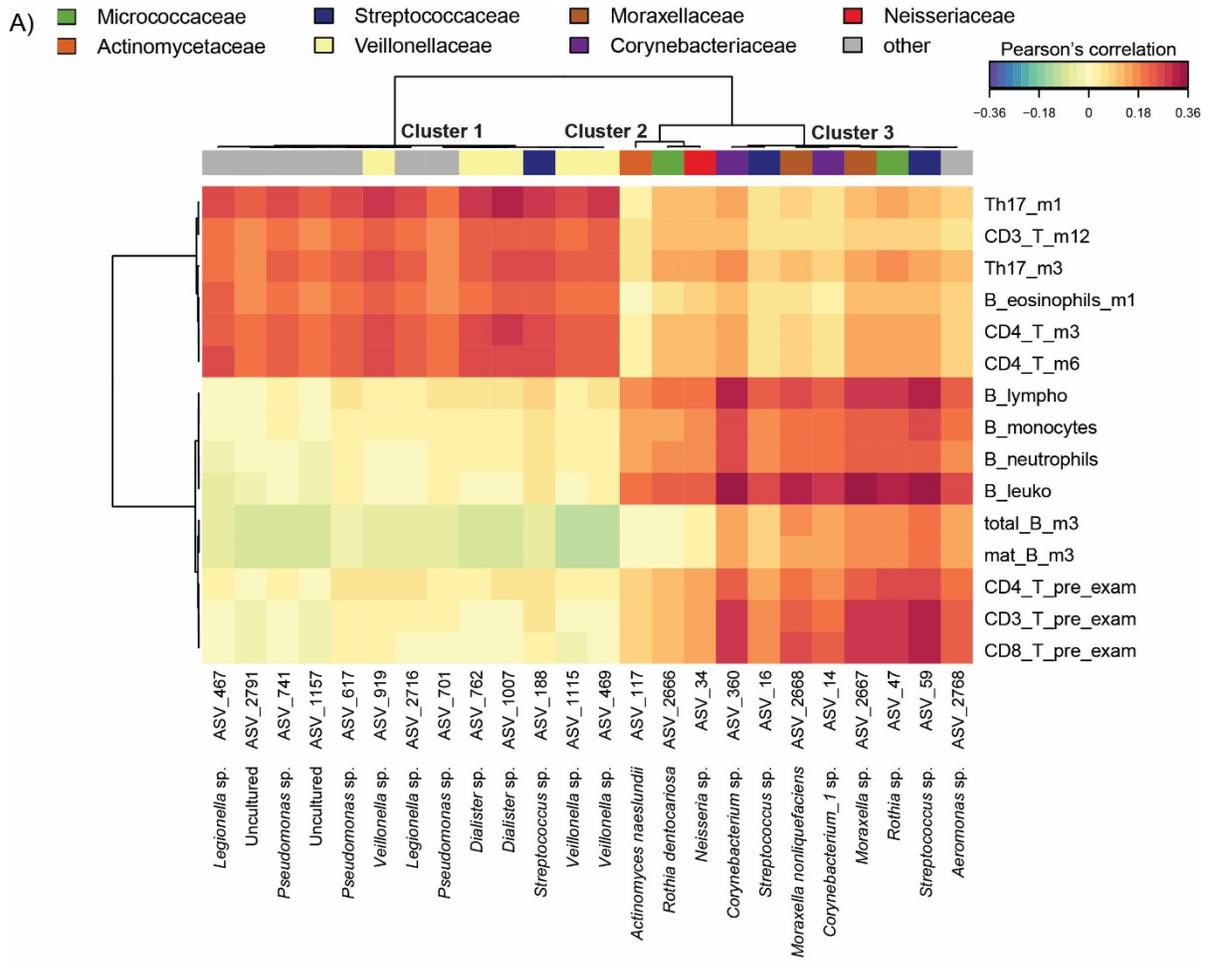
A)



B)



**Figure S5. Multivariate associations of the oral microbiota with immune and clinical parameters in HSCT.** A) Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis dimensions 1, 2, and 3, displaying pairwise correlations  $>0.2$ / $<-0.2$  between oral ASVs (bottom), and continuous immune and clinical parameters (right). Red indicated positive correlation, and blue indicates negative correlation, respectively. Based on the sPLS regression model, hierarchical clustering (clustering method: complete linkage, distance method: Pearson's correlation) was performed resulting in the three depicted clusters. B) Canonical correspondence analysis (CCpNA) relating oral microbial abundances (circles) to continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables with at least one correlation  $>0.2$ / $<-0.2$  in the sPLS analysis were included in the CCpNA. The triplot shows variables and ASVs with a score  $>0.3$ / $<-0.3$  on at least one the first three CCpNA axes, displayed on axis 1 versus 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval) correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. For visualization purposes, a focused section of the CCpNA triplot is shown. Abbreviations are described in Figure 6. Additional abbreviations: fungal, fungal infection; haploident, haploidentical donor; hemo, hemoglobin; leuko, leukocytes; lympho, lymphocytes; w1, week+1; w2, week+2; w3, week+3.



**Figure S6. Multivariate associations of the nasal microbiota with immune and clinical parameters in HSCT.** A) Clustered image map (CIM) based on sparse partial least squares (sPLS) regression analysis dimensions 1, 2, and 3, displaying pairwise correlations  $>0.2/ <-0.2$  between nasal ASVs (bottom), and continuous immune and clinical parameters (right). Red indicated positive correlation, and blue indicates negative correlation, respectively. Based on the sPLS regression model, hierarchical clustering (clustering method: complete linkage, distance method: Pearson's correlation) was performed resulting in the three depicted clusters. B) Canonical correspondence analysis (CCpNA) relating nasal microbial abundances (circles) to continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables with at least one correlation  $>0.2/ <-0.2$  in the sPLS analysis were included in the CCpNA. The triplot shows variables and ASVs with a score  $>0.3/ <-0.3$  on at least one the first three CCpNA axes, displayed on axis 1 versus 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval) correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. For visualization purposes, a focused section of the CCpNA triplot is shown. Abbreviations are described in Figures 6 and S5. Additional abbreviations: DonorMatch8, unrelated donor with 1 HLA mismatch; PB, peripheral blood.

### 1 Additional File 3 – Supplementary Discussion

2

3 We found that the patients for which we observed these associations were frequently treated with  
4 vancomycin, ciprofloxacin, and especially ceftazidime at time points early post-HSCT. This is in contrast  
5 to our previous study that revealed high *Lachnospiraceae* and *Ruminococcaceae* abundances and  
6 rapid B and NK cell reconstitution in the absence of vancomycin and ciprofloxacin treatment [1].  
7 Vancomycin and ciprofloxacin are among the broad-spectrum antibiotics that have previously been  
8 attributed a detrimental effect on the commensal microbiota, especially on *Clostridiales* [1–3]. The  
9 effect of ceftazidime however is controversially discussed. On the one hand, commensal sparing has  
10 previously been observed for cefepime, which belongs to the same antibiotic class as ceftazidime  
11 (cephalosporins) [4]. On the other hand, ceftazidime treatment could not rescue bacterial alpha  
12 diversity compared with e.g. vancomycin or ciprofloxacin treatment in a previous report [2]. In the  
13 present study, high abundances of clostridial microbiota members despite vancomycin and  
14 ciprofloxacin treatment might point to a potential beneficial effect of the additional treatment with  
15 ceftazidime. Moreover, future studies could assess if reconstitution in different NK cell subsets might  
16 differ with regards to time point and dependency on microbial composition and antimicrobial  
17 treatment.

18

19 We found that a specific *Parabacteroides distasonis* ASV (family *Tannerellaceae*) predicted moderate  
20 to severe aGvHD. In this context, it is interesting that our multivariate analyses indicated that a set of  
21 other *Parabacteroides* spp. ASVs was associated with high eosinophil counts from month +3 onwards.  
22 *Parabacteroides* members are providing the SCFA propionate which has previously been found to bind  
23 to the G 500 protein-coupled receptor GPR43 on eosinophils, although the exact mechanisms of  
24 GPR43-mediated immune modulation are yet to be investigated [5,6]. Interestingly, increased  
25 eosinophil numbers prior to and during aGvHD have been described before [7]. Therefore, an  
26 explanation for our prediction of moderate to severe aGvHD from high pre-HSCT *Parabacteroides* spp.  
27 abundances could involve a potential propionate-mediated activation of eosinophils, which might  
28 contribute to aGvHD. In contrast, propionate can also stimulate T<sub>reg</sub> cells, which in turn prevent T<sub>H</sub>17-  
29 induced inflammation involved in aGvHD [8,9]. We did not observe any correlation between  
30 *Parabacteroides* members and T cell counts in our study, however T cell subsets were measured after  
31 aGvHD and thus might be altered according to the inflammatory condition itself or its treatment.

32 In addition, a specific *Lachnospiraceae* ASV in the gut predicted subsequent moderate to severe  
33 aGvHD when highly abundant prior to HSCT. This family has previously been reported to be reduced  
34 prior to aGvHD (but post transplantation) [10,11]. However, here we examined pre-HSCT abundances  
35 to make predictions, i.e. we did not aim at elucidating *Lachnospiraceae* abundances concurrent to  
36 aGvHD. Particular members of this bacterial family are crucial producers of the SCFA butyrate, which  
37 has been suggested to prevent aGvHD directly by improving epithelial integrity in GvHD-target tissue,  
38 and/or indirectly by inducing T<sub>reg</sub> cells that reduce inflammation [9,12]. Of note, we did not observe  
39 any strong associations between bacterial abundances and inflammation (represented by CRP levels)  
40 at any body site in this study. It is of interest for future studies to reveal the exact mechanisms by  
41 which certain *Lachnospiraceae*, and specifically the predictive ASV we identified, might promote  
42 aGvHD, or whether their increase might be compensatory.

43

44 In addition to ASVs in the gut, we also identified oral and nasal cavity ASVs with pre-transplant  
45 abundances predicting subsequent aGvHD. For instance, in the oral cavity, an ASV affiliated with  
46 *Prevotella melaninogenica* strongly predicted aGvHD. This species is known to cause oral mucosal  
47 infections, and to be elevated in abundance in oral cancer [13,14]. Oral mucositis in allogeneic HSCT  
48 has been linked to elevated aGvHD risk [15,16]. It may be speculated that this could account for the  
49 observed association between aGvHD and high levels of this *Prevotella* ASV, although we did not  
50 investigate manifestations of mucositis in the present study. This highlights the need for direct  
51 evaluation of oral mucositis in this context.

52

53 Other microbial predictors of aGvHD in the oral and nasal cavity included members of the order  
54 *Actinomycetales* (*Actinomyces* spp., *Pseudopropionibacterium propionicum*, *Rothia* sp.). The order  
55 *Actinomycetales* comprises numerous commensals colonizing the oral and nasal cavities in healthy  
56 individuals, but can cause opportunistic infections in e.g. allo-HSCT patients [17]. It has been  
57 demonstrated that aGvHD increased the susceptibility to infections in allogeneic HSCT patients, but a  
58 predisposition to aGvHD due to a preceding infection has also been suggested for *Clostridium difficile*  
59 [18]. Whether actinomycosis might play a role in facilitating aGvHD development remains to be  
60 investigated. Interestingly, an increase of oral *Actinobacteria* in the gut at the time of neutrophil  
61 recovery has been found to be correlated with subsequent severe aGvHD [11]. One potential  
62 explanation involves the ability of *Actinomyces* spp. to drive biofilm formations, e.g. in periodontitis  
63 and during colorectal cancer [19].

64

65 Interestingly, among all ASVs significantly predicting aGvHD at any body site, we found only one ASV,  
66 namely a nasal *Rothia* sp. ASV, with high pre-HSCT abundances predicting that patients would be  
67 spared from subsequent aGvHD. *Rothia mucilaginosa* was among the species that, when present in  
68 the gut post-transplant, correlated with moderate to severe aGvHD in a previous study [11]. However,  
69 it has to date not been assessed in which way different *Rothia* species in the nasal cavity might be  
70 associated with aGvHD. This discrepancy again emphasizes the importance of high taxonomic  
71 resolution in this context, and the examination of the microbiota at different body sites.

72

73

74

## 75 **References**

76

- 77 1. Ingham AC, Kielsen K, Cilieborg MS, Lund O, Holmes S, Aarestrup FM, et al. Specific gut  
78 microbiome members are associated with distinct immune markers in pediatric allogeneic  
79 hematopoietic stem cell transplantation. *Microbiome* [Internet]. BioMed Central; 2019 [cited 2019  
80 Sep 18];7:131. Available from:  
81 <https://microbiomejournal.biomedcentral.com/articles/10.1186/s40168-019-0745-z>
- 82 2. Weber D, Hiergeist A, Weber M, Dettmer K, Wolff D, Hahn J, et al. Detrimental effect of broad-  
83 spectrum antibiotics on intestinal microbiome diversity in patients after allogeneic stem cell  
84 transplantation: Lack of commensal sparing antibiotics. *Clin Infect Dis* [Internet]. 2018 [cited 2018  
85 Sep 20]; Available from: <http://www.ncbi.nlm.nih.gov/pubmed/30124813>
- 86 3. Weber D, Jenq RR, Peled JU, Taur Y, Hiergeist A, Koestler J, et al. Microbiota Disruption Induced by  
87 Early Use of Broad Spectrum Antibiotics is an Independent Risk Factor of Outcome after Allogeneic  
88 Stem Cell Transplantation [Internet]. *Biol. Blood Marrow Transplant*. Elsevier Inc.; 2017. Available  
89 from: <http://linkinghub.elsevier.com/retrieve/pii/S1083879117302756>

- 90 4. Shono Y, Docampo MD, Peled JU, Perobelli SM, Velardi E, Tsai JJ, et al. Increased GVHD-related  
91 mortality with broad-spectrum antibiotic use after allogeneic hematopoietic stem cell  
92 transplantation in human patients and mice. *Sci Transl Med* [Internet]. 2016 [cited 2016 May  
93 23];8:339ra71-339ra71. Available from: <http://stm.sciencemag.org/content/8/339/339ra71>
- 94 5. Biagi E, Zama D, Nastasi C, Consolandi C, Fiori J, Rampelli S, et al. Gut microbiota trajectory in  
95 pediatric patients undergoing hematopoietic SCT. *Bone Marrow Transplant* [Internet]. Nature  
96 Publishing Group; 2015 [cited 2018 Jul 2];50:992–8. Available from:  
97 <http://www.nature.com/articles/bmt201516>
- 98 6. Maslowski KM, Vieira AT, Ng A, Kranich J, Sierro F, Yu D, et al. Regulation of inflammatory  
99 responses by gut microbiota and chemoattractant receptor GPR43. *Nature* [Internet]. NIH Public  
100 Access; 2009 [cited 2018 Oct 8];461:1282–6. Available from:  
101 <http://www.ncbi.nlm.nih.gov/pubmed/19865172>
- 102 7. Johnsson M, Cromvik J, Johansson J-E, Wennerås C, Vaht K. Eosinophils in the blood of  
103 hematopoietic stem cell transplanted patients are activated and have different molecular marker  
104 profiles in acute and chronic graft-versus-host disease. *Immunity, Inflamm Dis*. 2014;2:99–113.
- 105 8. Kverka M, Zakostelska Z, Klimesova K, Sokol D, Hudcovic T, Hrnčir T, et al. Oral administration of  
106 *Parabacteroides distasonis* antigens attenuates experimental murine colitis through modulation of  
107 immunity and microbiota composition. *Clin Exp Immunol* [Internet]. Wiley-Blackwell; 2011 [cited  
108 2018 Nov 19];163:250–9. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/21087444>
- 109 9. Arpaia N, Campbell C, Fan X, Dikiy S, Van Der Veeken J, Deroos P, et al. Metabolites produced by  
110 commensal bacteria promote peripheral regulatory T-cell generation. *Nature* [Internet]. Nature  
111 Publishing Group; 2013;504:451–5. Available from: <http://dx.doi.org/10.1038/nature12726>
- 112 10. Han L, Jin H, Zhou L, Zhang X, Fan Z, Dai M, et al. Intestinal Microbiota at Engraftment Influence  
113 Acute Graft-Versus-Host Disease via the Treg/Th17 Balance in Allo-HSCT Recipients. *Front Immunol*  
114 [Internet]. 2018 [cited 2018 May 17];9:669. Available from:  
115 <http://www.ncbi.nlm.nih.gov/pubmed/29740427>
- 116 11. Golob JL, Pergam SA, Srinivasan S, Fiedler TL, Liu C, Garcia K, et al. Stool Microbiota at Neutrophil  
117 Recovery Is Predictive for Severe Acute Graft vs Host Disease After Hematopoietic Cell  
118 Transplantation. *Clin Infect Dis* [Internet]. Oxford University Press; 2017 [cited 2018 Nov  
119 23];65:1984–91. Available from: <https://academic.oup.com/cid/article/65/12/1984/4085173>
- 120 12. Mathewson ND, Jenq R, Mathew A V, Koenigskecht M, Hanash A, Toubai T, et al. Gut  
121 microbiome-derived metabolites modulate intestinal epithelial cell damage and mitigate graft-  
122 versus-host disease. *Nat Immunol* [Internet]. NIH Public Access; 2016 [cited 2018 May 15];17:505–  
123 13. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/26998764>
- 124 13. Olczak-Kowalczyk D, Daszkiewicz M, Krasuska-Sławińska, Dembowska-Bagińska B, Gozdowski D,  
125 Daszkiewicz P, et al. Bacteria and *Candida* yeasts in inflammations of the oral mucosa in children  
126 with secondary immunodeficiency. *J Oral Pathol Med*. 2012;41:568–76.
- 127 14. Mager DL, Haffajee AD, Devlin PM, Norris CM, Posner MR, Goodson JM. The salivary microbiota  
128 as a diagnostic indicator of oral cancer: a descriptive, non-randomized study of cancer-free and oral  
129 squamous cell carcinoma subjects. *J Transl Med* [Internet]. BioMed Central; 2005 [cited 2018 Dec  
130 7];3:27. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/15987522>
- 131 15. Correa ME, Schubert MM, Storer B, Martin PJ, Flowers MED. Correlation Between Severity of  
132 Oral Mucositis and Both Overall and Gut Acute GvHD After High Intensity Conditioning  
133 Hematopoietic Stem Cell Transplantation. *Blood* [Internet]. 2012 [cited 2018 Nov 25];120. Available  
134 from: <http://www.bloodjournal.org/content/120/21/3057?sso-checked=true>
- 135 16. Chaudhry HM, Bruce AJ, Wolf RC, Litzow MR, Hogan WJ, Patnaik MS, et al. The Incidence and

- 136 Severity of Oral Mucositis among Allogeneic Hematopoietic Stem Cell Transplantation Patients: A  
137 Systematic Review. *Biol Blood Marrow Transplant* [Internet]. Elsevier; 2016 [cited 2018 Nov  
138 25];22:605–16. Available from:  
139 <https://www.sciencedirect.com/science/article/pii/S1083879115006394>
- 140 17. Barraco F, Labussière-Wallet H, Valour F, Ducastelle-Leprêtre S, Nicolini FE, Thomas X, et al.  
141 Actinomycosis after allogeneic hematopoietic stem cell transplantation despite penicillin  
142 prophylaxis. *Transpl Infect Dis*. 2016;18:595–600.
- 143 18. Chakrabarti S, Lees A, Jones S, Milligan D. Clostridium difficile infection in allogeneic stem cell  
144 transplant recipients is associated with severe graft-versus-host disease and non-relapse mortality.  
145 *Bone Marrow Transplant* [Internet]. Nature Publishing Group; 2000 [cited 2018 Nov 26];26:871–6.  
146 Available from: <http://www.nature.com/articles/1702627>
- 147 19. Flynn KJ, Baxter NT, Schloss PD. Metabolic and Community Synergy of Oral Bacteria in Colorectal  
148 Cancer. *mSphere* [Internet]. American Society for Microbiology (ASM); 2016 [cited 2018 Nov 26];1.  
149 Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27303740>
- 150

# Figures

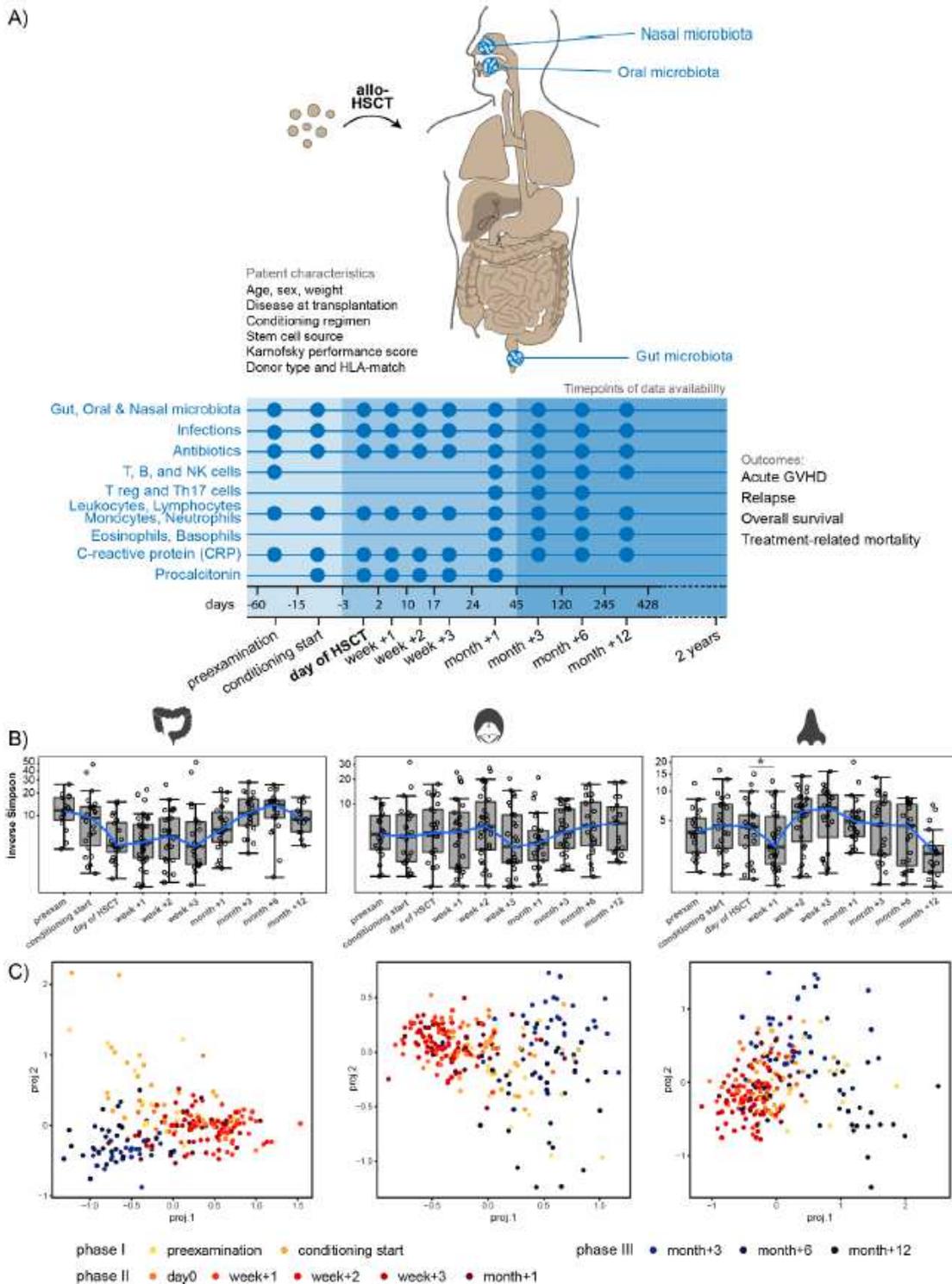


Figure 1

Monitoring gut, oral, and nasal microbiota and the host immune system in allogeneic hematopoietic stem cell transplantation (HSCT). A) Twenty-nine children were monitored before, at the time of, and immediately post allogeneic HSCT, as well as at late follow-up time points. Patients' baseline

characteristics, clinical outcomes, as well as immune cell counts, and inflammation and infection markers over time were monitored. Patient characteristics are described in detail in Table S1 (Additional File 1). Host immune system parameters were related to longitudinal dynamics of the gut, oral, and nasal microbiota that was assessed at the denoted time points. B) Bacterial alpha diversity before, at the time of, and after HSCT at each body site, displayed on a log<sub>10</sub> transformed y-axis for visualization purposes. Asterisks indicate significant differences in median inverse Simpson index between time points \* P < 0.05. C) Tree-based sparse linear discriminant (LDA) analyses by time point in relation to HSCT. For fecal samples, the positive LDA scores were observed for samples collected immediately post HSCT. For both oral and nasal samples, the positive LDA scores were observed for samples from before HSCT and from late follow up time points.

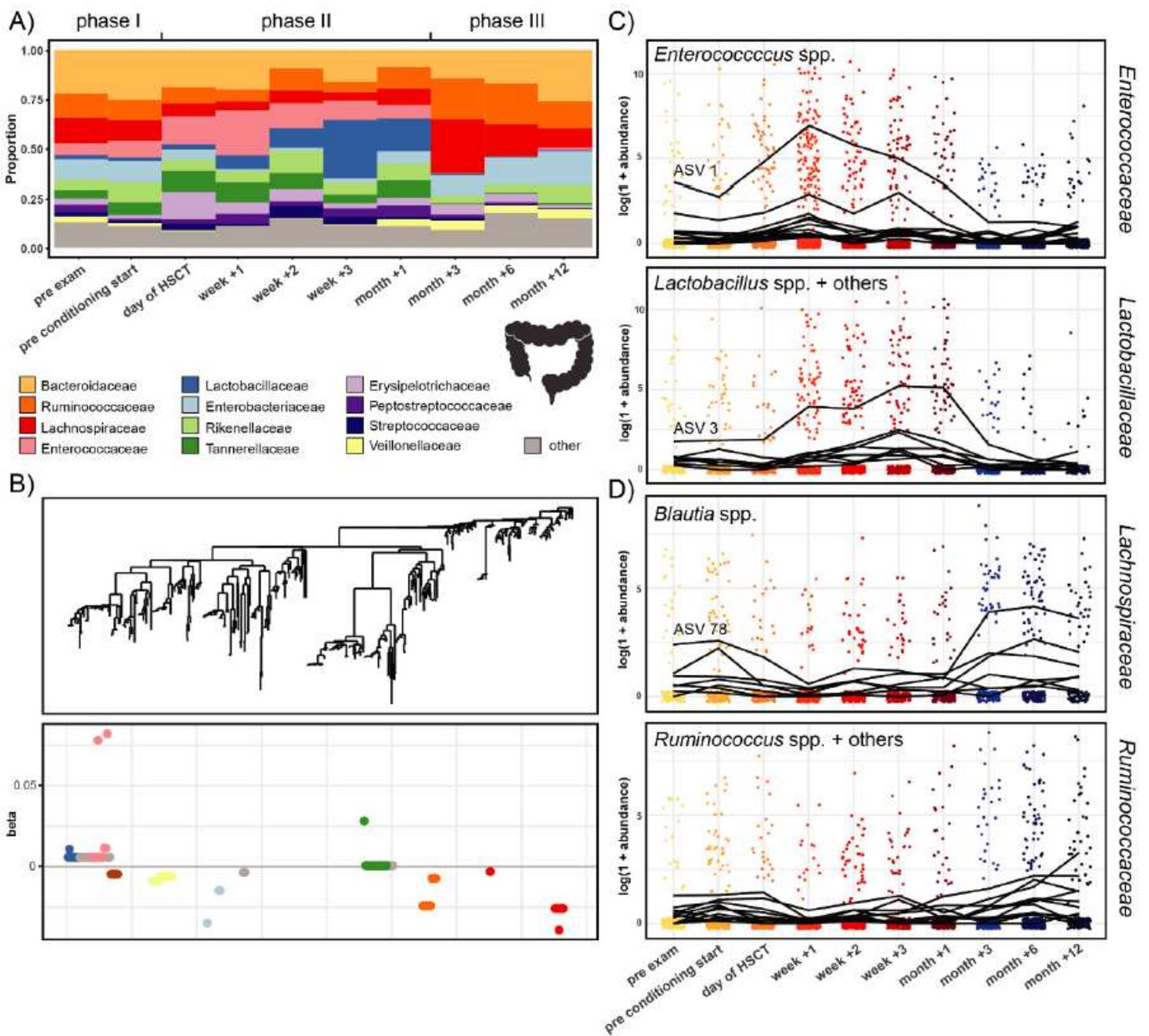


Figure 2

Temporal microbial community dynamics in the gut. A) Relative abundances over time of the 12 most abundant families in the gut. B) Tree-based sparse linear discriminant analysis (LDA). Coefficients of discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the phylogenetic tree. C) Trajectories of ASVs affiliated with the families Enterococcaceae and Lactobacillaceae, with increasing abundances after HSCT. The most abundant discriminating ASV for each family is indicated. D) Trajectories of ASVs affiliated with the families Lachnospiraceae and Ruminococcaceae, with decreasing abundances after HSCT and recovery at late follow-up time points. The most abundant discriminating ASV for *Blautia* spp. is indicated. Detailed taxonomic information and LDA-coefficients of the displayed ASVs are listed in Additional File 1: Table S2.

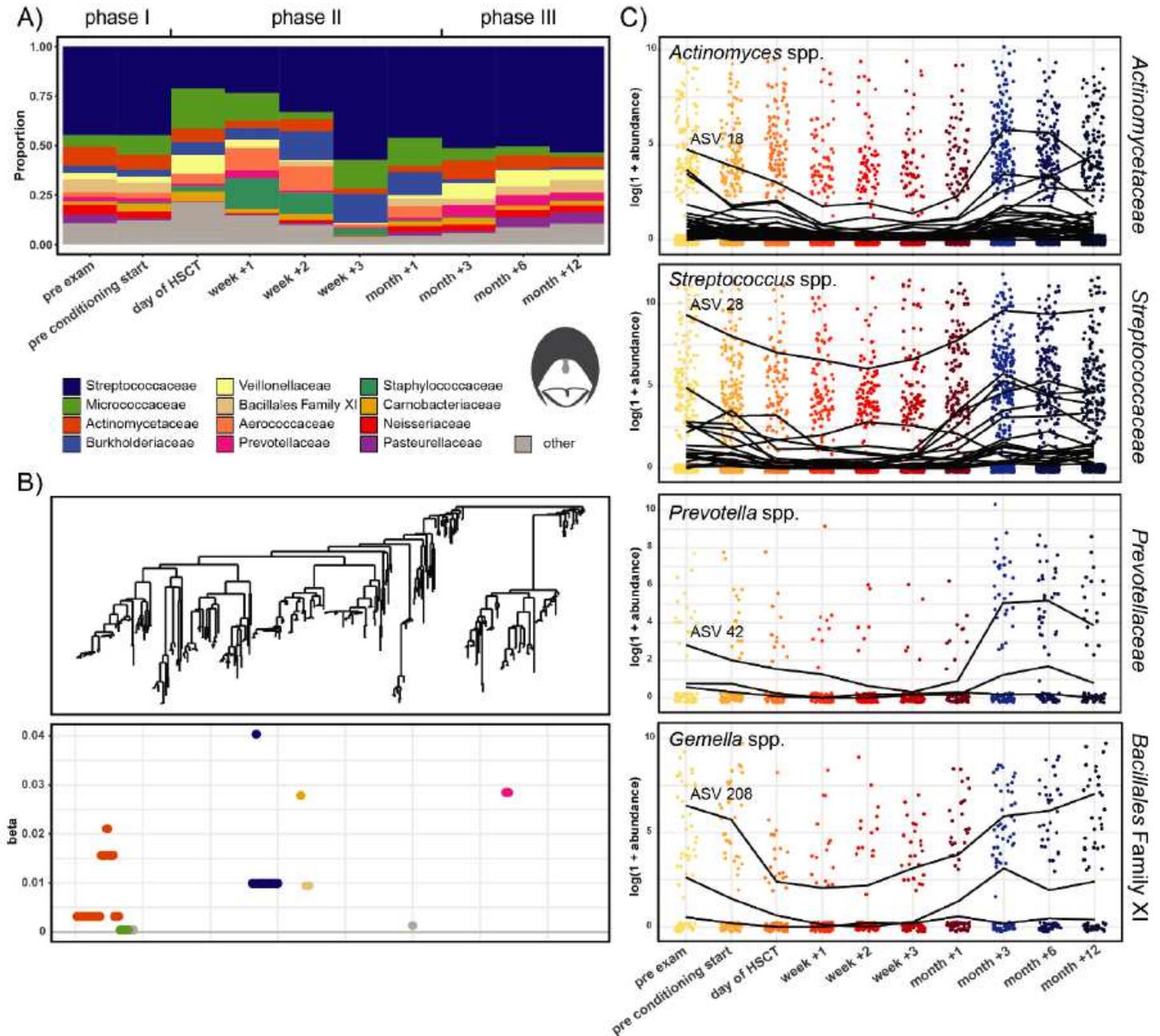
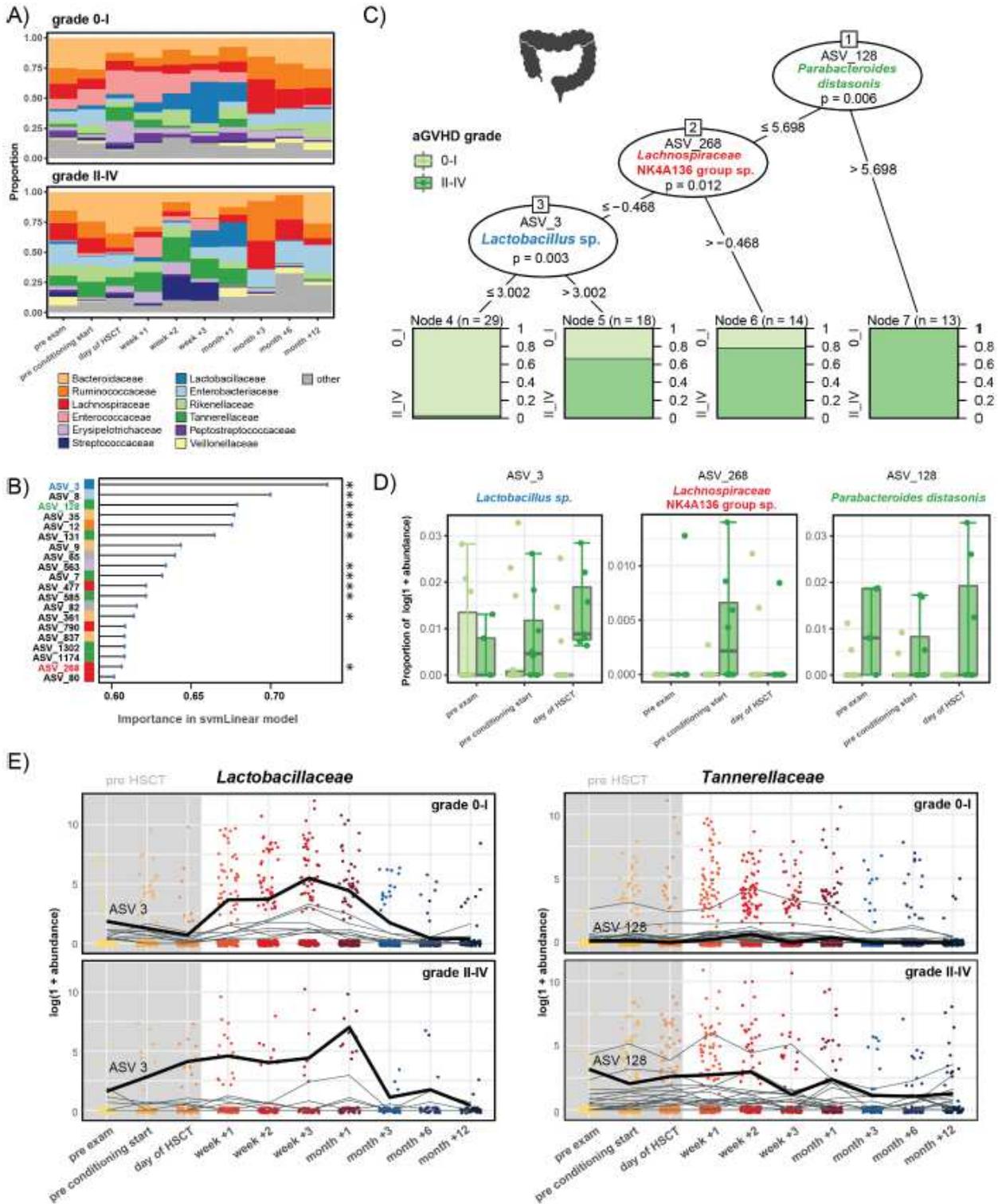


Figure 3

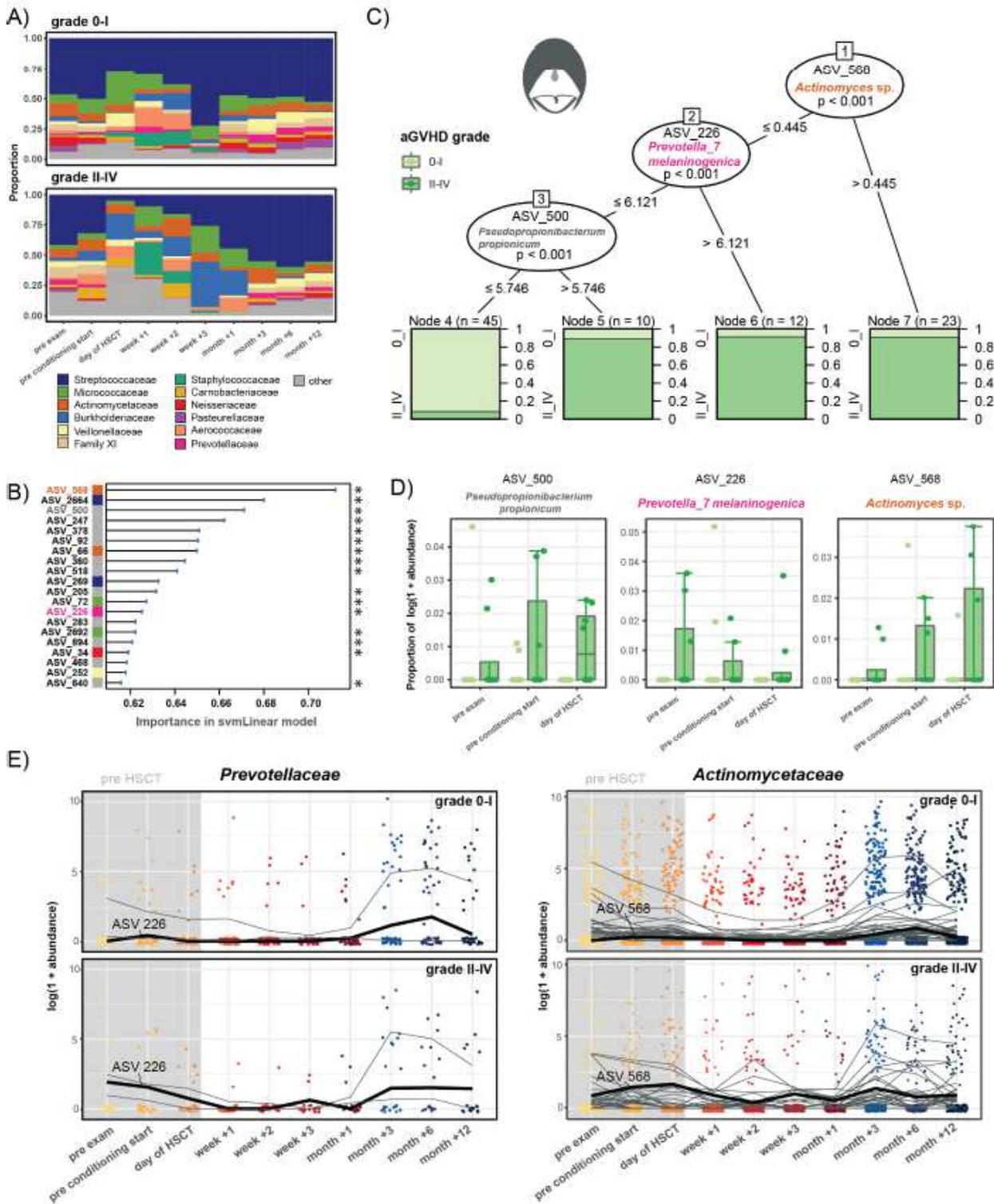
Temporal microbial community dynamics in the oral cavity. A) Relative abundances over time of the 12 most abundant families in the oral cavity. B) Tree-based sparse linear discriminant analysis (LDA). Coefficients of discriminating clades of ASVs on the first LDA axis, colored by taxonomic family, and plotted along the phylogenetic tree. C) Trajectories of ASVs affiliated with the families Actinomycetaceae, Streptococcaceae, Prevotellaceae, and Family XI (Class Bacillales), with decreasing abundances after HSCT and recovery at late follow-up time points. The most abundant discriminating ASV for each family is indicated. Detailed taxonomic information and LDA-coefficients of the displayed ASVs are listed in Additional File 1: Table S2.



**Figure 4**

Machine learning-based prediction of aGVHD severity from the pre-HSCT gut microbiota composition. A) Relative abundances of the 12 most abundant families over time in the gut in patients with aGVHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive gut ASVs identified by the svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted aGVHD (0-I versus II-IV)

from the abundances of gut ASVs pre- HSCT with 86% accuracy (95% CI: 65% to 97%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C) Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric regression for prediction of aGvHD. Numbers along the branches indicate split values of variance stabilized bacterial abundances. The terminal nodes show the proportion of samples originating from patients (n = number of samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depicting the log transformed relative abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared with grade II-IV patients. E) Trajectories of Lactobacillaceae and Tannerellaceae ASVs that were identified by tree-based sparse LDA, including ASV 3 and ASV 128 that were predictive for aGvHD (bold lines), in patients with aGvHD grade 0-I vs II-IV.



**Figure 5**

Machine learning-based prediction of aGvHD severity from the pre-HSCT oral microbiota composition. A) Relative abundances the 12 most abundant families over time in the oral cavity in patients with aGvHD grade 0-I versus II-IV. B) Importance plot of top 20 predictive oral ASVs identified by the svmLinear model with importance scores indicating the mean decrease in prediction accuracy in case the respective ASV would be excluded from the model. The final cross-validated svmLinear model predicted aGvHD (0-I

versus II-IV) from the abundances of oral ASVs pre-HSCT with 92% accuracy (95% CI: 73% to 99%). The ASVs that were also confirmed by Boruta feature selection are indicated with asterisk. C) Conditional inference tree (CTREE) displaying ASVs identified as significant split nodes by nonparametric regression for prediction of aGvHD. Numbers along the branches indicate split values of variance stabilized bacterial abundances. The terminal nodes show the proportion of samples originating from patients (n = number of represented samples) with aGvHD grade 0-I vs II-IV. D) Boxplots depict the log transformed relative abundances of the predictive ASVs at time points up to the transplantation in aGvHD grade 0-I compared with grade II-IV patients. E) Trajectories of Prevotellaceae and Actinomycetaceae ASVs that were identified by tree based sparse LDA, including ASV 226 and ASV 568 that were predictive for aGvHD (bold lines), in patients with aGvHD grade 0-I vs II-IV.



(clustering method: complete linkage, distance method: Pearson's correlation) was performed resulting in the three depicted clusters. B) Canonical correspondence analysis (CCpNA) relating gut microbial abundances (circles) to continuous (arrows) and categorical (+) immune and clinical parameters. ASVs and variables with at least one correlation  $>0.3$ / $<-0.3$  in the sPLS analysis were included in the CCpNA. The triplot shows variables and ASVs with a score  $>0.3$ / $<-0.3$  on at least one of the first three CCpNA axes, displayed on axis 1 versus 2 with samples depicted as triangles. The colored ellipses (depicted with 80% confidence interval) correspond to the clusters of ASVs identified by the sPLS-based hierarchical clustering. Abbreviations not mentioned in text: ATGmm, anti-thymocyte globulin; B\_, blood; BU, busulfan; CY, Cyclophosphamide; DonorMatch6, matched unrelated donor; FLU\_other, fludarabine combinations without thiotepa; GvHD.Prophylaxis1, treatment with cyclosporine; GvHD.Prophylaxis7, treatment with cyclosporine and methotrexate; immat\_B, immature B cells; K\_d100, Karnofsky score on day +100; K\_pre, Karnofsky score before HSCT; m1, month+1; m3, month+3; m6, month+6; m12, month+12; mat\_B, mature B cells; MEL, melphalan; total\_B, total B cells; P\_, plasma; parasitic, parasitic infection; pre\_cond, before conditioning start; pre\_exam, pre-examination; THIO, thiotepa; viral, viral infection; VP16, Etoposide.