# Microbial dynamics in newly diagnosed and treatment naïve IBD patients in the Mediterranean

**Authors:** Rausch, Dr. Philipp[1,2]*; Ellul, Dr. Sarah[3]*; Pisani, Dr. Anthea[4]; Bang, Dr. Corinna[1]; Tabone, Dr. Trevor[4]; Marantidis, Dr. Cordina Claire[5]; Zahra, Dr. Graziella[6]; Franke, Prof. Dr. Andre[1#] & Ellul, Dr. Pierre[4#]

[1] Institute of Clinical Molecular Biology, Christian-Albrechts-University of Kiel, Kiel, Germany

[2] Laboratory of Genomics and Molecular Biomedicine, Department of Biology, University of Copenhagen, Copenhagen, Denmark

[3] Division of Pediatric Surgery, Department of Surgery, Mater Dei Hospital, Malta,

[4] Division of Gastroenterology, Department of Medicine, Mater Dei Hospital, Malta

[5] Department of Microbiology, Mater Dei Hospital, Malta

[6] Molecular Diagnostics, Department of Pathology, Mater Dei Hospital, Malta

* equal contribution of PR and SE (joint first authors)

# corresponding authors: Pierre Ellul (pierre.ellul@gov.mt), Andre Franke (franke@mucosa.de)

## Supplemental materials and methods

**DNA Extraction:** DNA was extracted using the QIAamp DNA fast stool mini kit automated on the QIAcube (Qiagen, Hilden, Germany). Material is transferred to 0.70 mm Garnet Bead tubes (Dianova, Hamburg, Germany) filled with 1.1 ml ASL lysis buffer. Bead beating is performed using the a SpeedMill PLUS (Analytik Jena, Jena, Germany) for 45 s at 50 Hz. Samples are then heated to 95°C for 5 min with subsequent continuation of the manufacturer's protocol. DNA binds specifically to the QIAamp silica-gel membrane while contaminants pass through. PCR inhibitors are removed by the combined action of a unique adsorption resin and an optimized buffer. The approximate amount of DNA is between 10-40 µg per sample.

**Bacterial 16S rRNA Gene Sequencing:** Variable regions V1-V2 of the 16S rRNA gene are amplified using the primer pair 27F-338R in a dual-barcoding approach according to Caporaso *et al.* 2012 [1]. DNA is diluted 1:10 prior PCR, and 3 µl of this dilution are finally used for amplification. PCR-products are verified using the electrophoresis in agarose gel. PCR products are normalized using the SequalPrep Normalization Plate Kit (Thermo Fischer Scientific, Waltham, MA, USA), pooled in an equimolar fashion and sequenced on the Illumina MiSeq v3 2×300 bp (Illumina Inc., San Diego, CA, USA). Demultiplexing after sequencing was based on 0 mismatches in the barcode sequences via *bcl2fastq*.

**Quality control, classification, and binning of sequences:** Data processing was performed using the *DADA2* version 1.10 [2] workflow for big datasets (https://benjjneb.github.io/dada2/bigdata.html) resulting in abundance tables of **A**mplicon **S**equence **V**ariants (**ASV**s). All sequencing runs were handled separately for error correction, read merging, and combined chimera detection. ASVs underwent taxonomic annotation using the naïve Bayesian classifier implemented in DADA2 using the Ribosomal Database Project 16 release as a taxonomic classification database [3, 4]. ASV sequences were aligned via NAST-alignment to the SILVA core database and filtered for informative sites (constant gaps, constant bases) in *mothur* [5-8]. Phylogenetic tree construction on ASV alignment generated was carried out using *FastTree 2.1* using the CAT substitution model with Γ-correction and improved accuracy, employing more minimum evolution rounds for initial tree search [-spr 4], more exhaustive tree search [-mlacc 2], and a slower initial tree search [-slownni] [9].

**Statistical methods for microbiome analyses:** Microbiome data was rarefied to 11800 reads/sample to ensure comparable and sufficient coverage across samples during analysis (average Good's coverage 99.88% ± 0.001 SD).
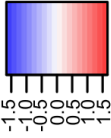
**Alpha diversity:** Phylogenetic measures of alpha diversity [Nearest Taxon Index (NTI), and Net Relatedness Index (NRI)] were derived using the *picante* package, based on 999 permutations against a null model preserving relative species richness within the communities [10, 11] [NRI=-1 × ($MPD_{observed}$-mean($MPD_{random}$))/SD($MPD_{random}$); NTI=-1 × ($MNTD_{observed}$-mean($MNTD_{random}$))/SD($MNTD_{random}$)]. Both metrics are phylogenetic effect sizes, for which positive values indicate phylogenetic clustering, values close to zero indicate neutral or random community assembly, and negative values indicate phylogenetic overdispersion, either over the whole phylogenetic tree (NRI) or across the closest related species/tips of the phylogenetic tree (NTI).

**Network analyses:** To generate co-abundance networks we employed the *SparCC* algorithm as implemented in *mothur* (50 samplings, 100 iterations, 10,000 permutations, *P*-value adjustment via

Benjamini-Hochberg procedure) was used to generate co-abundance networks based on ASV abundances within the study cohorts (shared among 10% samples) at a $P$-value cutoff of $P_{FDR} \leq 0.05$ and a correlation strength above $|R| \geq 0.25$ [12].

**Supplemental figures**

average Z-Value

-1.5 -1.0 -0.5 0.0 0.5 1.0 1.5

■ Proteobacteria    ■ Bacteroidetes    □ Firmicutes

A1    A2    A3

ASV 13 | G−Escherichia/Shigella uncl. | A1
ASV 16 | S−Alistipes shahii | A1
ASV 30 | F−Enterobacteriaceae uncl. | A1
ASV 31 | G−Escherichia/Shigella uncl. | A1
ASV 39 | G−Bacteroides uncl. | A1
ASV 42 | G−Bacteroides uncl. | A1
ASV 49 | G−Ruminococcus uncl. | A1
ASV 64 | G−Oscillibacter uncl. | A1
ASV 71 | S−Bacteroides cellulosilyticus | A1
ASV 77 | S−Haemophilus parainfluenzae | A1
ASV 93 | G−Phascolarctobacterium uncl. | A1
ASV 114 | S−Bacteroides fragilis | A1
ASV 118 | G−Escherichia/Shigella uncl. | A1
ASV 119 | S−Bacteroides cellulosilyticus | A1
ASV 124 | S−Roseburia intestinalis | A1
ASV 139 | F−Ruminococcaceae uncl. | A1
ASV 147 | S−Escherichia/Shigella coli | A1
ASV 157 | S−Bacteroides clarus | A1
ASV 166 | G−Alistipes uncl. | A1
ASV 183 | S−Escherichia/Shigella coli | A1
ASV 194 | F−Ruminococcaceae uncl. | A1
ASV 234 | G−Ruminococcus uncl. | A1
ASV 249 | G−Parasutterella uncl. | A1
ASV 253 | G−Oscillibacter uncl. | A1
ASV 258 | F−Desulfovibrionaceae uncl. | A1
ASV 394 | G−Oscillibacter uncl. | A1
ASV 595 | P−Firmicutes uncl. | A1
ASV 623 | S−Escherichia/Shigella coli | A1
ASV 647 | F−Ruminococcaceae uncl. | A1
ASV 808 | G−Streptococcus uncl. | A1
ASV 839 | G−Faecalibacterium uncl. | A1
ASV 893 | S−Morganella morganii | A1
ASV 18 | G−Bacteroides uncl. | A2
ASV 19 | F−Lachnospiraceae uncl. | A2
ASV 27 | G−Bacteroides uncl. | A2
ASV 40 | G−Paraprevotella uncl. | A2
ASV 79 | F−Lachnospiraceae uncl. | A2
ASV 117 | G−Alistipes uncl. | A2
ASV 176 | G−Dorea uncl. | A2
ASV 262 | G−Desulfovibrio uncl. | A2
ASV 310 | G−Bacteroides uncl. | A2
ASV 315 | F−Lachnospiraceae uncl. | A2
ASV 429 | G−Desulfovibrio uncl. | A2
ASV 662 | G−Faecalibacterium uncl. | A2
ASV 666 | G−Haemophilus uncl. | A2
ASV 768 | G−Faecalibacterium uncl. | A2
ASV 789 | G−Coprococcus uncl. | A2
ASV 850 | G−Faecalibacterium uncl. | A2
ASV 875 | G−Faecalibacterium uncl. | A2
ASV 903 | G−Coprococcus uncl. | A2
ASV 1113 | G−Bacteroides uncl. | A2
ASV 1446 | G−Clostridium XIVa uncl. | A2
ASV 1961 | G−Faecalibacterium uncl. | A2
ASV 2104 | F−Lachnospiraceae uncl. | A2
ASV 4 | G−Acidaminococcus uncl. | A3
ASV 5 | S−Bacteroides dorei | A3
ASV 8 | G−Dialister uncl. | A3
ASV 10 | S−Bacteroides uniformis | A3
ASV 14 | G−Alistipes uncl. | A3
ASV 50 | G−Bacteroides uncl. | A3
ASV 54 | G−Bacteroides uncl. | A3
ASV 83 | G−Bacteroides uncl. | A3
ASV 84 | G−Clostridium sensu stricto uncl. | A3
ASV 148 | G−Odoribacter uncl. | A3
ASV 178 | G−Holdemanella uncl. | A3
ASV 180 | G−Phascolarctobacterium uncl. | A3
ASV 192 | G−Bacteroides uncl. | A3
ASV 298 | G−Parabacteroides uncl. | A3
ASV 340 | G−Parabacteroides uncl. | A3
ASV 485 | G−Alistipes uncl. | A3
ASV 490 | G−Bacteroides uncl. | A3
ASV 569 | G−Bacteroides uncl. | A3
ASV 1016 | G−Faecalibacterium uncl. | A3
ASV 1411 | S−Morganella morganii | A3
ASV 1704 | G−Faecalibacterium uncl. | A3
ASV 2911 | G−Morganella uncl. | A3
ASV 51 | G−Barnesiella uncl. | A2
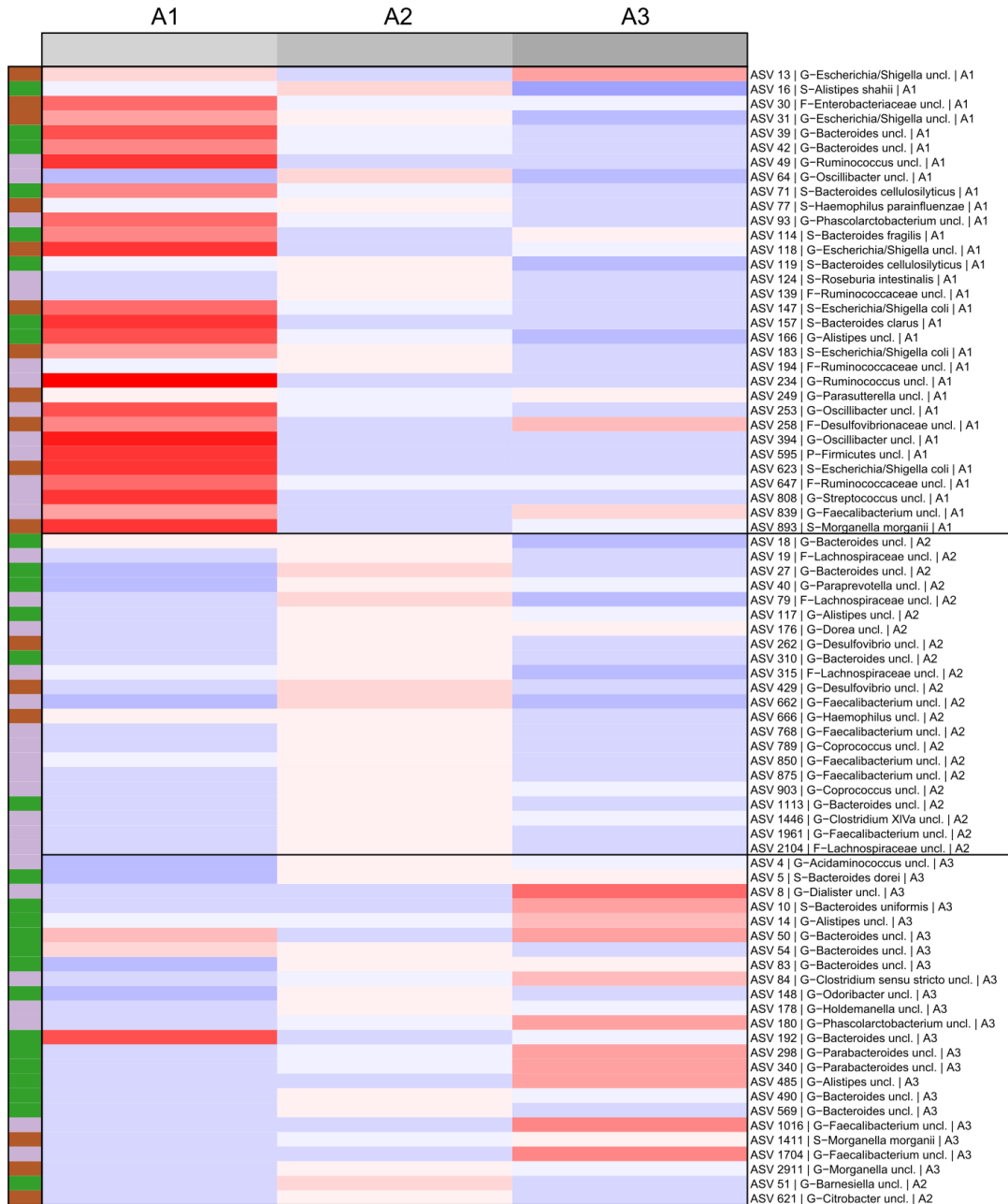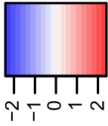ASV 621 | G−Citrobacter uncl. | A2

**Figure S1:** Heatmap visualizing significant differentially abundant ASVs in CD patients with respect to age subgroups following the Montreal classification (A1: <16 yrs., A2: 17-40 yrs., A3: >40 yrs., Table S3). Differential abundance was tested via *DESeq2* and only functions significant after p-value adjustment are displayed.
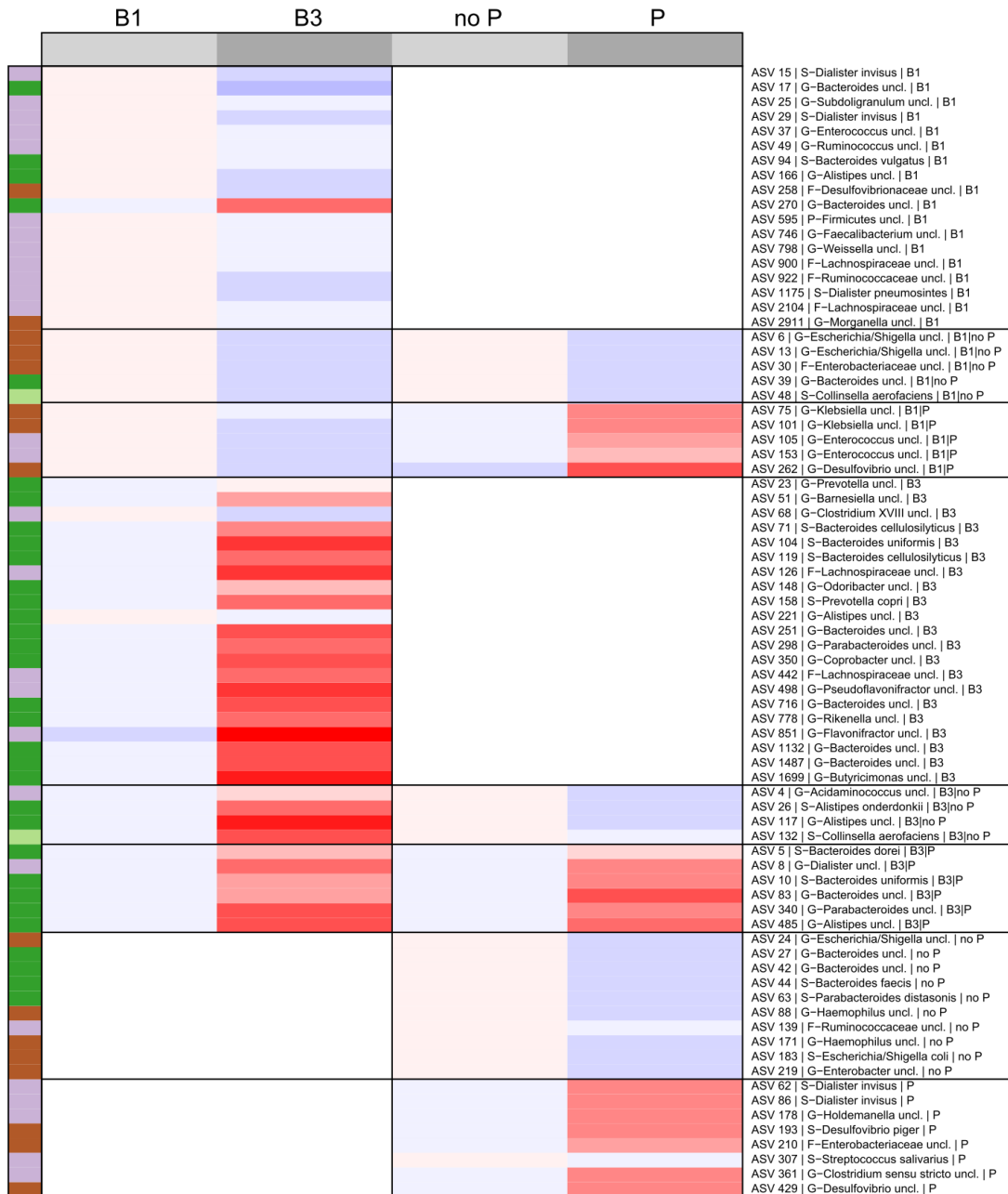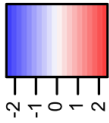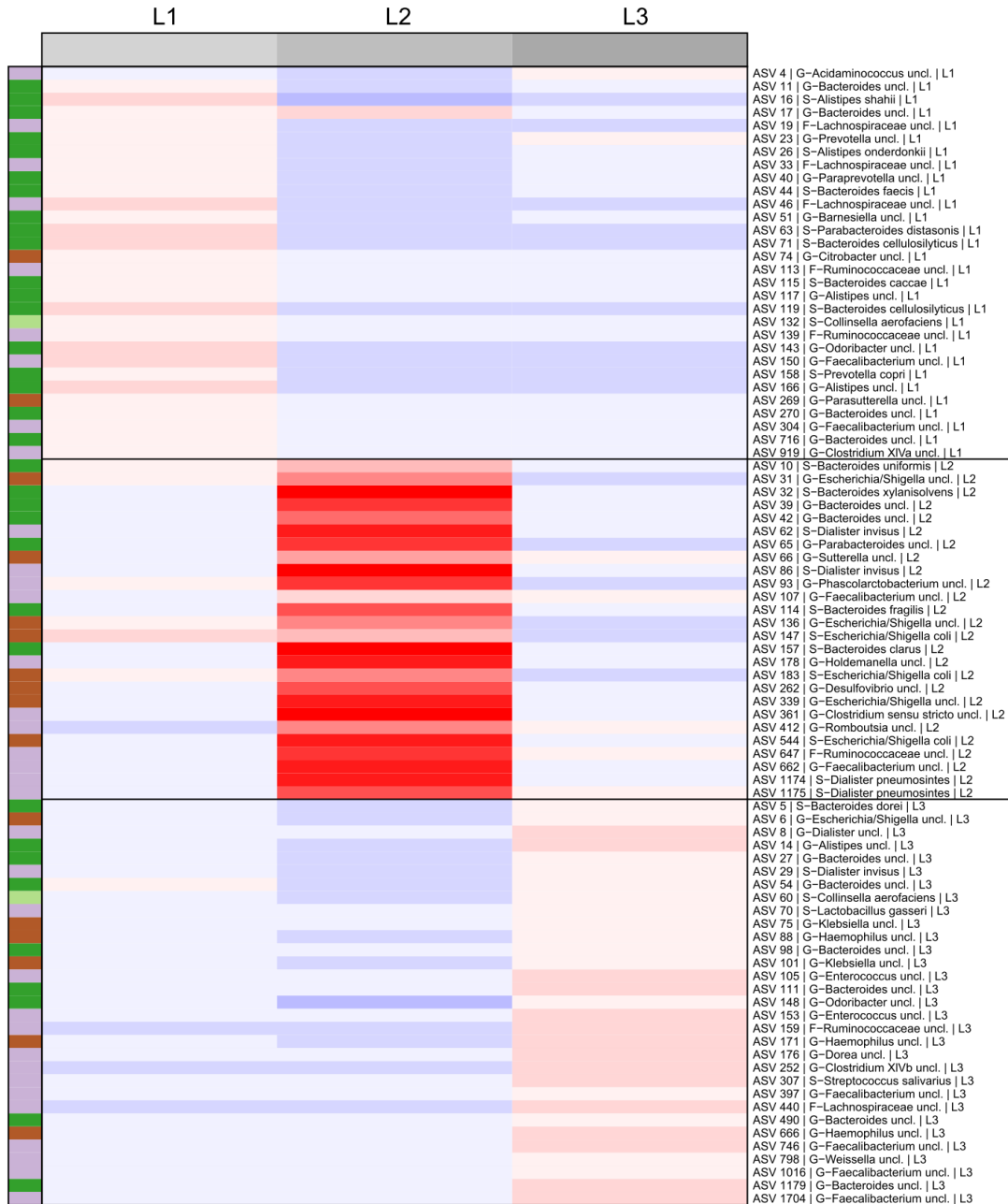
**Figure S2:** Heatmap visualizing significant differentially abundant ASVs in CD patients with respect to disease behavior subgroups following the Montreal classification (B1: non-stricturing/non-penetrating, B3: penetrating, Table S3), while using the classification of perianal disease as a separate subtype (P:

perianal disease manifestation, no-P: no perianal disease). Differential abundance was tested via DESeq2 and only functions significant after p-value adjustment are displayed.
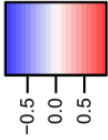
**Figure S3:** Heatmap visualizing significant differentially abundant ASVs in CD patients with respect to disease location subgroups following the Montreal classification (L1: ileal, L2: colonic, L3: ileocolonic, Table S3). Differential abundance was tested via DESeq2 and only functions significant after p-value adjustment are displayed.
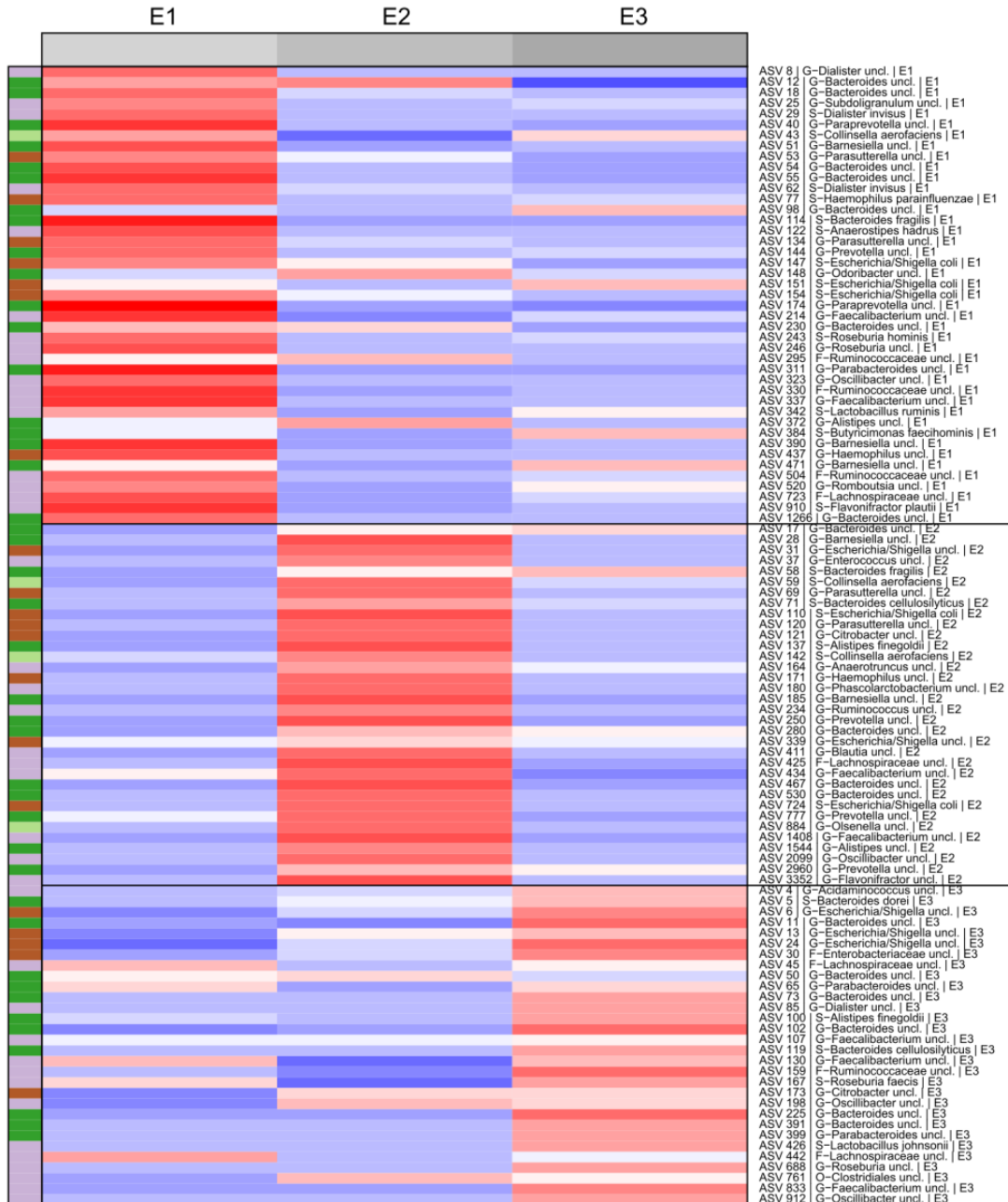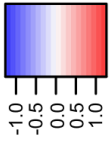
**Figure S4:** Heatmap visualizing significant differentially abundant ASVs in UC patients with respect to age subgroups following the Montreal classification (E1: Ulcerative proctitis, E2: left sided UC (distal UC), E3: extensive UC (pancolitis), Table S3). Differential abundance was tested via DESeq2 and only functions significant after p-value adjustment are displayed.

**Figure S5:** Heatmap visualizing significant differentially abundant ASVs in CD and UC patients with respect to prior antibiotic treatments (Table S5). Differential abundance was tested via DESeq2 and only functions significant after p-value adjustment are displayed.

**Figure S6:** Heatmap visualizing significant differentially abundant ASVs in CD and UC patients with respect to prior antibiotic treatments (Table S5). Differential abundance was tested via DESeq2 and only functions significant after p-value adjustment are displayed.

## Crohns disease

## Ulcerative colitis



**A** P=0.05165

**B** P=0.03937

**C** P=0.01244

**D** P=0.06129

**E** P=0.06541

**F** P=0.80043

**G** P=0.02285

**H** P=0.60042

**Figure S7:** Prediction of treatment probabilities (future treatment with mesalazine, azathioprine, anti-TNF, prednisolone) based on the underlying alpha diversity in CD and UC patients. As single treatment, only corticosteroid treatments are predictable from Simpson diversity (**A**: DF=1,29, Dev.=3.7872, $P$=0.05165, $P_{FDR}$=0.24177; **B**: DF=1,23, Dev.=4.2448, $P$=0.03937, $P_{FDR}$=0.61291), Chao1 species richness (**C**: DF=1,29, Dev.=6.2466, $P$=0.01244, $P_{FDR}$=0.20830; **D**: DF=1,23, Dev.=3.5021, $P$=0.06129, $P_{FDR}$=0.61291), NRI (**E**: DF=1,29, Dev.=3.3945, $P$=0.06541, $P_{FDR}$=0.24177; **F**: DF=1,23, Dev.=0.06391, $P$=0.8004, $P_{FDR}$=1.0000), and NTI (**G**: DF=1,29, Dev.=5.1799, $P$=0.02285, $P_{FDR}$=0.20830; **H**: DF=1,23, Dev.=0.27436, $P$=0.6004, $P_{FDR}$=1.0000). Ticks on the bottom of the graph indicate samples without treatment, ticks on the top indicate samples with future treatment.



**Figure S8:** Principle coordinate analyses of Jaccard distance visualizing the differences between the disease cohorts and including nominal significant factors and continuous variables (A), as well as separate betadiversity analyses of only CD- (B) and UC (C) patients. (D) Boxplots visualize the community

variability within each health group measured as distance to the centroid, which is lowest in healthy individuals.

**Figure S9:** Heatmap visualizing significant differentially abundant functions between healthy controls, CD and UC patients (Table S10). Functions are represented as the PiCRUST2 imputed abundance of single enzymes categories/EC categories (Enzyme Commission number). Differential abundance was tested via DESeq2 and only significantly different functions, after FDR adjustment, are displayed.



**Figure S10:** Global network descriptors for each health condition specific network (including subsamples of the control based networks), like **(A)** average transitivity, which describes the clustering and denseness of the respective network. (**B**) Degree assortativity describes preferential attachment of bacteria with similarly connected bacteria, while average degree and betweenness describe the **(C)**

average number of connections of nodes in the network and **(D)** the average of shortest paths in the networks.

**Figure S11:** Centrality measures of the 50 most important nodes of the respective importance measures (Betweenness, Eigenvalue centrality, PageRank). Red letters highlight network members with a higher importance than expected by chance, based on a Z-test against 10'000 randomized networks ($P \leq 0.05$). Highlighted are the associations of the respective ASVs as detected by differential abundance analysis as well ($P_{FDR} \leq 0.05$, Table S12).

## Degree

ASV 6-Escherichia/Shigella uncl. ASV 60-Collinsella aerofaciens ASV 351-Faecalibacterium uncl. ASV 903-Coprococcus uncl.
ASV 13-Escherichia/Shigella uncl. ASV 61-Alistipes uncl. ASV 354-Blautia uncl. ASV 906-Lachnospiraceae uncl.
ASV 22-Parasutterella uncl. ASV 95-Ruminococcaceae uncl. ASV 411-Blautia uncl. ASV 1069-Lachnospiraceae uncl.
ASV 30-Enterobacteriaceae uncl. ASV 99-Alistipes indistinctus ASV 696-Blautia uncl. ASV 1145-Clostridium XIVa uncl.
ASV 34-Parabacteroides uncl. ASV 167-Roseburia faecis ASV 762-Lachnospiraceae uncl. ASV 1331-Intestinimonas uncl.
ASV 43-Collinsella aerofaciens ASV 240-Blautia uncl. ASV 797-Coprococcus uncl. ASV 1461-Fusicatenibacter uncl.
ASV 47-Prevotella uncl. ASV 283-Dorea uncl. ASV 799-Butyricimonas uncl. ASV 2155-Butyricicoccus uncl.
ASV 48-Collinsella aerofaciens ASV 311-Parabacteroides uncl. ASV 846-Eggerthella lenta ASV 2604-Lachnospiraceae uncl.
ASV 4971-Candidatus Saccharibacteria uncl.
ASV 5581-Parvimonas micra

**UC**

ASV 368- Blautia uncl.
ASV 648- Eggerthella lenta

34

ASV 150- Faecalibacterium uncl.
ASV 387- Romboutsia uncl.

2  2

0

2

ASV 197- Blautia uncl.
ASV 412- Romboutsia uncl.

5  25

**Contr.**  **CD**

ASV 11-Bacteroides uncl.
ASV 87 -Flavonifractor plautii
ASV 206-Blautia uncl.
ASV 321-Clostridiales uncl.
ASV 855-Blautia uncl.

ASV_15-Dialister invisus   ASV_599-Veillonella uncl.
ASV_16-Alistipes shahii   ASV_615-Romboutsia uncl.
ASV_35-Parabacteroides distasonis   ASV_620-Butyricicoccus uncl.
ASV_64-Oscillibacter uncl.   ASV_643-Romboutsia uncl.
ASV_92-Romboutsia uncl.   ASV_650-Romboutsia uncl.
ASV_130-Faecalibacterium uncl.   ASV_713-Coprobacter uncl.
ASV_198-Oscillibacter uncl.   ASV_839-Faecalibacterium uncl.
ASV_249-Parasutterella uncl.   ASV_953-Ruminococcus2 uncl.
ASV_304-Faecalibacterium uncl.   ASV_1039-Romboutsia uncl.
ASV_339-Escherichia/Shigella uncl.   ASV_1096-Romboutsia uncl.
ASV_442-Lachnospiraceae uncl.   ASV_1152-Allisonella uncl.
ASV_478-Romboutsia uncl.   ASV_1347-Ruminococcaceae uncl.
ASV_2081-Bifidobacterium longum

## Betweenness

ASV 19-Lachnospiraceae uncl.   ASV 648-Eggerthella lenta
ASV 22-Parasutterella uncl.   ASV 657-Clostridium XIVa uncl.
ASV 43-Collinsella aerofaciens   ASV 774-Ruminococcus uncl.
ASV 48-Collinsella aerofaciens   ASV 797-Coprococcus uncl.
ASV 60-Collinsella aerofaciens   ASV 878-Blautia obeum
ASV 61-Alistipes uncl.   ASV 906-Lachnospiraceae uncl.
ASV 91-Faecalibacterium uncl.   ASV 1034-Flavonifractor uncl.
ASV 167-Roseburia faecis   ASV 1304-Blautia uncl.
ASV 176-Dorea uncl.   ASV 1331-Intestinimonas uncl.
ASV 227-Anaerostipes uncl.   ASV 1332-Pseudoflavonifractor uncl.
ASV 313-Faecalibacterium uncl.   ASV 1533-Butyricimonas uncl.
ASV 351-Faecalibacterium uncl.   ASV 2962-Pseudoflavonifractor uncl.
ASV 411-Blautia uncl.   ASV 4971-Candidatus Saccharibacteria uncl.

**UC**

ASV 240-Blautia uncl.
ASV 368-Blautia uncl.

26

ASV_7-Bacteroides massiliensis

2  1

0

ASV_87-Flavonifractor plautii
ASV_442-Lachnospiraceae uncl.

6  18

2

**Contr.**  **CD**

ASV_95-Ruminococcaceae uncl.
ASV_321-Clostridiales uncl. uncl.
ASV_412-Romboutsia uncl.
ASV_420-Coprococcus uncl.
ASV_855-Blautia uncl.
ASV_1480-Clostridium IV uncl.

ASV_35-Parabacteroides distasonis   ASV_387-Romboutsia uncl.
ASV_64-Oscillibacter uncl.   ASV_487-Firmicutes uncl.
ASV_150-Faecalibacterium uncl.   ASV_620-Butyricicoccus uncl.
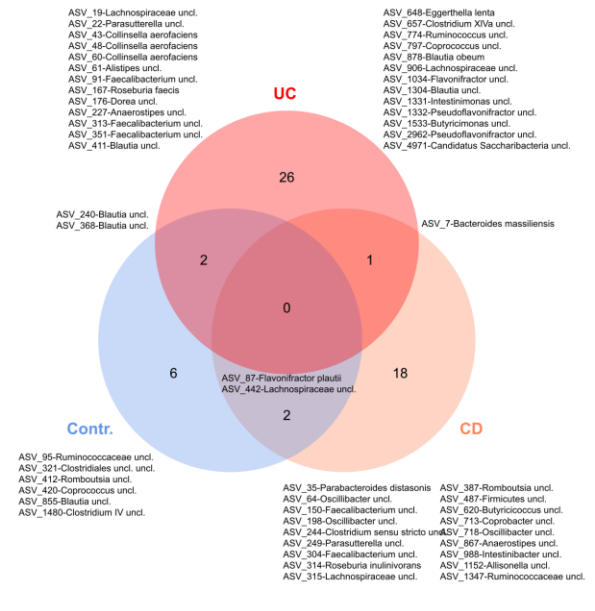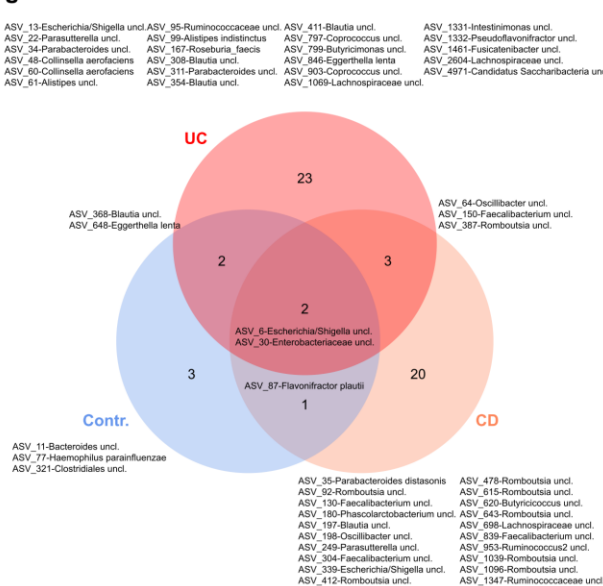ASV_198-Oscillibacter uncl.   ASV_713-Coprobacter uncl.
ASV_244-Clostridium sensu stricto uncl.   ASV_718-Oscillibacter uncl.
ASV_249-Parasutterella uncl.   ASV_867-Anaerostipes uncl.
ASV_304-Faecalibacterium uncl.   ASV_988-Intestinibacter uncl.
ASV_314-Roseburia inulinivorans   ASV_1152-Allisonella uncl.
ASV_315-Lachnospiraceae uncl.   ASV_1347-Ruminococcaceae uncl.

## PageRank

ASV_13-Escherichia/Shigella uncl.   ASV_95-Ruminococcaceae uncl.   ASV_411-Blautia uncl.   ASV_1331-Intestinimonas uncl.
ASV_22-Parasutterella uncl.   ASV_99-Alistipes indistinctus   ASV_799-Butyricimonas uncl.   ASV_1332-Pseudoflavonifractor uncl.
ASV_34-Parabacteroides uncl.   ASV_167-Roseburia_faecis   ASV_846-Eggerthella lenta   ASV_2604-Lachnospiraceae uncl.
ASV_48-Collinsella aerofaciens   ASV_308-Blautia uncl.   ASV_903-Coprococcus uncl.   ASV_4971-Candidatus Saccharibacteria uncl.
ASV_60-Collinsella aerofaciens   ASV_311-Parabacteroides uncl.   ASV_1069-Lachnospiraceae uncl.
ASV_61-Alistipes uncl.   ASV_354-Blautia uncl.

**UC**

ASV_368-Blautia uncl.
ASV_648-Eggerthella lenta

23

ASV_64-Oscillibacter uncl.
ASV_150-Faecalibacterium uncl.
ASV_387-Romboutsia uncl.

2  3

2

ASV_6-Escherichia/Shigella uncl.
ASV_30-Enterobacteriaceae uncl.

3  20

ASV_87-Flavonifractor plautii

1

**Contr.**  **CD**

ASV_11-Bacteroides uncl.
ASV_77-Haemophilus parainfluenzae
ASV_321-Clostridiales uncl.

ASV_35-Parabacteroides distasonis   ASV_478-Romboutsia uncl.
ASV_92-Romboutsia uncl.   ASV_615-Romboutsia uncl.
ASV_130-Faecalibacterium uncl.   ASV_620-Butyricicoccus uncl.
ASV_180-Phascolarctobacterium uncl.   ASV_643-Romboutsia uncl.
ASV_197-Blautia uncl.   ASV_698-Lachnospiraceae uncl.
ASV_198-Oscillibacter uncl.   ASV_839-Faecalibacterium uncl.
ASV_249-Parasutterella uncl.   ASV_953-Ruminococcus2 uncl.
ASV_304-Faecalibacterium uncl.   ASV_1039-Romboutsia uncl.
ASV_339-Escherichia/Shigella uncl.   ASV_1096-Romboutsia uncl.
ASV_412-Romboutsia uncl.   ASV_1347-Ruminococcaceae uncl.

## Eigencentrality

ASV_354-Blautia uncl.   ASV_884-Olsenella uncl.
ASV_368-Blautia uncl.   ASV_903-Coprococcus uncl.
ASV_411-Blautia uncl.   ASV_906-Lachnospiraceae uncl.
ASV_696-Blautia uncl.   ASV_929-Faecalibacterium uncl.
ASV_762-Lachnospiraceae uncl.   ASV_977-Romboutsia uncl.
ASV_789-Coprococcus uncl.   ASV_1069-Lachnospiraceae uncl.
ASV_846-Eggerthella_lenta   ASV_1145-Clostridium_XIVa uncl.
ASV_2604-Lachnospiraceae uncl.

**UC**

15

ASV_478-Romboutsia uncl.

0  1

1

ASV_387-Romboutsia uncl.

9  9

2

**Contr.**  **CD**

ASV_38-Haemophilus parainfluenzae
ASV_77-Haemophilus parainfluenzae
ASV_87-Flavonifractor plautii
ASV_88-Haemophilus uncl.
ASV_189-Haemophilus uncl.
ASV_197-Blautia uncl.
ASV_206-Blautia uncl.
ASV_220-Firmicutes uncl.
ASV_648-Eggerthella lenta

ASV_92-Romboutsia uncl.
ASV_412-Romboutsia uncl.

ASV_15-Dialister_invisus
ASV_339-Escherichia/Shigella uncl.
ASV_599-Veillonella uncl.
ASV_606-Blautia uncl.
ASV_615-Romboutsia uncl.
ASV_643-Romboutsia uncl.
ASV_650-Romboutsia uncl.
ASV_1039-Romboutsia uncl.
ASV_1096-Romboutsia uncl.

**Figure S12:** Overlapping and private ASVs with significant network positions among health condition specific networks. Centralities in the control cohort were based on the average of three independent

network calculations on three subsamples of the respective cohort to balance differences in cohort sizes (Table S12).

**Legend (top):**
- sign. important nodes in contr.
- sign. important nodes in UC
- sign. important nodes in CD
- insign. important nodes

**A** — Significant important nodes: Contr.↔CD / Contr.↔UC

**Node degree**

Cohort: CD / Control / UC

Contr. → CD/UC
Contr. + → CD/UC

Rank degree (0, 100, 200, 300, 400, 500)

**B** — Significant important nodes: UC↔CD

UC + → CD
UC – → CD

Rank degree (0, 100, 200, 300, 400, 500)

**C** — Node betweenness

Cohort: CD / Control / UC

Contr. → CD/UC
Contr. + → CD/UC

Rank betweenness (0, 100, 200, 300, 400, 500)

**D**

UC + → CD
UC – → CD

Rank betweenness (0, 100, 200, 300, 400, 500)

**E** — Node PageRank

Cohort: CD / Control / UC

Contr. → CD/UC
Contr. + → CD/UC

Rank PageRank (0, 100, 200, 300, 400, 500)

**F**

UC + → CD
UC – → CD

Rank PageRank (0, 100, 200, 300, 400, 500)

**G** — Node eigencentrality

Cohort: CD / Control / UC

Contr. → CD/UC
Contr. + → CD/UC

Rank eigencentrality (0, 100, 200, 300, 400, 500)

**H**

UC + → CD
UC – → CD

Rank eigencentrality (0, 100, 200, 300, 400, 500)

**Figure S13:** The parallel coordinate plots illustrate the change of positions of significantly central bacteria from one health condition to the other. Single bars visualize nodes ordered by their ranked importance within the respective networks, based on **(A, B)** node degree, **(C, D)** node betweenness, **(E, F)** PageRank index, and eigenvalue centrality **(G, H)**. The first bar shows taxa ranked by their importance in CD specific networks, the second bar taxa ranked by average node importance in networks based on the cohort, and the third bar taxa ranked by their importance in UC specific networks. The change in network position/importance is indicated by lines connecting the same ASVs/nodes between the respective networks/bars. Plots **A, C, E, G** display significant nodes and their respective position/rank in Contr. and CD networks, or Contr., and UC networks. Change of taxa, which increase in importance (smaller rank) when found in a diseased community as compared to the controls is shown in blue, while taxa decreasing in importance (increased rank) in a diseased community are shown in red. Plots **B, D, F, H** display significant nodes and their change in position between CD and UC, with increases of importance (decreasing rank) from CD to UC highlighted in blue, while increases of importance from UC to CD are highlighted in red. Centralities in the control cohort were based on the average of three independent network calculations on three subsamples of the respective cohort to balance differences in cohort sizes.

## Supplemental tables

**Table S1:** Results of differential ASV abundance tests via negative binomial generalized linear models (via *DESeq2,* LR test) with respect to IBD condition and smoking status, corrected for age and BMI. P-values were corrected for multiple testing via FDR-adjustment. ASVs are ordered by the cohort with their maximum abundance and overlaps with significant indicator species are shown. ASVs highlighted in grey are associated with more than one variable.

**Table S2:** Indicator species analysis abundance analyses of ASVs with respect to IBD- and smoking status ($P \leq 0.05$). ASVs are ordered by their association and overlaps with significant differentially abundant ASVs are highlighted.

**Table S3:** Differential abundance analyses of ASVs via *DESeq2* (LR test) with respect to IBD subtypes corrected for age and BMI ($P_{FDR} \leq 0.05$). ASVs are ordered by the cohort with their maximum abundance and overlaps with significant indicator species are shown. Subtypes follow the Montreal classification [13].

**Table S4:** Indicator species analysis abundance analyses of ASVs with respect to IBD subtypes ($P \leq 0.05$). ASVs are ordered by their association and overlaps with significant differentially abundant ASVs are highlighted. Subtypes follow the Montreal classification [13].

**Table S5:** Differential abundance analyses of ASVs via DESeq2 (LR test) with respect to medication with non-IBD related medication and prior antibiotic usage (within the last 6 weeks, within the last 6 months) corrected for age and BMI ($P_{FDR} \leq 0.05$). ASVs are ordered by the cohort with their maximum abundance and overlaps with significant indicator species are shown.

**Table S6:** Indicator species analysis abundance analyses of ASVs with respect to medication with non-IBD related medication and prior antibiotic usage (within the last 6 weeks, within the last 6 months) ($P \leq 0.05$). ASVs are ordered by their association and overlaps with significant differentially abundant ASVs are highlighted.

**Table S7:** Analysis of alpha diversity differences in CD and UC patients focusing on non-IBD related pharmaceutical treatments.

**Table S8:** Betadiversity analysis of CD and UC patients focusing on community variability with respect to anthropogenic and disease related variables, as estimated via a multivariate and permutative version of the Levene's test for homogeneity of variances.

**Table S9:** Betadiversity analysis of CD and UC patients focusing on anthropogenic and disease related variables (PERMANOVA).

**Table S10:** Analyses of betadiversity between CD and UC subtypes following the Montreal classification [13] (PERMANOVA).

**Table S11:** Differential abundance analyses of predicted functions (EC categories) via *DESeq2* (Wald test) with respect to health condition, corrected for age and BMI ($P_{FDR} \leq 0.05$). Functions are ordered by the

cohort with their maximum abundance. Each pairwise comparison is shown and only significant differences were evaluated for the direction of change as indicated by their maximum abundance.

**Table S12:** Differential abundance analyses of predicted pathways (MetaCyc pathways[14]) via *DESeq2* (Wald test) with respect to health condition, corrected for age and BMI ($P_{FDR} \leq 0.05$). Pathways are ordered by the cohort with their maximum abundance. Each pairwise comparison is shown and only significant differences were evaluated for the direction of change as indicated by their maximum abundance.

**Table S13:** Overview of significantly central ASVs in disease specific ASV correlation networks for healthy controls (average of 3 subsets), CD- and UC patients. The number of connections (degree), their position between different nodes (Betweenness) and their general importance (Eigenvector-centrality and PageRank) were used measures of centrality and their significance was determined by a one-sided Z-Test of the observed values against a null distribution of 10'000 permuted networks. Associations to different anthropometric characteristics, as detected via indicators species analyses ($P \leq 0.05$) and differential abundance analyses ($P_{FDR} \leq 0.05$), are included in the table. The color-code highlights in which networks ASVs are repeatedly important (■-important across all health conditions, ■-important only in control networks, ■-important only in the CD network, ■-important only in the UC network).

# Supplemental references

1       Caporaso JG, Lauber CL, Walters WA, et al. Ultra-high-throughput microbial community analysis on the Illumina HiSeq and MiSeq platforms. *ISME J*. 2012.

2       Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-resolution sample inference from Illumina amplicon data. *Nat Meth*. Brief Communication. 2016;13(7):581-583.

3       Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy. *Applied and environmental microbiology*. Research Support, U.S. Gov't, Non-P.H.S. 2007;73(16):5261-5267.

4       Cole JR, Wang Q, Cardenas E, et al. The Ribosomal Database Project: improved alignments and new tools for rRNA analysis. *Nucl Acids Res*. 2009;37(suppl_1):D141-145.

5       Schloss PD. A High-Throughput DNA Sequence Aligner for Microbial Ecology Studies. *PLoS One*. 2009;4(12).

6       Schloss PD, Westcott SL, Ryabin T, et al. Introducing mothur: Open Source, Platform-independent, Community-supported Software for Describing and Comparing Microbial Communities. *Applied and environmental microbiology*. Research Support, Non-U.S. Gov't. 2009;75(23):7537-7541.

7       Pruesse E, Quast C, Knittel K, et al. SILVA: a comprehensive online resource for quality checked and aligned ribosomal RNA sequence data compatible with ARB. *Nucl Acids Res*. 2007;35(21):7188-7196.

8       Quast C, Pruesse E, Yilmaz P, et al. The SILVA ribosomal RNA gene database project: improved data processing and web-based tools. *Nucleic Acids Res*. 2013;41(Database issue):D590-596.

9       Price MN, Dehal PS, Arkin AP. FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments. *PLoS One*. Research Support, U.S. Gov't, Non-P.H.S. 2010;5(3):e9490.

10      Gotelli NJ. Null model analysis of species co-occurrence patterns. *Ecology*. 2000;81(9):2606-2621.

11      Kembel SW, Cowan PD, Helmus MR, et al. Picante: R tools for integrating phylogenies and ecology. *Bioinformatics*. Research Support, Non-U.S. Gov't

Research Support, U.S. Gov't, Non-P.H.S. 2010;26(11):1463-1464.

12      Friedman J, Alm EJ. Inferring Correlation Networks from Genomic Survey Data. *PLoS Computational Biology*. 2012;8(9):e1002687.

13      Satsangi J, Silverberg MS, Vermeire S, Colombel JF. The Montreal classification of inflammatory bowel disease: controversies, consensus, and implications. *Gut*. 2006;55(6):749-753.

14      Karp PD, Riley M, Paley SM, Pellegrini-Toole A. The MetaCyc Database. *Nucleic Acids Res*. 2002;30(1):59-61.