**Supplementary Methods**

**Mucosa-attached microbiota analysis**

Profiles of mucosa-attached bacterial communities were generated from mucosal biopsy RNA (converted to cDNA using random hexamers, Qiagen, Germany). The 16S rRNA variable region V3-V4 was amplified using bacterial 16S rRNA gene-specific composite primers (319F and 806R) as described previously [1]. The primers included a 0 to a 12 nucleotide heterogeneity spacer to enhance sequence quality of low diversity regions and 12 nucleotide barcode to tag the sequences to samples. Mucosal cDNA from each individual was amplified in two separate PCR reactions using Phusion ® Hot start Flex 2X master Mix (New England Biolab, Germany). Pooled amplicon libraries were sequenced employing an Illumina MiSeq (2 x 300 bp), resulting in the generation of paired-end reads with an overlap of about 90 bases.

Match-paired forward and reverse sequence reads (fastq) were concatenated to single reads. Reads having any ambiguous base or more than 8 homopolymers were excluded from further analysis. Sequences were aligned against the mothur-curated [2] Silva alignment database and verified to have an alignment in the targeted region only. Subsequently, chimeric sequences were detected with the Uchime algorithm [3] and removed. Sequences were classified taxonomically with a confidence threshold of 80% using mothur-formatted Silva training sets (version: silva.nr_v119) and discarded if classified as unknown, Archaea, eukaryotes, chloroplast or mitochondria. A phylip-formatted distance matrix was computed from remaining quality aligned sequences. Sequences with at least 97% similarity were clustered into species Operational Taxonomic Units (OTUs) using the average neighbor joining algorithm within mothur. Consenus taxonomy of each OTUs was determined from Silva based taxonomical classification as stated above. Prior to downstream analysis, subsampling was performed to make samples comparable. Alpha diversity indices including observed number of OTUs, non-parametric Shannon and Simpson were calculated by mothur. Non parametric Kruskall-Wallis test was performed in *Graphpad* Prism (*GraphPad* Software, San Diego California USA) to test statistical significance of observed differences between the groups. Differences in overall mucosa attached transcriptionally active bacterial communities composition and structure between diagnosis, tissues and inflammation were tested using permutational multivariate analysis of variance (PerMANOVA). Principle coordinate analysis (PCoA) was performed in the PAST software [4] on OTUs incidence (Jaccard) and relative abundance (Bray Curtis) based distance matrices to visualize the differences of microbial communities under the influence of  disease type, tissue location, and inflammation status.

**Reconstruction of context-specific metabolic models**

To identify metabolic changes between inflamed and non-inflamed tissue, we reconstructed context-specific models for each tissue biopsy using the iMAT procedure [5]. iMAT requires the separation of genes into highly

expressed, moderately expressed and lowly expressed genes according to expression values. Thus, the FPKM-values for each sample were transformed using a Box-Cox-transformation [6] and genes were classified according to the mean $\mu$ and standard deviation $\sigma$ of the distribution of transformed FPKM-values. Genes were defined as lowly expressed if their expression value was smaller than $\mu - \sigma$, highly expressed if their expression was larger than $\mu + \sigma$ and moderately expressed otherwise. Context-specific metabolic networks were obtained based on the genome-scale metabolic reconstruction of humans, Recon 2.04 [7], using the iMAT-implementation of the COBRA-Toolbox [8]. For down-stream analyses, reactions were summarized into the metabolic subsystems to which they were associated by counting, for each sample, the number of reactions belonging to this subsystem, yielding a metabolic activity of each subsystem. Additionally, we added an artificial subsystem "All" for which the total number of reactions of the entire model was stored to test for significant differences in model size. Subsequently, the number of reactions in each subsystem was rank-normalized across all samples.

**Host transcriptome analysis**

Raw reads passing the chastity filter from Illumina were first pre-processed using cutadapt [9] and PrinSeq-lite [10] to remove adapter and low quality sequences. The reads were aligned to the human genome (hg19/CRCh37)Ensembl GRCh37 (hg19) reference genome) using TopHat2 [11], while the corresponding GTF annotation file was obtained from the Ensembl database (Homo_sapiens.GRCh37.63.gtf).

Differential gene expression levels of the transcripts quantified by HTSeq [12] were analyzed using the Bioconductor package DESeq2 [13] which is capable of handling multifactorial experimental designs, in this case considering tissue type (sigmoidal colon versus terminal ileum), inflammatory stage (inflamed versus not-inflamed), and diagnosis status (Crohn's disease, terminal ileum inflamed versus non-inflamed and ulcerative colitis, sigmoidal colon inflamed versus non-inflamed). DESeq2 was then used to perform a likelihood ratio test modeling tissue type, inflammatory stage or diagnostic status to identify genes differentially expressed in the respective experimental subgroups. Genes with an adjusted p-value below 0.01 were considered significantly differentially expressed, representing a conservative and stringent approach.

To gain insight into the nature of the genes differentially expressed in each analysis (tissue type, inflammatory stage and diagnosis status), KEGG pathway [14] enrichment analysis was performed for the top 150 up- and down-regulated genes (by p-value) separately using the InnateDB database (www.innatedb.com) [15].

**Transcription factor binding site analysis**

Over-represented transcription factor binding sites (TFBS) in the promoter region of the top 150 differentially expressed genes (up- and down- regulated) were identified employing innateDB, which incorporates predicted transcription factor binding site data from the CisRED database (www.cisred.org).

**Analysis of splicing events**

The essential criteria for a given event to be considered as alternative splicing [16]:

- For exon skipping (Ex): (i) ≥10 actual reads mapping to the sum of exclusion exon-exon junctions (EEJs), or (ii) ≥10 actual reads mapping to one of the two inclusion EEJs and ≥5 actual reads mapped to the other inclusion EEJ.

- For intron retention (IR): (i) ≥10 actual reads mapping to the sum of skipping EEJs, or (ii) ≥10 actual reads mapping to one of the two inclusion and exon-intron junctions (EIJs) and ≥5 actual reads to the other inclusion EIJ.

- For alternative 3' (ALTA) and 5' (ALTD) splice sites: ≥10 actual reads mapping to the sum of all EEJs involved in the specific event.

An additional filtering step was performed on retained introns that discards intron retention events with a binomial p value above 0.05 across all samples [17]. Alternative splicing events with sufficient read coverage and expression value (cRPKM) ≥2 across all samples were subjected to an analysis of differentially splicing. Splicing events with an absolute Δ (PSI/PIR/PSU) above 10 were considered differentially regulated in each pair-wise comparison. Furthermore, for those regulated and not regulated (up and down) splicing events, we calculated the Spearman correlation coefficients between their PSI/PIR/PSU and the gene expression of the corresponding genes. KEGG pathway [14] enrichment analysis was performed for all alternative splicing events with an absolute Δ(PSI/PIR/PSU) greater than 10 employing innateDB [15], while p-values ≤ 0.05 were considered significant.

**Identification of differentially regulated metabolic pathways**

To infer differentially regulated pathways between inflamed and non-inflamed tissue, we used unbalanced type-II repeated measures analysis of variance (ANOVA), implemented in the R-package 'car' [18] with 'lmer' from the 'lme4'-package [19] for fitting linear mixed-effects models. The outcome (dependent variable) was 'subsystem activity' and 'inflammatory status' (inflamed, non-inflamed), 'disease type' (healthy control, disease control, Ulcerative colitis, Crohn's disease), gender (female, male) and 'tissue' (terminal ileum and sigmoidal colon) were treated as fixed independent factors. The ID of the patient from which each sample originated was treated as error term nested within 'tissue' (i.e. samples from the same patient originated always from the two tissues). All interactions up to four-way were included into the model. We did not consider 4 samples derived from the same patient at two different time points, since this could not be explicitly considered. Normality was tested using a Shapiro-Wilks test and subsystems discarded for p<0.01. Interactions were iteratively reduced by starting with the

highest order of interaction. Main effects or lower order interactions were only removed, if no higher order interaction in the model included this interaction or a main effect. Model reduction continued until no influence with p>0.05 remained or until the main effect 'inflammatory status' was removed. Reported p values correspond to the p value of the interaction between the main effect 'inflammatory status' and 'subsystem activity'. P values were corrected for multiple testing by controlling the false-discovery rate according to Benjamini-Hochberg [20].

**Analysis of metabolic network coherence**

The metabolic network coherence was computed by mapping significant gene expression differences of metabolic genes onto a gene-centric projection of a metabolic network derived from the Recon2 metabolic model [7]. The set of metabolic genes with significant expression differences, together with the links in the gene-centric metabolic network form an effective network, which shows, how connected the expression differences between two conditions are from a metabolic perspective. The percentage of connected genes in this effective network (i.e., genes with non-zero network degree) serves as a quantifier of this connectedness. The metabolic network coherence M is the z-score of this quantifier with respect to randomized gene expression data. A value of M = 1 thus means that the percentage of connected metabolic genes derived from the expression differences is one standard deviation higher than the same quantity in random data. The method was employed as previously described [21], while applying the modifications required for human gene expression data [22]. To integrate the networks into the interaction between host transcriptome, splicing patterns and microbiome, differentially expressed genes were forwarded to the network analysis (fold change threshold -1.5 and +1.5).

To contextualize the individual nodes obtained in these metabolic networks, genes from these networks have been used as a basis to determine the number of occurrence of metabolic functional categories based on the Recon 2 model.

**References for supplementary methods**

1   Fadrosh DW, Ma B, Gajer P, *et al.* An improved dual-indexing approach for multiplexed 16S rRNA gene sequencing on the Illumina MiSeq platform. *Microbiome* 2014;**2**:6. doi:10.1186/2049-2618-2-6

2   Schloss PD, Westcott SL, Ryabin T, *et al.* Introducing mothur: Open-Source, Platform-Independent, Community-Supported Software for Describing and Comparing Microbial Communities. *Appl Env Microbiol* 2009;**75**:7537–41. doi:10.1128/AEM.01541-09

3   Edgar RC, Haas BJ, Clemente JC, *et al.* UCHIME improves sensitivity and speed of chimera detection. *Bioinforma Oxf Engl* 2011;**27**:2194–200. doi:10.1093/bioinformatics/btr381

4   Hammer O, Harper D, Ryan P. PAST: Paleontological Statistics Software Package for Education and Data Analysis. *Palaentologigia Electron* 2001;**4**:9.

5  Zur H, Ruppin E, Shlomi T. iMAT: an integrative metabolic analysis tool. *Bioinforma Oxf Engl* 2010;**26**:3140–2. doi:10.1093/bioinformatics/btq602

6  Box, G.E.P., Cox, D.R. An Analysis of Transformations on JSTOR. *J R Stat Soc Ser B Methodol* 1964;**26**:211–52.

7  Thiele I, Swainston N, Fleming RMT, *et al.* A community-driven global reconstruction of human metabolism. *Nat Biotechnol* 2013;**31**:419–25. doi:10.1038/nbt.2488

8  Schellenberger J, Que R, Fleming RMT, *et al.* Quantitative prediction of cellular metabolism with constraint-based models: the COBRA Toolbox v2.0. *Nat Protoc* 2011;**6**:1290–307. doi:10.1038/nprot.2011.308

9  Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* 2011;**17**:10–2. doi:10.14806/ej.17.1.200

10     Schmieder R, Edwards R. Quality control and preprocessing of metagenomic datasets. *Bioinforma Oxf Engl* 2011;**27**:863–4. doi:10.1093/bioinformatics/btr026

11     Trapnell C, Roberts A, Goff L, *et al.* Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* 2012;**7**:562–78. doi:10.1038/nprot.2012.016

12     Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput sequencing data. *Bioinforma Oxf Engl* 2015;**31**:166–9. doi:10.1093/bioinformatics/btu638

13     Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* 2014;**15**:550. doi:10.1186/s13059-014-0550-8

14     Kanehisa M, Goto S, Kawashima S, *et al.* The KEGG resource for deciphering the genome. *Nucleic Acids Res* 2004;**32**:D277-280. doi:10.1093/nar/gkh063

15     Breuer K, Foroushani AK, Laird MR, *et al.* InnateDB: systems biology of innate immunity and beyond--recent updates and continuing curation. *Nucleic Acids Res* 2013;**41**:D1228-1233. doi:10.1093/nar/gks1147

16     Irimia M, Weatheritt RJ, Ellis JD, *et al.* A highly conserved program of neuronal microexons is misregulated in autistic brains. *Cell* 2014;**159**:1511–23. doi:10.1016/j.cell.2014.11.035

17     Braunschweig U, Barbosa-Morais NL, Pan Q, *et al.* Widespread intron retention in mammals functionally tunes transcriptomes. *Genome Res* 2014;**24**:1774–86. doi:10.1101/gr.177790.114

18     Fox, J., Weisberg, S. An R Companion to Applied Regression | SAGE Publications Ltd. *Sage Publ* Published Online First: 2010.https://uk.sagepub.com/en-gb/eur/an-r-companion-to-applied-regression/book233899 (accessed 11 Dec2015).

19     Bates D, Mächler M, Bolker B, *et al.* Fitting Linear Mixed-Effects Models using lme4. *ArXiv14065823 Stat* Published Online First: 23 June 2014.http://arxiv.org/abs/1406.5823 (accessed 11 Dec2015).

20     Benjamini Y, Hochberg Y. A practical and powerful approach to multiple testing. *J Roy Stat Soc Ser B* 1995;**57**:289–300.

21     Sonnenschein N, Geertz M, Muskhelishvili G, *et al.* Analog regulation of metabolic demand. *BMC Syst Biol* 2011;**5**:40. doi:10.1186/1752-0509-5-40

22     Sonnenschein N, Golib Dzib JF, Lesne A, *et al.* A network perspective on metabolic inconsistency. *BMC Syst Biol* 2012;**6**:41. doi:10.1186/1752-0509-6-41