

# Broliden\_5325

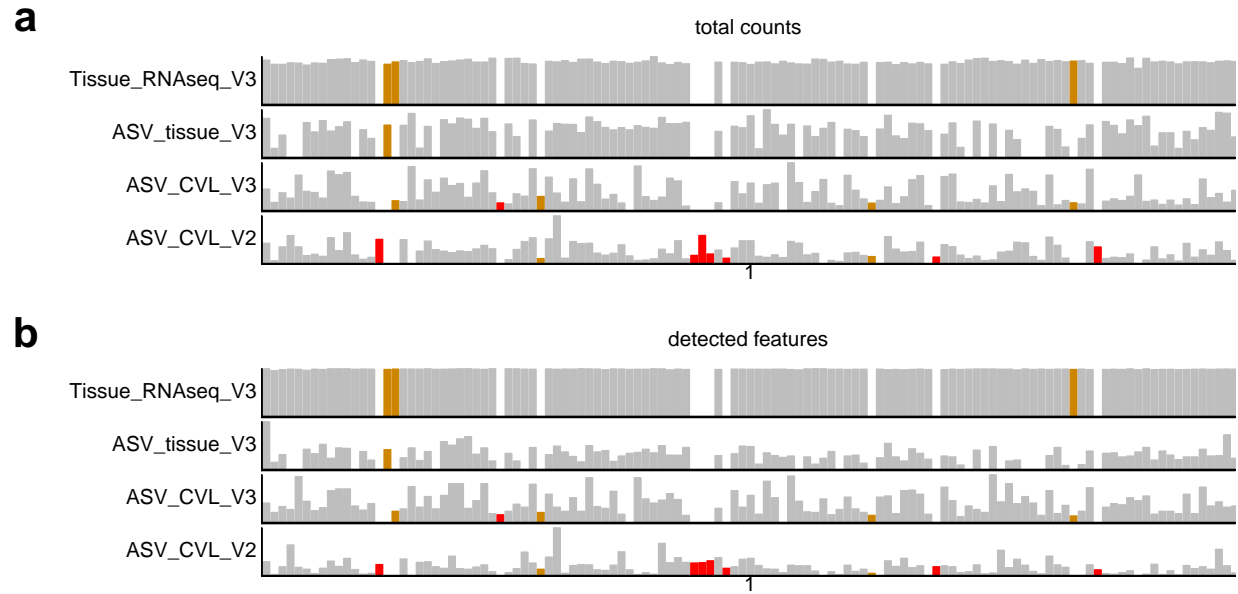
28 October, 2020

## Contents

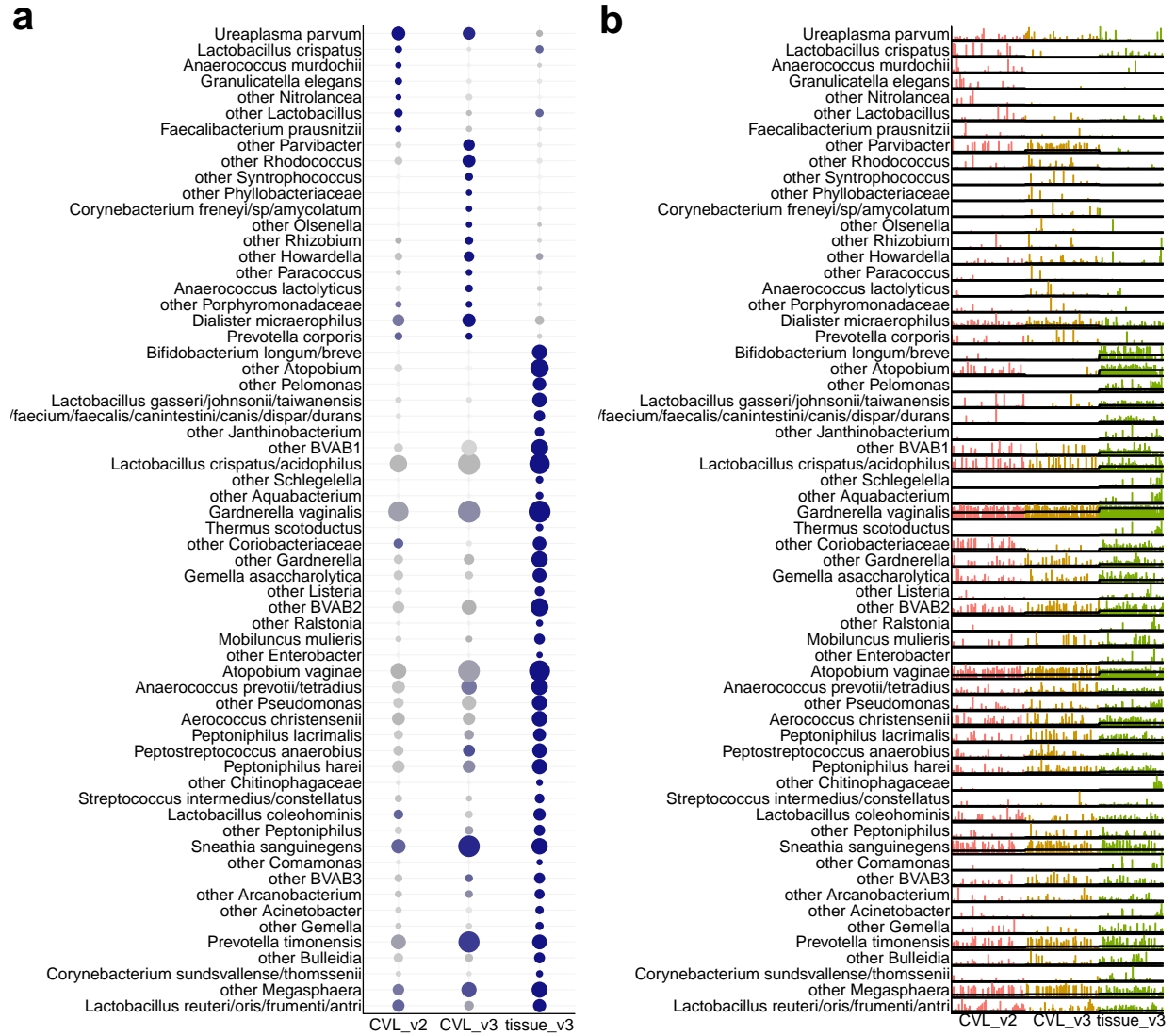
<b>Loading data and metadata</b>	<b>1</b>
<b>Calculate QC metrics</b>	<b>1</b>
<b>Computing differential expression across microbiome datasets</b>	<b>4</b>
#Load libraries and other scripts	
#Defining some variables for the analysis	

## Loading data and metadata

## Calculate QC metrics

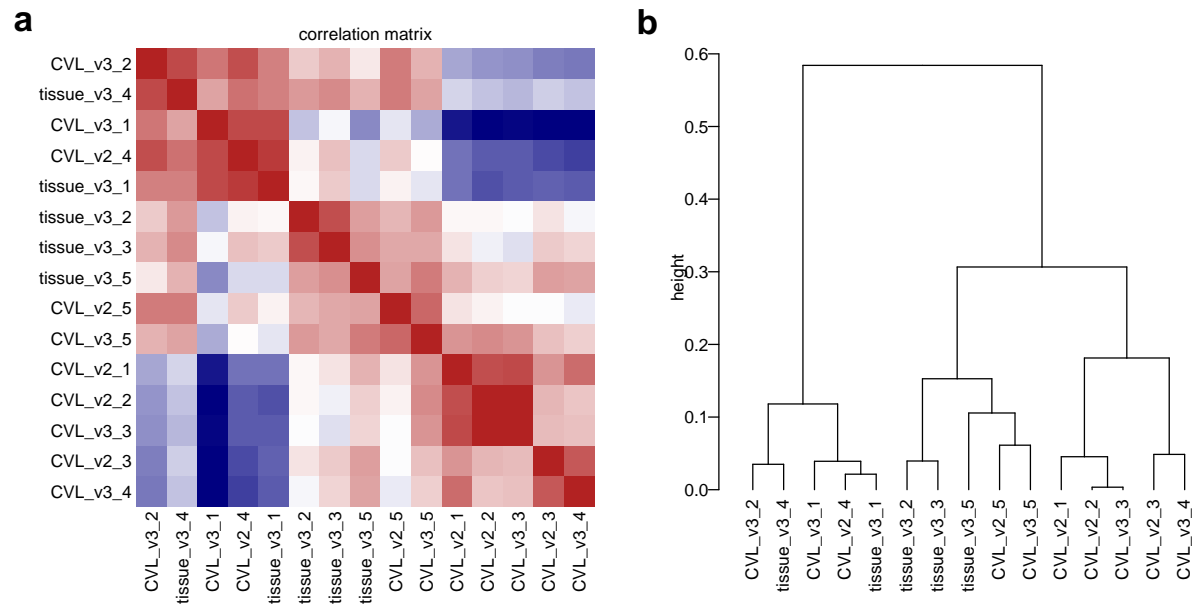


**Figure 1.** Comparative barplot for the **a)** total counts and **b)** number of non-zero detected features (genes / bacteria) for each of the sequencing datasets. Samples are ordered alphabetically according to the patient ID. Samples that are present in exactly two datasets are shown in orange (P015,P030,P054,P055,P056,P058,P084,P104). Samples that are present in exactly two datasets are shown in red (P016,P017,P035,P076,P101).



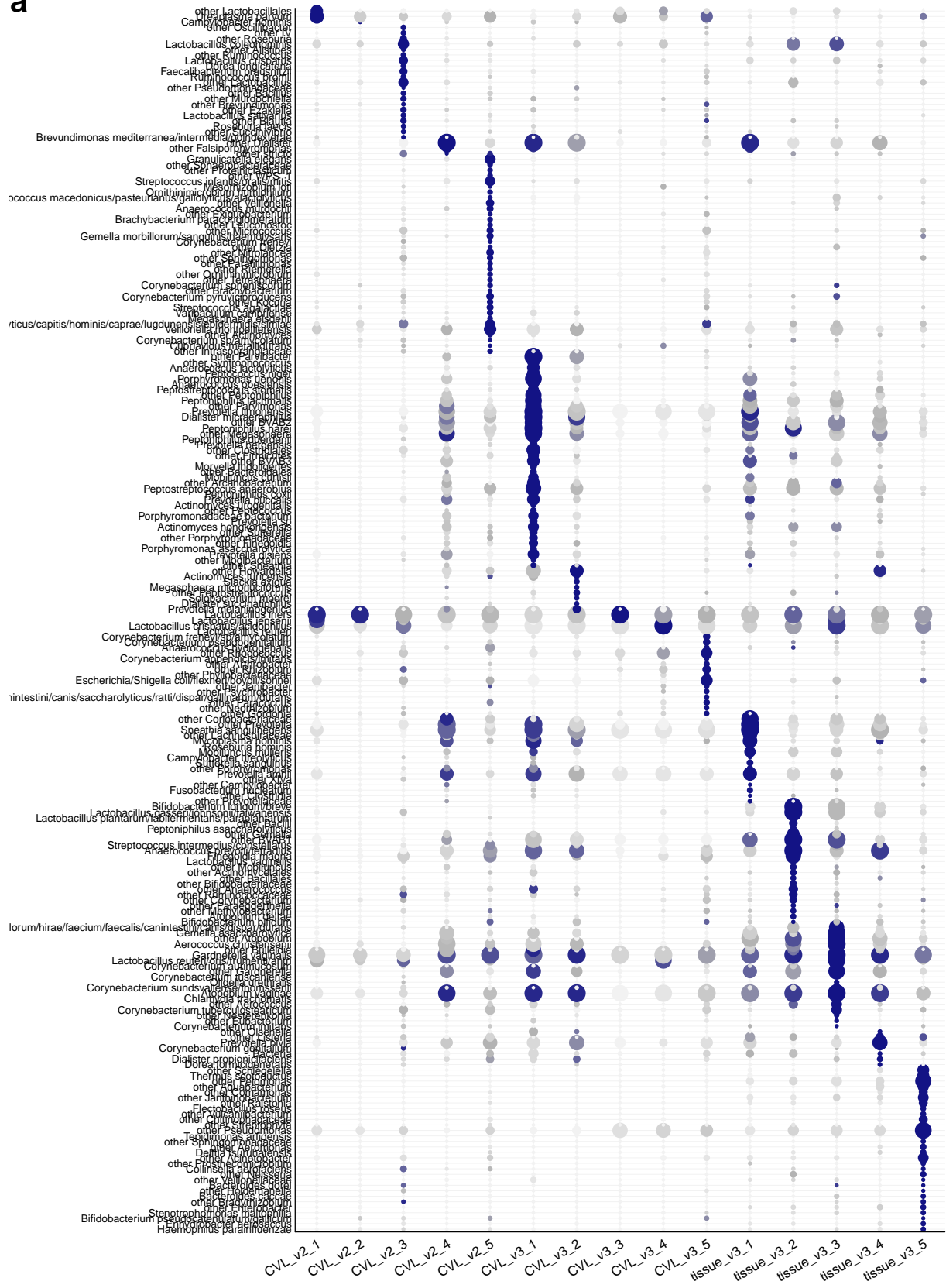
**Figure 2.** Differential bacterial abundance across microbiome datasets. The results are shown both as a) Dot plots and b) barplots. Bacteria with log2FC above 0.25 and p-value below 0.01 were considered significant and were sorted by the highest expression.

# Computing differential expression across microbiome datasets



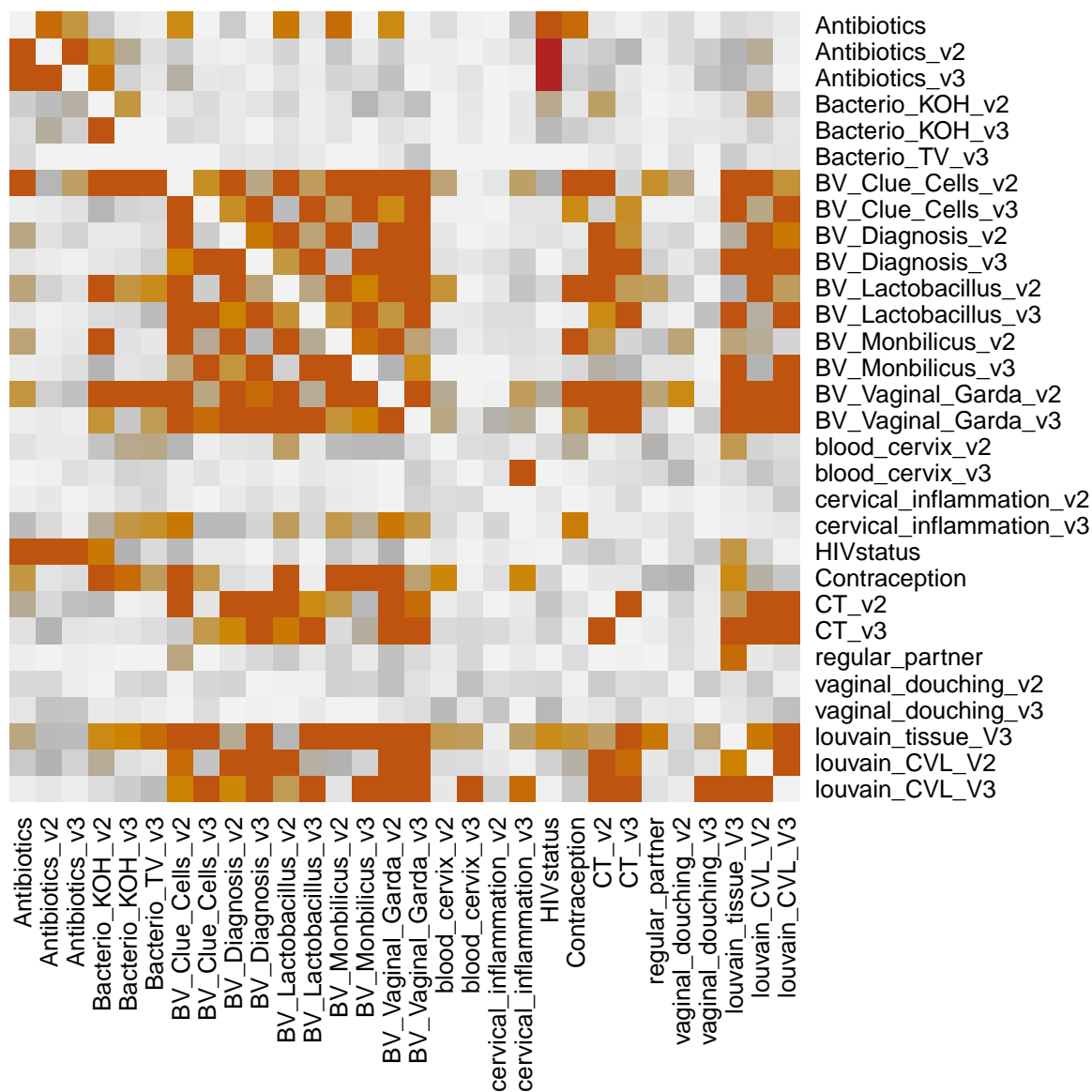
**Figure 3.** Comparisson among patient groups across datasets. **(a)** Correlation matrix across sample groups. **(b)** Hierarchical clustering of sample groups.

**a**

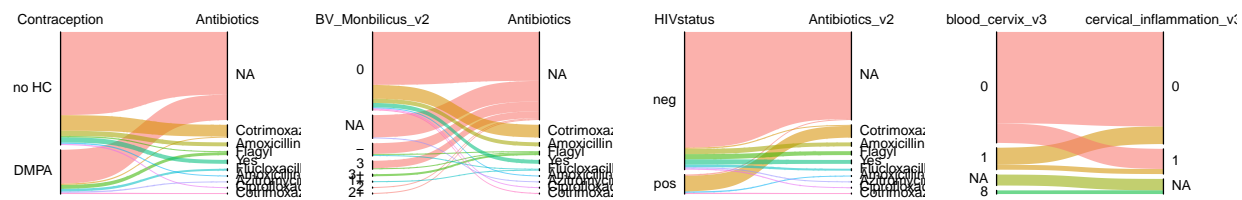


**Figure 4.** Differential bacterial abundance across all groups and all microbiome datasets. The results are shown both as **a)** Dot plots and **b)** barplots. Bacteria with  $\log_2FC$  above 0.25 and p-value below 0.01 were considered significant and were sorted by the highest expression.

**Figure 5.** Comparison of microbiome datasets, showing only the significant bacteria. Samples are ordered by the CVL3 groupings. The colors represent their respective bacterial groupings for each dataset.

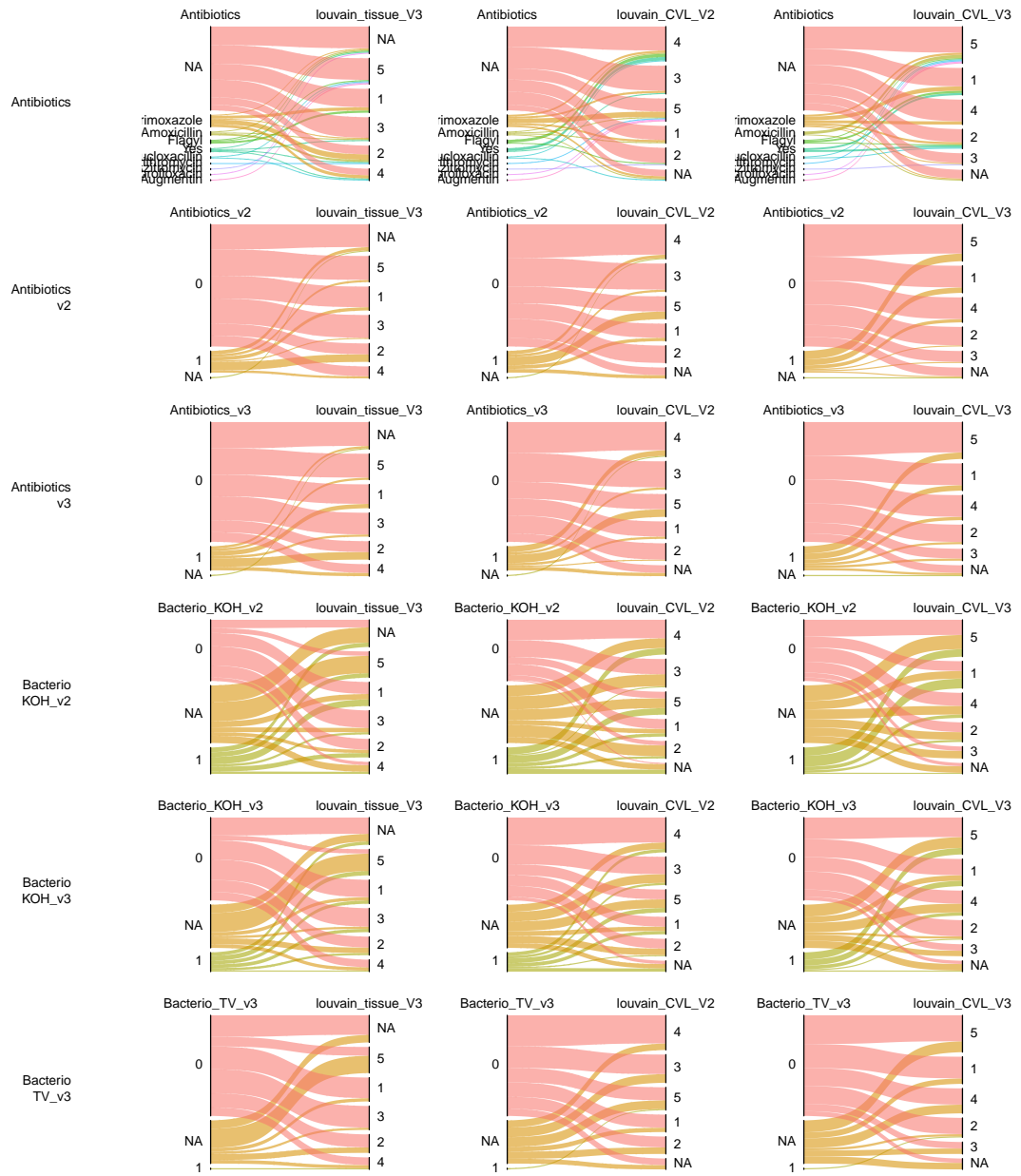


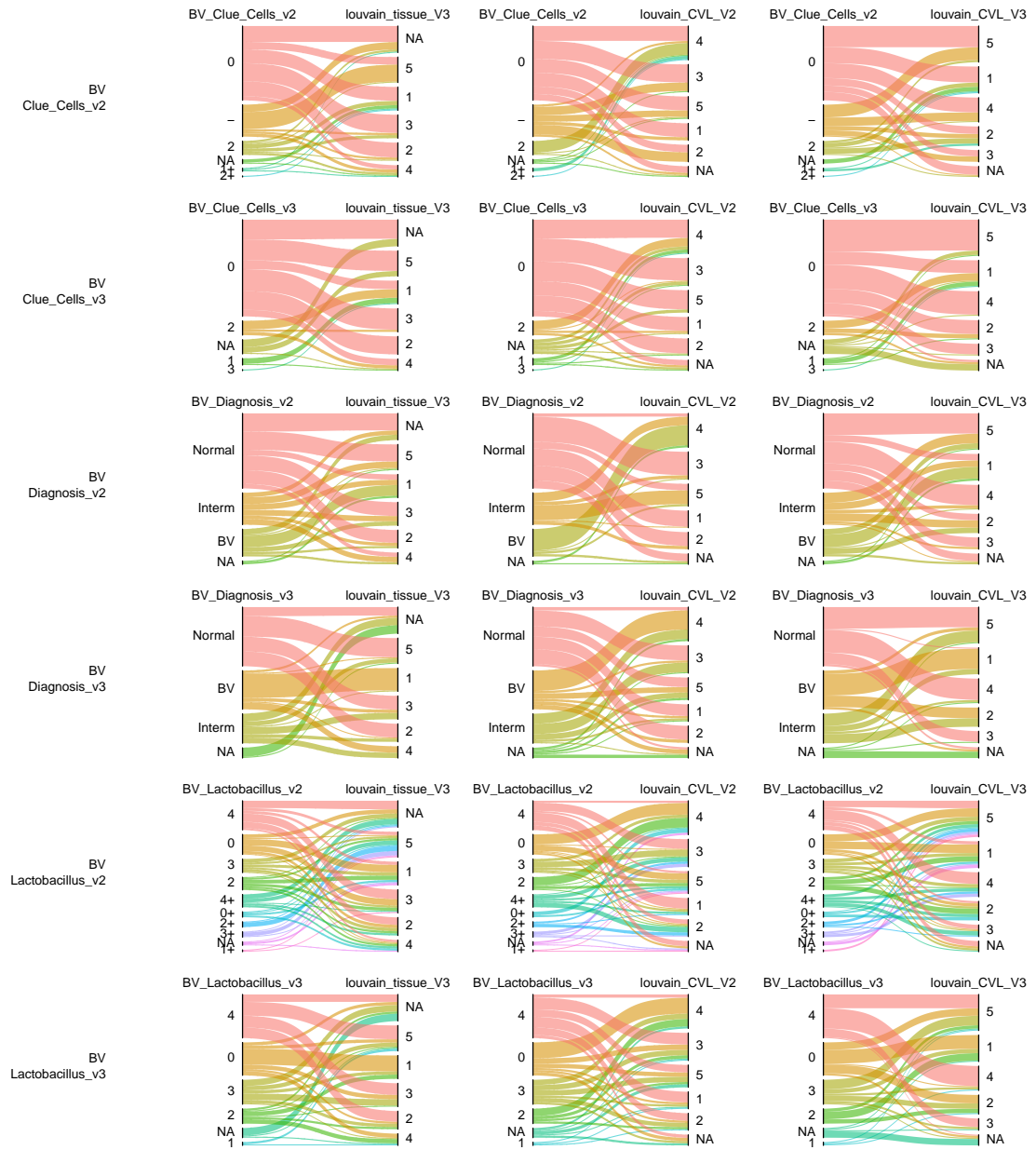
**Figure 6.** Association analysis across several patient categorical metadata parameters, including patient groupings annotations from microbiome.

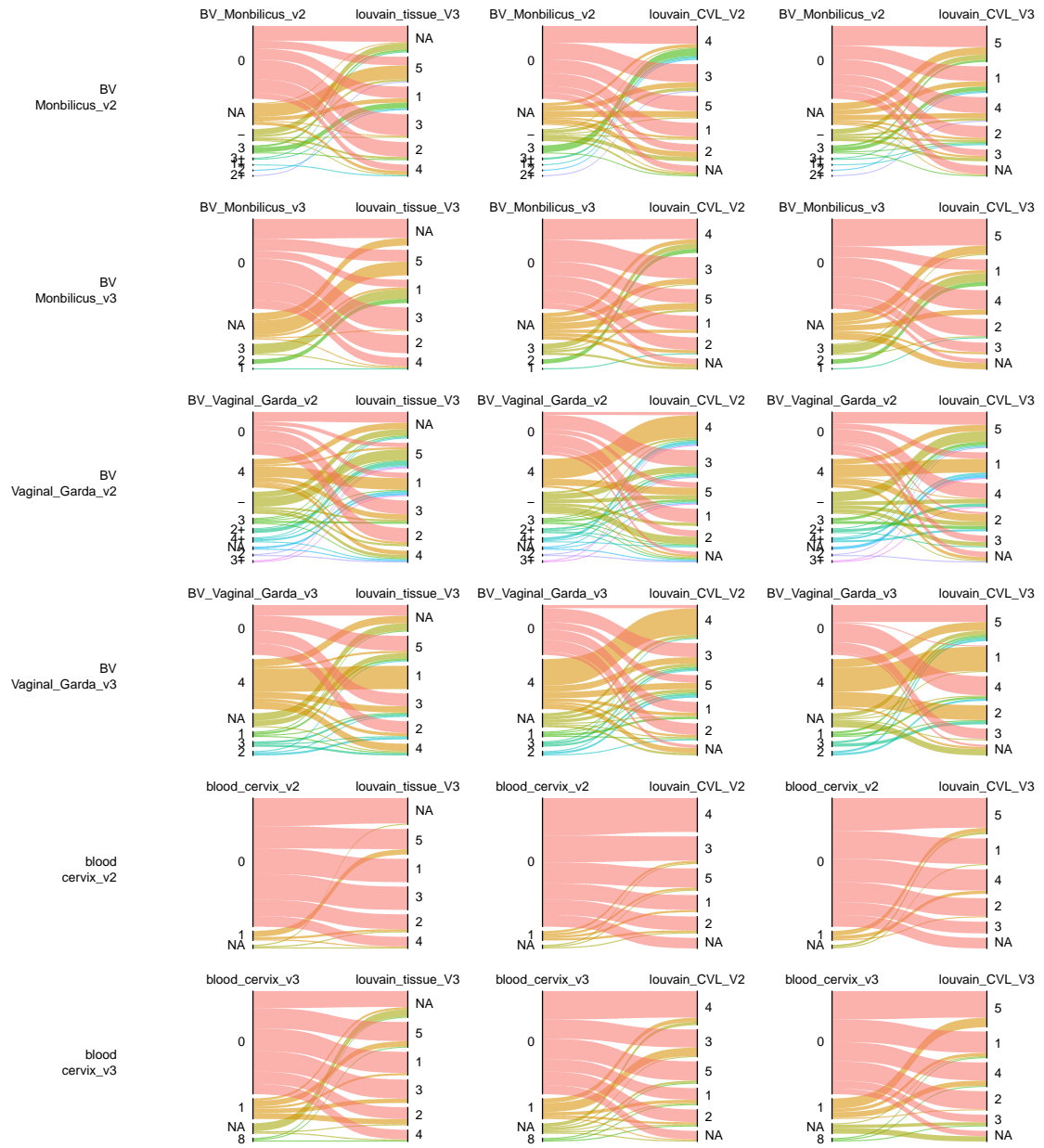


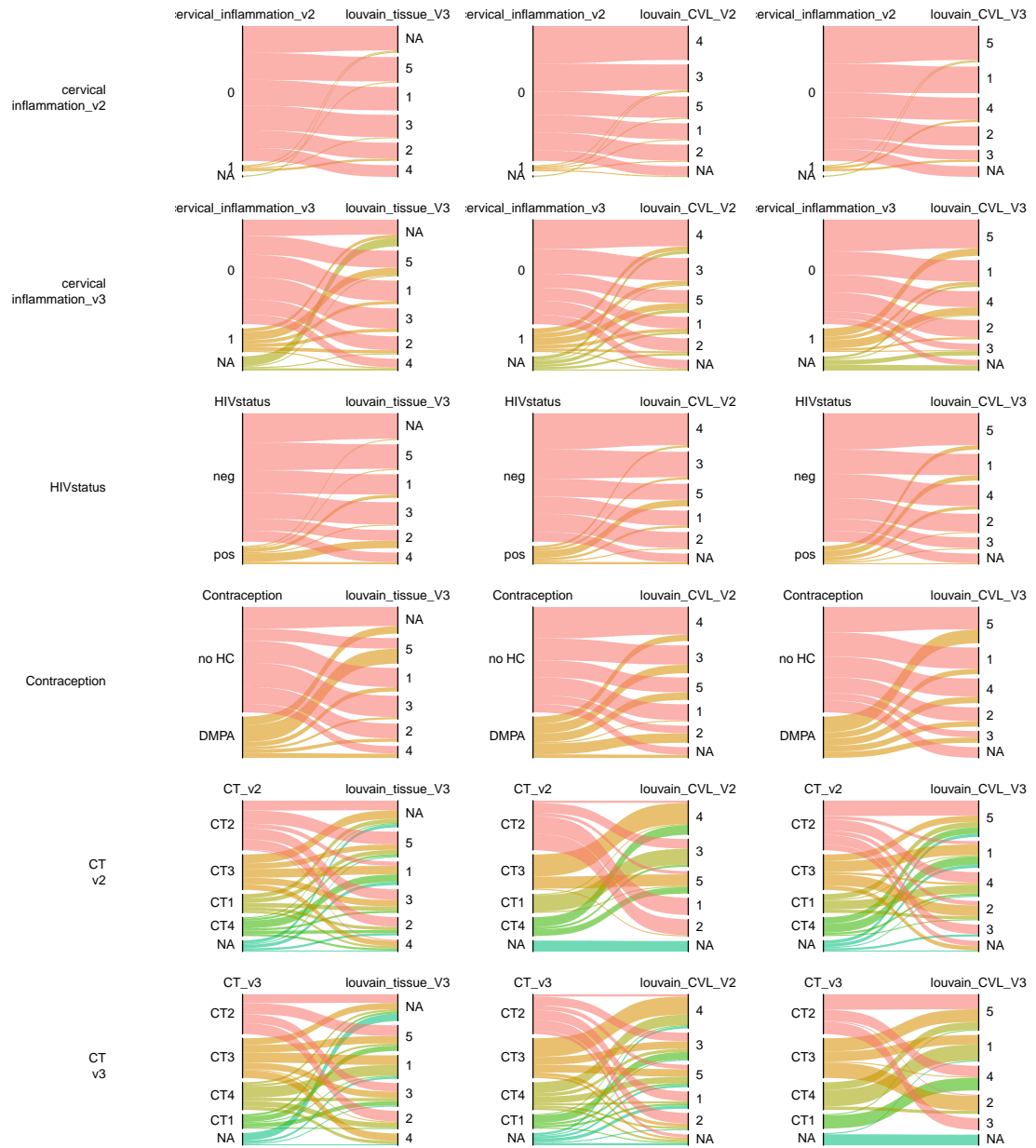
**Figure 7.** A few examples of significant association between metadata parameters shown as sankey plots.

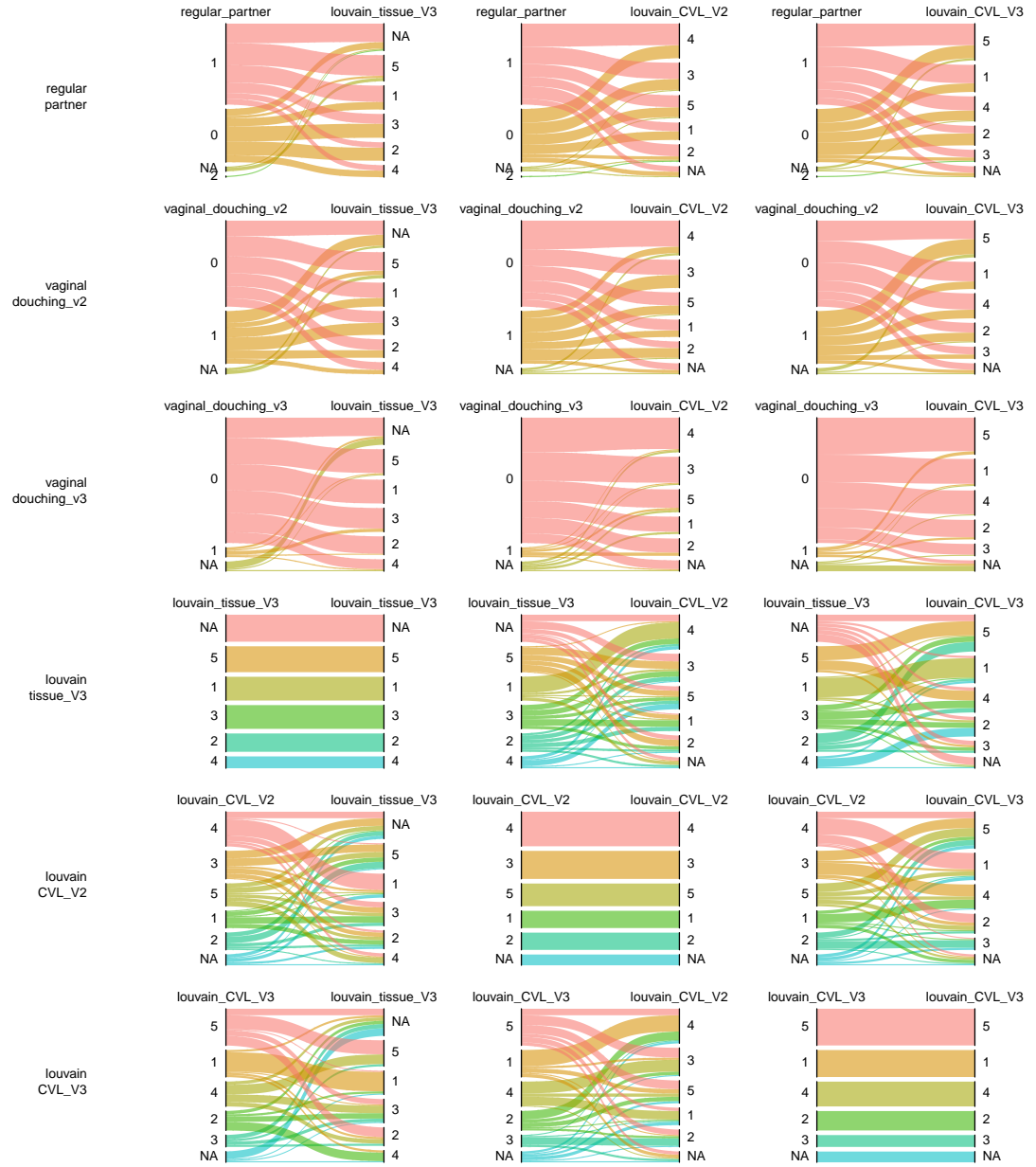










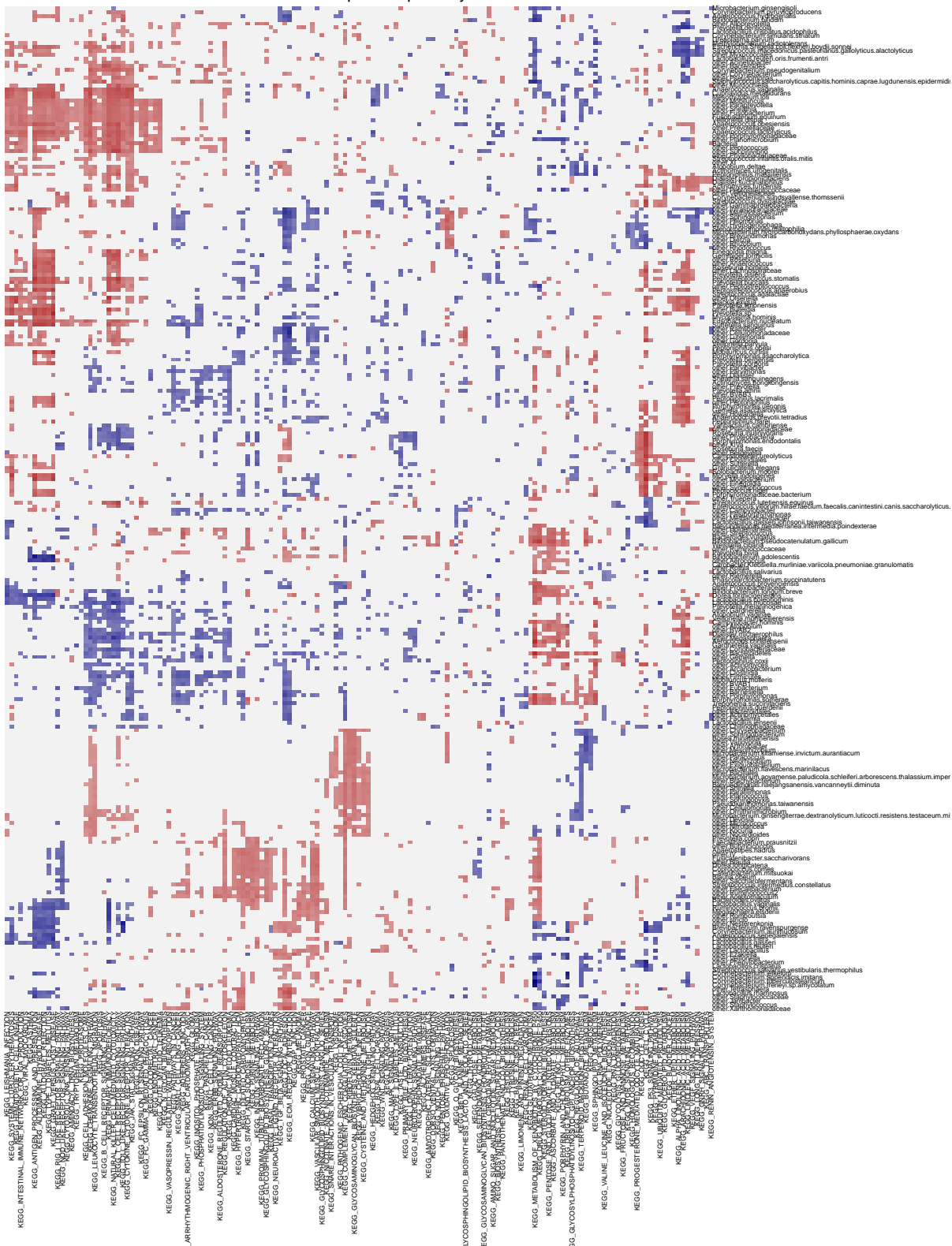


**Figure 8.** Sankey plots for all tested associations between the patient groups identified in the microbiome datasets.

ASV\_tissue\_V3\_normalized\_batch\_corrected

ASV\_CVL\_V3\_normalized\_batch\_corrected

## tissue RNAseq KEGG pathways

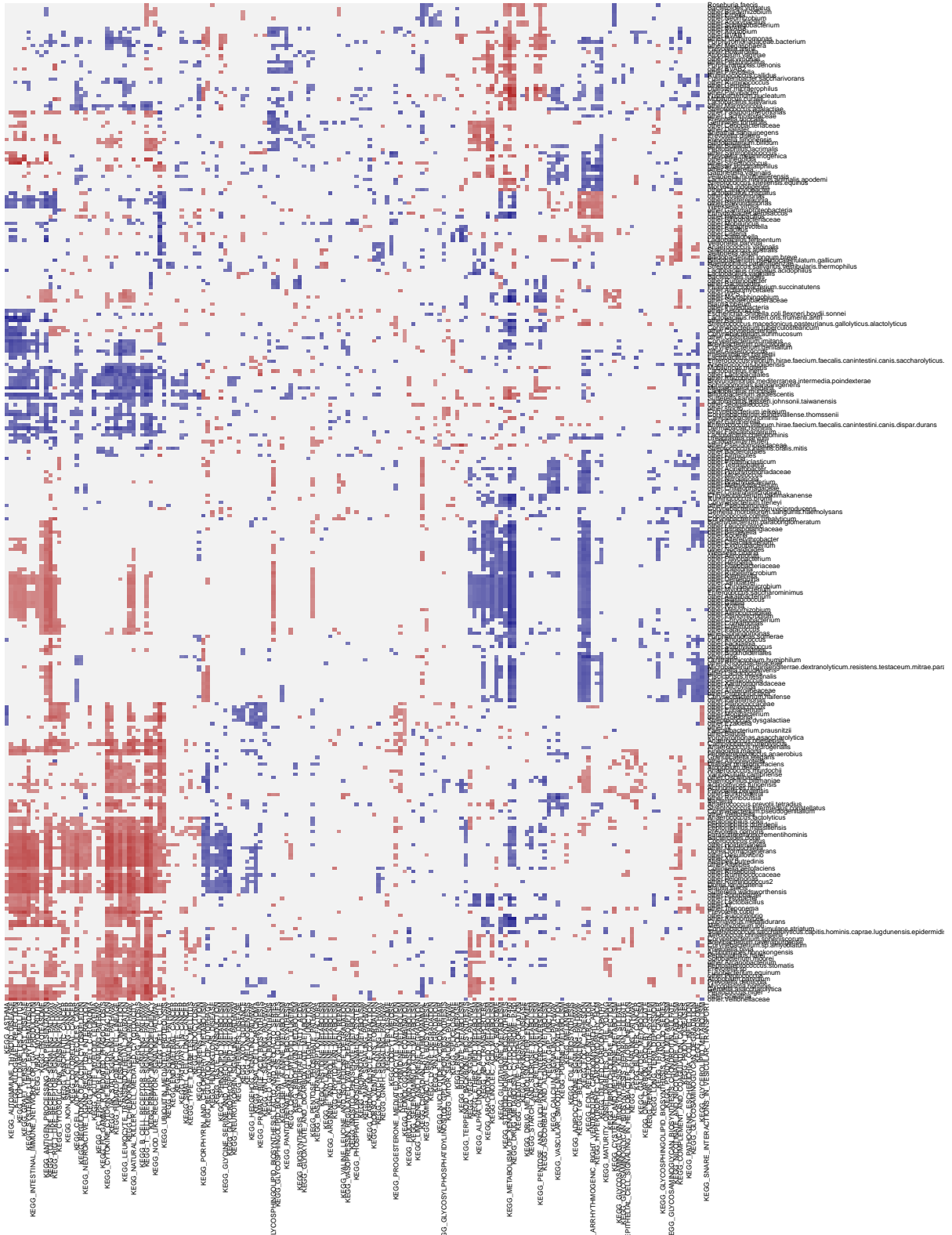




**C**

tissue RNAseq KEGG pathways

ASV\_CVL\_V2\_normalized\_NOT\_batch\_corrected





of the top 5000 highly variable genes from the RNAseq dataset, generating a correlation matrix between bacteria and genes. Then for each bacteria, we rank genes based on their correlation to that bacteria and perform gene set enrichment analysis (GSEA) using the KEGG gene annotation database. This, in turn, will result in a matrix associating every bacteria with every KEGG process in the tissue. The heatmap shows the normalized enrichment score (NES). Only enrichments with pvalue below 0.05 are shown. Bacteria and pathways significant in less than 10 pathways and bacteria, respectively, were omitted.