# ADVANZ4 Exploratory Report

## Introduction

This report contains the analysis of microbiome data including Alpha-Diversity, Ordination, Hierarchical Clustering and detection of differentially abundant taxa. The report was generated using the following parameters:

- Diversity slot type (**taxa__slot**): igc
- Taxonomic classification system (**taxa__slot**): metaphlan
- Taxonomic levels (**tax__level**): Species
- Metadata variables (**metadata__vars**): group, risk_group, center, gender, ethnic_group, CD4diff_48, CD8diff_48, CD4after_48, CD8after_48, CD4, CD8, CD8_CD38_DR, CRP, IL6, TNFa, sCD14, time_point
- Group variable (**group__var**): group

The following table shows a summary of selected metadata variables with descriptive statistics.

| Characteristic | N | **Overall**, N = 271[1] | **DTG**, N = 144[1] | **RTVr**, N = 127[1] | p-value[2] |
|---|---|---|---|---|---|
| risk_group | 259 | | | | 0.12 |
| hts | | 109 (42%) | 50 (36%) | 59 (48%) | |
| msm | | 145 (56%) | 85 (62%) | 60 (49%) | |
| pwid | | 5 (1.9%) | 2 (1.5%) | 3 (2.5%) | |
| Unknown | | 12 | 7 | 5 | |
| center | 271 | | | | |
| bellvitge | | 51 (19%) | 35 (24%) | 16 (13%) | |
| clinic | | 120 (44%) | 57 (40%) | 63 (50%) | |
| hgtp | | 17 (6.3%) | 9 (6.2%) | 8 (6.3%) | |
| mataro | | 9 (3.3%) | 0 (0%) | 9 (7.1%) | |
| sant__pau | | 23 (8.5%) | 14 (9.7%) | 9 (7.1%) | |
| vall__hebron | | 51 (19%) | 29 (20%) | 22 (17%) | |
| Unknown | | 0 | 0 | 0 | |
| gender | 271 | | | | 0.9 |
| female | | 35 (13%) | 19 (13%) | 16 (13%) | |
| male | | 236 (87%) | 125 (87%) | 111 (87%) | |
| Unknown | | 0 | 0 | 0 | |
| ethnic_group | 267 | | | | 0.005 |
| asian | | 4 (1.5%) | 4 (2.8%) | 0 (0%) | |
| black | | 18 (6.7%) | 7 (4.9%) | 11 (8.9%) | |
| caucassian | | 106 (40%) | 66 (46%) | 40 (33%) | |
| hispanic | | 105 (39%) | 56 (39%) | 49 (40%) | |
| other | | 34 (13%) | 11 (7.6%) | 23 (19%) | |
| Unknown | | 4 | 0 | 4 | |
| CD4diff_48 | 205 | | | | 0.006 |
| <50 | | 7 (3.4%) | 0 (0%) | 7 (7.0%) | |

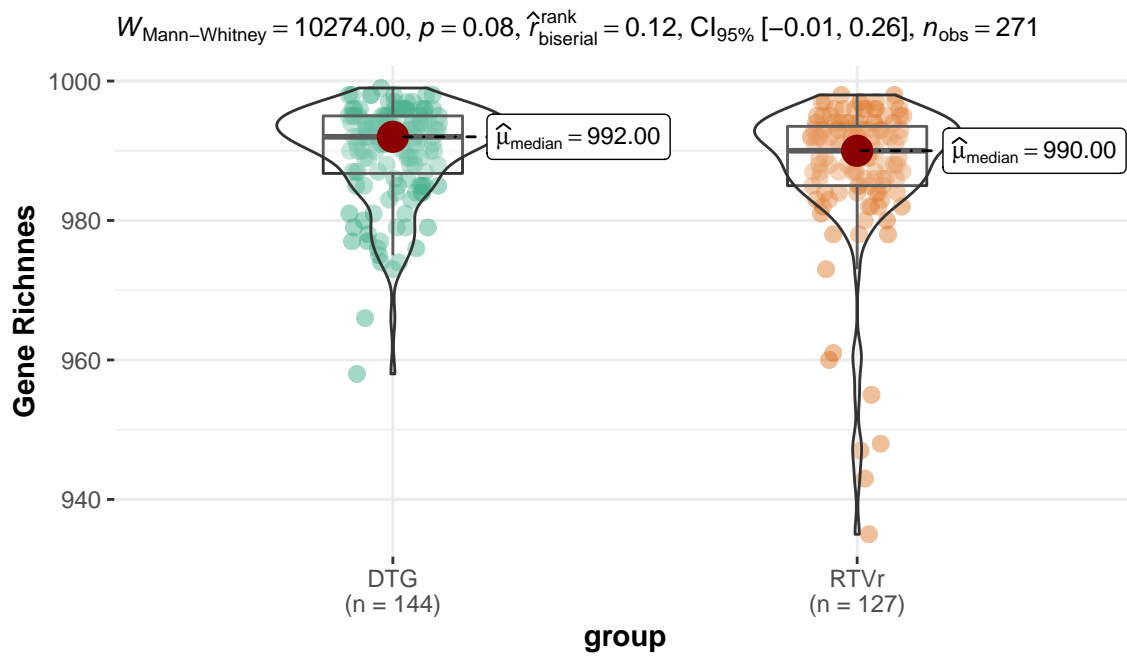| | | | | | |
|---|---|---|---|---|---|
| >50 | | 198 (97%) | 105 (100%) | 93 (93%) | |
| Unknown | | 66 | 39 | 27 | |
| CD8diff__48 | 205 | | | | 0.024 |
| <50 | | 69 (34%) | 43 (41%) | 26 (26%) | |
| >50 | | 136 (66%) | 62 (59%) | 74 (74%) | |
| Unknown | | 66 | 39 | 27 | |
| CD4after__48 | 205 | | | | <0.001 |
| high | | 15 (7.3%) | 4 (3.8%) | 11 (11%) | |
| low | | 90 (44%) | 37 (35%) | 53 (53%) | |
| mid | | 100 (49%) | 64 (61%) | 36 (36%) | |
| Unknown | | 66 | 39 | 27 | |
| CD8after__48 | 205 | | | | 0.048 |
| high | | 171 (83%) | 93 (89%) | 78 (78%) | |
| low | | 3 (1.5%) | 0 (0%) | 3 (3.0%) | |
| mid | | 31 (15%) | 12 (11%) | 19 (19%) | |
| Unknown | | 66 | 39 | 27 | |
| CD4 | 201 | 113 (40, 237) | 140 (60, 245) | 89 (30, 216) | 0.049 |
| Unknown | | 70 | 38 | 32 | |
| CD8 | 201 | 697 (462, 1,132) | 725 (504, 1,168) | 688 (414, 1,053) | 0.3 |
| Unknown | | 70 | 38 | 32 | |
| CD8_CD38_DR | 151 | 31 (19, 48) | 31 (17, 48) | 31 (20, 47) | 0.8 |
| Unknown | | 120 | 67 | 53 | |
| CRP | 133 | 0.18 (0.09, 0.51) | 0.15 (0.08, 0.42) | 0.20 (0.10, 0.59) | 0.14 |
| Unknown | | 138 | 77 | 61 | |
| IL6 | 128 | 8 (2, 20) | 8 (2, 23) | 7 (2, 17) | 0.7 |
| Unknown | | 143 | 77 | 66 | |
| TNFa | 152 | 12 (9, 17) | 12 (8, 19) | 13 (9, 16) | 0.8 |
| Unknown | | 119 | 66 | 53 | |
| sCD14 | 151 | 2,204 (1,720, 2,986) | 1,976 (1,620, 2,836) | 2,250 (1,885, 3,139) | 0.046 |
| Unknown | | 120 | 67 | 53 | |
| time_point | 271 | | | | 0.9 |
| 0 | | 81 (30%) | 40 (28%) | 41 (32%) | |
| 24 | | 51 (19%) | 29 (20%) | 22 (17%) | |
| 48 | | 71 (26%) | 38 (26%) | 33 (26%) | |
| 96 | | 68 (25%) | 37 (26%) | 31 (24%) | |
| Unknown | | 0 | 0 | 0 | |

[1]n (%); Median (IQR)

[2]Fisher's exact test; Pearson's Chi-squared test; Wilcoxon rank sum test
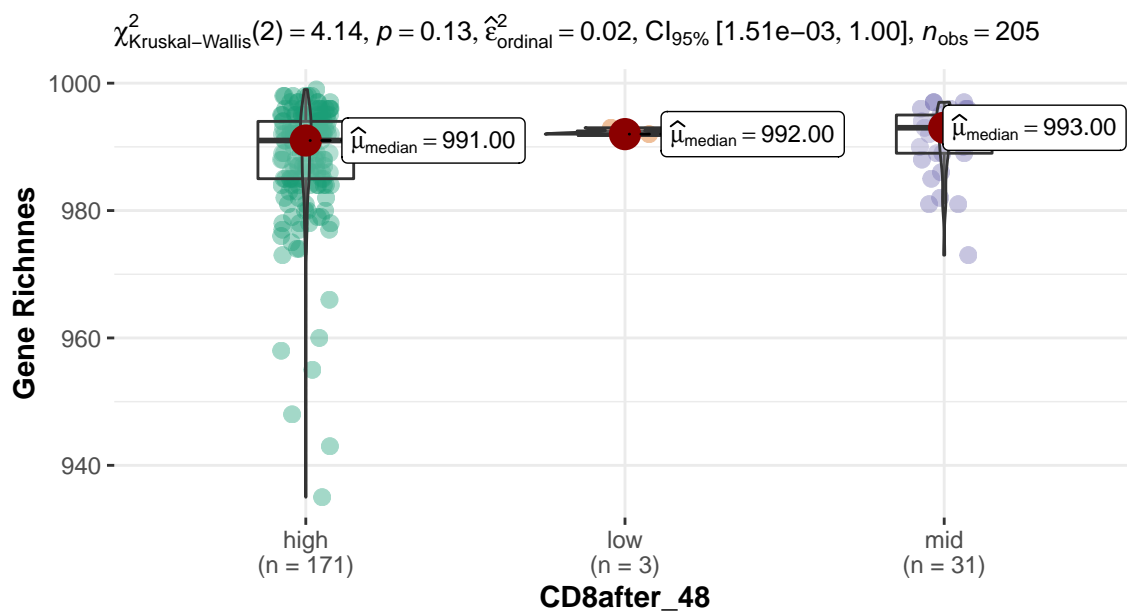
# Quality Control

The objective of this section is to check if there is any association between the total number of sequencing reads and the different levels of each categorical variable. Only the two most significant associations are shown.

# Alpha Diversity

Representation of the increase in the Gene Richness in relation to the total number of mapped reads for each of the samples. Dashed line represents the quantile at a probability of 2%.



## Gene Richness by categorical variables

This section shows the genetic wealth for each of the different levels of each of the categorical variables. Plots were produced using the ggstatsplot package. The upper text presents information on inferential statistics and the bottom one provides information about Bayesian hypothesis-testing and estimation.
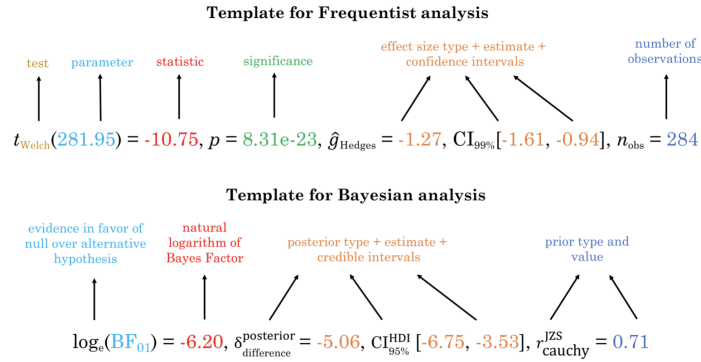


Figure 1: Stats structure

```
#> $group
```

$W_{\text{Mann–Whitney}} = 10274.00, p = 0.08, \hat{r}^{\text{rank}}_{\text{biserial}} = 0.12, \text{CI}_{95\%} \, [-0.01, 0.26], n_{\text{obs}} = 271$



$\widehat{\mu}_{\text{median}} = 992.00$

$\widehat{\mu}_{\text{median}} = 990.00$

DTG
(n = 144)

RTVr
(n = 127)

**group**

```
#>
#> $CD8after_48
```

$\chi^2_{\text{Kruskal–Wallis}}(2) = 4.14, p = 0.13, \hat{\varepsilon}^2_{\text{ordinal}} = 0.02, \text{CI}_{95\%} \, [1.51\text{e}{-}03, 1.00], n_{\text{obs}} = 205$



$\widehat{\mu}_{\text{median}} = 991.00$

$\widehat{\mu}_{\text{median}} = 992.00$

$\widehat{\mu}_{\text{median}} = 993.00$

high
(n = 171)

low
(n = 3)

mid
(n = 31)

**CD8after_48**

Pairwise test: **Dunn test**, Comparisons shown: **only significant**

# Gene Richness by numeric variables

In this section, Gene Richness was correlated with numeric metadata variables. Plots were produced using the ggstatsplot package. The upper text presents information on inferential statistics and the bottom one provides information about Bayesian hypothesis-testing and estimation.
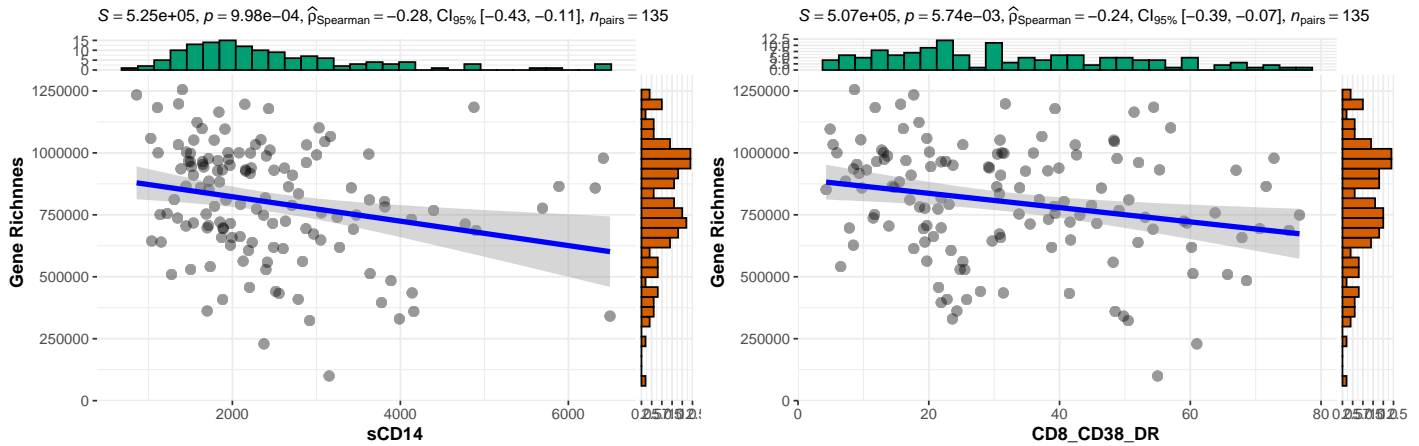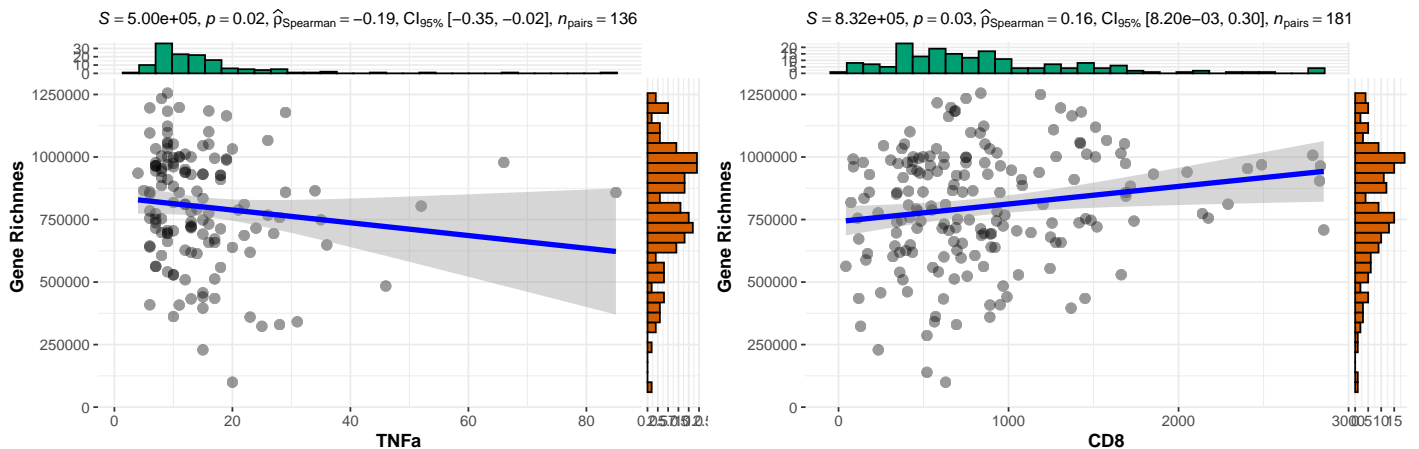
**Template for Frequentist analysis**

test   parameter   statistic   significance   effect size type + estimate + confidence intervals   number of observations

$t_{\text{Welch}}(281.95) = -10.75$, $p = 8.31\mathrm{e}{-23}$, $\widehat{g}_{\text{Hedges}} = -1.27$, $\text{CI}_{99\%}[-1.61, -0.94]$, $n_{\text{obs}} = 284$

**Template for Bayesian analysis**

evidence in favor of null over alternative hypothesis   natural logarithm of Bayes Factor   posterior type + estimate + credible intervals   prior type and value

$\log_{e}(\text{BF}_{01}) = -6.20$, $\delta_{\text{difference}}^{\text{posterior}} = -5.06$, $\text{CI}_{95\%}^{\text{HDI}}[-6.75, -3.53]$, $r_{\text{cauchy}}^{\text{JZS}} = 0.71$

Figure 2: Stats structure
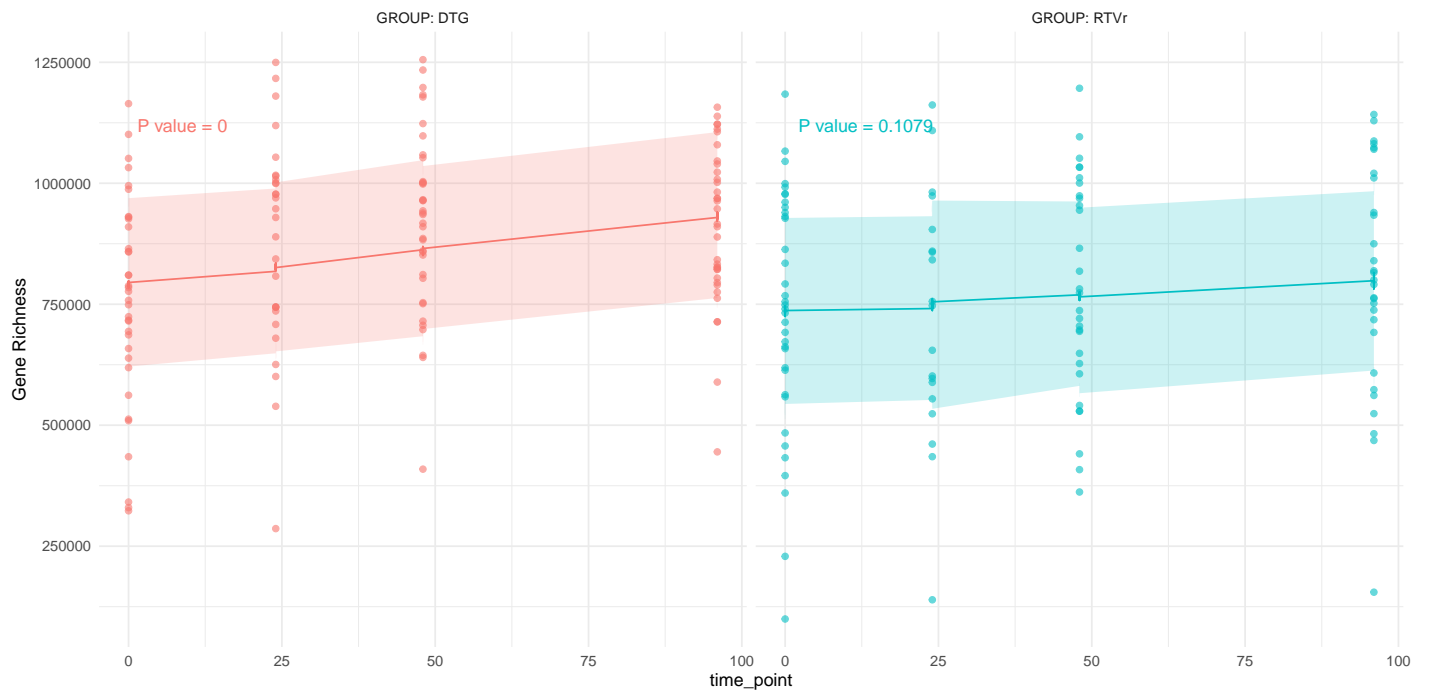
```
#> [[1]]
```



```
#>
#> [[2]]
```

## Gene Richness by longitudinal variable

In this section, Gene Richness was correlated with longitudinal metadata variables.

```
#> $time_point
```

# Taxa description based on Non-Metric Multidimendional Scaling (NMDS) ordering

Bar plot showing the relative abundance of different taxa in each sample. Sample axis order was determined using Non-metric Multidimensional Scaling (NMDS). The Bray method was used for distance calculations on Shotgun data, for 16s data Wunifrac distance was used. In the upper part of the graph are shown the distribution of the values of the variables selected from the metadata.

```
#> $Genus
```

Relative Abundance (%)

```
#>
#> $Species
```
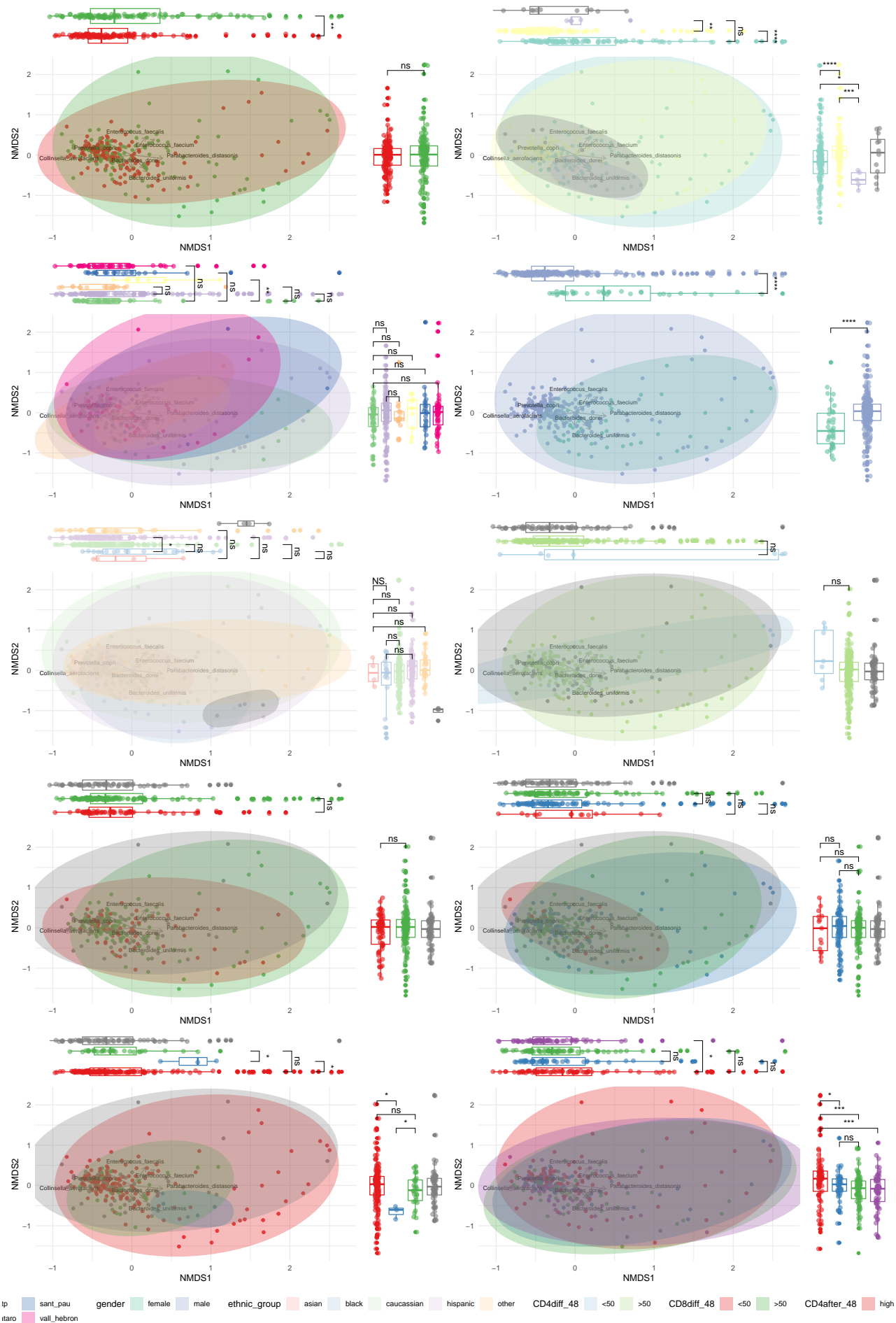
# Hierarchical Clustering Analysis

Heatmap showing the abundance of different taxas in each sample. Sample order was determined using ward.D2 hierarchical clustering. The categorical and numeric variables present in the mre object were used for sample annotation.

# Ordination Analysis (Non-metric multidimensional scaling).

## Ordination by categorical variables

Non-metric multidimensional scaling plot of categorical metadata variables and microbial community compositions. One plot for each categorical metadata. NMDS analysis within the vegan package of R software package based on dissimilarities calculated using the Bray-Curtis (Shotgun data) or WUnifrac (16s data) index of bacterial communities composition for the relative abundance of each OTU in relation to the categorical variables.
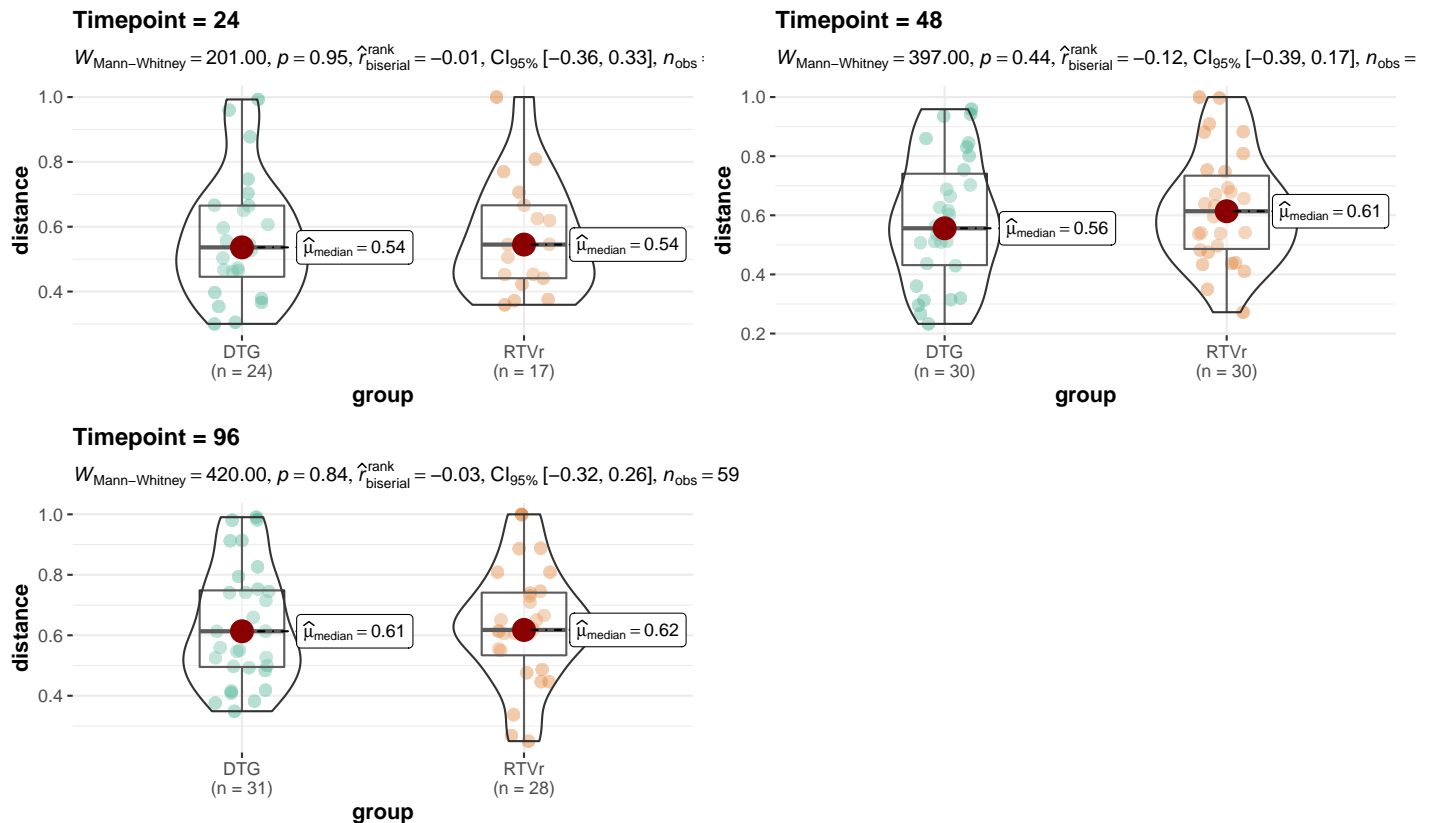
The statistics of the marginal boxplots were calculated using the ANOVA test. The Permutational Multivariate Analysis of Variance Using Distance Matrices (PERMANOVA) was computed using the `vegan::adonis()` function. The bottom table shows the results of the PERMANOVA for each categorical variable.
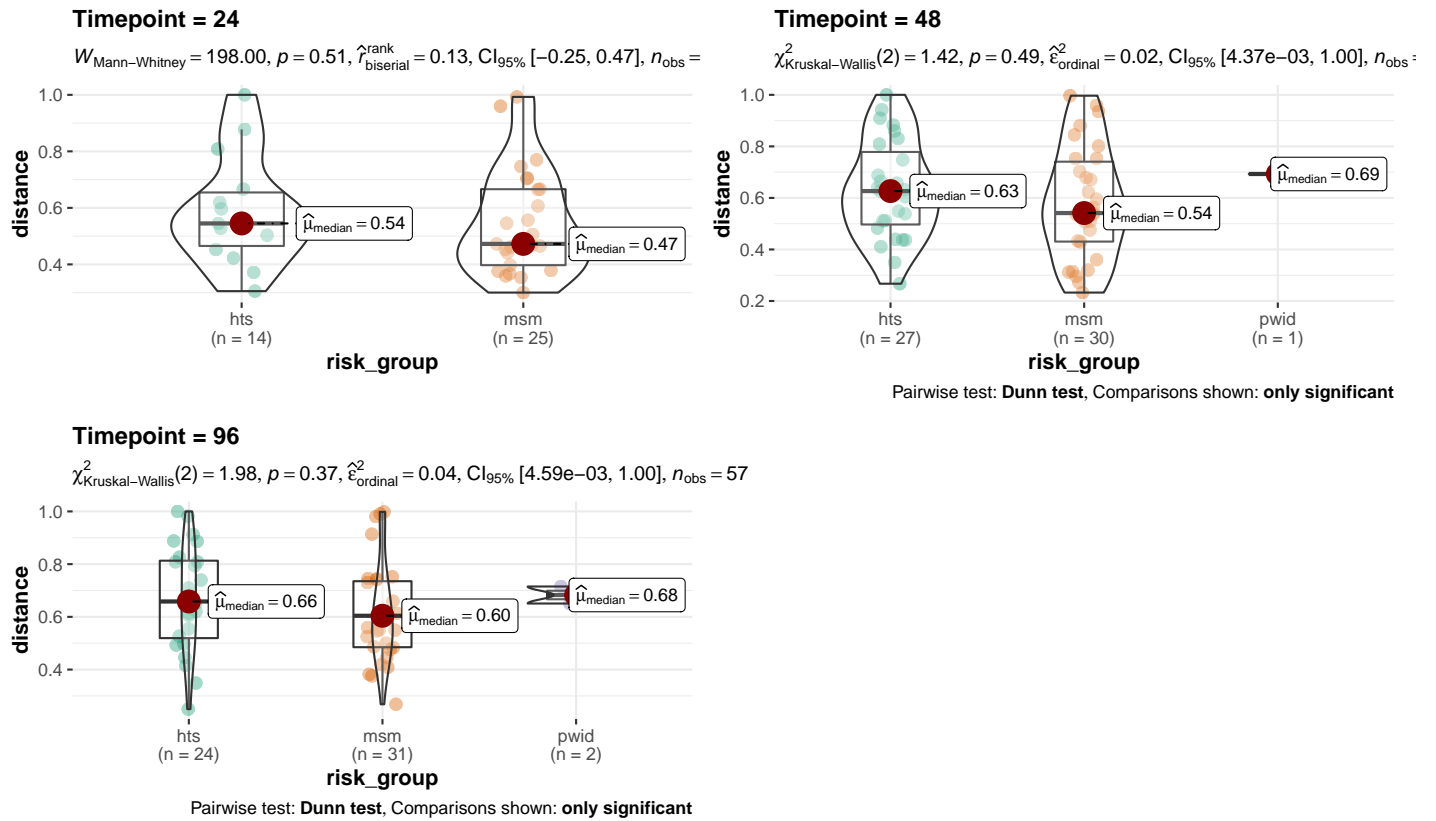
| id | Df | SumsOfSqs | MeanSqs | F.Model | R2 | Pr..F. |
|---|---|---|---|---|---|---|
| group | 1 | 0.9565124 | 0.9565124 | 3.409143 | 0.012514788 | 0.001 |
| risk_group | 2 | 2.8339413 | 1.4169706 | 5.165222 | 0.038788070 | 0.001 |
| center | 5 | 2.2744560 | 0.4548912 | 1.625573 | 0.029758458 | 0.002 |
| gender | 1 | 2.7318254 | 2.7318254 | 9.971147 | 0.035742574 | 0.001 |
| ethnic_group | 4 | 1.8604200 | 0.4651050 | 1.684405 | 0.025071370 | 0.001 |
| CD4diff_48 | 1 | 0.7980095 | 0.7980095 | 2.833267 | 0.013764863 | 0.004 |
| CD8diff_48 | 1 | 0.3653768 | 0.3653768 | 1.287498 | 0.006302382 | 0.170 |
| CD4after_48 | 2 | 0.7747815 | 0.3873907 | 1.368068 | 0.013364202 | 0.078 |
| CD8after_48 | 2 | 1.5293299 | 0.7646650 | 2.736507 | 0.026379405 | 0.002 |
| time_point | 3 | 1.2042928 | 0.4014309 | 1.424795 | 0.015756690 | 0.026 |
| record_id | 94 | 46.4281638 | 0.4939166 | 2.897412 | 0.607455395 | 0.001 |
| cluster | 1 | 8.2019180 | 8.2019180 | 32.337087 | 0.107312005 | 0.001 |

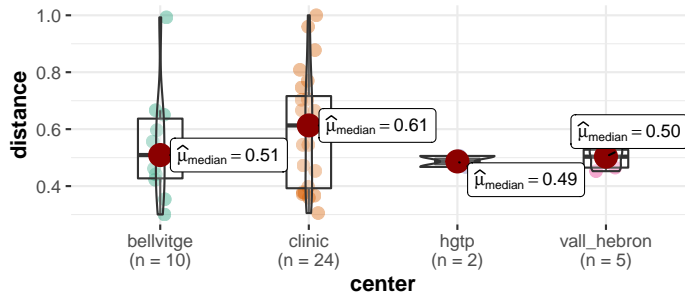**Changes in Beta diversity along longitudinal variable by each categorical variable**

```
#> $group
```

```
#>
#> $risk_group
```

**Timepoint = 24**

$W_{\text{Mann–Whitney}} = 198.00$, $p = 0.51$, $\hat{r}^{\text{rank}}_{\text{biserial}} = 0.13$, $CI_{95\%}$ [−0.25, 0.47], $n_{\text{obs}} =$

$\widehat{\mu}_{\text{median}} = 0.54$

$\widehat{\mu}_{\text{median}} = 0.47$

hts
(n = 14)

msm
(n = 25)

**risk_group**

**Timepoint = 48**

$\chi^2_{\text{Kruskal–Wallis}}(2) = 1.42$, $p = 0.49$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.02$, $CI_{95\%}$ [4.37e−03, 1.00], $n_{\text{obs}} =$

$\widehat{\mu}_{\text{median}} = 0.63$

$\widehat{\mu}_{\text{median}} = 0.54$

$\widehat{\mu}_{\text{median}} = 0.69$

hts
(n = 27)

msm
(n = 30)

pwid
(n = 1)

**risk_group**

Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 96**

$\chi^2_{\text{Kruskal–Wallis}}(2) = 1.98$, $p = 0.37$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.04$, $CI_{95\%}$ [4.59e−03, 1.00], $n_{\text{obs}} = 57$

$\widehat{\mu}_{\text{median}} = 0.66$

$\widehat{\mu}_{\text{median}} = 0.60$

$\widehat{\mu}_{\text{median}} = 0.68$

hts
(n = 24)

msm
(n = 31)

pwid
(n = 2)

**risk_group**

Pairwise test: **Dunn test**, Comparisons shown: **only significant**

```
#>
#> $center
```

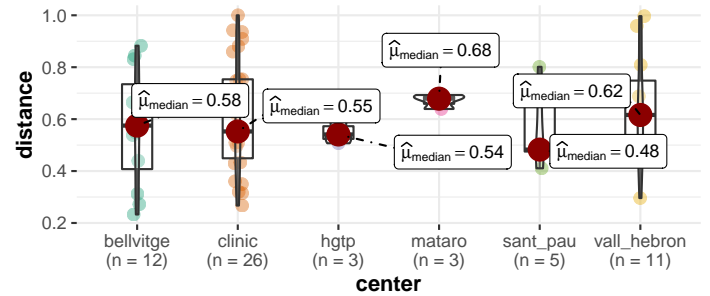**Timepoint = 24**

$\chi^2_{\text{Kruskal-Wallis}}(3) = 1.47$, $p = 0.69$, $\hat{\epsilon}^2_{\text{ordinal}} = 0.04$, CI$_{95\%}$ [5.39e−03, 1.00], $n_{\text{obs}}$



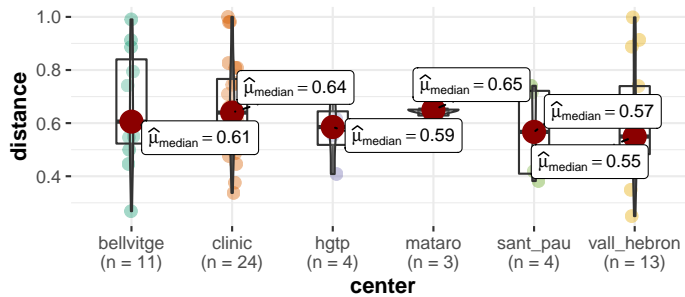Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 48**

$\chi^2_{\text{Kruskal-Wallis}}(5) = 1.76$, $p = 0.88$, $\hat{\epsilon}^2_{\text{ordinal}} = 0.03$, CI$_{95\%}$ [0.03, 1.00], $n_{\text{obs}} = 60$
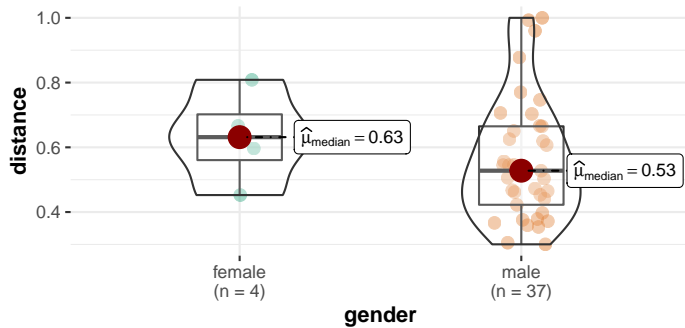


Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 96**

$\chi^2_{\text{Kruskal-Wallis}}(5) = 2.03$, $p = 0.84$, $\hat{\epsilon}^2_{\text{ordinal}} = 0.04$, CI$_{95\%}$ [0.01, 1.00], $n_{\text{obs}} = 59$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**
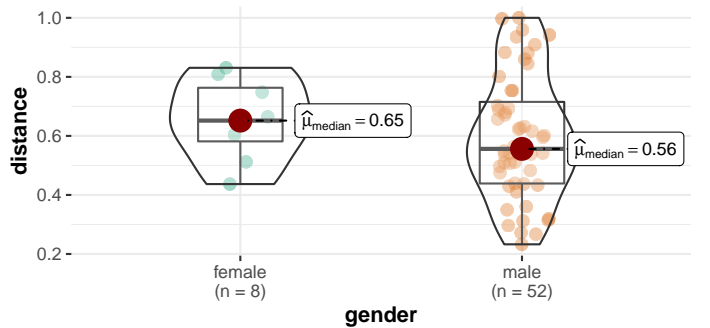
```
#>
#> $gender
```

**Timepoint = 24**

$W_{\text{Mann-Whitney}} = 96.00$, $p = 0.34$, $\hat{r}^{\text{rank}}_{\text{biserial}} = 0.30$, CI$_{95\%}$ [−0.29, 0.72], $n_{\text{obs}} = 4$
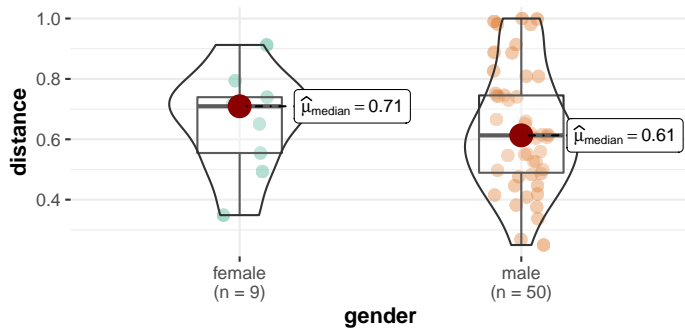


**Timepoint = 48**

$W_{\text{Mann-Whitney}} = 251.00$, $p = 0.36$, $\hat{r}^{\text{rank}}_{\text{biserial}} = 0.21$, CI$_{95\%}$ [−0.22, 0.57], $n_{\text{obs}} = 6$
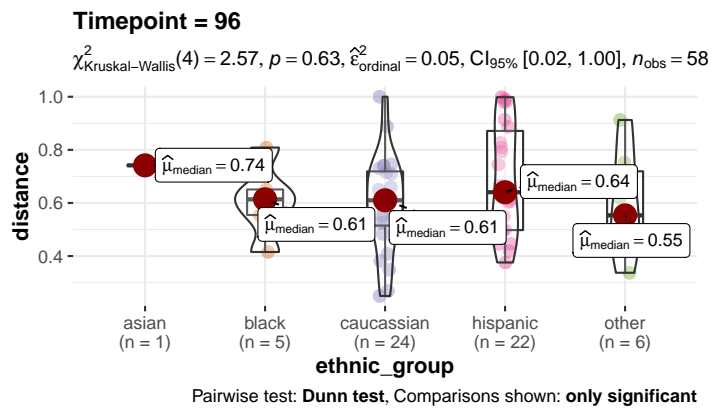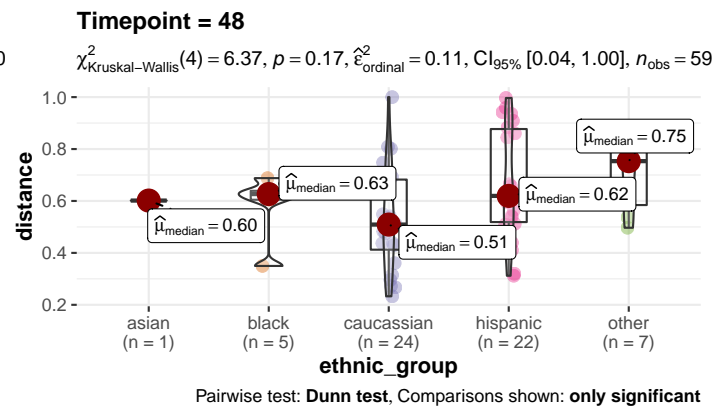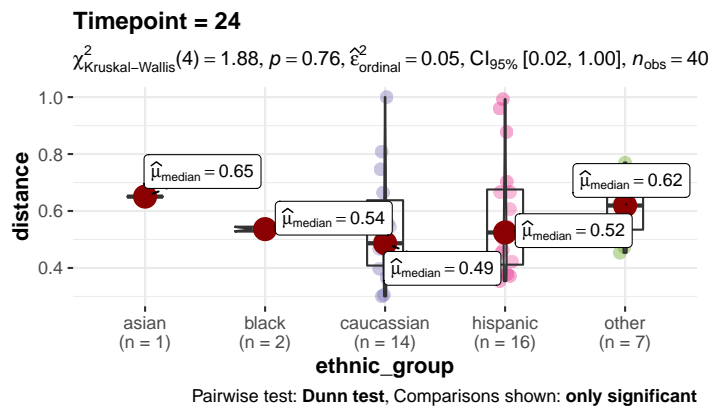


**Timepoint = 96**

$W_{\text{Mann-Whitney}} = 250.00$, $p = 0.61$, $\hat{r}^{\text{rank}}_{\text{biserial}} = 0.11$, CI$_{95\%}$ [−0.29, 0.48], $n_{\text{obs}} = 59$
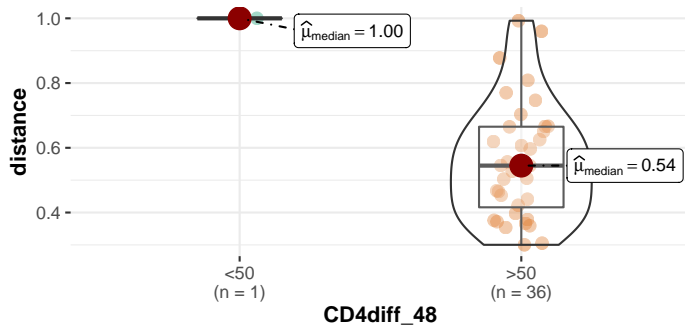


```
#>
```

```
#> $ethnic_group
```

**Timepoint = 24**

$\chi^2_{\text{Kruskal-Wallis}}(4) = 1.88$, $p = 0.76$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.05$, $\text{CI}_{95\%}$ [0.02, 1.00], $n_{\text{obs}} = 40$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 48**

$\chi^2_{\text{Kruskal-Wallis}}(4) = 6.37$, $p = 0.17$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.11$, $\text{CI}_{95\%}$ [0.04, 1.00], $n_{\text{obs}} = 59$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 96**

$\chi^2_{\text{Kruskal-Wallis}}(4) = 2.57$, $p = 0.63$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.05$, $\text{CI}_{95\%}$ [0.02, 1.00], $n_{\text{obs}} = 58$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**
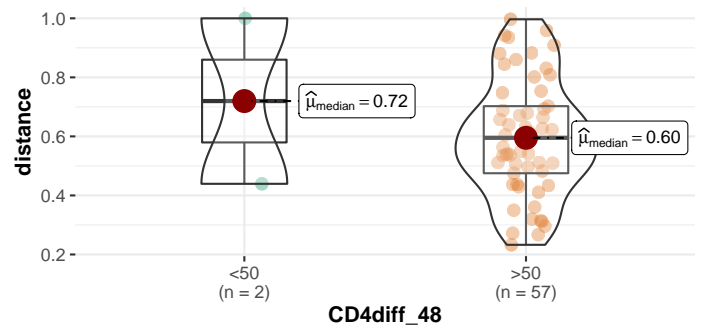
```
#>
#> $CD4diff_48
```

**Timepoint = 24**

$W_{\text{Mann-Whitney}} = 36.00$, $p = 0.10$, $\hat{r}_{\text{biserial}}^{\text{rank}} = 1.00$, $\text{CI}_{95\%}$ [1.00, 1.00], $n_{\text{obs}} = 37$
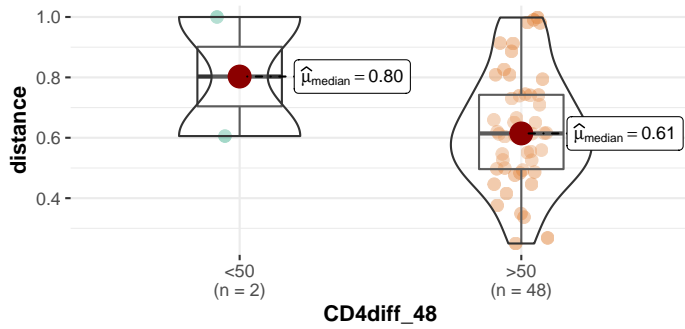


$\widehat{\mu}_{\text{median}} = 1.00$

$\widehat{\mu}_{\text{median}} = 0.54$

distance

<50
(n = 1)

>50
(n = 36)

**CD4diff_48**

**Timepoint = 48**

$W_{\text{Mann-Whitney}} = 71.00$, $p = 0.57$, $\hat{r}_{\text{biserial}}^{\text{rank}} = 0.25$, $\text{CI}_{95\%}$ [−0.52, 0.79], $n_{\text{obs}} = 59$



$\widehat{\mu}_{\text{median}} = 0.72$

$\widehat{\mu}_{\text{median}} = 0.60$

distance

<50
(n = 2)

>50
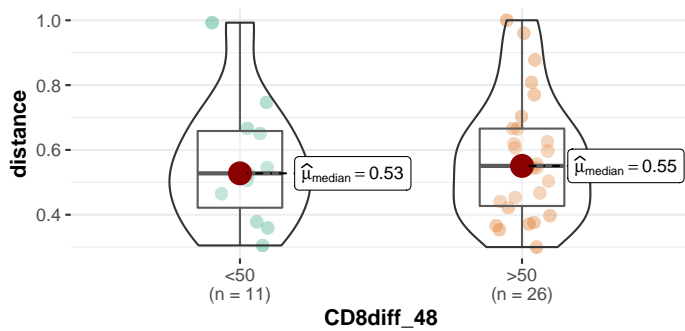(n = 57)

**CD4diff_48**

**Timepoint = 96**

$W_{\text{Mann-Whitney}} = 69.00$, $p = 0.31$, $\hat{r}_{\text{biserial}}^{\text{rank}} = 0.44$, $\text{CI}_{95\%}$ [−0.34, 0.86], $n_{\text{obs}} = 50$



$\widehat{\mu}_{\text{median}} = 0.80$

$\widehat{\mu}_{\text{median}} = 0.61$

distance

<50
(n = 2)

>50
(n = 48)

**CD4diff_48**
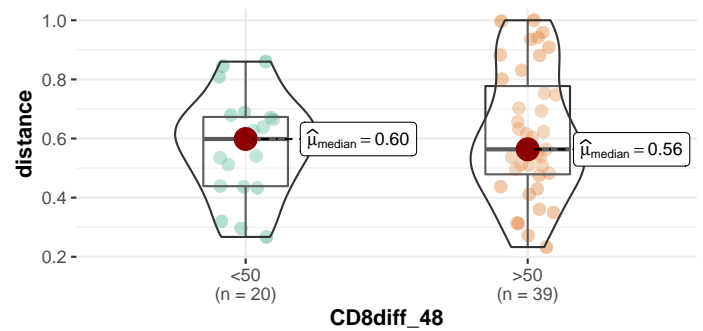
```
#>
#> $CD8diff_48
```

**Timepoint = 24**

$W_{\text{Mann-Whitney}} = 136.00$, $p = 0.83$, $\hat{r}_{\text{biserial}}^{\text{rank}} = -0.05$, $\text{CI}_{95\%}$ [−0.43, 0.35], $n_{\text{obs}}$
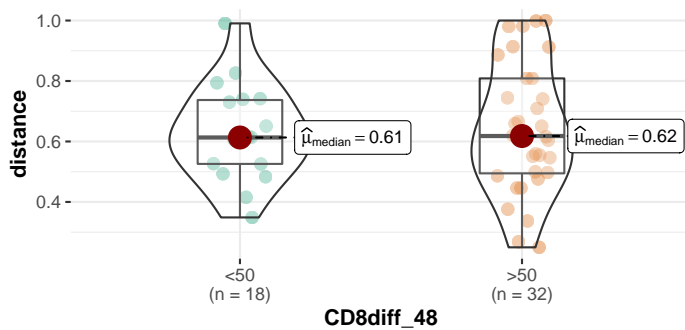


$\widehat{\mu}_{\text{median}} = 0.53$

$\widehat{\mu}_{\text{median}} = 0.55$

distance

<50
(n = 11)

>50
(n = 26)

**CD8diff_48**

**Timepoint = 48**

$W_{\text{Mann-Whitney}} = 358.00$, $p = 0.61$, $\hat{r}_{\text{biserial}}^{\text{rank}} = -0.08$, $\text{CI}_{95\%}$ [−0.38, 0.23], $n_{\text{obs}} =$



$\widehat{\mu}_{\text{median}} = 0.60$

$\widehat{\mu}_{\text{median}} = 0.56$

distance
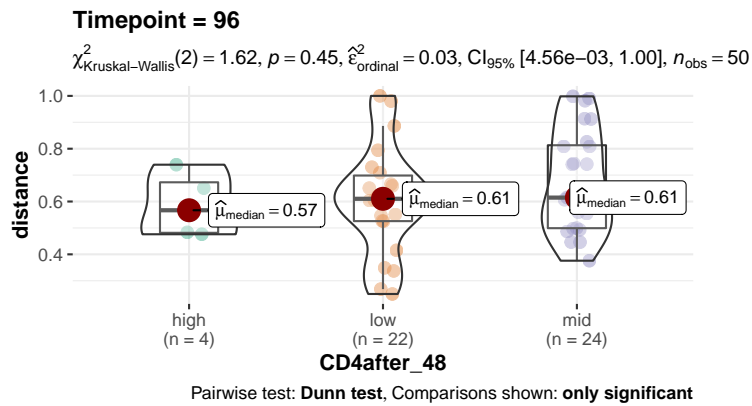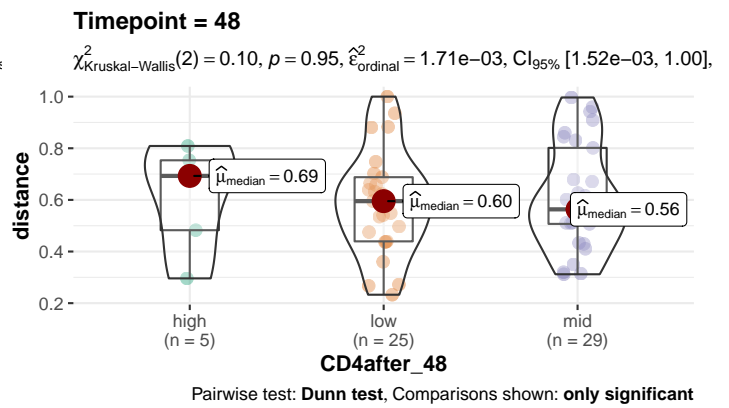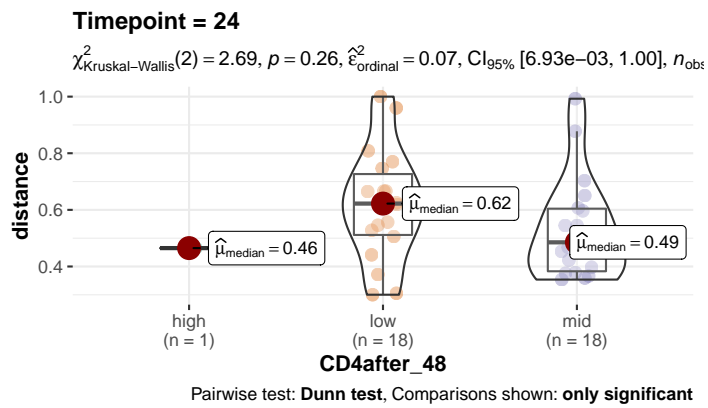
<50
(n = 20)

>50
(n = 39)

**CD8diff_48**

**Timepoint = 96**

$W_{\text{Mann-Whitney}} = 277.00$, $p = 0.83$, $\hat{r}_{\text{biserial}}^{\text{rank}} = -0.04$, $\text{CI}_{95\%}$ [−0.36, 0.29], $n_{\text{obs}} = 50$



$\widehat{\mu}_{\text{median}} = 0.61$

$\widehat{\mu}_{\text{median}} = 0.62$

distance

<50
(n = 18)

>50
(n = 32)

**CD8diff_48**

```
#>
```

```
#> $CD4after_48
```

**Timepoint = 24**

$\chi^2_{\text{Kruskal-Wallis}}(2) = 2.69$, $p = 0.26$, $\widehat{\varepsilon}^2_{\text{ordinal}} = 0.07$, $\text{CI}_{95\%}$ [6.93e–03, 1.00], $n_{\text{obs}}$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 48**

$\chi^2_{\text{Kruskal-Wallis}}(2) = 0.10$, $p = 0.95$, $\widehat{\varepsilon}^2_{\text{ordinal}} = 1.71\text{e–03}$, $\text{CI}_{95\%}$ [1.52e–03, 1.00],



Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 96**

$\chi^2_{\text{Kruskal-Wallis}}(2) = 1.62$, $p = 0.45$, $\widehat{\varepsilon}^2_{\text{ordinal}} = 0.03$, $\text{CI}_{95\%}$ [4.56e–03, 1.00], $n_{\text{obs}} = 50$



Pairwise test: **Dunn test**, Comparisons shown: **only significant**

```
#>
#> $CD8after_48
```

**Timepoint = 24**

$W_{\text{Mann–Whitney}} = 75.00$, $p = 0.47$, $\hat{r}^{\text{rank}}_{\text{biserial}} = -0.19$, $\text{CI}_{95\%}$ [−0.61, 0.31], $n_{\text{obs}} =$



$\hat{\mu}_{\text{median}} = 0.54$

$\hat{\mu}_{\text{median}} = 0.58$

high
(n = 31)

mid
(n = 6)

**CD8after_48**

**Timepoint = 48**

$\chi^2_{\text{Kruskal–Wallis}}(2) = 1.14$, $p = 0.57$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.02$, $\text{CI}_{95\%}$ [0.01, 1.00], $n_{\text{obs}} = 59$



$\hat{\mu}_{\text{median}} = 0.81$

$\hat{\mu}_{\text{median}} = 0.55$

$\hat{\mu}_{\text{median}} = 0.60$

high
(n = 49)

low
(n = 1)

mid
(n = 9)

**CD8after_48**

Pairwise test: **Dunn test**, Comparisons shown: **only significant**

**Timepoint = 96**

$\chi^2_{\text{Kruskal–Wallis}}(2) = 0.84$, $p = 0.66$, $\hat{\varepsilon}^2_{\text{ordinal}} = 0.02$, $\text{CI}_{95\%}$ [2.85e−03, 1.00], $n_{\text{obs}} = 50$



$\hat{\mu}_{\text{median}} = 0.74$

$\hat{\mu}_{\text{median}} = 0.61$

$\hat{\mu}_{\text{median}} = 0.65$

high
(n = 42)

low
(n = 1)

mid
(n = 7)

**CD8after_48**

Pairwise test: **Dunn test**, Comparisons shown: **only significant**

## Ordination by numeric variables

Correlation analysis between NMDS components (NMDS1 and NMDS2) and numeric variables presents in the mre object. Plots were produced using the ggstatsplot package. The upper text presents information on inferential statistics and the bottom one provides information about Bayesian hypothesis-testing and estimation.
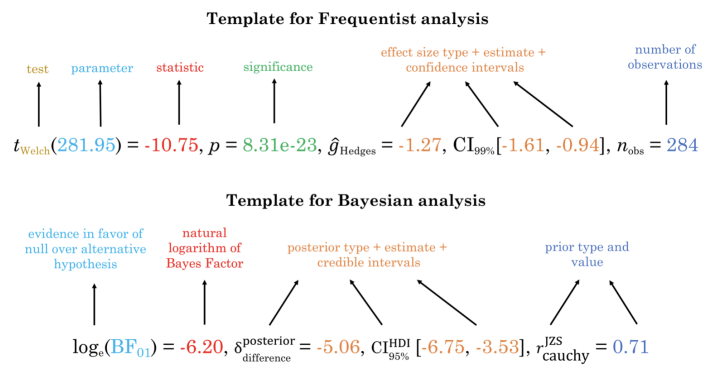
**Template for Frequentist analysis**

test    parameter    statistic    significance    effect size type + estimate + confidence intervals    number of observations

$t_{\text{Welch}}(281.95) = -10.75$, $p = 8.31\text{e-}23$, $\hat{g}_{\text{Hedges}} = -1.27$, $\text{CI}_{99\%}[-1.61, -0.94]$, $n_{\text{obs}} = 284$

**Template for Bayesian analysis**

evidence in favor of null over alternative hypothesis    natural logarithm of Bayes Factor    posterior type + estimate + credible intervals    prior type and value

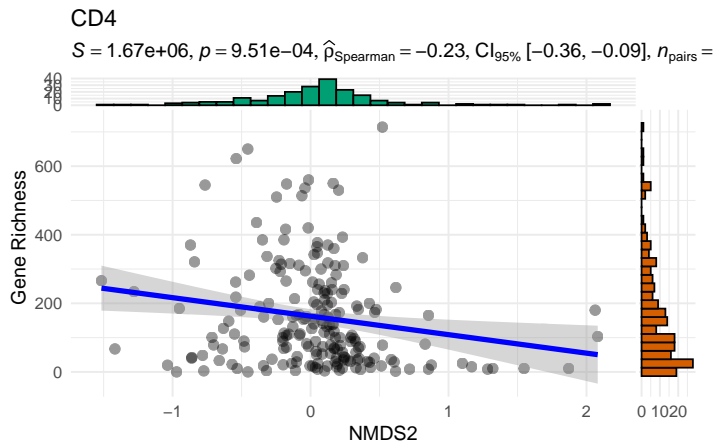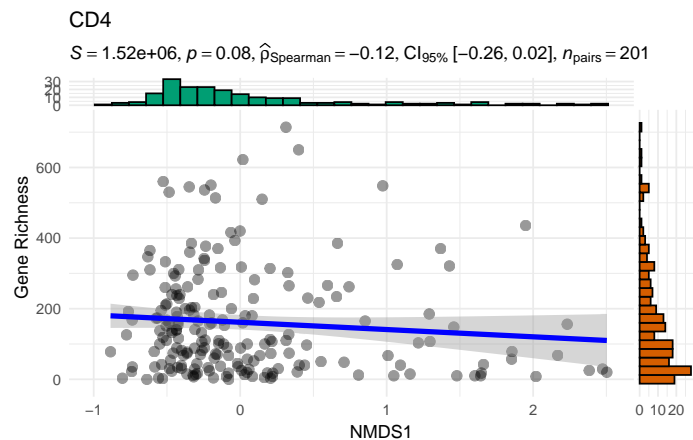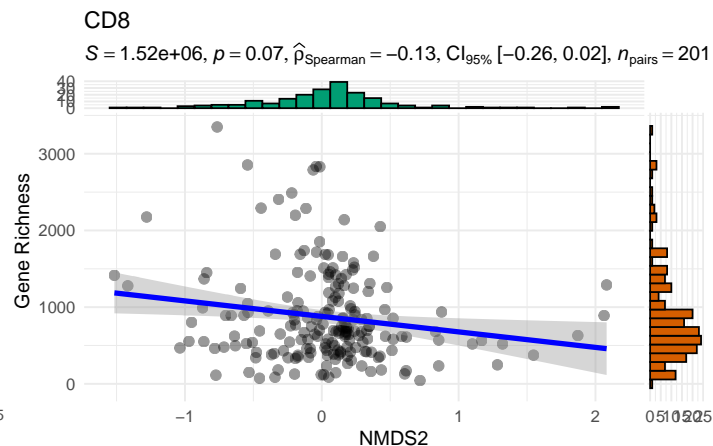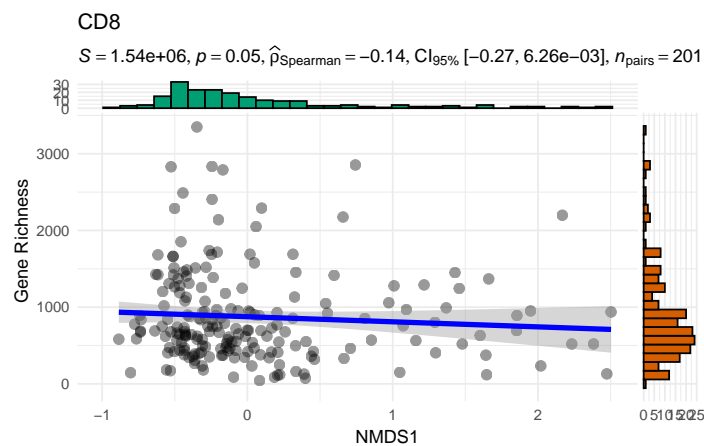$\log_e(\text{BF}_{01}) = -6.20$, $\delta^{\text{posterior}}_{\text{difference}} = -5.06$, $\text{CI}^{\text{HDI}}_{95\%}$ [-6.75, -3.53], $r^{\text{JZS}}_{\text{cauchy}} = 0.71$

Figure 3: Stats structure

```
#>  $CD4
```

**CD4**

$S = 1.52\mathrm{e}{+}06$, $p = 0.08$, $\widehat{\rho}_{\mathrm{Spearman}} = -0.12$, $\mathrm{CI}_{95\%}\ [-0.26, 0.02]$, $n_{\mathrm{pairs}} = 201$

**CD4**

$S = 1.67\mathrm{e}{+}06$, $p = 9.51\mathrm{e}{-}04$, $\widehat{\rho}_{\mathrm{Spearman}} = -0.23$, $\mathrm{CI}_{95\%}\ [-0.36, -0.09]$, $n_{\mathrm{pairs}} =$

```
#>
#> $CD8
```

**CD8**

$S = 1.54\mathrm{e}{+}06$, $p = 0.05$, $\widehat{\rho}_{\mathrm{Spearman}} = -0.14$, $\mathrm{CI}_{95\%}\ [-0.27, 6.26\mathrm{e}{-}03]$, $n_{\mathrm{pairs}} = 201$

**CD8**

$S = 1.52\mathrm{e}{+}06$, $p = 0.07$, $\widehat{\rho}_{\mathrm{Spearman}} = -0.13$, $\mathrm{CI}_{95\%}\ [-0.26, 0.02]$, $n_{\mathrm{pairs}} = 201$

```
#>
#> $CD8_CD38_DR
```

**CD8_CD38_DR**

$S = 5.03\mathrm{e}{+}05$, $p = 0.13$, $\widehat{\rho}_{\mathrm{Spearman}} = 0.12$, $\mathrm{CI}_{95\%}\ [-0.04, 0.28]$, $n_{\mathrm{pairs}} = 151$

**CD8_CD38_DR**

$S = 4.33\mathrm{e}{+}05$, $p = 2.36\mathrm{e}{-}03$, $\widehat{\rho}_{\mathrm{Spearman}} = 0.25$, $\mathrm{CI}_{95\%}\ [0.08, 0.39]$, $n_{\mathrm{pairs}} = 151$

```
#>
#> $CRP
```

## CRP

$S = 3.85e{+}05$, $p = 0.83$, $\widehat{\rho}_{\text{Spearman}} = 0.02$, $\text{CI}_{95\%}$ [$-0.16$, $0.19$], $n_{\text{pairs}} = 133$

## CRP

$S = 2.50e{+}05$, $p = 1.77e{-}05$, $\widehat{\rho}_{\text{Spearman}} = 0.36$, $\text{CI}_{95\%}$ [$0.20$, $0.51$], $n_{\text{pairs}} = 133$

```
#>
#> $IL6
```

## IL6

$S = 3.28e{+}05$, $p = 0.48$, $\widehat{\rho}_{\text{Spearman}} = 0.06$, $\text{CI}_{95\%}$ [$-0.12$, $0.24$], $n_{\text{pairs}} = 128$

## IL6

$S = 2.88e{+}05$, $p = 0.05$, $\widehat{\rho}_{\text{Spearman}} = 0.18$, $\text{CI}_{95\%}$ [$-1.71e{-}03$, $0.34$], $n_{\text{pairs}} = 12$

```
#>
#> $TNFa
```

## TNFa

$S = 5.17e{+}05$, $p = 0.16$, $\widehat{\rho}_{\text{Spearman}} = 0.12$, $\text{CI}_{95\%}$ [$-0.05$, $0.27$], $n_{\text{pairs}} = 152$

## TNFa

$S = 4.25e{+}05$, $p = 6.58e{-}04$, $\widehat{\rho}_{\text{Spearman}} = 0.27$, $\text{CI}_{95\%}$ [$0.11$, $0.42$], $n_{\text{pairs}} = 152$

```
#>
#> $sCD14
```

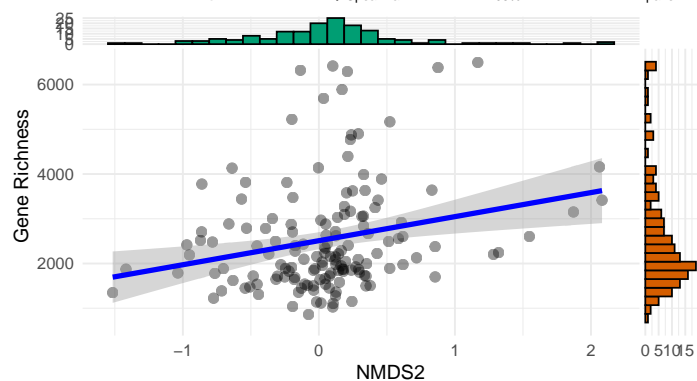sCD14

$S = 4.73e+05$, $p = 0.03$, $\widehat{\rho}_{Spearman} = 0.17$, $CI_{95\%}$ [0.01, 0.33], $n_{pairs} = 151$

sCD14

$S = 4.21e+05$, $p = 9.78e{-}04$, $\widehat{\rho}_{Spearman} = 0.27$, $CI_{95\%}$ [0.11, 0.41], $n_{pairs} = 15$
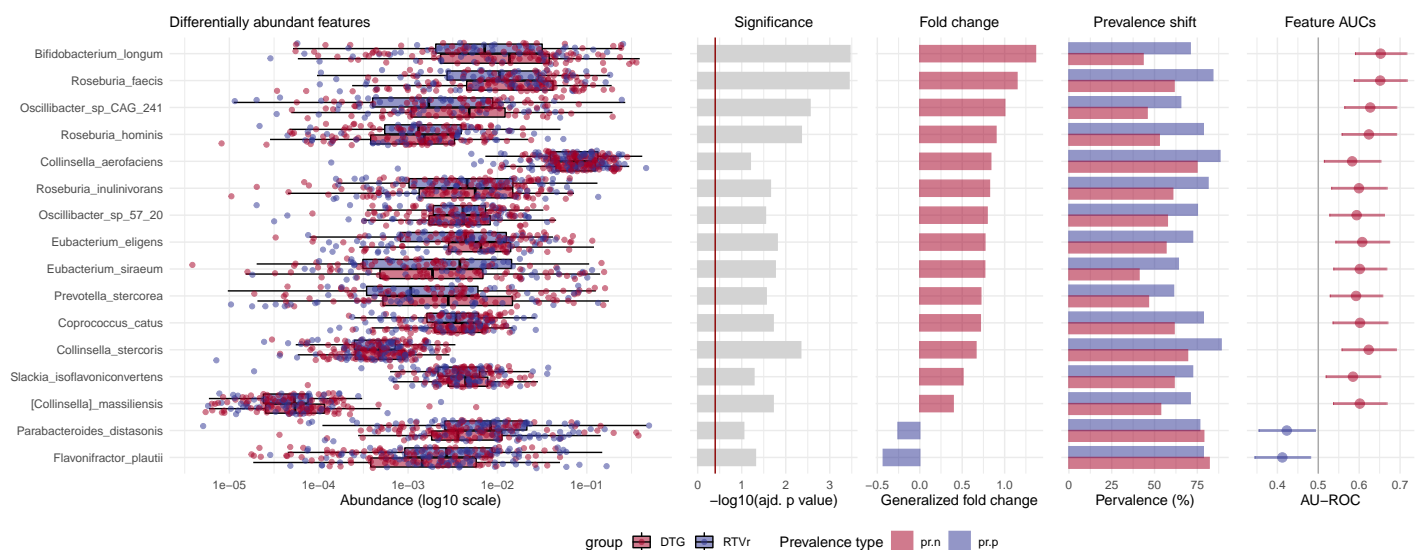
# Differential abundance

## Statistical Inference of Associations between Microbial Communities And host phenoTypes (SIAMCAT).

Detection of changes in community composition that are associated with metadata variables using LASSO logistic regression modeling. For this purpose, the abundance matrix was relativized and expressed in times ones. Feature selection was performed with "prevalence" method removing those features with low prevalence across samples (relative abundance cutoff value set as 0.5). Feature normalization was performed using centred log-ratio ("log.clr") transformation. Finally, for visualization purposes, only those associations with an fdr < 0.05 (Default value) were considered.

```
#> $group
```



```
#>
#> $gender
```