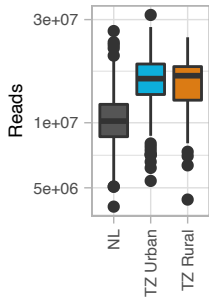
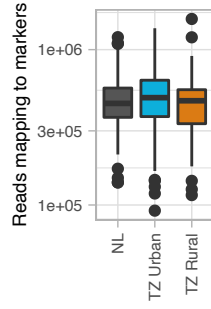
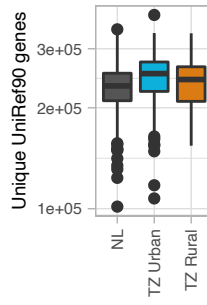
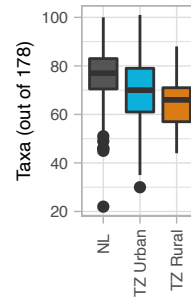
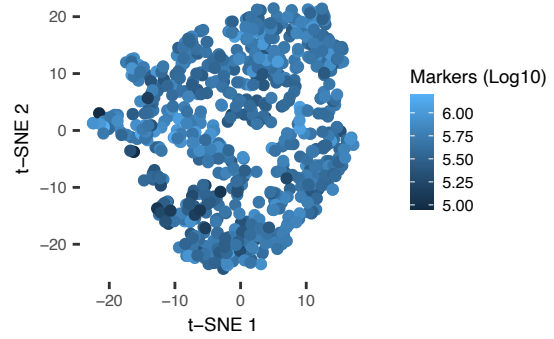
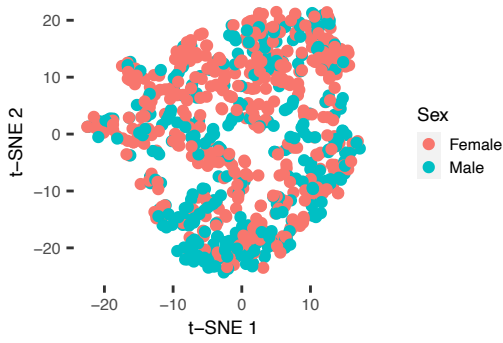
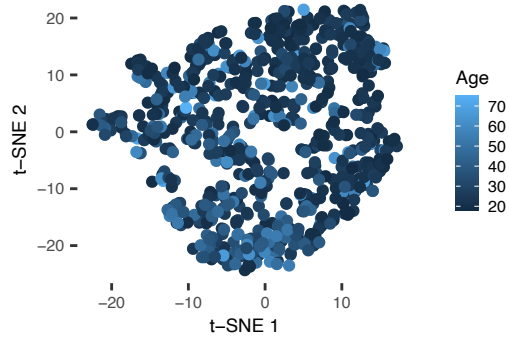
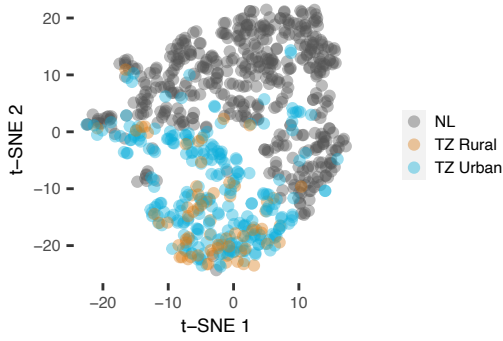
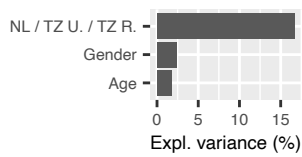
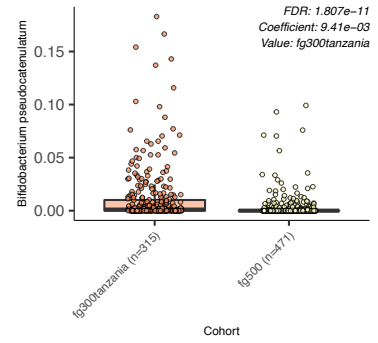
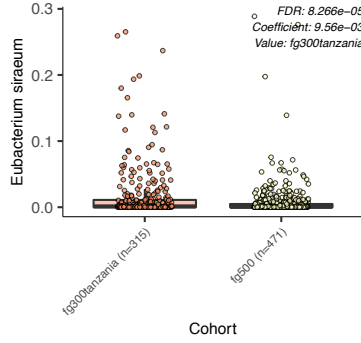
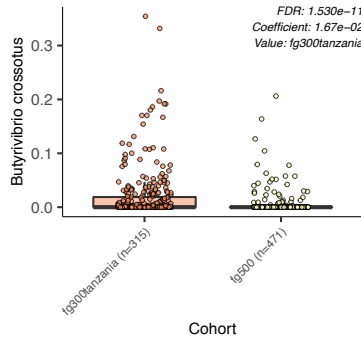
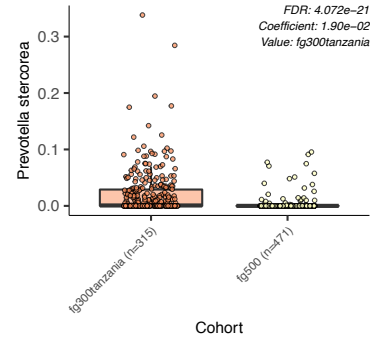
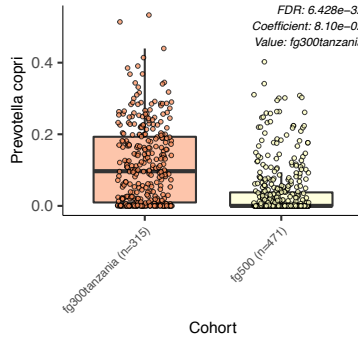
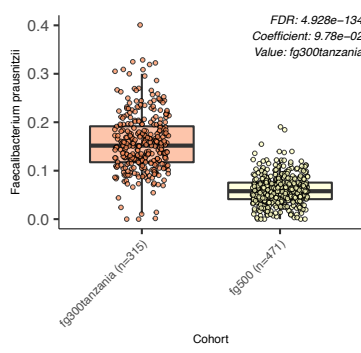
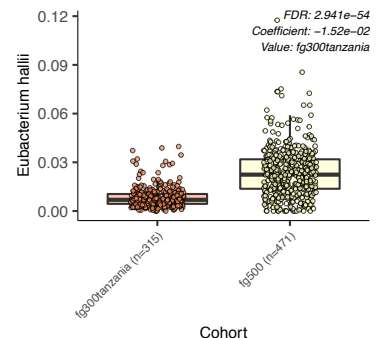
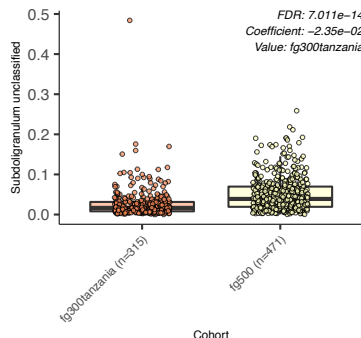
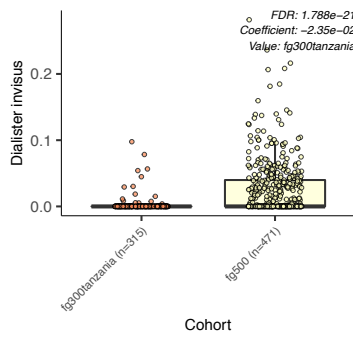
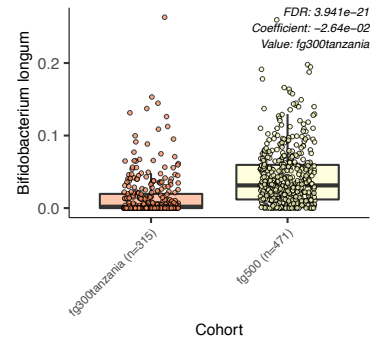
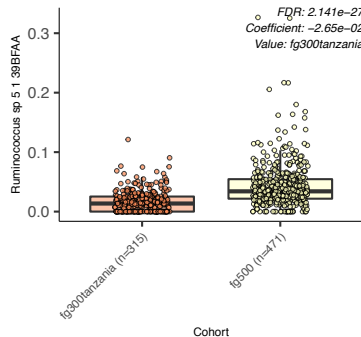
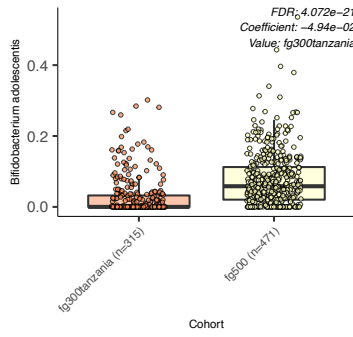


Gut microbiome-mediated metabolism effects on immunity in rural and urban populations

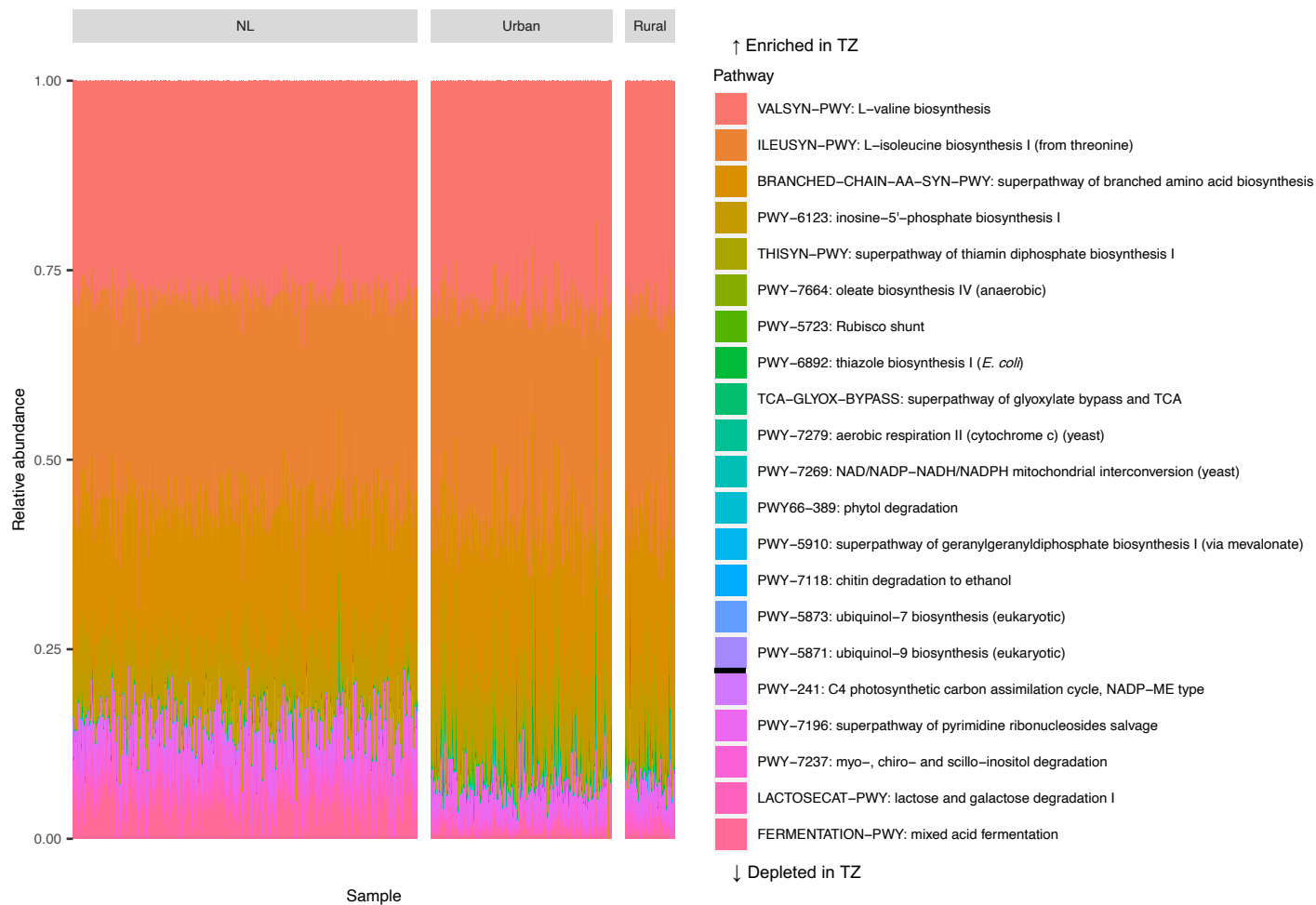
Supplementary Information

a**b****c****d****e****f**

Supplementary Figure 1. Quality control and metagenomic profiling statistics between NL (n=471 biologically independent samples) TZ rural (n = 68) and TZ urban (n = 247) cohorts, with samples yielding at least 4 million reads. a) Number of reads after adapter and host DNA removal. The NL cohorts uses Illumina HiSeq 2000 platform (101 bp paired-end reads) while the TZ cohort was analyzed on Illumina NovaSeq 6000 (151bp paired-end). **b)** Reads mapping to marker genes in Metaphlan taxonomic profiling. **c)** Number of identified UniRef90 protein families. **d)** Number of detected species per subject among 176 species common to the two cohorts. Box plots show the median (center), first (Q1) and third quartile (Q3, box). Lower and upper whiskers represent $Q1 + 1.5 \text{ IQR}$ and $Q3 + \text{IQR}$, respectively. IQR, inter-quartile range. Data points outside whiskers' bounds are shown individually. **e)** t-SNE projection of microbial composition labeled by residency, sex, age and number of marker genes from the Metaphlan reference. **f)** PERMANOVA analysis of microbial compositions, represented by relative abundances of 176 common taxa. Bars show explained variance for residency, sex and age, all of which were significant at $P < 0.001$.

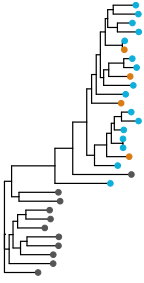
a**b**

Supplementary Figure 2. Differences with between the Tanzanian (n = 315 biologically independent samples) and 500FG cohort (n = 471 biologically independent samples). The differential test is performed using Maaslin2 with cohort as a fixed effect and log-transformed relative abundance data. Out of total 104 significant species (FDR < 5%), the top six **a)** enriched and **b)** depleted species sorted by regression coefficient are displayed. Box plots show the median (center), first (Q1) and third quartile (Q3, box). Lower and upper whiskers represent $Q1 + 1.5 \times IQR$ and $Q3 + IQR$, respectively. IQR, inter-quartile range. Data points outside whiskers' bounds are shown individually.

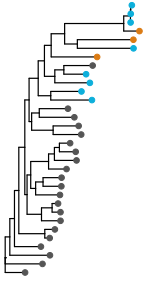


Supplementary Figure 3. Functional diversity in metagenomic pathways between cohorts. Reads mapping to MetaCyc pathways were quantified using pathway copy numbers using Humann2. The differential relative abundance of pathways was analyzed using Maaslin2, with residency (NL, TZ urban, TZ rural), age and sex as fixed effects and cohort (NL or TZ) as a random effect. Differentially abundant pathways were selected subject to $FDR < 0.05$ and presence in at least 20 samples. The bar plots show relative abundance within 21 selected pathways for each subject. Black line in the legend separates pathways enriched in TZ (16) and those enriched in NL (5).

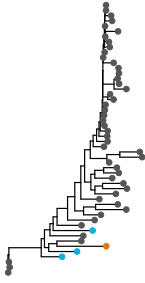
Lactobacillus ruminis



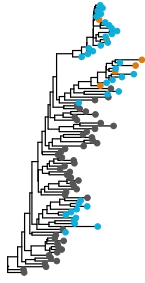
Clostridium sp L2 50



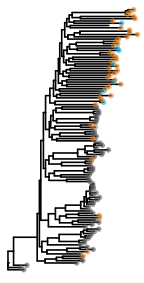
Bacteroides uniformis



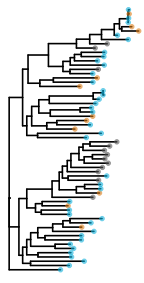
Ruminococcus lactaris



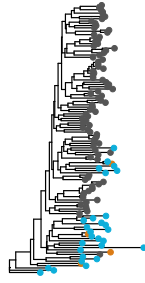
Ruminococcus bromii



Butyrivibrio crossotus



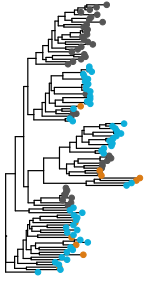
Coproccoccus sp ART55 1



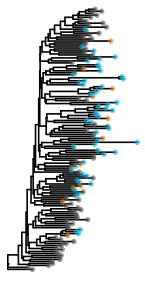
Phascolarctobacterium



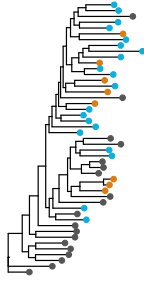
Eubacterium siraeum



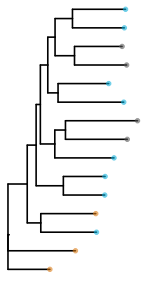
Roseburia intestinalis



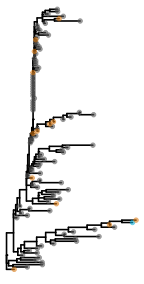
Streptococcus salivarius



Roseburia hominis



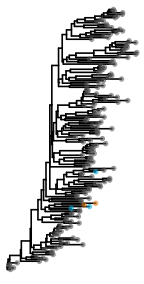
Ruminococcus torques



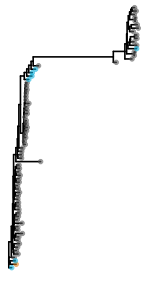
Akkermansia



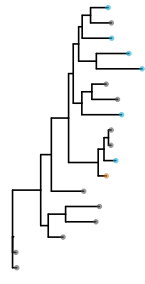
Eubacterium hallii



Methanobrevibacter smithii



Bacteroides stercoris



Ruminococcus obeum



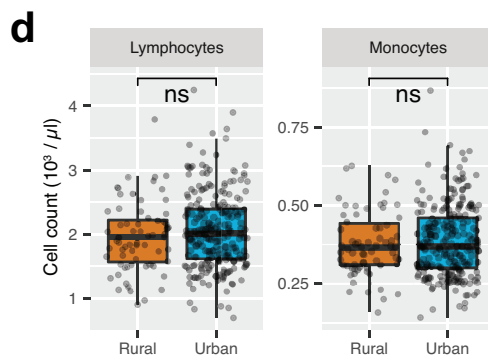
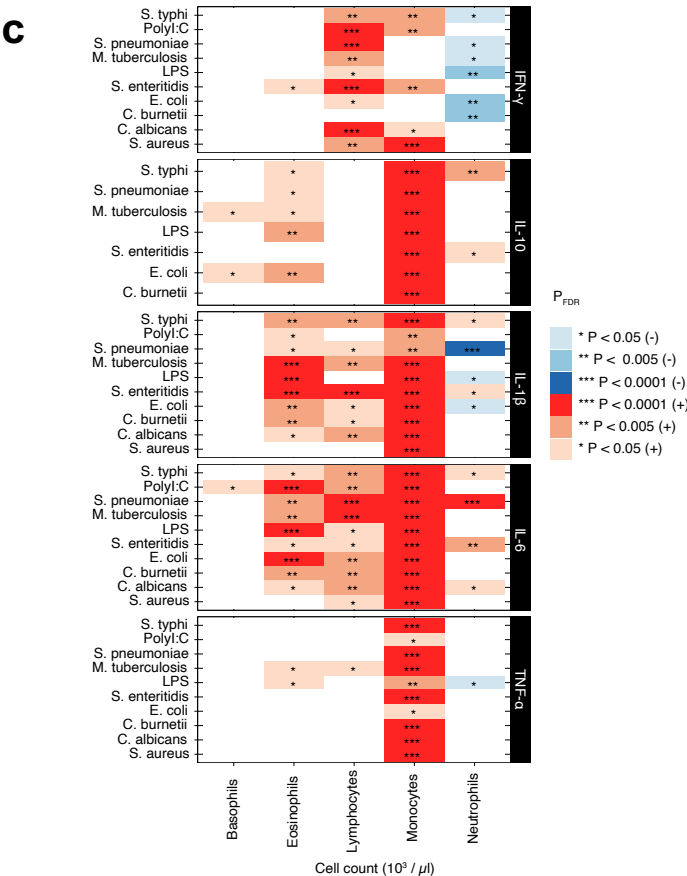
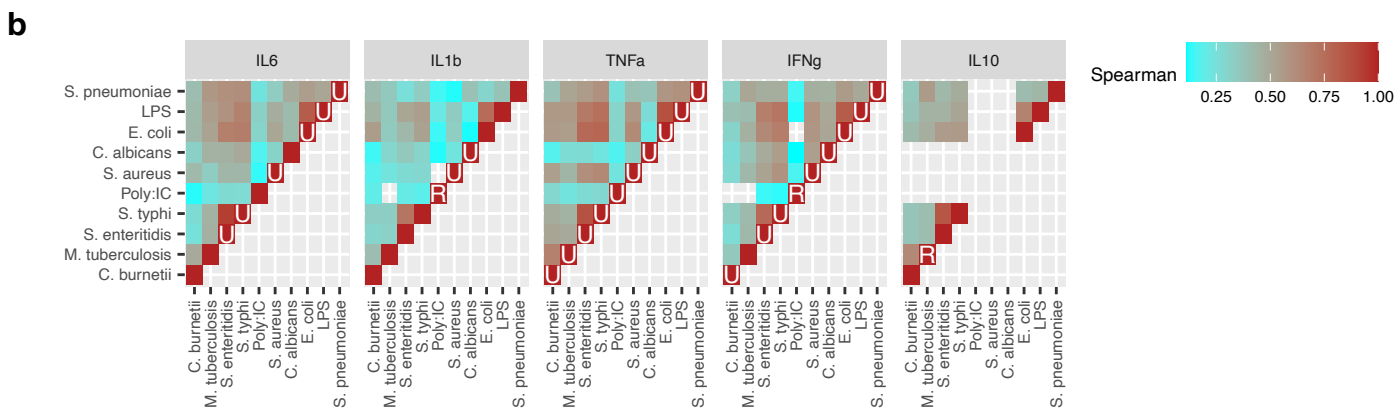
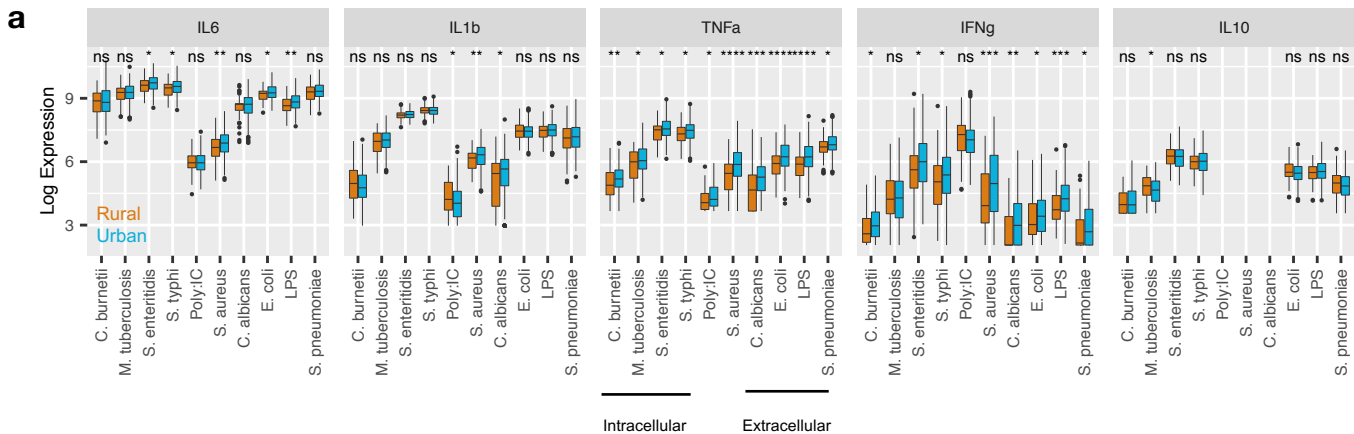
Dorea longicatena



Residency

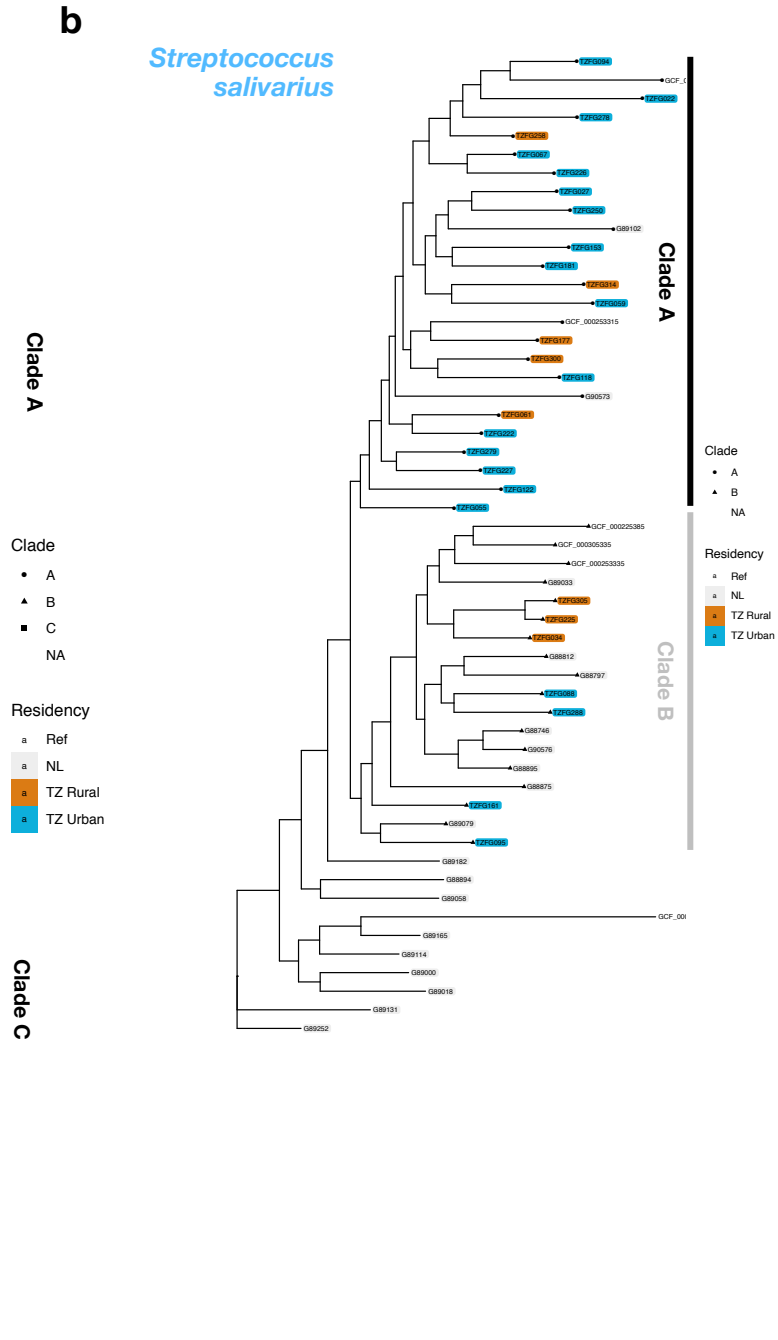
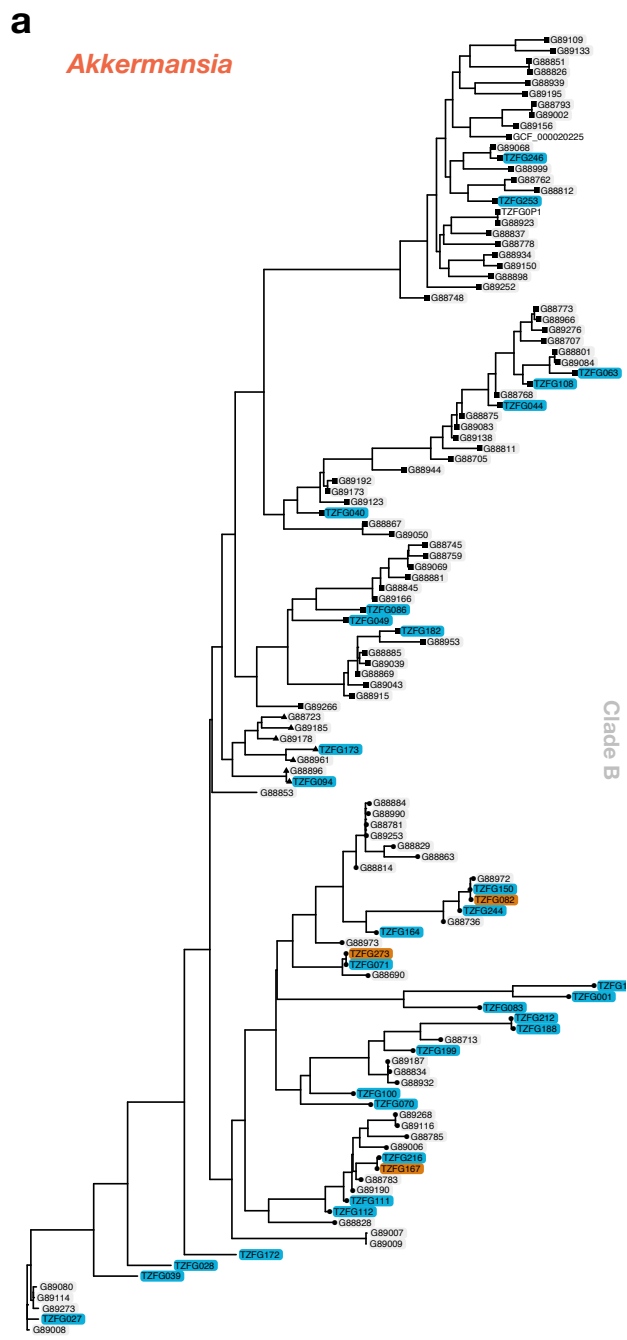
- NL
- TZ Rural
- TZ Urban

Supplementary Figure 4. Strain analysis of NL, TZ Urban and TZ Rural cohorts with Strainphlan. Significance strain divergence driven by residency for each taxon was estimated using the PERMANOVA analysis using the underlying phylogenetic distances matrices ($P < 0.05$, minimum samples = 10, presence in all three residency areas).



Supplementary Figure 5. Cytokine expression in rural and urban TZ cohorts. **a)** Box plots of cytokine expression levels upon ex vivo stimulations for TZ rural (n=68) and TZ urban (n=239) cohorts. The difference between urban and rural sample levels is assessed by two-sided Wilcoxon Rank Sum test (*P <= 0.05, **P <= 0.01, ***P <= 0.001, ****P <= 0.0001). Box plots show the median (center), first (Q1) and third quartile (Q3, box). Lower and upper whiskers represent $Q1 + 1.5 \text{ IQR}$ and $Q3 + \text{IQR}$, respectively. IQR, inter-quartile range. Data points outside whiskers' bounds are shown individually. **b)** Significant Spearman correlations ($P < 0.05$) between samples in cytokine responses to all pairs of stimuli. Diagonal entries display significantly increased levels in urban (U) or rural (R) samples (Wilcoxon Rank Sum test, $P < 0.05$). **c)** Spearman correlation between counts of basophils, lymphocytes monocytes and neutrophils versus all measured cytokine responses. Blue and red respectively denote significant negative and positive correlation coefficients after FDR correction. **d)** Comparison of relative monocyte and lymphocyte counts ($10^3 / \mu\text{l}$) in urban (n=239) and rural (n=68) whole blood samples. Definition of box plot as in a. Significance assessed with two-sided Wilcoxon Rank Sum test.

Supplementary Figure 6. Overlapping cytokine responses in the 500 FG cohort. a) Spearman correlation plots between cytokine responses for different stimuli. **b)** Log-linear model of immunomodulatory applied to the 500FG cohort (same as with TZ cohort; Methods).



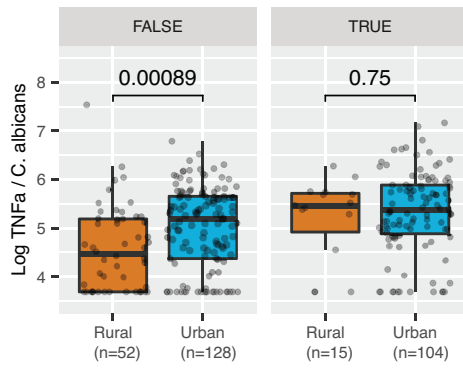
c

	effect	p-value
Residency_Area Urban	-0.03	0.606
<i>Akkermansia</i> Clade B	0.11	0.183
<i>Akkermansia</i> Clade C	0.08	0.070

d

	effect	pvalue
Residency Area Urban	0.262	3.28e-08 ***
<i>S. salivarius</i> Clade B	0.262	1.13e-07 ***

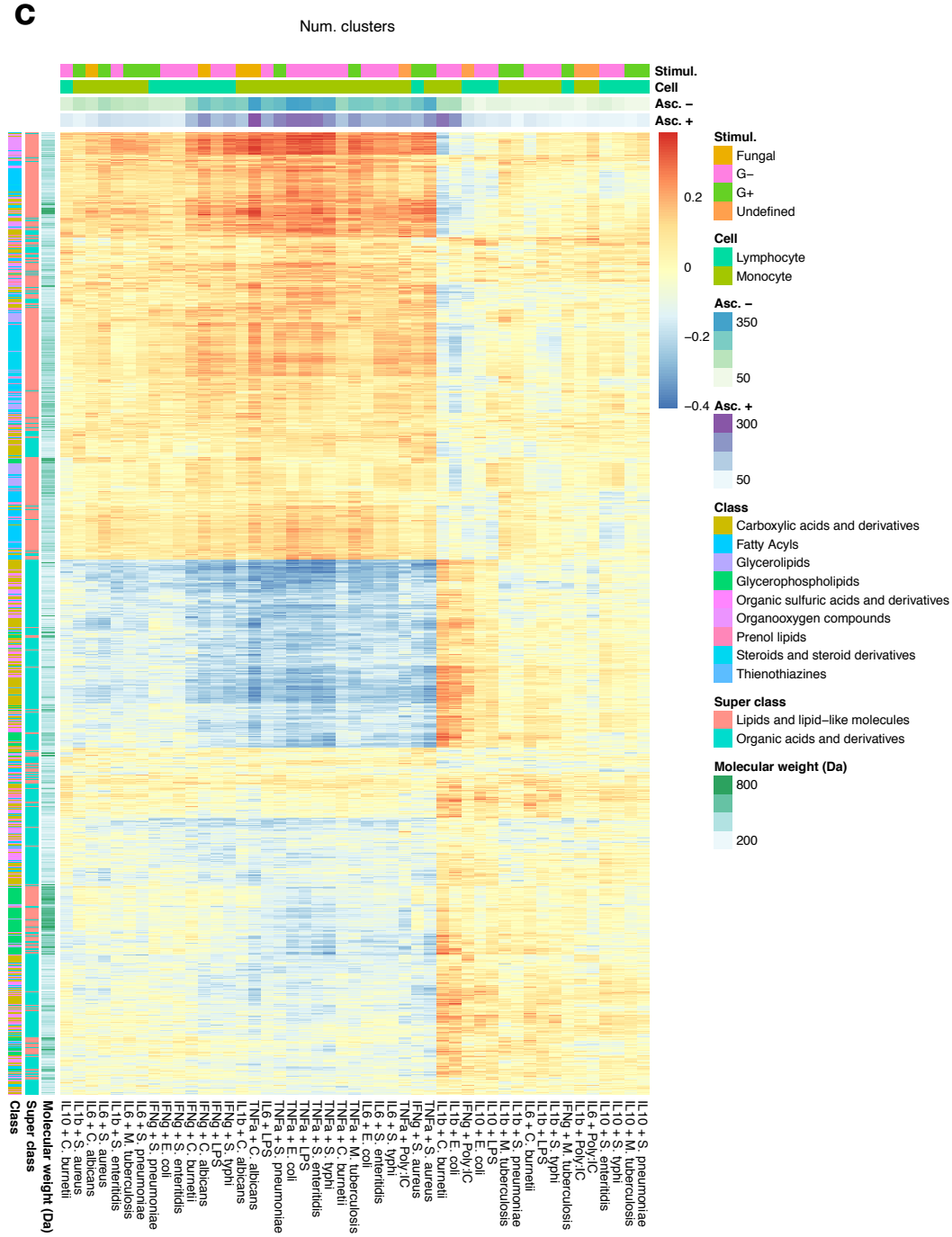
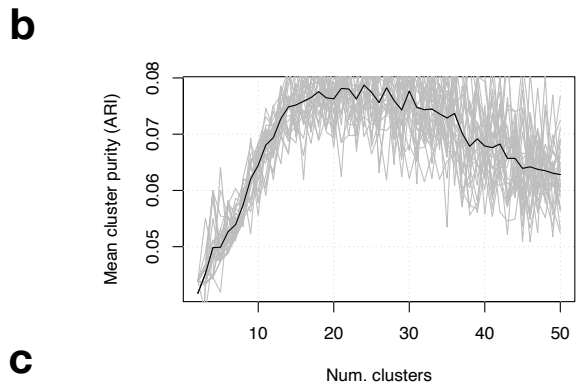
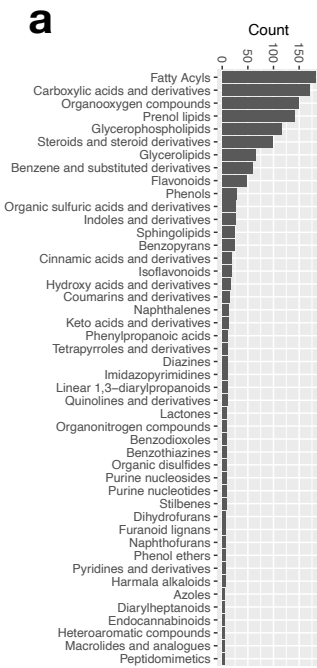
Supplementary Figure 7. Clade analysis for two immunomodulatory taxa, *Akkermansia* and *S. Salivarius*. **a-b)** Phylogenetic trees of samples from urban, rural (TZ) and Dutch samples obtained with StrainPhlan. Clades are determined such that the ordering within them retains a significant split in rural-urban-TZ (Kruskal-Wallis test). **c-d)** Clade effects on cytokine responses using the ANOVA model of the form $\text{Log Cytokine Expression} \sim \text{Age} + \text{Gender} + \text{Residency} + \text{Clade}$, where the observation including the taxa are used.

a*Akkermansia***b**

Ascomycota



Supplementary Figure 8. Effect on expression of TNF- α stimulated with *C. albicans* in rural and urban samples. Expression based on detection of **a)** Akkermansia and **b)** Ascomycota (two-tailed Wilcoxon rank sum test, n denotes the number of biologically independent samples in each category). Box plots show the median (center), first (Q1) and third quartile (Q3, box). Lower and upper whiskers represent $Q1 + 1.5 \text{ IQR}$ and $Q3 + \text{IQR}$, respectively. IQR, inter-quartile range. Data points outside whiskers' bounds are shown individually.



Supplementary Figure 9. Clustering and quality control of serum metabolomics data. A total of 1,607 m/z peaks identified by MS/MS were identified and matched to at least one molecule in the Human Metabolome Database (HMDB) **a)** The distribution of molecular classes for measured molecules. **b)** Number of clusters (28) is selected such that the clusters show maximal correspondence with the molecular classes as measured by the Adjusted Random Index (ARI). The clustering is repeated 50 times for each number of clusters of K with different random initializations. **c)** Significant spearman correlations between cytokine responses and metabolites (Spearman correlation, PFDR < 0.05).

Supplementary Figure 10. Arginine biosynthesis pathway detection in microbial species. a, b) Copy numbers of arginine biosynthesis pathway II (MetaCyc, ARGSYNBSUB), detected in the metagenome of negative, neutral and positive species. **c)** Detected enzymes, metabolites and associations in arginine biosynthesis pathway (KEGG KO00220). Colored nodes represent detected microbial enzymes (squares) and compounds (circles), with significant correlations with cytokine responses in red (positive) and blue (negative). Enzymes detected in stool microbiomes are shown as colored squares, with contribution from negative (blue) or positive species (red) quantified as fraction of total RPKM in the cohort.