

SUPPLEMENTARY TEXT 3

COMPLETE DESCRIPTION OF METHODS USED FOR THE BIOINFORMATIC STATISTICAL ANALYSIS OF THE MICROBIOME DATA

Pre-processing of amplicon reads

The FLASH program was used to join the paired-end reads (25). The data was barcode-corrected and quality filtered using the QIIME package; followed by clustering of reads into Operational Taxonomic Units (OTUs) (97% identity threshold) using USEARCH Clustering algorithm; followed by chimeric removal (26, 27). The taxonomic classification of the representative sequences for each OTU was performed using both the RDP classifier (genus level: 0.8 confidence threshold) and the SPINGO classifier (species level: 0.7 confidence threshold) (28, 29).

Multi-variate analysis of dietary profiles and taxonomic profiles

Multivariate analyses using Principal Coordinate Analysis (PCoA) were performed using the ade4 package of the R programming interface, using Spearman distances of the individual sample profiles as well as the across time point changes (final-baseline). To test the significance of the between-country variation of the baseline dietary and microbiome profiles, Permutational Analysis of Variance (PERMANOVA) was performed on the PCoA objects using the adonis function of the vegan R package. Procrustes analysis was performed to quantify the relationships between the baseline diet and microbiome profiles using the procrustes function of the vegan package. The Shannon diversities of the samples were obtained using the diversity function of the vegan R package.

Machine Learning-based identification of microbiome taxa associated with the dietary intervention

The Machine learning based Random Forest (RF) approach (implemented in the randomForest package of R) was used to identify microbiome taxa significantly associated with NU-AGE FBDG adherence scores. We first divided individuals into three equal tertiles, namely 'High Adherence', 'Medium Adherence' and 'Low Adherence' in decreasing order of the change in adherence across time-points and the samples from each into two cohorts corresponding to the baseline and final time-points. Two separate models were created for the baseline and the final time points. The performance of the models was measured by calculating the correlation

between the actual and the predicted food scores obtained using the models. The RF approach provided the feature score importance scores for each microbiome component (OTUs) (indicating the extent of association of these with the dietary adherence scores). For identifying the most-predictive features, iterative random forest models ($n=100$, sample subset-size=100) with varying number of features (selected in decreasing order of their feature importance scores) were obtained using the randomForest package (two-fold cross validation) and their performances compared. Finally, to identify the OTUs associated with dietary adherence, a Reboot approach (using Spearman correlations) was used to identify OTUs that were significantly associated with adherence scores with an FDR corrected P-value $< 1e-5$ (30). OTUs positively and negatively associated with diet were classified as DietPositive and DietNegative, respectively. A pictorial representation of the workflow adopted for this entire step is provided in **Supplementary figure 1**.

Overview: iBBiG is based on the detection profile of the taxonomic units (in this case, the Operational Taxonomic Units (OTUs)). It then utilizes an iterative, heuristic, genetic-algorithm based methodology to identify modules of taxa within a microbial community that tend to show strong co-occurrence relationships across a given population of microbiomes. The primary advantage of this strategy is its flexibility, as it allows identification of over-lapping modules such that certain taxonomic units can be part of multiple modules. Such a partitioning strategy makes more biological sense as certain taxa (or species) can be part of multiple guilds because of their functional versatility or may be functionally specialized (i.e. belonging to specific guilds).

Method: For identifying modules within the gut microbiome, we used the iterative **Binary Bi**-clustering of **Gene**-sets (iBBiG) approach (38). Rather than profiling abundances or proportions, iBBiG investigates the detection profile of the taxonomic units or OTUs. Subsequently, an iterative, heuristic, genetic-algorithm based methodology is used to identify taxonomic modules that tend to show strong co-occurrence relationships across a given population of microbiomes. For performing the iBBiG based clustering, we used the iBBiG function available within the Bioconductor package of R. While OTUs belonging to the different modules were then classified based on their clustering patterns, samples were classified based on the occurrence of the different iBBiG modules within them. The taxonomic compositional pattern of each module was then obtained by collating the RDP-based genus classification of each OTU and subsequently rank-normalizing these based on the abundance

of each genus (in terms of the number of OTUs) across a module. To associate the modules with frailty, we first obtained the frailty status of each individual at each time-point (0: Non Frail; 1: Pre Frail; 2: Frail). Subsequently based on the changes across time-points, individuals across the cohorts were classified as ‘Reduced Frailty’, ‘No Change’ and ‘Increased Frailty’. The representation of each of the modules were obtained at both the time-points for each of three groups of individuals. The occurrence changes of each module (the number of samples in which a module is present at follow-up divided by the number of samples the module is present in at the baseline) were computed for each group. The log fold changes in these ratios in the Reduced Frailty with respect to the Increased frailty groups would provide the enrichment or depletion of the modules in individuals with reduced frailty as compared to those showing an increase in frailty across time-points. A positive change would indicate enrichment, and a negative value would indicate depletion. To compare the patterns across modules X and Y, Chi-square tests (using the `chisq.test` function of R) were then performed on the contingency tables containing four values, namely occurrence at baseline and follow-up of reduced frailty and occurrence at baseline and follow-up of increased frailty, corresponding to the two modules. To check for the significance of the differences of the occurrences across modules in terms of their diet association, we obtained number of times a module was present in the list of DietPositive and the DietNegative OTUs, and subsequently compared them using the Fishers’ Exact test (`fisher.test` function of R).

Associating dietary adherence and microbiome changes with frailty and inflammation

For associating the abundances of the adherence associated marker OTUs with the different measures of frailty, cognitive function and cytokine profiles, we computed Spearman correlations using the `corr.test` function of the `psych` package in R (along with the Benjamini-Hochberg corrected p-values).

To account for various confounders, we used Partial Correlations (`partial.r` and the `corr.p` functions of the `psych` R package). Partial correlations measure the strength and the direction of the association between two variables considering the effect of confounding variable (s). Partial Correlations are like multiple regressions with confounders but not limited to specific distributions of the response and predictor variables. Further, one can compute rank-based non-parametric measures of association like the Spearman rho (which we have used in this study), after considering the confounding effect of other factors like adherence scores or age/BMI/gender.

Computation of Microbiome Indices

A pictorial representation of the methodology for this purpose is described in **Supplementary figure 2**. This scoring scheme ‘rewards’ samples with higher abundances of Marker OTUs with increasingly positive association with adherence scores and taxes those which have higher abundances of Marker OTUs with negative associations with adherence scores.

For each sample, the diet-modulated microbiome score was computed using the following formula:

$\sum_{\text{across all marker OTUs}} (\text{OTU correlation with Diet adherence scores}) * \text{Abundance of the OTU}$

To avoid over-fitting, leave-one out strategy was applied where for computing the microbiome index for a given sample, the sample was not considered while calculating the OTU correlations (with Diet Adherence scores).

Obtaining Inferred Microbial Metabolite Profiles based on Species Abundance Profiles

Literature annotated Species-to-Metabolite consumption/production associations were already available as part of the Virtual Metabolic Human database as well as those obtained in a recent meta-analysis by Sung *et al* (32, 33). These were parsed to create a present/absence information map of around 300 metabolite production and consumption profiles in greater than 900 species in a 0 (absent) and 1 (present) notation. Given the SPINGO-based species abundance profile, from the 16S amplicon data, the inferred metabolite profile was then obtained as an inner product of the species abundance profile and the species-to-metabolite map.

Generation of co-occurrence networks and computation of centrality measures

We used the Reboot Approach for generating the inter-microbial co-occurrence/co-inhibition networks (30) (described in **Supplementary text 4**). The co-occurrence networks obtained were visualized using Cytoscape (34). For any network, two different centrality measures were calculated for the nodes, namely degree centrality and betweenness centrality using the igraph R package. The relative co-occurrence propensities between any two groups of taxa were calculated as the log of the number of positive edges divided by the number of negative edges.

Given any two features (in this case, the OTUs), the Reboot approach computes the association between the two features using two different distributions of association measures obtained using repeated iterations as described below(52). The association measure can be any score, like the Pearson correlation, Spearman correlation, the Regression coefficients, or even the effect size measures. The first distribution (bootstrap distribution) was obtained by taking the repeated sub-samples of randomly selected observations and then computing the

association between the two features. This profiled the association values across an entire observation landscape, thereby removing biases which could be present because of specific samples. The second distribution (null distribution) was obtained by performing an equal number of iterations, where in each iteration, a fixed set of values (which in this case was 50%) are swapped across samples for both the features. The profiles were then re-normalized and the associations computed for the two features. The distribution of the values obtained in the two distributions were then compared using any comparative tests (which in this case was Mann-Whitney). The p-values thus obtained were then False Discovery Rate (FDR) corrected (Benjamini-Hochberg) and those pairs of features having FDR-corrected associations of less than $1e-5$ (threshold used in this study) were inferred to be significant and an edge drawn between them in the network. The directionality of the association was taken as the sign of the median value of the bootstrap distribution. While pairs of features with significant positive associations were used to create the co-occurrence network, those with negative associations were used to create the co-inhibition network.

*Please refer to the main document for the corresponding reference numbers.