

## Article

# MicrobiomeGWAS: A Tool for Identifying Host Genetic Variants Associated with Microbiome Composition

Xing Hua<sup>1,2</sup>, Lei Song<sup>1</sup>, Guoqin Yu<sup>3</sup>, Emily Vogtmann<sup>1</sup>, James J. Goedert<sup>1</sup>, Christian C. Abnet<sup>1</sup>, Maria Teresa Landi<sup>1</sup> and Jianxin Shi<sup>1,\*</sup>

<sup>1</sup> Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institute of Health, Rockville, MD 20850, USA; xhua2@fredhutch.org (X.H.); lei.song@nih.gov (L.S.); emily.vogtmann@nih.gov (E.V.); jamesgoedert@gmail.com (J.J.G.); abnetc@mail.nih.gov (C.C.A.); landim@mail.nih.gov (M.T.L.)

<sup>2</sup> Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA

<sup>3</sup> Molecular Genetics and Genomics Branch, Center for Scientific Review, National Institute of Health, Bethesda, MD 20817, USA; guoqin.yu@nih.gov

\* Correspondence: jianxin.shi@nih.gov; Tel.: +1-240-276-7419

**Abstract:** The microbiome is the collection of all microbial genes and can be investigated by sequencing highly variable regions of 16S ribosomal RNA (rRNA) genes. Evidence suggests that environmental factors and host genetics may interact to impact human microbiome composition. Identifying host genetic variants associated with human microbiome composition not only provides clues for characterizing microbiome variation but also helps to elucidate biological mechanisms of genetic associations, prioritize genetic variants, and improve genetic risk prediction. Since a microbiota functions as a community, it is best characterized by  $\beta$  diversity; that is, a pairwise distance matrix. We develop a statistical framework and a computationally efficient software package, microbiomeGWAS, for identifying host genetic variants associated with microbiome  $\beta$  diversity with or without interacting with an environmental factor. We show that the score statistics have positive skewness and kurtosis due to the dependent nature of the pairwise data, which makes  $p$ -value approximations based on asymptotic distributions unacceptably liberal. By correcting for skewness and kurtosis, we develop accurate  $p$ -value approximations, whose accuracy was verified by extensive simulations. We exemplify our methods by analyzing a set of 147 genotyped subjects with 16S rRNA microbiome profiles from non-malignant lung tissues. Correcting for skewness and kurtosis eliminated the dramatic deviation in the quantile–quantile plots. We provided preliminary evidence that six established lung cancer risk SNPs were collectively associated with microbiome composition for both unweighted ( $p = 0.0032$ ) and weighted ( $p = 0.011$ ) UniFrac distance matrices. In summary, our methods will facilitate analyzing large-scale genome-wide association studies of the human microbiome.

**Keywords:** microbiome; genome-wide association study; gene–environment interaction; host genetics; tail probabilities; skewness and kurtosis



**Citation:** Hua, X.; Song, L.; Yu, G.; Vogtmann, E.; Goedert, J.J.; Abnet, C.C.; Landi, M.T.; Shi, J. MicrobiomeGWAS: A Tool for Identifying Host Genetic Variants Associated with Microbiome Composition. *Genes* **2022**, *13*, 1224. <https://doi.org/10.3390/genes13071224>

Academic Editor: Yi-Juan Hu

Received: 2 June 2022

Accepted: 1 July 2022

Published: 9 July 2022

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2022 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

The human body is colonized by bacteria, viruses, and other microbes that exceed the number of human cells by at least 10-fold and that exceed the number of human genes by at least 100-fold. The relationship between a person and his or her microbial population, termed the microbiota, is generally mutualistic. The microbiota may promote human health by inhibiting infection by pathogens, conditioning the immune system, synthesizing and digesting nutrients, and maintaining overall homeostasis. The microbiome, which is the collection of all microbial genes, can be investigated through massively parallel, next-generation DNA sequencing technologies. By amplifying and sequencing highly variable

regions of 16S ribosomal RNA genes that are present in all eubacteria, cost-effective and informative microbiome profiles down to the genus level are obtained.

The human microbiome has been associated with diseases, including obesity [1], inflammatory bowel disease (IBD) [2], colorectal cancer [3], and breast cancer [4]. Thus, identifying factors that have a sustained impact on the microbiome is fundamental for elucidating its role in health conditions and for developing treatment strategies. Increasing evidence suggests that microbiome composition at a specific site of the human body is impacted by environmental factors [5,6], host genetics [7,8], and possibly by their interactions. In the mouse, quantitative trait loci (QTL) studies have identified loci contributing to the variation in the gut microbiome using linkage analysis [9,10]. Recently, Goodrich et al. [11] systematically investigated the heritability of the human gut microbiome by comparing monozygotic twins to dizygotic twins and found substantial heritability in different microbiome metrics, suggesting the important role of host genetics on gut microbiome diversity. Associations between individual host genetic variants and microbiome taxa abundances have also begun to emerge in other human samples [7,8,12]. These studies suggest that genome-wide association studies (GWAS) have great potential to identify host genetic variants associated with microbiome diversity.

GWAS of complex human diseases have identified many risk SNPs; however, the biological mechanisms are largely unknown for the majority of the risk SNPs. QTL studies of intermediate traits, e.g., gene expression [13,14], DNA methylation [15,16], chromatin structure [17,18], and metabolite production [19,20], have provided useful insights into the biological mechanisms of the GWAS findings. The human microbiome at a specific body site is another important and informative intermediate trait for interpreting GWAS signals. Knights et al. [8] reported that a risk SNP for IBD located in *NOD2* was associated with the relative abundance of *Enterobacteriaceae* in the human gut microbiome. Tong et al. [7] show that a loss-of-function allele in *FUT2* that increases the risk of developing Crohn's Disease (CD) may modulate the energy metabolism of the gut microbiome. In both examples, the microbiome is a potential intermediate for explaining the association between risk SNPs and disease risks, although a formal mediation analysis is required based on samples with genotype, microbiome, and disease status data. Moreover, identifying microbiome-associated host genetic variants has the potential to prioritize SNPs for discovery and to improve the performance of polygenetic risk prediction.

Three types of microbiome metrics can be derived as phenotypes for GWAS analysis. First, for each taxon at a specified taxonomic level (phylum, class, order, family, genus, and species), we calculate the relative abundance (RA) of the taxon as the ratio of the number of sequencing reads assigned to the taxon to the total number of sequencing reads. In 16S ribosomal RNA sequence profiles, approximately 100–200 taxa with average RAs  $\geq 0.1\%$  (from the phylum level to the genus level) across samples are abundant enough for QTL analysis. One can perform a Poisson regression to examine the association between the RA of each taxon and each SNP. Significant associations are identified using Bonferroni correction ( $p < 5 \times 10^{-8} / 200 = 2.5 \times 10^{-10}$ ) or by controlling FDR at an appropriate level. Second, multiple  $\alpha$ -diversity metrics [21] can be calculated to reflect the richness (e.g., number of unique taxa) and evenness of each microbiome community after a procedure called rarefaction, which eliminates the dependence between the estimated  $\alpha$  diversity and the variable total number of sequence reads across subjects. Once the  $\alpha$ -diversity metrics are derived, one may perform standard GWAS with  $\alpha$  diversity as the phenotype using linear regression.

Because a microbiota functions as a community, the most important analysis for a microbiome GWAS may be by assessing the complete structure of the community by using a pairwise microbiome distance matrix (or  $\beta$  diversity) of the microbial community. Microbiome distances can be defined in different ways, based on using phylogenetic tree information or each taxon's abundance information. Bray–Curtis dissimilarity [22] quantifies the difference between two microbiome communities using the abundance information of specific taxa. UniFrac [23–25] is another widely used distance metric. Unlike

the Bray–Curtis dissimilarity metric, UniFrac compares microbiome communities by using information on the relative relatedness of each taxon, specifically by phylogenetic distance (branch lengths on a phylogenetic tree). UniFrac has two variants: the weighted UniFrac [24], which accounts for the taxa abundance information, and the unweighted UniFrac [23], which only models the information of presence or absence. Recently, a generalized UniFrac distance metric [26] was developed to automatically appreciate the advantages of weighted and unweighted UniFrac metrics and was shown to provide better statistical power to detect associations between human health conditions and microbiome communities. GWAS based on a microbiome distance matrix aims to identify the host SNPs associated with microbiome composition. This has been done frequently by fitting non-parametric multivariate models [27]. This approach requires permutations to assess significance [28], which is computationally prohibitive, particularly when evaluating  $p$ -values less than  $5 \times 10^{-8}$ —the standard GWAS  $p$ -value threshold—or even lower when testing multiple-diversity matrices. In a recent microbiome GWAS, the computation is prohibitive even using a moment matching method based on the F-statistic.

Intuitively, the microbiome distances tend to be smaller for pairs of subjects with similar genotypic values at the associated SNP. In addition, it is also of great interest to identify host SNPs that interact with an environmental factor to affect microbiome composition. Importantly,  $\beta$  diversity is temporally more stable compared with RA of taxa and  $\alpha$ -diversity metrics based on the data from the Human Microbiome Project [29], suggesting a smaller power loss for a GWAS due to temporal variability. To our knowledge, no statistical methods or software packages have been designed to efficiently analyze microbiome GWAS data using distance matrices as phenotypes.

In this paper, we develop a statistical framework and a computationally efficient package, microbiomeGWAS, for analyzing microbiome GWAS data. Our package allows the detection of host SNPs with the main effect or interaction with an environment factor; i.e., host SNPs interacting with an environment factor to affect the microbiome composition. We calculate the variance of the score statistics by appropriately considering the dependence of the pairwise distances. Importantly, we show that the score statistics have positive skewness and kurtosis due to the dependence in pairwise distances, which makes the approximation of small  $p$ -values based on the asymptotic distribution too liberal, which easily yields false positive associations. Resampling methods, e.g., bootstrap or permutation, are computationally prohibitive for accurately approximating small  $p$ -values. We propose to improve the tail probability approximation by correcting for skewness and kurtosis of the score statistics. Numerical investigations demonstrate that our method provides a very accurate approximation, even for  $p = 5 \times 10^{-8}$ . MicrobiomeGWAS runs very efficiently, taking 36 min for analyzing main effects and 69 min for analyzing both main and interaction effects for a study with 2000 subjects and 500,000 SNPs, using a single core. MicrobiomeGWAS is available at <https://github.com/lncibb/microbiomeGWAS> [30], accessed on 30 May 2022.

We illustrate our methods by applying microbiomeGWAS to non-malignant lung tissue samples ( $N = 147$ ) in the Environment And Genetics in Lung cancer Etiology (EAGLE) study [31,32]. Because smoking may alter microbiome composition, we tested both the main effect and gene–smoking interaction effect. When  $p$ -values were calculated based on asymptotic distributions, the quantile–quantile (QQ) plots strongly deviated from the uniform distribution. Nine loci also achieved genome-wide significance based on asymptotic approximations. Correcting for skewness and kurtosis eliminated the inflation and also the genome-wide significance of these loci. However, we provide evidence that the established lung cancer risk SNPs are associated with lung microbiome composition.

## 2. Material and Methods

### 2.1. A Score Statistic for Testing Main Effect

Suppose that we have a set of  $N$  subjects genotyped with SNP arrays. For notational simplicity, we consider only one SNP with a minor allele frequency (MAF) denoted as  $f$ .

Our interest centers on testing whether the genotype of the SNP is associated with microbiome composition. Let  $g_n = 0, 1, 2$  represent the number of the minor alleles for subject  $n$ . We assume that the 16S rRNA gene of microbiota from a target site (e.g., gut) has been sequenced for these samples. Let  $d_{ij}$  be the microbiome distance between subject  $i$  and subject  $j$  and  $D$  be the distance matrix.

Intuitively, if the SNP is associated with the microbiome composition, the microbiome distances tend to be smaller for subject pairs with similar genotypic values, as is illustrated in Figure 1. For  $N$  subjects,  $N(N - 1)/2$  pairs can be divided into three groups with genetic distance 0, 1, and 2. For example, a pair of subjects with genotype (AA, AA) or (BB, BB) has genetic distance 0; a pair of subjects with genotype (AA, BB) or (BB, AA) has genetic distance 2; all other pairs have genetic distance 1. Apparently, we expect the microbiome distance to be positively correlated with genetic distance for subject pairs.

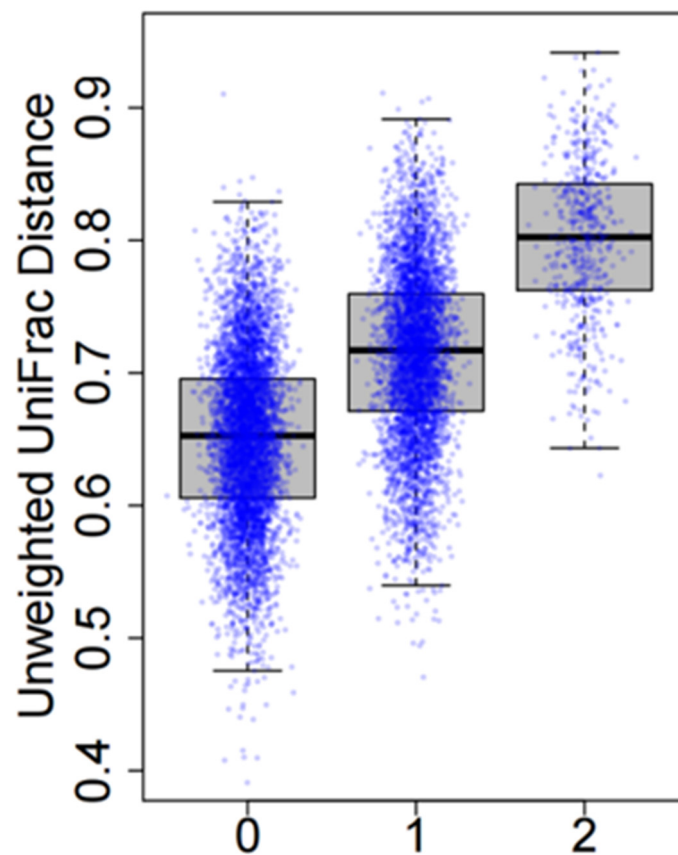


Figure 1. Microbiome distances are positively correlated with genetic distances at an associated SNP.

We define  $G_{ij} = |g_i - g_j|$  as the genetic distance for a pair of subjects  $(i, j)$ . We assume  $d_{ij} = \alpha + \beta_M G_{ij} + \varepsilon_{ij}$ . A score statistic for testing  $H_0 : \beta_M = 0$  (main effect) vs.  $\beta_M > 0$  is derived as:

$$S_M = \sum_{i < j} d'_{ij} G_{ij} \quad \text{with} \quad d'_{ij} = d_{ij} - \frac{2}{N(N-1)} \sum_{k < l} d_{kl}. \tag{1}$$

The variance  $Var_0(S_M|D)$  under  $H_0 : \beta_M = 0$  is calculated by considering the dependence in  $(G_{ij}, G_{kl})$  and conditioning on the distance matrix  $D$ . Briefly, we have  $Var_0(S_M|D) = \sum_{i < j, k < l} d'_{ij} d'_{kl} Cov(G_{ij}, G_{kl})$ . When  $(i, j, k, l)$  are distinct,  $G_{ij}$  and  $G_{kl}$  are independent; i.e.,  $Cov(G_{ij}, G_{kl}) = 0$ . Some algebra leads to

$$Var_0(S_M|D) = \frac{N(N-1)}{2} Var(G_{ij}) \mu_2 + N(N-1)(N-2) Cov(G_{ij}, G_{ik}) \mu_3 \tag{2}$$

where

$$\mu_2 = \frac{2}{N(N-1)} \sum_{i < j} (d'_{ij})^2 \tag{3}$$

and

$$\mu_3 = \frac{2}{N(N-1)(N-2)} \sum_{i < j < k} (d'_{ij}d'_{ik} + d'_{ij}d'_{jk} + d'_{ik}d'_{jk}) \tag{4}$$

The details for calculating  $Var(G_{ij})$  and  $Cov(G_{ij}, G_{ik})$  are in Appendix A. The normalized statistic  $Z_M = S_M / \sqrt{Var_0(S_M|D)} \sim N(0, 1)$  under  $H_0$  asymptotically.

In analyses of real data, we typically have to adjust for covariates, including demographic variables and principal component analysis (PCA) scores derived based on genotypes, to eliminate potential population stratification. Given a distance matrix  $D$  and  $v$  covariates  $(x_{i1}, \dots, x_{iv})$ , we perform distance-based redundancy analyses using function *capscale* in the *vegan* package [33]. The residual matrix  $D'$ , extracted using the *residuals* function in the *vegan* package [33], is now adjusted for these potential confounding factors and can be used for genetic analysis.

### 2.2. A Score Statistic for Testing Gene–Environment Interaction

Let  $E_i$  denote an environmental variable. Define  $\Delta_{ij} = |g_i E_i - g_j E_j|$ . We extend the statistical framework to detect the SNP–environment interaction by assuming  $d_{ij} = \alpha + \beta_M G_{ij} + \beta_E |E_i - E_j| + \beta_I \Delta_{ij} + \varepsilon_{ij}$ , where  $\beta_M$  denotes the main genetic effect,  $\beta_I$  denote the additive gene–environment effect, and  $\beta_E$  denotes the main effect of the environmental factor. We consider testing the null hypothesis that the SNP is not associated with microbiome composition either directly or by interacting with  $E$ , i.e.  $H_0 : \beta_M = \beta_I = 0$ . The alternative hypothesis is  $H_1 : \beta_M > 0$  or  $\beta_I > 0$ .

We estimate  $\beta_E$  and  $\alpha$  under  $H_0$  and calculate  $d'_{ij} = d_{ij} - \hat{\alpha} - \hat{\beta}_E |E_i - E_j|$ . Let  $D' = (d'_{ij})$  be the residual matrix. The scores evaluated under  $H_0$  are  $S_M = \sum_{i < j} d'_{ij} G_{ij}$  for  $\beta_M$  and  $S_I = \sum_{i < j} d'_{ij} \Delta_{ij}$  for  $\beta_I$ . Similar to (2), we derive the variance  $Var_0(S_I|D')$  by accounting for the dependence in  $(\Delta_{ij}, \Delta_{kl})$ :

$$Var_0(S_I|D') = \frac{N(N-1)}{2} Var(\Delta_{ij}) \mu_2 + N(N-1)(N-2) Cov(\Delta_{ij}, \Delta_{ik}) \mu_3 \tag{5}$$

Let  $Z_M = S_M / \sqrt{Var_0(S_M|D')}$  and  $Z_I = S_I / \sqrt{Var_0(S_I|D')}$ . Asymptotically,  $Z_M \sim N(0, 1)$  and  $Z_I \sim N(0, 1)$  under  $H_0$ . In Appendix B, we derive

$$Cov_0(S_M, S_I|D') = \frac{N(N-1)}{2} Cov(G_{ij}, \Delta_{ij}) \mu_2 + N(N-1)(N-2) Cov(G_{ij}, \Delta_{ik}) \mu_3 \tag{6}$$

Let  $\rho = Cor_0(Z_M, Z_I|D')$  be the correlation between the two statistics. Asymptotically,  $(Z_M, Z_I)$  follows a bivariate normal distribution with a correlation matrix  $\Omega = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ . In Appendix C, we derive a statistic for jointly testing  $H_0 : \beta_M = \beta_I = 0$  vs.  $H_1 : \beta_M > 0$  or  $\beta_I > 0$ . Briefly, the 2D plane is partitioned to four parts (Figure 2). The joint statistic is derived as

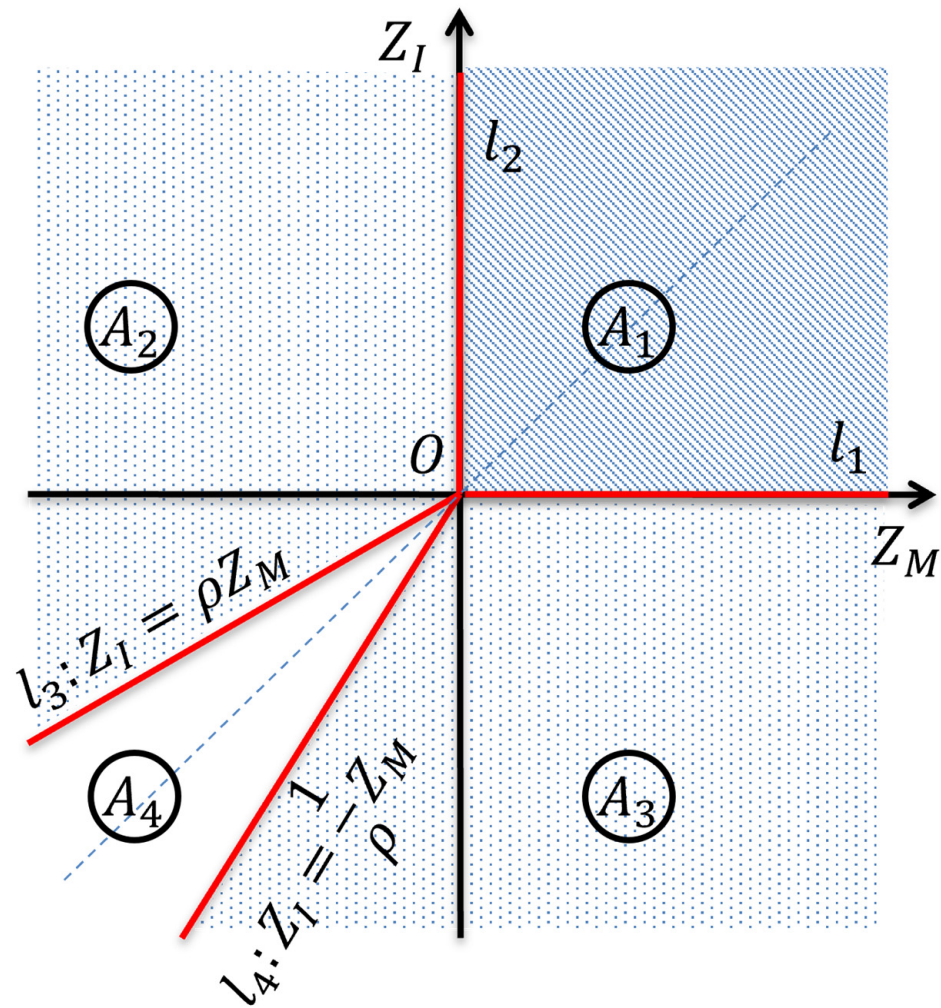
$$Q = \begin{cases} (Z_M, Z_I) \Omega^{-1} (Z_M, Z_I)^T & (Z_M, Z_I) \in A_1 \\ (w_1 Z_M + w_2 Z_I)^2 & (Z_M, Z_I) \in A_2 \\ (w_2 Z_M + w_1 Z_I)^2 & (Z_M, Z_I) \in A_3 \\ 0 & (Z_M, Z_I) \in A_4 \end{cases} \tag{7}$$

where  $w_1 = (\theta - 1/\theta)/2$ ,  $w_2 = (\theta + 1/\theta)/2$  and  $\theta = \sqrt{(1-\rho)/(1+\rho)}$ . The asymptotic  $p$ -value is calculated as

$$P(Q > b^2) = q_1 P(\chi^2_2 > b^2) + q_2 P(N(0, 1) > b) + q_3 P(N(0, 1) > b), \tag{8}$$



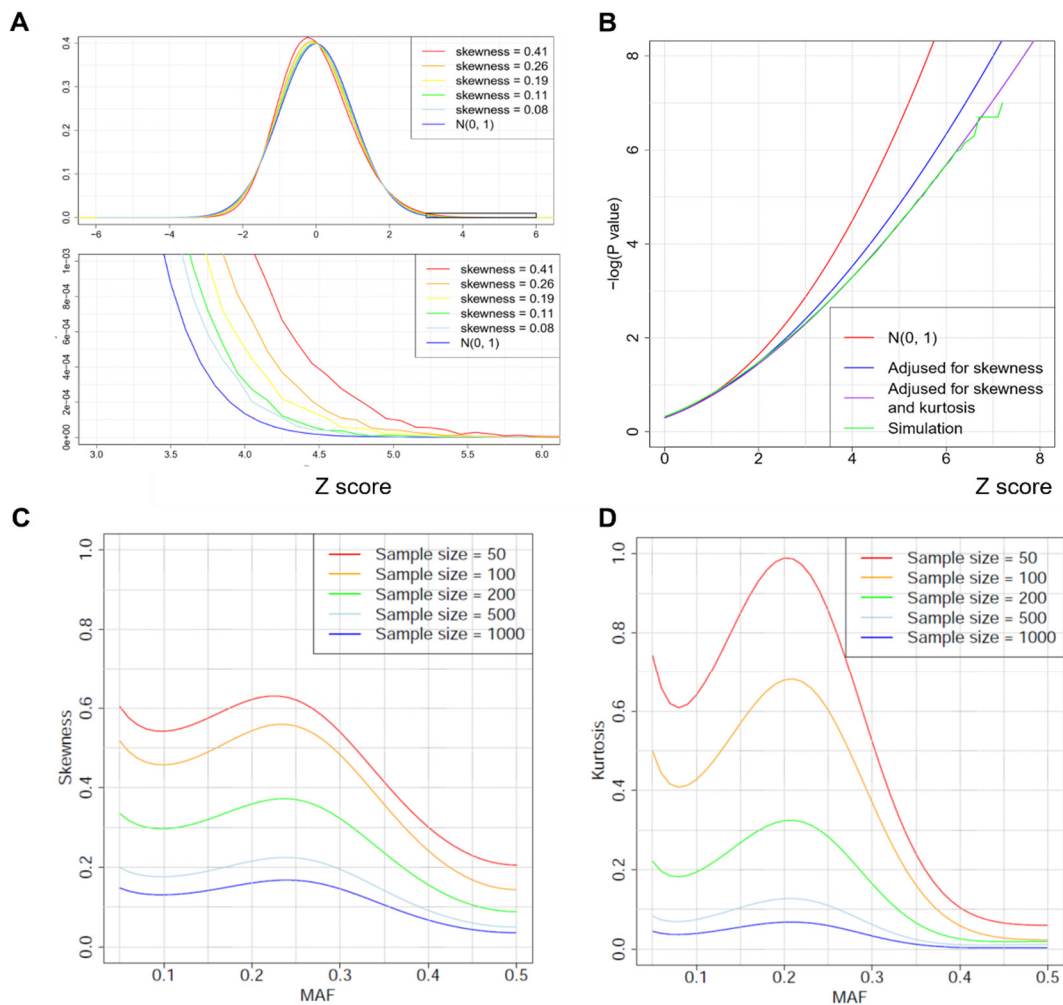
where  $q_i = P((Z_M, Z_I) \in A_i)$ .



**Figure 2.** Define the joint test for testing  $H_0 : \beta_M = \beta_I = 0$  vs.  $\beta_M > 0$  or  $\beta_I > 0$ . We assume that  $Z_M \sim N(0,1)$ ,  $Z_I \sim N(0,1)$  and  $cor(Z_M, Z_I) = \rho$  under  $H_0$ . Details are in Appendix C.

### 2.3. Improved *p*-Value Approximations by Correcting for Skewness and Kurtosis

Theoretic investigation suggests that the score statistics  $Z_M$  and  $Z_I$  have a positive skewness, which makes the tail probability approximations based on the asymptotic distribution  $N(0,1)$  unacceptably liberal (Figure 3A,B). In a numeric example with skewness  $\gamma = 0.2$ ,  $P(Z > 5) = 2.9 \times 10^{-7}$  based on  $N(0, 1)$ , which is approximately two orders of magnitude more significant than  $p = 3.9 \times 10^{-5}$  based on  $10^8$  permutations. The significance inflation becomes worse for smaller *p*-values and larger skewness  $\gamma$ . Similar but more tedious calculations suggest that both statistics have positive kurtosis, making the approximation based on  $N(0, 1)$  even worse. One possible solution is to approximate tail probabilities using permutations or bootstrap. However, these resampling methods are computationally prohibitive for testing millions of common SNPs in a large-scale study.



**Figure 3.** Correcting tail probabilities for skewness and kurtosis. **(A)** The standard normal distribution  $N(0,1)$  and an approximately normal distribution with positive skewness. The skewness has big impact when calculating the tail probability  $P(Z > b)$  for a large value of  $b$ . **(B)** Numerical evaluation of tail probability approximation for  $Z_M$ . We used the unweighted UniFrac distance matrix of 500 samples from the American Gut Project (AGP). For each value of  $b (> 0)$ , we calculated  $p$ -values  $P(Z_M > b)$  based on  $N(0,1)$ , skewness correction, both skewness and kurtosis correction, and  $10^8$  simulations. **(C)** Skewness depends on minor allele frequency (MAF) of SNPs and the sample size of the study, calculated based on the weighted UniFrac distance matrix in AGP data. **(D)** Kurtosis depends on MAF of SNPs and the sample size, calculated based on the weighted UniFrac distance matrix in the AGP data.

To address this problem, we calculated the skewness  $\gamma$  and kurtosis  $\kappa$  of the score statistics under  $H_0$  (Appendix D). We propose to improve the tail probability approximation  $P_0(Z > b)$  by correcting for the skewness and kurtosis, following the skewness correction in linkage analysis [34,35]. Technical details are provided in Appendix E. Correcting for both skewness and kurtosis leads to an approximation

$$P_0(Z > b) \approx e^{-b\zeta_1 + (1+\sigma_1^2)\zeta_1^2/2 + \gamma\zeta_1^3/6 + \kappa\zeta_1^4/24} \Phi(-\sigma_1\zeta_1) \tag{9}$$

where  $\zeta_1$  satisfies  $\zeta + \gamma\zeta^2/2 + \kappa\zeta^3/6 = b$ ,  $\sigma_1^2 = 1 + \gamma\zeta_1 + \kappa\zeta_1^2/2$  and  $\Phi(\cdot)$  is the cumulative distribution function of  $N(0,1)$ . Correcting for skewness but ignoring kurtosis (i.e., assuming  $\kappa = 0$ ) leads to an approximation

$$P_0(Z > b) \approx e^{-b\zeta_2 + (1+\sigma_2^2)\zeta_2^2/2 + \gamma\zeta_2^3/6} \Phi(-\sigma_2\zeta_2) \tag{10}$$

where  $\xi_2 = (\sqrt{1 + 2\gamma b} - 1)/\gamma$ ,  $\sigma_2^2 = 1 + \gamma\xi_2$ . Numerical results presented in Figure 3B demonstrate that (9) works very well.

Given the distance matrix  $\mathbf{D}$ ,  $\gamma_M \propto 1/N^{1/2}$ ,  $\gamma_I \propto 1/N^{1/2}$ ,  $\kappa_M \propto 1/N$  and  $\kappa_I \propto 1/N$  (Appendix D). Thus, skewness decays much more slowly with sample size  $N$  than kurtosis (Figure 3C,D). Thus, even for a large study with thousands of samples, correcting for skewness is necessary for accurately evaluating the tail probabilities. Importantly, both skewness and kurtosis highly depend on the MAF, suggesting that the impact of skewness and kurtosis is different across SNPs with a different MAF. Numerical studies (Figure 3C,D) show that skewness and kurtosis are minimized when MAF = 0.5 and maximized when MAF  $\approx$  0.2–0.3.

Finally, we discuss how to approximate the tail probability of  $Q$  in (7) for testing  $H_0 : \beta_M = \beta_I = 0$  by correcting for non-normality in  $Z_M$  and  $Z_I$ . When  $(Z_M, Z_I) \in A_2$  (or  $A_3$ ), we calculate the skewness  $E(w_1Z_M + w_2Z_I)^3$  and the kurtosis  $E(w_1Z_M + w_2Z_I)^4 - 3$  and use (9) to approximate  $P(w_1Z_M + w_2Z_I > b)$ . When  $(Z_M, Z_I) \in A_1$ , we first approximate their marginal  $p$ -values as  $p_M$  and  $p_I$  by (9), and then calculate the normal quantile  $z_M = \Phi(1 - p_M)$  and  $z_I = \Phi(1 - p_I)$ . Because the correction primarily impacts the tails of the distributions, the correlation between the two statistics will remain roughly unchanged; i.e.,  $cor_0(Z_M, Z_I) \approx cor_0(z_M, z_I)$ . Thus, when  $(Z_M, Z_I) \in A_1$ , the tail probability is approximated as  $P(\chi_2^2 > (z_M, z_I)\Omega^{-1}(z_M, z_I)')$ .

### 3. Results

#### 3.1. Simulation Results

The main purpose of simulations was to investigate the type-I error of  $Z_M$  (for testing the main genetic effect),  $Z_I$  (for detecting SNP–environment interactions), and  $Q$  (for detecting either the main genetic effect or SNP–environment effect or both). Simulations were performed under different combinations of sample size, MAF, and microbiome distance matrices. To make the simulations realistic, we used an unweighted distance matrix of the fecal microbiome samples with the 16S rRNA V4 region sequences from the American Gut Project (AGP) [36]. The OTU table, rarefied to 10,000 sequence reads per sample, as well as the metadata were downloaded from the AGP website. Samples with less than 10,000 sequence reads were excluded from the analysis. The weighted and the unweighted UniFrac distance matrices were generated in the Quantitative Insights Into Microbial Ecology [21] (QIIME) pipeline. Because antibiotics may substantially change the microbiome composition to generate outliers that may distort the null distribution, we excluded samples with self-reported history of antibiotic usage within one month. After quality control, 1879 subjects remained for analysis. In the simulations, we randomly selected  $N$  samples for a given sample size  $N$ .

For each setting, the type-I error rates were evaluated based on  $10^8$  simulations under  $H_0$ . For the interaction test and the joint test, the binary environment factor had a frequency of 50% and was simulated independent of the SNP. The type-I error rates are summarized in Table 1 for the weighted UniFrac distance matrix. The skewness and kurtosis are reported in Figure 3C,D. The statistics adjusted for skewness and kurtosis have accurate type-I error rates while the statistics without adjustment have unacceptably high type-I error rates. As the sample size increases, the impact of skewness and kurtosis decreases. However, even for a study with  $N = 1000$ , the type-I error rates are still seriously inflated. The results for the unweighted UniFrac distance matrix and for MAF = 0.5 are reported in Table S1.

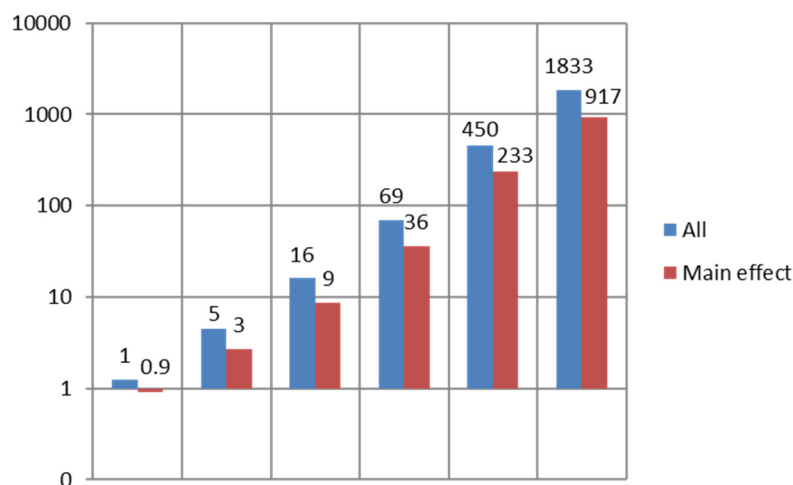


**Table 1.** Type-I error rates estimated based on  $10^8$  simulations. Minor allele frequency = 20%. Simulations were based on the weighted UniFrac distance matrix of the gut microbiome data from the American Gut Project. Reported are the type-I error inflation factor. A value greater than 1 indicates an inflated type-I error.

	N	$\alpha = 10^{-3}$	$Z_M$			$Z_I$			Q	
			$10^{-5}$	$10^{-7}$	$10^{-3}$	$10^{-5}$	$10^{-7}$	$10^{-3}$	$10^{-5}$	$10^{-7}$
Asymptotic approximation	100	5.5	51.6	610.0	4.7	36.1	342.8	7.3	80.9	1148.0
	200	3.7	23.0	187.3	3.1	15.8	105.5	4.6	33.0	316.7
	500	2.4	9.4	45.2	2.1	6.7	25.5	2.8	11.9	64.1
	1000	2.0	5.7	21.3	1.8	4.4	14.0	2.2	6.9	28.5
Adjusted for skewness and kurtosis	100	1.0	1.2	0.7	1.0	1.1	0.6	1.0	1.5	2.0
	200	1.0	1.1	1.0	1.0	1.1	0.7	0.9	1.3	1.8
	500	1.0	1.1	1.3	1.0	1.0	0.9	0.9	1.0	1.7
	1000	1.0	1.0	1.2	1.0	1.0	0.8	0.9	1.0	1.1

### 3.2. Software Implementation, Memory Requirement, and Computational Complexity

We implemented our algorithms in a software package, microbiomeGWAS, which is freely available at <https://github.com/lscibb/microbiomeGWAS> [30], accessed on 30 May 2022. MicrobiomeGWAS requires three sets of files: a microbiome distance matrix file, a set of PLINK binary files for GWAS genotypes, and a set of covariates. MicrobiomeGWAS processes one SNP at a time and does not load all genotype data into memory; thus, it requires only memory for storing the distance matrix. Variance, skewness, and kurtosis can be partitioned into two parts related with the microbiome distance matrix and the MAF of the SNP separately; thus, we can quickly calculate these quantities for a predefined grid of MAFs. The overall computational complexity is about  $O(N^2M)$ , where  $N$  is the sample size and  $M$  is the number of SNPs. Figure 4 reports the computation time on a Linux server using a single core. For a study with 10,000 subjects, it takes approximately 15 h for analyzing the main effect and approximately 30 h for analyzing both the main and interaction effects for 0.5 million variants. As a comparison, in a recent microbiome GWAS [37], to analyze  $7 \times 10^{-6}$  variants for the main effect and  $n = 3382$  subjects in the SHIP-TREND cohort [37], it would take 61 years using one CPU and 94 days using one graph-processing unit for parallel computation. Moreover, their analytic pipeline could not jointly analyze all 8956 subjects from five cohorts because of the computational burden; instead, they performed a stepwise search that may cause power loss.



**Figure 4.** Computation time for a microbiome GWAS with 500,000 SNPs. “Main”: computation time for testing main effect only. “All”: computation time for testing main effect, interaction and the joint null hypothesis  $H_0 : \beta_M = 0, \beta_I = 0$ .

### 3.3. GWAS of Microbiome Diversity in Adjacent Normal Lung Tissues

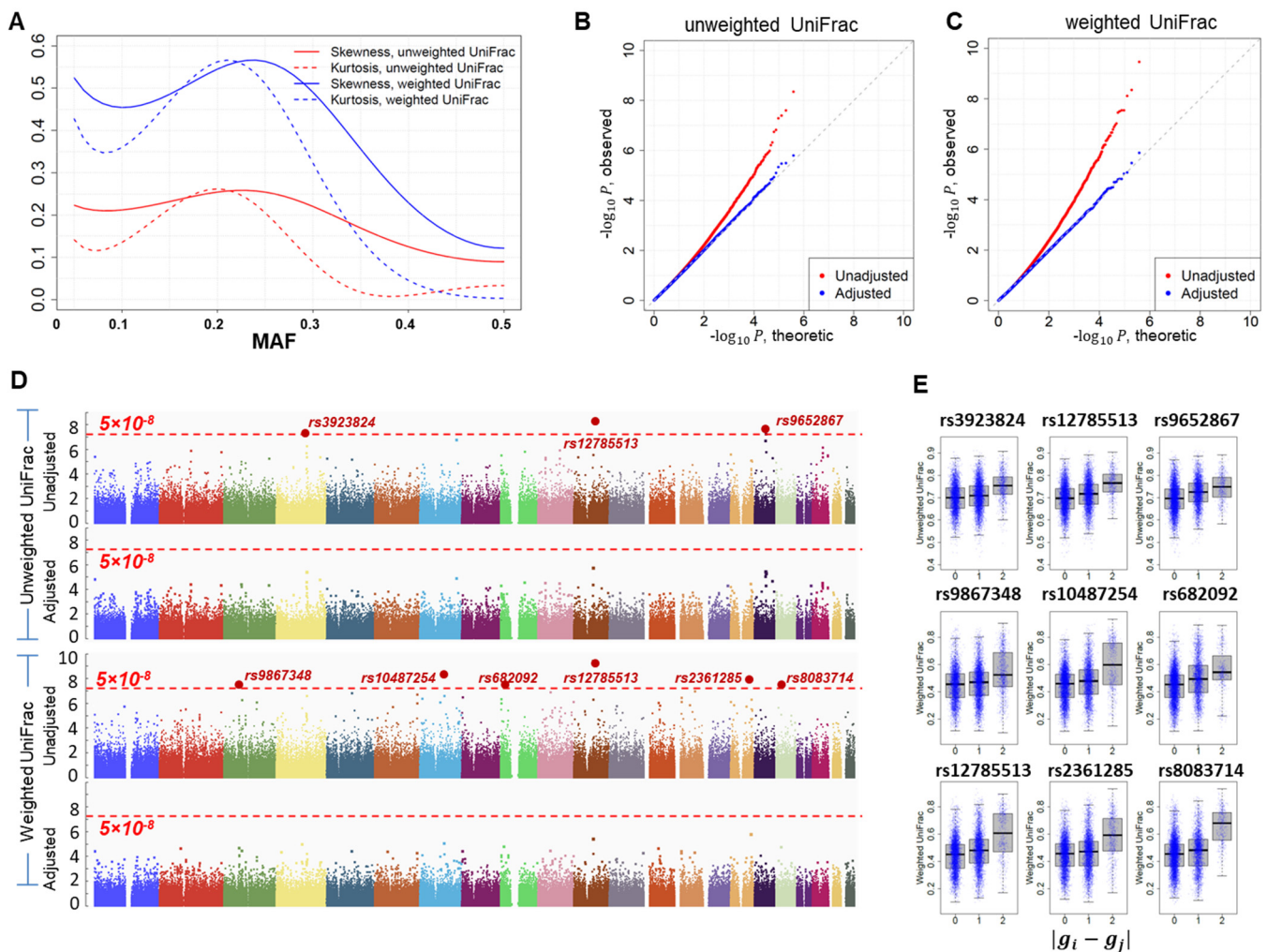
We applied our methods to a set of lung cancer patients of Italian ancestry in the EAGLE [31] study. All subjects have germline genome-wide SNPs [32] and 16S rRNA microbiome data (V3-V4 region, Illumina MiSeq, 300 paired-end) in histologically normal lung tissues from these patients. Here, the histologically normal lung tissues were 1–5 cm from the tumor tissue. We performed a series of quality control steps to filter out low-quality sequence reads: average quality score <20 over 30 bp windows, less than 60% similarity to the Greengenes [38] reference, or identified as chimera reads using UCHIME [39]. Sequence reads were then processed by QIIME [21] to produce the relative abundances (RA) of taxa, two  $\alpha$ -diversity metrics (observed number of species and Shannon's index), and  $\beta$  diversity metrics (unweighted and weighted UniFrac distances) rarified to 1000 reads. We included 147 subjects with at least 1000 high-quality sequence reads for genetic association analysis.

Out of the 147 subjects, 78 are current smokers, 8 never smoked, and 61 are former smokers. Because of the small number of never smokers, we merged never and former smokers as non-current smokers. All of the genetic association analyses were adjusted for sex, age, smoking status, and the top three PCA scores derived based on genome-wide SNPs. Here, the top three PCA scores were selected for controlling population stratification because the other PCA scores were unassociated with the distance matrices. We included 383,263 common SNPs with  $MAF \geq 10\%$  because rarer SNPs were expected to have no statistical power given the current sample size. We first performed GWAS analysis using PLINK [40] to identify the SNPs associated with taxa with an average RA greater than 0.1% or two  $\alpha$ -diversity metrics. We did not detect genome-wide significant associations with either the main effects or gene–smoking interactions.

Next, we performed GWAS analysis using unweighted and weighted UniFrac distance matrices as a representation of eubacteria  $\beta$  diversity. The results for testing the main effects are reported in Figure 5. Results for testing the joint effects (main effect and SNP by smoking status interaction) are reported in Figure S1. Because of the small sample size, we observed large values of skewness and kurtosis, with the magnitude varying with the MAF of the SNPs (Figure 5A). The score statistics based on the weighted UniFrac distance matrix had a much larger skewness and kurtosis than did the unweighted UniFrac matrix. Figure 5B,C report the quantile–quantile (QQ) plot of the logarithm of the association  $p$ -values for the unweighted and weighted UniFrac distance matrices, respectively. For each distance matrix, we produced QQ plots for  $p$ -values based on the asymptotic approximation and for  $p$ -values adjusted for skewness and kurtosis. For both distance matrices, the QQ plots before adjustment strongly deviated from the expected uniform distribution. Our adjustment eliminated the deviation. In addition, consistent with the observation that the skewness and kurtosis were larger for the weighted UniFrac distance matrix, the QQ plot deviated more for the analysis based on the weighted UniFrac distance. Note that the skewness and kurtosis only affect the tail probabilities; thus, the inflation of the QQ plot is not reflected by the genomic control lambda value [41], calculated as the median of the  $p$ -values. In fact,  $\lambda \approx 1$  for all four QQ plots.

Without correcting for skewness and kurtosis, we identified three and six loci achieving genome-wide significance ( $p < 5 \times 10^{-8}$ ) for the unweighted and weighted UniFrac distance matrices, respectively (Figure 5D). After correcting for skewness and kurtosis, no locus remained genome-wide significant (Figure 5D), which was verified by  $10^8$  permutations. Importantly, skewness and kurtosis had a dramatic effect on tail probabilities. Here, we use SNP rs12785513 as an example, which was identified as the top SNP in both analyses. In the unweighted UniFrac analysis,  $p = 4.4 \times 10^{-9}$  without adjustment and  $p = 1.6 \times 10^{-6}$  after adjustment, a 364-fold inflation. The inflation was even larger for weighted UniFrac analysis because of larger skewness and kurtosis (Figure 5A). In fact,  $p = 3.4 \times 10^{-10}$  without adjustment and  $p = 3.5 \times 10^{-6}$  after adjustment, a 1000-fold inflation. Although these SNPs were not significant genome-wide, they were the top SNPs from the current study. Thus, we report box-plots for each of these nine SNPs (Figure 5E). As expected, in all box plots, microbiome distances tend to be larger in subject pairs with

greater genetic distance at these SNPs. These associations remain to be replicated in studies with larger sample sizes.



**Figure 5.** Results of analyzing the microbiome GWAS data of 147 adjacent normal lung tissues in the EAGLE study. (A) Skewness and kurtosis for the main effect test using the unweighted and the weighted UniFrac distance matrices. (B) Quantile–quantile (QQ) plot for association  $p$ -values using the unweighted UniFrac distance matrix. “Adjusted”:  $p$ -values were corrected for skewness and kurtosis. “Unadjusted”:  $p$ -values were approximated based on the asymptotic distribution  $N(0, 1)$ . (C) Quantile–quantile (QQ) plot for association  $p$ -values using the weighted UniFrac distance matrix. (D) Manhattan plots based on the unweighted or the weighted UniFrac distance matrices. (E) Box plots for the top nine loci in microbiome GWAS analysis. Subject pairs are classified into three groups according to the genetic distance  $|g_i - g_j|$  at the SNP. The  $y$ -coordinate is the microbiome distance.

Finally, we concentrated on the six common SNPs in four genomic regions reported to be associated with lung cancer risk in GWAS of European subjects: rs2036534 and rs1051730 at 15q25.1 [42–45] (*CHRNA5–CHRNA3–CHRNA4*), rs2736100 and rs401681 at locus 5p15.33 [31,46] (*TERT/CLPTM1L*), rs6489769 [47] at 12p13.3 (*RAD52*), and rs1333040 at 9p21.3 [48] (*CDKN2A/CDKN2B*). The SNPs at 15q25.1 and 5p15.33 have the largest effect sizes for lung cancer risk based on the meta-analysis from the Transdisciplinary Research in Cancer of the Lung (TRICL) consortium [48]: OR = 1.32 for rs1051730, OR = 1.26 for rs2036534, OR = 1.13 for rs2736100, and OR = 1.14 for rs401681. Rs3131379 at locus 6p21.33 [46] (*BAT3/MSH5*) was excluded because the MAF = 7.5%. No SNPs were significantly associated with taxa RAs or  $\alpha$ -diversity metrics after correcting for multiple testing. However, association analysis based on the UniFrac distance matrices provided evidence

that these SNPs may be associated with the lung microbiota (Table 2). These SNPs were independent except that rs2036534 and rs1051730 at 15q25.1 were weakly correlated with  $R^2 = 0.15$ . A test combining six  $z$ -scores ( $Z_M$ ) and adjusting for the weak correlation yielded overall  $p$ -values of 0.0033 and 0.011 for the unweighted and the weighted UniFrac distance matrices, respectively. These results suggest that lung cancer risk SNPs were enriched for genetic association with the composition of the lung microbiome. The results for testing interactions and joint effects are reported in Table S2.

**Table 2.** Association  $p$ -values between lung cancer risk SNPs and microbiome composition in the EAGLE data.

SNP	Chr	Annotated Genes	Unweighted UniFrac	Weighted UniFrac
rs2036534	15q25.1	<i>CHRNA3/4/5</i>	0.425	0.167
rs1051730	15q25.1	<i>CHRNA3/4/5</i>	0.020	0.401
rs2736100	5p15.33	<i>TERT</i>	0.089	0.267
rs401681	5p15.33	<i>CLPTM1L</i>	0.056	0.005
rs6489769	12p13.3	<i>RAD52</i>	0.197	0.329
rs1333040	9p21.3	<i>CDKN2A/B</i>	0.249	0.224
Overall test			0.0032	0.011

#### 4. Discussion

We developed a software package, microbiomeGWAS, for identifying host genetic variants associated with microbiome composition. MicrobiomeGWAS can test both the main effect and SNP–environment interactions. Importantly, we found that the score statistics had positive skewness and kurtosis and that the tail probabilities evaluated based on asymptotic approximations were very liberal. We addressed this problem by explicitly adjusting for skewness and kurtosis. MicrobiomeGWAS runs very efficiently and takes only 36 min for testing main effects and 69 min for testing joint effects in a GWAS with 2000 subjects and 500,000 markers. Other statistical methods exist for testing the association of microbiome distance matrices. PERMANOVA [27] is an extension of multivariate analysis of variance to a matrix of pairwise distances and relies on permutations to evaluate significance. MiRKAT [49], a recently proposed method based on kernel regression, takes hours for evaluating one association for 2000 subjects. Neither is computationally feasible for analyzing a large-scale GWAS of a microbiome. Recently, an asymptotic distribution was proposed to approximate the  $p$ -value for the PERMANOVA pseudo-F statistic [50]; however, whether it is sufficiently accurate for very small  $p$ -values ( $p < 5 \times 10^{-8}$ , for GWAS) remains to be investigated.

Interactions of host genetic susceptibility with the microbiome have been postulated for many conditions, including inflammatory bowel diseases [51,52], autoimmune and rheumatic diseases [53–56], diabetes [57], and cancer, especially of the colon [58]. All models of these host–microbiome interactions also note the critical role of environmental factors, including diet, smoking, drugs, and antibiotics and other medications [59]. Although based on a very small initial sample set, the suggestive associations that we found between the six known lung cancer risk SNPs and the microbiome of adjacent normal lung tissue samples, including effects of cigarette smoking, provide preliminary evidence that our microbiomeGWAS method is likely to be a useful tool for generating data that will unravel host–microbiome interactions with high confidence.

We are working on two extensions for microbiomeGWAS: (1) jointly testing additive and dominant effects; and (2) testing genetic associations using many microbiome distance matrices. We have assumed an additive effect model (Figure 1); however, several top SNPs in the EAGLE data suggest a dominant effect (e.g., rs8083714 in Figure 5E). Thus, a statistic for jointly testing the additive and dominant effects might be powerful for this scenario. The second extension is motivated by the fact that the power to detect associations depends heavily on the choice of distance matrix. The recently developed generalized UniFrac [26] (gUniFrac) defines a series of distance matrices to reflect the different emphases of using taxa relative abundance information. gUniFrac has been shown to have a robust power



for association studies [26]. Extending microbiomeGWAS to gUniFrac, however, requires solving two problems. First, the computational complexity is proportional to the number of distance matrices analyzed for associations, which can be addressed by implementing the algorithms using multithreading technology. Second, we need to derive accurate analytic approximations to the association  $p$ -values by correcting for the multiple testing introduced by many distance matrices. MiRKAT [49] has an option for using gUniFrac; however, intensive permutations are required to evaluate  $p$ -values.

In summary, GWAS of the microbiome of each body site has the potential to help one understand microbiome variation, to elucidate the biological mechanisms of genetic associations, to improve the power of identifying novel disease-associated genetic variants, and to improve the performance of genetic risk prediction. We expect our methods and software to be useful for large-scale GWAS of the human microbiome.

**Supplementary Materials:** The following supporting information can be downloaded at: <https://www.mdpi.com/article/10.3390/genes13071224/s1>, Table S1: Type-I error rates estimated based on  $10^8$  simulations. Table S2: Association P-values between lung cancer risk SNPs and microbiome composition in the EAGLE data. Figure S1: Quantile-quantile (QQ) plot for association  $p$ -values testing the joint effects (main effect and SNP by smoking interaction) using the unweighted UniFrac distance matrices. Figure S2: Derivation of the likelihood ratio statistic  $Q$  in (7) and (8). Figure S3: Calculations related with genetic dependence.

**Author Contributions:** Conceptualization, X.H. and J.S.; methodology, X.H. and J.S.; software, X.H., L.S. and J.S.; data analysis, X.H., L.S. and J.S.; investigation, all authors; EAGLE data resource, M.T.L.; writing—original draft preparation, X.H. and J.S.; writing—review and editing, all authors. All authors have read and agreed to the published version of the manuscript.

**Funding:** This research was supported by the NIH Intramural Research Program.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Informed consent was obtained from all subjects in the study.

**Data Availability Statement:** The genetic data for the EAGLE study can be accessed from dbGap with accession number phs000093.v2.p2. The American Gut Project data used for simulations can be obtained from <https://github.com/biocompare/American-Gut> (accessed on 25 May 2022).

**Acknowledgments:** This study utilized the high-performance computational capabilities of the Biowulf Linux cluster at the National Institutes of Health, Bethesda, MD. (<http://biowulf.nih.gov> (accessed on 25 May 2022)).

**Conflicts of Interest:** The authors declare no conflict of interest.

## Appendix A. Calculating Calculating $Var(G_{ij})$ , $Cov(G_{ij}, G_{ik})$ , $Var(\Delta_{ij})$ and $Cov(\Delta_{ij}, \Delta_{ik})$

We first calculate  $E(G_{ij})$ ,  $Var(G_{ij})$ , and  $Cov(G_{ij}, G_{ik})$ . Let  $p_t = P(g_i = t)$  with  $p_0, p_1, p_2 \geq 0$  and  $p_0 + p_1 + p_2 = 1$ . We can also assume the Hardy–Weinberg equilibrium and characterize the probabilities as the allele frequency:  $p_0 = (1 - f)^2$ ,  $p_1 = 2f(1 - f)$  and  $p_2 = f^2$ . Some algebra leads to

$$E(G_{ij}) = E|g_i - g_j| = \sum_{m,n \in \{0, 1, 2\}} p_m p_n |m - n| = 2p_0 p_1 + 2p_1 p_2 + 4p_0 p_2 \quad (A1)$$

$$Var(G_{ij}) = E(G_{ij}^2) - E(G_{ij})^2 = (2p_0 p_1 + 2p_1 p_2 + 8p_0 p_2) - (2p_0 p_1 + 2p_1 p_2 + 4p_0 p_2)^2 \quad (A2)$$

$$Cov(G_{ij}, G_{ik}) = p_1(1 - p_1) + 4p_0 p_2(1 + p_1) - (2p_0 p_1 + 2p_1 p_2 + 4p_0 p_2)^2 \quad (A3)$$

Now consider  $\Delta_{ij} = |g_i E_i - g_j E_j|$ . When  $E_i$  is binary,  $g_i E_i = 0, 1$  or  $2$ . Let  $p'_t = P(g_i E_i = t)$ . Then,  $E(\Delta_{ij})$ ,  $Var(\Delta_{ij})$ , and  $Cov(\Delta_{ij}, \Delta_{ik})$  can be calculated similarly using (A1)–(A3).

**Appendix B. Calculating  $\rho = Cor_0(Z_M, Z_I|D')$**

Let  $G'_{ij} = G_{ij} - EG_{ij}$  and  $\Delta'_{ij} = \Delta_{ij} - E\Delta_{ij}$ . We first calculate covariance under  $H_0$ :

$$Cov_0(S_M, S_I|D') = Cov_0\left(\sum_{i<j} d'_{ij}G'_{ij}, \sum_{m<n} d'_{mn}\Delta'_{mn}\right) = \sum_{i<j, m<n} d'_{ij}d'_{mn}Cov(G_{ij}, \Delta_{mn}).$$

When  $(i, j, m, n)$  are distinct,  $Cov(G_{ij}, \Delta_{mn}) = 0$ . Some algebra leads to

$$Cov_0(S_M, S_I|D') = \binom{N}{2}Cov(G_{ij}, \Delta_{ij})\mu_2 + 6\binom{N}{3}Cov(G_{ij}, \Delta_{ik})\mu_3 \tag{A4}$$

with  $\mu_2$  and  $\mu_3$  specified in (3) and (4). Combining (2), (5) and (A4), we have

$$\rho = \frac{Cov_0(S_M, S_I|D')}{\sqrt{Var_0(S_M|D')Var_0(S_I|D')}} \xrightarrow{N \rightarrow \infty} \frac{Cov(G_{ij}, \Delta_{ik})}{\sqrt{Cov(G_{ij}, G_{ik})Cov(\Delta_{ij}, \Delta_{ik})}} \tag{A5}$$

Equation (A5) suggests that the correlation is asymptotically independent of the microbiome distance matrix. In the real-data analyses, we found that (A5) was very accurate when sample size  $N \geq 50$ . The details of calculating  $Cov(G_{ij}, \Delta_{ij})$  and  $Cov(G_{ij}, \Delta_{ik})$  are provided in Supplemental Data.

**Appendix C. A Statistic for Testing  $H_0 : \beta_M = \beta_I = 0$  vs.  $H_1 : \beta_M > 0$  or  $\beta_I > 0$**

Denote  $Z = (Z_M, Z_I)^T$ . Under  $H_0$ ,  $Z \sim N(0, \Sigma)$  with  $\Sigma = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ . Let  $\xi_M = E_1Z_M \geq 0$  and  $\xi_I = E_1Z_I \geq 0$  be the non-centrality parameter of the two score statistics. Apparently, the original testing problem is equivalent for testing  $H_0 : \xi_M = \xi_I = 0$  vs.  $H_1 : \xi_M > 0$  or  $\xi_I > 0$ . Given the observed values  $(Z_M, Z_I)$ , the likelihood ratio statistic is simplified as

$$Q = Z^T \Sigma^{-1} Z - (Z - \xi)^T \Sigma^{-1} (Z - \xi) \tag{A6}$$

where  $\xi = (\xi_M, \xi_I)^T = \operatorname{arginf}_{\xi_M \geq 0, \xi_I \geq 0} Q$  (Figure S2A).

To simplify the optimization problem in (A6), we perform a linear transformation:  $Y^T = Z^T \Sigma^{-\frac{1}{2}}$  and  $v^T = \xi^T \Sigma^{-\frac{1}{2}}$ , where

$$\Sigma^{-\frac{1}{2}} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ -1 & 1 \end{pmatrix} \begin{pmatrix} 1/\sqrt{1-\rho} & 0 \\ 0 & 1/\sqrt{1+\rho} \end{pmatrix} \tag{A7}$$

Under this transformation,  $Q = Y^T Y - (Y - v)^T (Y - v)$  and can be interpreted as the difference of the square of two distances (Figure S2B). The original parameter space  $\{(\xi_M, \xi_I) : \xi_M \geq 0, \xi_I \geq 0\}$  is now transformed to  $\{(v_1, v_2) : v_2 \geq \theta v_1, v_2 \geq -\theta v_1\}$  with  $\theta = \sqrt{(1-\rho)/(1+\rho)}$ . Thus, the new parameter space is bounded by two lines represented by  $v_2 \geq \theta v_1$  and  $v_2 \geq -\theta v_1$ . We partition the 2D plane into four parts (see Figure S2B), identify  $v = \operatorname{arginf}_{v \in A_1} (Y - v)^T (Y - v)$  and calculate  $Q$ :

$$Q = \begin{cases} Y_1^2 + Y_2^2 & (Y_1, Y_2) \in A_1 \\ (Y_2 - Y_1/\theta)^2 / (1 + \theta^{-2}) & (Y_1, Y_2) \in A_2 \\ (Y_2 + Y_1/\theta)^2 / (1 + \theta^{-2}) & (Y_1, Y_2) \in A_3 \\ 0 & (Y_1, Y_2) \in A_4 \end{cases} \tag{A8}$$

We now perform an inverse transformation using matrix

$$\Sigma^{\frac{1}{2}} = \left[ \frac{1}{\sqrt{2}} \begin{pmatrix} \sqrt{1-\rho} & 0 \\ 0 & \sqrt{1+\rho} \end{pmatrix} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} \right] \tag{A9}$$

to return to the original parameter space. The four areas  $\{A_1, A_2, A_3, A_4\}$  under the original space are in Figure 2 and Figure S2C.

Tedious calculations show that  $(Y_2 + Y_1/\theta)^2 / (1 + \theta^{-2}) = (w_2 Z_M + w_1 Z_I)^2$  with  $w_1 = (\theta - 1/\theta)/2$  and  $w_2 = (\theta + 1/\theta)/2$ . Similarly,  $(Y_2 - Y_1/\theta)^2 / (1 + \theta^{-2}) = (w_1 Z_M + w_2 Z_I)^2$ . This proves (7). In addition,  $w_1 Z_M + w_2 Z_I \geq 0$  and  $w_1^2 + 2\rho w_1 w_2 + w_2^2 = 1$ ; thus,  $P\{(w_1 Z_M + w_2 Z_I)^2 > b^2\} = P\{w_1 Z_M + w_2 Z_I > b\} = P\{N(0,1) > b\}$ . This proves (8). The probabilities in (8) could also be calculated from Figure S2B:  $q_1 = 1/2 - (\arctan\theta)/\pi$ ,  $q_2 = q_3 = 1/4$ .

### Appendix D. Calculating Skewness and Kurtosis under $H_0$

By definition,  $\gamma = E_0(S_M^3 | D') / Var_0^{3/2}(S_M | D')$  and  $\kappa = E_0(S_M^4 | D') / Var_0^2(S_M | D') - 3$ . We first calculate  $E_0(S_M^3 | D')$ . Let  $G'_{ij} = G_{ij} - EG_{ij}$ . We have

$$E_0(S_M^3 | D') = E_0\left(\sum_{i < j} d'_{ij} G'_{ij}\right)^3 = \sum_{i < j, m < n, s < t} d'_{ij} d'_{mn} d'_{st} EG'_{ij} G'_{mn} G'_{st}.$$

Figure S3A lists all combinations of  $(i, j, m, n, s, t)$  with  $EG'_{ij} G'_{mn} G'_{st} \neq 0$ ; then

$$E_0(S_M^3 | D') = \binom{N}{2} \mu_4 EG_{ij}^4 + \binom{N}{3} (\mu_5 EG_{ij}^2 G'_{ik} + \mu_6 EG_{ij} G'_{jk} G'_{ik}) + \binom{N}{4} (\mu_7 EG_{ij} G'_{jk} G'_{kl} + \mu_8 EG_{ij} G'_{ik} G'_{il}),$$

where  $(\mu_4, \mu_5, \mu_6, \mu_7, \mu_8)$  are provided in Supplemental Data. Similarly,

$$E_0(S_M^4 | D') = E_0\left(\sum_{i < j} d'_{ij} G'_{ij}\right)^4 = \sum_{i < j, m < n, s < t, x < y} d'_{ij} d'_{mn} d'_{st} d'_{xy} EG'_{ij} G'_{mn} G'_{st} G'_{xy}.$$

Figure S3B lists combinations of  $(i, j, m, n, s, t, x, y)$  with  $EG'_{ij} G'_{mn} G'_{st} G'_{xy} \neq 0$ . Thus,

$$E_0(S_M^4 | D) = \binom{N}{2} \mu_9 EG_{ij}^4 + \binom{N}{3} (\mu_{10} EG_{ij}^3 G'_{ik} + \mu_{11} EG_{ij}^2 G'_{ik}^2 + \mu_{12} EG_{ij}^2 G'_{jk} G'_{ik}) + \binom{N}{4} (\mu_{13} EG_{ij}^2 G'_{jk} G'_{kl} + \mu_{14} EG_{ij} G'_{jk}^2 G'_{kl} + \mu_{15} EG_{ij}^2 G'_{ik} G'_{il} + \mu_{16} EG_{ij} G'_{jk} G'_{ik} G'_{il} + \mu_{17} EG_{ij} G'_{jk} G'_{kl} G'_{il} + \mu_{18} EG_{ij}^2 G'_{kl}^2) + \binom{N}{5} (\mu_{19} EG_{ij} G'_{jk} G'_{kl} G'_{lm} + \mu_{20} EG_{ij} G'_{ik} G'_{il} G'_{im} + \mu_{21} EG_{ij} G'_{ik} G'_{il} G'_{lm} + \mu_{22} EG_{ij} G'_{ik} G'_{lm}^2) + \binom{N}{6} \mu_{23} EG_{ij} G'_{ik} G'_{lm} G'_{ln}$$

The constants  $(\mu_9, \dots, \mu_{23})$  are dependent on  $D$  and are provided in Supplemental Data. Note that  $Var_0(S_M | D') \sim O(N^3)$ ,  $E_0(S_M^3 | D') \sim O(N^4)$ ; thus,  $\gamma \sim O(1/\sqrt{N})$ . Similarly, we can prove  $\kappa \sim O(1/N)$ .

### Appendix E. Improve $p$ -Value Approximations by Adjusting for Skewness and Kurtosis

We assume that  $E_0 Z = 0$ ,  $Var_0 Z = 1$ ,  $\gamma = E_0 Z^3$  and  $\kappa = E_0 Z^4 - 3$  under the original probability measure  $P_0$ . The tail probability  $P_0(Z > b)$  for a large value of  $b$  is sensitive to the non-normality of  $Z$ , characterized by  $\gamma$  and  $\kappa$ . We define a new probability measure by embedding to the exponential probability density

$$dP_\xi = \exp(\xi Z - \phi(\xi)) dP_0 \tag{A10}$$

where  $\phi(\xi) = \log E_0 \exp(\xi Z)$  is the log moment generating function. Note that  $\gamma = \phi'''(0)$  and  $\kappa = \phi''''(0)$ . Because  $E_0(Z) = 0$  and  $Var_0(Z) = 1$ , Taylor's expansion leads to  $\phi(\xi) \approx \xi^2/2 + \gamma \xi^3/6 + \kappa \xi^4/24$ . Under  $P_\xi$ , we have

$$E_\xi Z = \int Z dP_\xi = \phi'(\xi) \approx \xi + \frac{\gamma}{2} \xi^2 + \frac{\kappa}{6} \xi^3 \tag{A11}$$

and

$$\text{Var}_{\xi} Z = \phi''(\xi) \approx 1 + \gamma\xi + \frac{\kappa}{2}\xi^2 \quad (\text{A12})$$

We choose  $\xi$  such that  $E_{\xi} Z \approx b$  by numerically solving an equation

$$\xi + \frac{\gamma}{2}\xi^2 + \frac{\kappa}{6}\xi^3 = b \quad (\text{A13})$$

Under the probability measure  $P_{\xi}$ ,  $Z \sim N(b, \sigma^2)$  approximately with  $\sigma^2 = 1 + \gamma\xi + \kappa\xi^2/2$  in (A12). By the likelihood ratio identity and (A10), we have

$$P_0(Z > b) = E_0 I_{Z>b} = E_{\xi} \frac{dP_0}{dP_{\xi}} I_{Z>b} = E_{\xi} e^{\phi(\xi) - \xi Z} I_{Z>b} = e^{\phi(\xi)} E_{\xi} e^{-\xi Z} I_{Z>b} \quad (\text{A14})$$

Note that  $e^{-\xi Z}$  decays very fast when  $Z$  increases. Thus, the integral  $E_{\xi} e^{-\xi Z} I_{Z>b}$  does not heavily depend on the tail distribution of  $Z$ . Assuming  $Z \sim N(b, \sigma^2)$  under  $P_{\xi}$ , we can verify that

$$E_{\xi} e^{-\xi Z} I_{Z>b} = e^{-b\xi + \frac{\sigma^2 \xi^2}{2}} \Phi(-\sigma\xi) \quad (\text{A15})$$

Combining (A14) and (A15) gives  $P_0(Z > b) \approx \exp(\phi(\xi) - b\xi + \sigma^2 \xi^2 / 2) \Phi(-\sigma\xi)$ , which is further approximated as

$$P_0(Z > b) \approx \exp\left(-b\xi + \frac{1 + \sigma^2}{2}\xi^2 + \frac{\gamma}{6}\xi^3 + \frac{\kappa}{24}\xi^4\right) \Phi(-\sigma\xi),$$

because  $\phi(\xi) \approx \xi^2/2 + \gamma\xi^3/6 + \kappa\xi^4/24$  based on the Taylor expansion. This proves (9). If we correct skewness but assume kurtosis  $\kappa = 0$ , then  $\phi(\xi) \approx \xi^2/2 + \gamma\xi^3/6$ . We recalculate  $\xi$  by setting  $\kappa = 0$  in (A13) to derive  $\xi = \left(\sqrt{1 + 2\gamma b} - 1\right) / \gamma$ . This proves (10).

## References

- Turnbaugh, P.J.; Hamady, M.; Yatsunencko, T.; Cantarel, B.L.; Duncan, A.; Ley, R.E.; Sogin, M.L.; Jones, W.J.; Roe, B.A.; Affourtit, J.P.; et al. A core gut microbiome in obese and lean twins. *Nature* **2009**, *457*, 480–484. [[CrossRef](#)] [[PubMed](#)]
- Morgan, X.C.; Tickle, T.L.; Sokol, H.; Gevers, D.; Devaney, K.L.; Ward, D.V.; Reyes, J.A.; Shah, S.A.; LeLeiko, N.; Snapper, S.B.; et al. Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol.* **2012**, *13*, R79. [[CrossRef](#)]
- Ahn, J.; Sinha, R.; Pei, Z.; Dominianni, C.; Wu, J.; Shi, J.; Goedert, J.J.; Hayes, R.B.; Yang, L. Human gut microbiome and risk for colorectal cancer. *J. Natl. Cancer Inst.* **2013**, *105*, 1907–1911. [[CrossRef](#)] [[PubMed](#)]
- Goedert, J.J.; Jones, G.; Hua, X.; Xu, X.; Yu, G.; Flores, R.; Falk, R.T.; Gail, M.H.; Shi, J.; Ravel, J.; et al. Investigation of the Association Between the Fecal Microbiota and Breast Cancer in Postmenopausal Women: A Population-Based Case-Control Pilot Study. *J. Natl. Cancer Inst.* **2015**, *107*, djv147. [[CrossRef](#)] [[PubMed](#)]
- Lax, S.; Smith, D.P.; Hampton-Marcell, J.; Owens, S.M.; Handley, K.M.; Scott, N.M.; Gibbons, S.M.; Larsen, P.; Shogan, B.D.; Weiss, S.; et al. Longitudinal analysis of microbial interaction between humans and the indoor environment. *Science* **2014**, *345*, 1048–1052. [[CrossRef](#)]
- Wu, G.D.; Chen, J.; Hoffmann, C.; Bittinger, K.; Chen, Y.Y.; Keilbaugh, S.A.; Bewtra, M.; Knights, D.; Walters, W.A.; Knight, R.; et al. Linking long-term dietary patterns with gut microbial enterotypes. *Science* **2011**, *334*, 105–108. [[CrossRef](#)] [[PubMed](#)]
- Tong, M.; McHardy, I.; Ruegger, P.; Goudarzi, M.; Kashyap, P.C.; Haritunians, T.; Li, X.; Graeber, T.G.; Schwager, E.; Huttenhower, C.; et al. Reprogramming of gut microbiome energy metabolism by the FUT2 Crohn's disease risk polymorphism. *ISME J.* **2014**, *8*, 2193–2206. [[CrossRef](#)] [[PubMed](#)]
- Knights, D.; Silverberg, M.S.; Weersma, R.K.; Gevers, D.; Dijkstra, G.; Huang, H.; Tyler, A.D.; van Sommeren, S.; Imhann, F.; Stempak, J.M.; et al. Complex host genetics influence the microbiome in inflammatory bowel disease. *Genome Med.* **2014**, *6*, 107. [[CrossRef](#)]
- McKnite, A.M.; Perez-Munoz, M.E.; Lu, L.; Williams, E.G.; Brewer, S.; Andreux, P.A.; Bastiaansen, J.W.M.; Wang, X.; Kachman, S.D.; Auwerx, J.; et al. Murine Gut Microbiota Is Defined by Host Genetics and Modulates Variation of Metabolic Traits. *PLoS ONE* **2012**, *7*, e39191. [[CrossRef](#)]
- Benson, A.K.; Kelly, S.A.; Legge, R.; Ma, F.; Low, S.J.; Kim, J.; Zhang, M.; Oh, P.L.; Nehrenberg, D.; Hua, K.; et al. Individuality in gut microbiota composition is a complex polygenic trait shaped by multiple environmental and host genetic factors. *Proc. Natl. Acad. Sci. USA* **2010**, *107*, 18933–18938. [[CrossRef](#)]
- Goodrich, J.K.; Waters, J.L.; Poole, A.C.; Sutter, J.L.; Koren, O.; Blekhan, R.; Beaumont, M.; Van Treuren, W.; Knight, R.; Bell, J.T.; et al. Human genetics shape the gut microbiome. *Cell* **2014**, *159*, 789–799. [[CrossRef](#)]
- Davenport, E.R.; Cusanovich, D.A.; Michelini, K.; Barreiro, L.B.; Ober, C.; Gilad, Y. Genome-Wide Association Studies of the Human Gut Microbiota. *PLoS ONE* **2015**, *10*, e0140301. [[CrossRef](#)]



13. GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis: Multitissue gene regulation in humans. *Science* **2015**, *348*, 648–660. [[CrossRef](#)] [[PubMed](#)]
14. Battle, A.; Mostafavi, S.; Zhu, X.; Potash, J.B.; Weissman, M.M.; McCormick, C.; Haudenschild, C.D.; Beckman, K.B.; Shi, J.; Mei, R.; et al. Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome Res.* **2014**, *24*, 14–24. [[CrossRef](#)] [[PubMed](#)]
15. Bell, J.T.; Pai, A.A.; Pickrell, J.K.; Gaffney, D.J.; Pique-Regi, R.; Degner, J.F.; Gilad, Y.; Pritchard, J.K. DNA methylation patterns associate with genetic and gene expression variation in HapMap cell lines. *Genome Biol.* **2011**, *12*, R10. [[CrossRef](#)]
16. Shi, J.; Marconett, C.N.; Duan, J.; Hyland, P.L.; Li, P.; Wang, Z.; Wheeler, W.; Zhou, B.; Campan, M.; Lee, D.S.; et al. Characterizing the genetic basis of methylome diversity in histologically normal human lung tissue. *Nat. Commun.* **2014**, *5*, 3365. [[CrossRef](#)]
17. McVicker, G.; van de Geijn, B.; Degner, J.F.; Cain, C.E.; Banovich, N.E.; Raj, A.; Wellen, N.; Myrthil, M.; Gilad, Y.; Pritchard, J.K. Identification of Genetic Variants That Affect Histone Modifications in Human Cells. *Science* **2013**, *342*, 747–749. [[CrossRef](#)]
18. Kilpinen, H.; Waszak, S.M.; Gschwind, A.R.; Raghav, S.K.; Witwicki, R.M.; Orioli, A.; Migliavacca, E.; Wiederkehr, M.; Gutierrez-Arcelus, M.; Panousis, N.I.; et al. Coordinated Effects of Sequence Variation on DNA Binding, Chromatin Structure, and Transcription. *Science* **2013**, *342*, 744–747. [[CrossRef](#)]
19. Suhre, K.; Wallaschowski, H.; Raffler, J.; Friedrich, N.; Haring, R.; Michael, K.; Wasner, C.; Krebs, A.; Kronenberg, F.; Chang, D.; et al. A genome-wide association study of metabolic traits in human urine. *Nat. Genet.* **2011**, *43*, 565–569. [[CrossRef](#)]
20. Sabatti, C.; Service, S.K.; Hartikainen, A.L.; Pouta, A.; Ripatti, S.; Brodsky, J.; Jones, C.G.; Zaitlen, N.A.; Varilo, T.; Kaakinen, M.; et al. Genome-wide association analysis of metabolic traits in a birth cohort from a founder population. *Nat. Genet.* **2009**, *41*, 35–46. [[CrossRef](#)]
21. Caporaso, J.G.; Kuczynski, J.; Stombaugh, J.; Bittinger, K.; Bushman, F.D.; Costello, E.K.; Fierer, N.; Pena, A.G.; Goodrich, J.K.; Gordon, J.I.; et al. QIIME allows analysis of high-throughput community sequencing data. *Nat. Methods* **2010**, *7*, 335–336. [[CrossRef](#)] [[PubMed](#)]
22. Bray, J.R.; Curtis, J.T. An ordination of upland forest communities of southern Wisconsin. *Ecol. Monogr.* **1957**, *27*, 325–349. [[CrossRef](#)]
23. Lozupone, C.; Knight, R. UniFrac: A new phylogenetic method for comparing microbial communities. *Appl. Environ. Microbiol.* **2005**, *71*, 8228–8235. [[CrossRef](#)] [[PubMed](#)]
24. Lozupone, C.A.; Hamady, M.; Kelley, S.T.; Knight, R. Quantitative and qualitative  $\beta$  diversity measures lead to different insights into factors that structure microbial communities. *Appl. Environ. Microbiol.* **2007**, *73*, 1576–1585. [[CrossRef](#)] [[PubMed](#)]
25. Lozupone, C.; Hamady, M.; Knight, R. UniFrac—an online tool for comparing microbial community diversity in a phylogenetic context. *BMC Bioinform.* **2006**, *7*, 371. [[CrossRef](#)]
26. Chen, J.; Bittinger, K.; Charlson, E.S.; Hoffmann, C.; Lewis, J.; Wu, G.D.; Collman, R.G.; Bushman, F.D.; Li, H. Associating microbiome composition with environmental covariates using generalized UniFrac distances. *Bioinformatics* **2012**, *28*, 2106–2113. [[CrossRef](#)]
27. Anderson, M.J. A new method for non-parametric multivariate analysis of variance. *Austral. Ecol.* **2001**, *26*, 32–46.
28. Wang, J.; Thingholm, L.B.; Skieceviciene, J.; Rausch, P.; Kummel, M.; Hov, J.R.; Degenhardt, F.; Heinsen, F.A.; Ruhlemann, M.C.; Szymczak, S.; et al. Genome-wide association analysis identifies variation in vitamin D receptor and other host factors influencing the gut microbiota. *Nat. Genet.* **2016**, *48*, 1396–1406. [[CrossRef](#)]
29. Gevers, D.; Knight, R.; Petrosino, J.F.; Huang, K.; McGuire, A.L.; Birren, B.W.; Nelson, K.E.; White, O.; Methe, B.A.; Huttenhower, C. The Human Microbiome Project: A community resource for the healthy human microbiome. *PLoS Biol.* **2012**, *10*, e1001377. [[CrossRef](#)]
30. MicrobiomeGWAS. Available online: <https://github.com/lsncibb/microbiomeGWAS> (accessed on 30 May 2022).
31. Landi, M.T.; Consonni, D.; Rotunno, M.; Bergen, A.W.; Goldstein, A.M.; Lubin, J.H.; Goldin, L.; Alavanja, M.; Morgan, G.; Subar, A.F.; et al. Environment And Genetics in Lung cancer Etiology (EAGLE) study: An integrative population-based case-control study of lung cancer. *BMC Public Health* **2008**, *8*, 203. [[CrossRef](#)]
32. Landi, M.T.; Chatterjee, N.; Yu, K.; Goldin, L.R.; Goldstein, A.M.; Rotunno, M.; Mirabello, L.; Jacobs, K.; Wheeler, W.; Yeager, M.; et al. A genome-wide association study of lung cancer identifies a region of chromosome 5p15 associated with risk for adenocarcinoma. *Am. J. Hum. Genet.* **2009**, *85*, 679–691. [[CrossRef](#)]
33. Oksanen, J.; Blanchet, F.G.; Friendly, M.; Kindt, R.; Legendre, P.; McGlinn, D.; Minchin, P.R.; O’Hara, R.B.; Simpson, G.L.; Solymos, P.; et al. vegan: Community Ecology Package. R Package Version 2.5-7. 2020. Available online: <https://CRAN.R-project.org/package=vegan> (accessed on 30 May 2022).
34. Tu, I.P.; Siegmund, D.O. The maximum of a function of a Markov chain and application to linkage analysis. *Adv. Appl. Probab.* **1999**, *31*, 510–531. [[CrossRef](#)]
35. Siegmund, D. *Sequential Analysis: Tests and Confidence Intervals*; Springer: New York, NY, USA, 1985.
36. McDonald, D.; Hyde, E.; Debelius, J.W.; Morton, J.T.; Gonzalez, A.; Ackermann, G.; Aksenov, A.A.; Behsaz, B.; Brennan, C.; Chen, Y.; et al. American Gut: An Open Platform for Citizen Science Microbiome Research. *mSystems* **2018**, *3*, e00031-18. [[CrossRef](#)]
37. Rühlemann, M.C.; Hermes, B.M.; Bang, C.; Doms, S.; Moitinho-Silva, L.; Thingholm, L.B.; Frost, F.; Degenhardt, F.; Wittig, M.; Kässens, J.; et al. Genome-wide association study in 8956 German individuals identifies influence of ABO histo-blood groups on gut microbiome. *Nat. Genet.* **2021**, *53*, 147–155. [[CrossRef](#)]

38. DeSantis, T.Z.; Dubosarskiy, I.; Murray, S.R.; Andersen, G.L. Comprehensive aligned sequence construction for automated design of effective probes (CASCADE-P) using 16S rDNA. *Bioinformatics* **2003**, *19*, 1461–1468. [[CrossRef](#)]
39. Edgar, R.C.; Haas, B.J.; Clemente, J.C.; Quince, C.; Knight, R. UCHIME improves sensitivity and speed of chimera detection. *Bioinformatics* **2011**, *27*, 2194–2200. [[CrossRef](#)]
40. Purcell, S.; Neale, B.; Todd-Brown, K.; Thomas, L.; Ferreira, M.A.; Bender, D.; Maller, J.; Sklar, P.; de Bakker, P.I.; Daly, M.J.; et al. PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **2007**, *81*, 559–575. [[CrossRef](#)]
41. Devlin, B.; Roeder, K. Genomic control for association studies. *Biometrics* **1999**, *55*, 997–1004. [[CrossRef](#)]
42. Hung, R.J.; McKay, J.D.; Gaborieau, V.; Boffetta, P.; Hashibe, M.; Zaridze, D.; Mukeria, A.; Szeszenia-Dabrowska, N.; Lissowska, J.; Rudnai, P.; et al. A susceptibility locus for lung cancer maps to nicotinic acetylcholine receptor subunit genes on 15q25. *Nature* **2008**, *452*, 633–637. [[CrossRef](#)]
43. McKay, J.D.; Hung, R.J.; Gaborieau, V.; Boffetta, P.; Chabrier, A.; Byrnes, G.; Zaridze, D.; Mukeria, A.; Szeszenia-Dabrowska, N.; Lissowska, J.; et al. Lung cancer susceptibility locus at 5p15.33. *Nat. Genet.* **2008**, *40*, 1404–1406. [[CrossRef](#)]
44. Thorgeirsson, T.E.; Geller, F.; Sulem, P.; Rafnar, T.; Wiste, A.; Magnusson, K.P.; Manolescu, A.; Thorleifsson, G.; Stefansson, H.; Ingason, A.; et al. A variant associated with nicotine dependence, lung cancer and peripheral arterial disease. *Nature* **2008**, *452*, 638–642. [[CrossRef](#)] [[PubMed](#)]
45. Amos, C.I.; Wu, X.; Broderick, P.; Gorlov, I.P.; Gu, J.; Eisen, T.; Dong, Q.; Zhang, Q.; Gu, X.; Vijayakrishnan, J.; et al. Genome-wide association scan of tag SNPs identifies a susceptibility locus for lung cancer at 15q25.1. *Nat. Genet.* **2008**, *40*, 616–622. [[CrossRef](#)] [[PubMed](#)]
46. Wang, Y.; Broderick, P.; Webb, E.; Wu, X.; Vijayakrishnan, J.; Matakidou, A.; Qureshi, M.; Dong, Q.; Gu, X.; Chen, W.V.; et al. Common 5p15.33 and 6p21.33 variants influence lung cancer risk. *Nat. Genet.* **2008**, *40*, 1407–1409. [[CrossRef](#)]
47. Shi, J.; Chatterjee, N.; Rotunno, M.; Wang, Y.; Pesatori, A.C.; Consonni, D.; Li, P.; Wheeler, W.; Broderick, P.; Henrion, M.; et al. Inherited variation at chromosome 12p13.33, including RAD52, influences the risk of squamous cell lung carcinoma. *Cancer Discov.* **2012**, *2*, 131–139. [[CrossRef](#)]
48. Timofeeva, M.N.; Hung, R.J.; Rafnar, T.; Christiani, D.C.; Field, J.K.; Bickeboller, H.; Risch, A.; McKay, J.D.; Wang, Y.; Dai, J.; et al. Influence of common genetic variation on lung cancer risk: Meta-analysis of 14 900 cases and 29 485 controls. *Hum. Mol. Genet.* **2012**, *21*, 4980–4995. [[CrossRef](#)]
49. Zhao, N.; Chen, J.; Carroll, I.M.; Ringel-Kulka, T.; Epstein, M.P.; Zhou, H.; Zhou, J.J.; Ringel, Y.; Li, H.; Wu, M.C. Testing in Microbiome-Profiling Studies with MiRKAT, the Microbiome Regression-Based Kernel Association Test. *Am. J. Hum. Genet.* **2015**, *96*, 797–807. [[CrossRef](#)]
50. Chen, J.; Zhang, X. D-MANOVA: Fast distance-based multivariate analysis of variance for large-scale microbiome association studies. *Bioinformatics* **2021**, *38*, 286–288. [[CrossRef](#)]
51. Leone, V.A.; Cham, C.M.; Chang, E.B. Diet, gut microbes, and genetics in immune function: Can we leverage our current knowledge to achieve better outcomes in inflammatory bowel diseases? *Curr. Opin. Immunol.* **2014**, *31*, 16–23. [[CrossRef](#)]
52. Huang, H.; Vangay, P.; McKinlay, C.E.; Knights, D. Multi-omics analysis of inflammatory bowel disease. *Immunol. Lett.* **2014**, *162*, 62–68. [[CrossRef](#)] [[PubMed](#)]
53. Troncone, R.; Discepolo, V. Celiac disease and autoimmunity. *J. Pediatr. Gastroenterol. Nutr.* **2014**, *59* (Suppl. S1), S9–S11. [[CrossRef](#)]
54. Yeoh, N.; Burton, J.P.; Suppiah, P.; Reid, G.; Stebbings, S. The role of the microbiome in rheumatic diseases. *Curr. Rheumatol. Rep.* **2013**, *15*, 314. [[CrossRef](#)]
55. Sparks, J.A.; Costenbader, K.H. Genetics, environment, and gene-environment interactions in the development of systemic rheumatic diseases. *Rheum. Dis. Clin. N. Am.* **2014**, *40*, 637–657. [[CrossRef](#)]
56. Smith, J.A. Update on ankylosing spondylitis: Current concepts in pathogenesis. *Curr. Allergy Asthma Rep.* **2015**, *15*, 489. [[CrossRef](#)]
57. Nielsen, D.S.; Krych, L.; Buschard, K.; Hansen, C.H.; Hansen, A.K. Beyond genetics. Influence of dietary factors and gut microbiota on type 1 diabetes. *FEBS Lett.* **2014**, *588*, 4234–4243. [[CrossRef](#)]
58. Birt, D.F.; Phillips, G.J. Diet, genes, and microbes: Complexities of colon cancer prevention. *Toxicol. Pathol.* **2014**, *42*, 182–188. [[CrossRef](#)]
59. Marietta, E.; Rishi, A.; Taneja, V. Immunogenetic control of the intestinal microbiota. *Immunology* **2015**, *145*, 313–322. [[CrossRef](#)] [[PubMed](#)]