

1 **Microbial Risk Score for Capturing Microbial Characteristics, Integrating Multi-omics**
2 **Data, and Predicting Disease Risk**

3
4 Chan Wang¹, Leopoldo N. Segal², Jiyuan Hu¹, Boyan Zhou¹, Richard Hayes³, Jiyoung Ahn³, Huilin Li^{1*}

5 ¹Division of Biostatistics, Department of Population Health, New York University Grossman School of
6 Medicine, New York, 10016, NY, USA

7 ²Division of Pulmonary and Critical Care Medicine, New York University Grossman School of Medicine,
8 New York, 10017, NY, USA

9 ³Division of Epidemiology, Department of Population Health, New York University Grossman School of
10 Medicine, New York, 10016, NY, USA

11 *Correspondence: Huilin.Li@nyulangone.org

12
13 Emails: Chan Wang: Chan.Wang@nyulangone.org, Leopoldo N. Segal: Leopoldo.Segal@nyumc.org,

14 Jiyuan Hu: Jiyuan.Hu@nyulangone.org, Boyan Zhou: Boyan.Zhou@nyulangone.org, Richard Hayes:

15 Richard.B.Hayes@nyulangone.org, Jiyoung Ahn: Jiyoung.Ahn@nyulangone.org, Huilin Li:

16 Huilin.Li@nyulangone.org

17 **Abstract**

18 **Background:** With the rapid accumulation of microbiome-wide association studies, a great amount of
19 microbiome data are available to study the microbiome's role in human disease and advance the
20 microbiome's potential use for disease prediction. However, the unique features of microbiome data hinder
21 its utility for disease prediction.

22 **Methods:** Motivated from the polygenic risk score framework, we propose a microbial risk score (MRS)
23 framework to aggregate the complicated microbial profile into a summarized risk score that can be used to
24 measure and predict disease susceptibility. Specifically, the MRS algorithm involves two steps: 1)
25 identifying a sub-community consisting of the signature microbial taxa associated with disease, and 2)
26 integrating the identified microbial taxa into a continuous score. The first step is carried out using the
27 existing sophisticated microbial association tests and pruning and thresholding method in the discovery
28 samples. The second step constructs a community-based MRS by calculating alpha diversity on the
29 identified sub-community in the validation samples. Moreover, we propose a multi-omics data integration
30 method by jointly modeling the proposed MRS and other risk scores constructed from other omics data in
31 disease prediction.

32 **Results:** Through three comprehensive real data analyses using the NYU Langone Health COVID-19
33 cohort, the gut microbiome health index (GMHI) multi-study cohort, and a large type 1 diabetes cohort
34 separately, we exhibit and evaluate the utility of the proposed MRS framework for disease prediction and
35 multi-omics data integration. In addition, the disease-specific MRSs for colorectal adenoma, colorectal
36 cancer, Crohn's disease, and rheumatoid arthritis based on the relative abundances of 5, 6, 12, and 6
37 microbial taxa respectively are created and validated using the GMHI multi-study cohort. Especially,
38 Crohn's disease MRS achieves AUCs of 0.88 ([0.85-0.91]) and 0.86 ([0.78-0.95]) in the discovery and
39 validation cohorts, respectively.

40 **Conclusions:** The proposed MRS framework sheds light on the utility of the microbiome data for disease
41 prediction and multi-omics integration, and provides great potential in understanding the microbiome's role
42 in disease diagnosis and prognosis.

43 **Keywords:** Alpha diversity; Disease prediction; Microbiome-wide association studies; Microbial risk score;
44 Multi-omics data integration; Sub-community

45

46 **Background**

47 Recent microbiome-wide association studies (MWASs) have uncovered that microbiome plays a crucial
48 role in human health and disease [1-4], with linkage of microbiota dysbiosis to a variety of complex diseases,
49 including diabetes, cardiovascular and mental disease, and cancer [5-12]. These studies provide great
50 opportunities to study microbiome's role in disease prediction, which, however, is challenging because of
51 its unique data structure.

52 Rapid advances in high-throughput sequencing technologies identify diverse microorganisms in a single
53 sample by targeted sequencing of their unique 16S rRNA gene, or shotgun sequencing of the collective
54 genomes of all microbes. For 16S rRNA sequencing data, QIIME 2 [13] is commonly used to assign the
55 sequencing reads to amplicon sequence variants or clustered operational taxonomic units based on the
56 similarity of sequences. For shotgun sequencing data, MetaPhlAn [14] or StrainPhlAn [15] can be used to
57 map the sequencing reads to species/strains against a reduced set of clade-specific marker sequences. Either
58 method produces the count or relative abundance table which typically contain hundreds to thousands of
59 taxonomic or functional features, i.e. microbiome data are high-dimensional, especially compared to the
60 available number of samples in most existing studies. In addition, these feature tables are usually sparse
61 with excessive zero counts, compositional with a sum constrained to a constant, and heterogeneous with a
62 phylogenetic tree structure to reveal the evolutionary relationship among the taxa. How to deal with these

63 unique characteristics of microbiome data and effectively utilize them in predicting disease risk is
64 challenging and needs comprehensive explorations and validations.

65 Polygenic risk score (PRS), a continuous score of an individual's genetic liability to a complex disease or
66 phenotype, has become more routine and powerful in current genomic research [16, 17]. PRS aggregates
67 the results from genome-wide association studies (GWASs) and is defined as the sum of risk alleles linked
68 to a phenotype of interest weighted by the corresponding effect sizes. The construction of PRS involves
69 two key steps: determining the risk alleles and their effect sizes based on discovery samples or published
70 GWASs, and calculating the PRS for each subject in the target population. The PRS framework motivates
71 us to construct a similar microbial risk score (MRS) to summarize the disease-specific microbial profile in
72 the increasing large-scale population-based microbiome studies [11, 18, 19] and to investigate its potential
73 in disease prediction. However, microbiome's unique community features make MWASs differ from
74 GWASs. First, the microbiota is a complex ecosystem, whose dynamics are driven by the interactions
75 among microbes and between microbes and their host. The link between this complex ecosystem and
76 disease process often involves interwoven mechanisms [20]. Further, the microbiota is composed of various
77 sub-communities related to different traits [21, 22], and its influence on disease development may act at the
78 community rather than the single-microbe level [23]. Thus, it is less informative or efficient to simply define
79 MRS as the weighted sum of the relative abundances of the associated microbes. Instead, we propose a
80 community-based MRS by calculating alpha diversity on a sub-community with member taxa identified as
81 being associated with the study trait. Alpha diversity is the diversity in a single ecosystem or sample with
82 respect to its richness, evenness, or both characteristics [24, 25]. Several indices, including Observed,
83 Simpson, Shannon, and Faith's phylogenetic diversity (PD), have been extensively used to characterize
84 microbial community. With the NYU Langone Health (NYULH) COVID-19 cohort [26] and the gut
85 microbiome health index (GMHI) multi-study cohort [27], we propose and validate a few MRSs on
86 COVID-19, colorectal adenoma (CA), colorectal cancer (CC), Crohn's disease (CD), and rheumatoid
87 arthritis (RA) to exhibit the utility of the proposed MRS framework.

88 With the recent advances in the next-generation sequencing and mass spectrometry, there is a growing need
89 for the ability to merge biological features to study an ecosystem as a whole. Aspects such as the
90 metagenome, metatranscriptome, host genome, host gene expression, and metabolome each provides a
91 snapshot of one level of regulation in a system. The proposed MRS framework provides a simple and
92 interpretable approach to integrate the microbial profiles with other biological omics data and elucidate the
93 microbial interactions with other omics datasets in the disease prediction. We use the NYULH COVID-19
94 cohort, which characterized the lung microbiome in a large prospective cohort of critically ill patients with
95 SARS-CoV-2 infection who required invasive mechanical ventilation, to illustrate, evaluate, and validate
96 the proposed MRS and its integrations with other omics data in the prediction for COVID-19 mortality. In
97 addition, we elucidate the join effect of MRS and PRS on T1D risk stratification using the Environmental
98 Determinants of Diabetes in the Young (TEDDY) study (<https://teddy.epi.usf.edu/>) [28-30].

99 **Methods**

100 **MRS framework**

101 **MRS workflow.** We propose a microbial risk score framework to convert the high-dimensional
102 microbiome data into a summarized risk score that can measure and predict disease susceptibility. As
103 illustrated in Figure 1, with the ready-for-downstream-analysis microbial data, the microbial risk score
104 algorithm involves two key steps: 1) to identify a sub-community consisting of the signature microbial taxa
105 associated with disease, and 2) to integrate the identified microbial taxa into a continuous score.

106 **Microbial signature identification.** We propose to employ the existing sophisticated microbial association
107 tests [1, 3, 4, 31-33] to identify microbial taxa associated with disease using the discovery samples. Great
108 amount of abundance-based methods examining the difference of microbial abundance directly, which is
109 also called differential abundance (DA) analysis [31-39] have been proposed recently. Based on the results
110 in two recent benchmarking works [32, 33], ANCOM-BC (Analysis of compositions of microbiomes with
111 bias correction) [31] is one of the top-performing methods and has been widely used in microbiome research.

112 ANCOM-BC [31] models the observed abundances using an offset-based log-linear model, in which the
113 offset term is sample-specific to account for sampling fraction. We use it as the default microbial association
114 test to identify the candidate taxa in the first step of our microbial risk score algorithm. Considering
115 developing novel differential abundance test is still an active area of research, in the Discussion section, we
116 discuss the performance of the proposed MRS framework with other DA tests.

117 In addition to the above-mentioned statistical methods, a variety of machine learning (ML) techniques have
118 been applied in microbiome studies for microbial feature selection, biomarker identification, disease
119 prediction and classification, as recently reviewed [40]. As an example, Gou et al. [41] defined an MRS
120 with the microbiome features selected by the Light Gradient Boosting Machine method [42] and examined
121 its association with type 2 diabetes (T2D) as well as T2D-related traits. Despite the visible contributions of
122 characterizing the microbial profiles and uncovering the relationship between microbiome and disease, the
123 applications of ML methods including traditional methods and deep learning techniques in the microbiome
124 studies share several drawbacks [40]. One is that most ML methods input all available microbial features
125 into the model to determine the final output solely based on algorithms, without considering the inherent
126 structure of microbiome data, such as compositionality and zero inflation. Another unavoidable drawback
127 of ML methods is the model instability in the relatively small-scale biomedical human studies [43]. Because
128 the nature of ML algorithms is to learn the pattern by training the data, they usually require a large sample
129 size to reach stable results, especially for the algorithms involving various parameters or various layers that
130 need to be trained via cross-validation (CV). Given these common pitfalls and relatively small sample size
131 in biomedical studies due to the high cost of patients' in-person visit, sample collection and sequencing,
132 ML's application in microbiome research may provide inexplicable results and even lead to the loss of
133 statistical power. With the NYULH COVID-19 cohort example, we illustrate the inefficient utility of ML
134 methods in analyzing the microbiome data compared to the proposed MRS method. The details are reported
135 in the Results section.

136 **Sub-community determination.** Pruning and thresholding (P+T) method is a heuristic approach
137 commonly used in PRS studies for identification of genetic variants based on an empirically determined p -
138 value threshold [44]. We propose to use P+T method to determine the final candidate microbial taxa in
139 discovery cohort. Specifically, we calculated a series of MRSs proposed below using the nested sets of
140 microbial taxa with the increasingly relaxed significance thresholds. We set the final threshold at the value
141 that produced the largest area under the receiver operating characteristic (ROC) curve (AUC). All the taxa
142 whose p -values are less than the final threshold form a disease or trait specific sub-community. If there is
143 only one dataset available, CV will be used to determine the sub-community along with P+T method. More
144 details are provided in the Results section.

145 **MRS calculation.** We propose an MRS, denoted by MRS_{α} , which is defined as the alpha diversity of the
146 sub-community consisting of the identified candidate taxa. Alpha diversity is the diversity in a single
147 ecosystem or sample with respect to its richness, evenness, or both characteristics [24, 25]. The core concept
148 of alpha diversity index in biology is to find the effective number of elements of a system to measure its
149 complexity or diversity [45]. Note that multiple alpha diversity indices are available. Some measure species
150 richness such as observed index, Chao1, and ACE. PD is a phylogenetic metric which is defined as the sum
151 of the lengths of all those branches on the tree that span the members of the set. Simpson index is a
152 dominance index as it gives more weight to the common or dominant species and does not account for
153 species richness. While Shannon index is an information statistic index (entropy) which accounts for both
154 species richness and its evenness in a community or sample, and it has a unique ability to weigh taxa by
155 their frequency, without disproportionately favoring either rare or common elements. As the most popular
156 and accepted index for diversity [46], we adopt Shannon index in the proposed MRS_{α} . Other indices are
157 also investigated in the Discussion section and included in the MRS framework (MRS R package).

158 Suppose there are n samples (each sample represents one ecosystem or microbial community) and Q taxa.

159 Let M_{ij} be the relative abundance of the j th taxon in the i th sample with the constraint $\sum_{ij=1}^Q M_{ij} = 1, i =$

160 $1, \dots, n$, and $j = 1, \dots, Q$. Assume p ($< Q$) taxa are identified as a sub-community to construct MRSs.

161 Without loss of generality, we assume that the first p taxa are the identified candidate taxa.

162 For the i th sample, its MRS_{α} is calculated as $MRS_{\alpha}^i = \sum_{j=1}^p \tilde{M}_{ij} \ln(\tilde{M}_{ij})$, where \tilde{M}_{ij} is relative abundance

163 of the j th identified candidate taxon within the sub-community for the i th sample ($\tilde{M}_{ij} = \frac{M_{ij}}{\sum_{j=1}^p M_{ij}}$).

164 MRS_{α} is constructed based on the Shannon index [24, 25] without the negative sign, so that the smaller is

165 MRS_{α} , the healthier is the microbial community [47]. As a comparison, we also derive a standard MRS as

166 an analogy to PRS, denoted by MRS_S . It is a (weighted) sum of relative abundances of the identified

167 candidate taxa as $MRS_S^i = \sum_{j=1}^p w_j M_{ij}$, where w_j is the weight for the j th taxon. We propose two sets of

168 weights: all weights are equal to 1 (denoted by MRS_{unwS}^i); and the weights are the effect sizes estimated

169 from the training or discovery samples by certain microbial association method (denoted by MRS_{wS}^i).

170 Noticeably, MRS_{α} integrates p identified taxa as a community by measuring its diversity. While, MRS_S

171 focuses on the additive effect of the identified taxa and doesn't account for the microbial community

172 feature.

173 **Validation.** The proposed MRSs need to be validated either by external validation or internal validation.

174 Since the GMHI multi-study cohort [27] has independent discovery and validation cohorts, the MRSs are

175 created using the discovery cohort and validated using the validation cohort. For the NYULH COVID-19

176 [26] and TEDDY studies [28-30], due to the lack of independent additional samples, we employ CV to

177 perform independent internal validation.

178 **Risk score-based multi-omics data integration**

179 Note that the proposed MRS summarizes a complex microbial profile into a quantifiable score, which

180 provides a fast and flexible way to integrate microbiome data with other omics data to better predict disease

181 risk. Both the NYULH COVID-19 and TEDDY studies contain not only microbial profile data, but also

182 other omics data. We propose to jointly model MRS and other risk scores built on other omics data to

183 further improve the performance of disease prediction. In the COVID study, on one hand, the enrichment
184 of SARS-CoV-2 and some oral commensals in the lower-airway microbiota are associated with poor
185 outcome, and on the other hand, host lower-airway immune phenotypes reveal a failure of adaptive and
186 innate immune response to SARS-CoV-2 among deceased subjects. Jointly modeling these omics profiles
187 can improve the predictive accuracy of mortality. For the TEDDY study, since that genotype data in the
188 regions containing autoimmunity and inflammatory response genes are available, one can build a PRS for
189 each subject using the existing PRS algorithms [48-50]. By combining the PRS and the proposed MRS, we
190 can jointly model the association of genetic and environmental risk in T1D prediction.

191 **Prediction performance evaluation**

192 With the constructed risk scores from various omics data, one can employ a logistic regression model for
193 the prediction of disease status(binary outcome), or a Cox proportional-hazards model [51] for the
194 prediction of disease onset (survival outcome). Prediction performance can be evaluated by AUC for
195 binary outcome or by hazard ratio (HR) for survival outcome. The additive model can be used to integrate
196 multiple risk scores in these two regression models. The interaction terms between scores can be explored
197 further for risk stratification [52], as illustrated in the TEDDY study in Result section.

198 **NYULH COVID-19 cohort**

199 The NYULH COVID-19 cohort [26] includes a subset of 142 patients with COVID-19, at the NYULH
200 Manhattan campus from March 3 to June 18, 2020, who required invasive mechanical ventilation and
201 underwent bronchoscopy for airway clearance and/or tracheostomy. Among all patients, 108 (76%)
202 survived hospitalization and 34 (24%) died. The study has collected and processed lower-airway samples
203 and performed: a) metagenomic sequencing for bacterial, fungal and DNA viral genomes; and b)
204 metatranscriptome assays for viral, bacterial, fungal, and human transcriptomes and the RNA virome. In
205 addition, comprehensive demographic, longitudinal clinical, and treatment data are available.

206 **GMHI multi-study cohort**

207 An integrated dataset of 4,347 human stool metagenomics samples (cross-sectional) from 34 published
208 studies (discovery cohort) and an independent dataset of 679 samples from 9 additional studies (validation
209 cohort) are publicly available [27]. Both cohorts consist of healthy subjects and patients with various
210 diseases. Using these two cohorts, Gupta et al. [27] introduced and validated the gut microbiome health
211 index (GMHI) to quantify the likelihood of disease presence based on subject's gut microbiome data. In
212 both cohorts, they pooled samples from different disease conditions together into one nonhealthy group,
213 and the proposed GMHI exclusively identifies the difference of microbiome profile between healthy and
214 non-healthy samples. After the pre-processing and quality control, there are 2,636 healthy and 1,711
215 nonhealthy samples in the discovery cohort and 118 healthy and 561 nonhealthy samples in the validation
216 cohort respectively. Among nonhealthy samples, discovery and validation cohorts both have samples from
217 patients with CA, CC, CD, and RA. Sample sizes are shown in Table S1. For microbiome data, there are
218 313 species and 576 species in the discovery and validation cohorts respectively available for analysis.
219 More details are described in Gupta et al. [27].

220 **TEDDY cohort**

221 TEDDY is a large-scale prospective study designed to identify the genetic and environmental triggers that
222 cause childhood T1D [28-30]. Children with high genetic risk for islet autoimmunity or T1D were enrolled
223 and multiple biomarkers were assessed longitudinally for prediction of T1D development. A total of 12,005
224 fecal samples from 903 children, collected from 3 to 46 months of age, were characterized by 16S rRNA
225 sequencing. Of this cohort, 114 children were ascertained to T1D by year 5 [29]. The findings in the
226 previous TEDDY publications [53, 54] focus exclusively on the microbiome profiles, and suggest that the
227 gut microbiome data may have the potential to predict the progression of T1D. In addition to microbiome
228 data and metadata, the TEDDY cohort also includes genomic, longitudinal metabolomic, and host
229 transcriptomic data which together provide opportunity to explore the integrated information from multiple
230 aspects on the pathogenesis of T1D through the multi-omics analysis.

231 **Results**

232 **Evaluation and validation of MRS framework**

233 **NYULH COVID-19 cohort.** With the same quality control, sequencing process, and filtering criteria
234 described in Sulaiman et al. [26], we analyzed data from 118 patients (28 Deceased and 118 Alive) who
235 had all metagenome, metatranscriptome, and host transcriptome samples. We included 374 taxa in
236 metagenome, 1,149 taxa in metatranscriptome, and 14,697 genes in host transcriptome data for our
237 analyses. We used the binary outcome (Deceased vs. Alive) to illustrate the predictive performance of
238 MRS here.

239 Figure 2 presents the optimal p -value thresholds (0.42, 0.38, and 0.02) used to identify the associated taxa
240 in MRSs (MRS_{α} , MRS_{wS} , and MRS_{unwS} , respectively) using the metagenomic data. The optimal thresholds
241 were determined by P+T method as described in the sub-section “Sub-community determination” using the
242 leave-one-out CV. With the optimal p -value cutoffs, the community-based MRS_{α} has the best performance
243 in predicting deceased/alive status (AUC=0.74), compared to two summation-based standard MRSs:
244 MRS_{wS} (AUC=0.72) and MRS_{unwS} (AUC=0.70). This reflects that analyzing the microbial profile as a
245 community can characterize more microbial information and work better than analyzing microbes
246 individually. Additionally, MRS_{wS} performs better than MRS_{unwS} , as expected, since MRS_{wS} incorporates
247 the strength of the association effects of taxa on the outcome, as well as the microbial relative abundances,
248 while MRS_{unwS} is just the summation of the microbial relative abundances from the selected taxa.

249 Figure S1 shows prediction performance for various ML algorithms which have been commonly applied in
250 microbiome research [40]. The leave-one-out CV was used for the predictions and the predicted probability
251 for deceased/alive status was used for ROC analysis. All ML algorithms have lower AUCs than the
252 proposed MRS_{α} . Among these ML algorithms, the ML algorithms based on regularization (Figure S1A) all
253 perform better with higher AUCs, compared to the ML algorithms that have various tuning parameters or
254 layers (Figure S1B). Elastic-net logistic regression and penalized discriminant analysis (regression-based)

255 algorithms have the best prediction performance. On the other hand, ML algorithms were also applied to
256 select the candidate taxa used for the construction of MRS_{α} based on the variable importance. The top K
257 features were determined based on leave-one-out CV. Take the elastic-net logistic regression which has the
258 best prediction above for example, the top 30 taxa were ultimately selected to construct MRS_{α} with the
259 AUC being the largest based on CV, and its AUC for deceased/alive status prediction is 0.66, which is 11%
260 lower than the AUC of the above MRS_{α} . The efficiency of ML algorithms is evidently limited due to the
261 small sample size and not being able to take care of the unique features of microbiome data, such as
262 compositionality and zero inflation.

263 In addition, we checked the prediction performance of the alpha diversity indices on the whole microbial
264 community in terms of AUC. Table S2 reports the AUC values for six common alpha diversity indices in
265 predicting alive and deceased status. All alpha diversity indices have similar prediction performance, with
266 AUC being 0.50 to 0.53, which are much poorer than the proposed MRS_{α} . Comparisons between MRS_{α}
267 and alpha diversity indices underline the significance of identification of the associated taxa in the microbial
268 risk score framework, which condenses the signal by excluding the non-associated taxa and provides full
269 potential for the proposed MRS to measure and predict disease susceptibility.

270 **GMHI multi-study cohort.** With the discovery and validation cohorts [27], we evaluated and validated
271 the proposed MRS_{α} in terms of predictive performance. Specifically, for CA, CC, CD, and RA diseases,
272 respectively, we performed ANCOM-BC to identify candidate species that were differentially abundant
273 between samples from healthy subjects and patients with this disease in the discovery cohort, constructed
274 disease-specific MRS_{α} based on the identified species, and performed the independent validation of
275 disease-specific MRS_{α} using samples from healthy subjects and patients with the disease in the validation
276 cohort.

277 Figure 3A presents that AUC values and 95% confidence intervals for MRS_{α} s to predict healthy and 4
278 different diseases in discovery and validation cohorts, respectively. Overall, MRS_{α} s achieve great

279 predictive performance in both discovery (AUCs: 0.60-0.88) and validation (AUCs: 0.68-0.86) cohorts.
280 Notably three MRS_{α} s (healthy vs. CA, healthy vs. CC, and healthy vs. RA) have higher AUCs in
281 validation cohort, compared to discovery cohort. Among these four disease-specific MRS_{α} s, MRS_{α}
282 specific for CD disease has the best predictive performance (AUC=0.88 in discovery and AUC=0.86 in
283 validation). In addition, different MRS_{α} s are constructed by different identified taxa. 5, 6, 12, and 6 taxa
284 are used for constructions of MRS_{α} s for CA, CC, CD, and RA, respectively (Figure 3B; Table S3).
285 Several taxa contribute multiple MRS_{α} s, for example, species *Bifidobacterium angulatum* is involved for
286 constructions of MRS_{α} s for CA, CC, and RA (Table S3). On the other hand, 21 taxa are disease-specific
287 and exclusively used in one MRS_{α} (Table S3). They are differentially abundant in Healthy, CA, CC, CD,
288 and RA samples (Tables S4 and S5). This demonstrates that the proposed MRS framework powerfully
289 improves disease prediction by incorporating the disease-specific microbial profile. This feature makes
290 the proposed MRS framework more crucial in practice, as most research studies aim to identify the
291 microbial taxa specifically playing a role in a certain disease, rather than the generalized disease-
292 associated microbial taxa.

293 Similar to disease-specific MRS, we also assessed the MRS framework that distinguishes two disease
294 groups, as well as healthy and nonhealthy conditions defined as in the original study [27] in the discovery
295 and validation cohorts, respectively. Figure S2 presents the AUC values and 95% confidence intervals for
296 MRS_{α} s to classify any two diseases of CA, CC, CD, and RA, and healthy and nonhealthy conditions in
297 discovery and validation cohorts, respectively. Table S3 correspondingly reports which taxa are involved
298 for these MRS_{α} calculations, respectively. Again, the MRS framework achieves notable performance. For
299 example, discovery cohort has AUCs of 0.91 and 0.89, meanwhile, validation cohort has AUCs of 0.84
300 and 0.84, to distinguish CD from RA and CC, respectively. Validation cohort has a relatively lower AUC
301 for classifying CA and RA, due to the small sample size. In terms of healthy vs. nonhealthy prediction,
302 MRS_{α} achieves consistently competitive performance but with much fewer species, whose AUCs are 0.7
303 and 0.71 in discovery and validation cohorts, respectively, compared to GMHI whose AUCs are 0.7 and

304 0.74 in discovery and validation cohorts, respectively. And the identified 6 species for MRS_{α} construction
305 is a subset of 50 microbial species used in GMHI [27].

306 **Results of risk score-based multi-omics data integration**

307 **NYULH COVID-19 cohort.** In addition to metagenome data, the NYULH COVID-19 cohort has
308 metatranscriptome and host transcriptome data. In the following, we present how to integrate metagenomic,
309 metatranscriptomic, and host transcriptomic datasets using the proposed MRS_{α} and the evaluation of
310 different methods. For the metatranscriptomic data, we employed the same MRS algorithm as we described
311 in the Methods section, in terms of determining the p -value cutoff, identifying candidate taxa, and
312 constructing the microbial risk score, to construct its MRS_{α} . In order to differentiate various MRS_{α} s, we
313 denoted the MRS_{α} using the metagenomic and metatranscriptomic data by DNA_MRS_{α} and RNA_MRS_{α} ,
314 respectively in the rest of manuscript. For the transcriptomic data, we employed DESeq2 [36] to evaluate
315 the association effects of genes on the deceased/alive status, determined the p -value cutoff based on the
316 P+T method, and identified the candidate genes by AUC evaluation. Then we defined the weighted sum of
317 log-transformed counts of the selected candidate genes for each sample as the risk score (denoted as Host),
318 with the weight being 1 if the corresponding logarithmic fold change estimate from DESeq2 was positive,
319 otherwise -1. Computational details are reported in Section S1. Figure 4A shows that the risk scores based
320 on metagenomic, metatranscriptomic, and host transcriptomic data separately have the AUC values of 0.74,
321 0.69, and 0.63, respectively, in terms of predicting deceased/alive status. Furthermore, the combinations of
322 risk scores from different datasets can obviously improve the predictive performance (Figure 4B) of
323 mortality. The combinations of any two datasets have comparable AUC values and perform similarly. As
324 expected, the integration of all three datasets ($DNA_MRS_{\alpha} + RNA_MRS_{\alpha} + Host$) has the highest AUC of
325 0.85, which yields at least a 15% increase in AUC compared to DNA_MRS_{α} , RNA_MRS_{α} , or Host alone.
326 In Figure 5, comparing the risk scores between the alive and deceased groups, the deceased group always
327 has a significantly higher average risk score than the alive group, no matter the score was constructed based
328 on a single omics dataset or the integration of different omics datasets (p -values<0.05).

329 Figure 6 presents the 2D or 3D scatterplots of risk scores from metagenomics, metatranscriptomic, and host
330 transcriptomic data. The subjects were first classified into “High risk” and “Low risk” groups by each risk
331 score’s mean. We next checked how well these risk classifications can be used to predict disease status by
332 reporting the classification metrics [55]: sensitivity, specificity, accuracy, and F1 score in Table 1.
333 Specifically, the predicted values for the subjects labeled as “High risk” by two risk scores (in Figures 6A-
334 C) or by all three risk scores (in Figure 6D) datasets are “Deceased”, and the predicted values for the
335 subjects labeled as “Low risk” by two risk scores (in Figures 6A-C) or by all three risk scores (in Figure
336 6D) datasets are “Alive”. From Table 1, we can see that among the combinations of two risk scores for
337 classification, the combination of metagenomic and host transcriptomic risk scores has the highest
338 sensitivity, accuracy and F1 score, but is still inferior to the combination of all three omics risk scores,
339 which identifies the mortality status with 86% sensitivity, 91% specificity, 88% accuracy, and an F1 score
340 of 0.89. In this real study, from different angles, including the AUC in Figure 4, the scatterplots of risk
341 scores in Figure 7, and the test results in Table 1, we show that combining risk scores from metagenomics,
342 metatranscriptomic, and host transcriptomic data increases the predictive accuracy for COVID-19 mortality.
343 Table S6 reports the included features in the metagenomic, metatranscriptomic, and host transcriptomic
344 risk scores separately. The feature importance was determined by the selection proportion among all CV
345 iterations. For the host transcriptomic data, the fold change between deceased and alive was used to
346 determine the feature importance when the selection proportions were the same. Here we take the top 50
347 features in each data as an illustration to investigate the correlation networks among these three datasets.
348 Figures 7, S3, and S4 show the paired correlation heatmaps among the selected metagenomic,
349 metatranscriptomic, and host transcriptomic features in the alive and deceased groups, respectively. Notably,
350 the alive and deceased groups have different correlation patterns among these top 50 features from any two
351 datasets. Specifically, the metagenomics features tend to have stronger correlations with the host
352 transcriptomic and metatranscriptomic features in the deceased group, compared to the alive group; and the

353 metatranscriptomic features tend to have more negative correlations with the host transcriptome in the alive
354 group.

355 Note that the results reported in this section are different from those in Sulaiman et al. [26] in which the
356 main goal was to reveal the scientific findings and the Cox proportional-hazards model [51] was employed
357 to identify the candidate taxa and genes associated with the time to death. In this paper, we formally
358 introduce the MRS concept and propose it as a general method with the detailed instruction on how to
359 construct MRS. As a validation of the proposed method, the results presented above based on the binary
360 outcome (Deceased vs. Alive) agree with the previous scientific conclusions [26]. Table 2 reports the hazard
361 ratios of all risk scores constructed in this paper and their combinations on the time to death based on the
362 Cox proportional-hazards model. All risk scores are significantly associated with the time to death. As we
363 found in [26], metatranscriptomic data alone, or combined with the other two datasets, always has a higher
364 hazard of death, because it involves SARS-CoV-2 viral, which is a key risk factor on the COVID-19
365 mortality.

366 Overall, these results highlight that the proposed community-based MRS_{α} , can characterize and summarize
367 the microbial profiles effectively and provide a flexible way to integrate microbiome data with other omics
368 data. Integrations of risk scores from different omics data further improves the predictive performance on
369 the alive/deceased status in the NYULH COVID-19 study.

370 **TEDDY study.** Although the TEDDY cohort includes both genome and microbiome data, the previous
371 microbiome research on TEDDY study [53, 54] focused exclusively on the microbiome profiles and only
372 identified very few microbial signatures associated with T1D. Given the fact that T1D is a multifactorial
373 disease caused by both genetic and environmental factors and the children enrolled in the TEDDY study
374 all have high genetic risk for T1D development (they have at least one of nine HLA DR-DQ genotypes
375 associated with high risk for T1D) [29], we here propose a new angle to employ the proposed MRS along
376 with the existing PRS for T1D to investigate the combined effect of microbial profile and host genetic
377 profile on T1D risk prediction. Specifically, we analyzed 551 TEDDY subjects who have both

378 microbiome data and genotype data; 75 of them developed T1D. Using the available genotype data and
379 the PRS algorithm which has the robust and superior prediction performance on T1D [48, 49], we built
380 the PRS for subjects. We used the microbial samples that were collected at the time point most close to
381 month 30 when microbiome profile got stable and the largest sample size was available, to build MRS_{α} to
382 predict T1D status independently. The practice of MRS calculations are the same as those used in the
383 NYULH COVID-19 study.

384 Figure 8A compares the AUCs for predicting T1D based on the individual risk scores and the combination
385 of PRS and MRS_{α} , and Figures 8B-D show the Kaplan–Meier survival curve comparisons between high
386 and low risk group identified by PRS, MRS_{α} and PRS + MRS_{α} respectively. Specifically, subjects whose
387 risk scores are above the third quartile are defined as high risk, others as low risk. Although the predictive
388 models considered in Figure 8A have only modest predictive ability in the TEDDY cohort (AUC range:
389 [0.58, 0.63]), we found that integrating PRS and MRS_{α} scores is more useful in stratifying the subjects into
390 high and low risk groups for T1D development (Figure 8D) than the PRS (Figure 8B) or MRS_{α} (Figure 8C)
391 alone, which indicates that the potential genetic-microbial interaction effect on the T1D progression. These
392 results exhibit the utility of modeling multi-omics risk scores to identify the high risk populations who can
393 benefit from more targeted interventions.

394 **Discussion**

395 With the recent proliferation of large-scale microbial association studies, we propose a two-step novel
396 microbial risk score framework to aggregate the high-dimensional microbiome profile into a summarized
397 risk score and apply it in disease prediction. Specifically, we first identify the associated taxa based on the
398 recommended microbial association tests by two recent benchmarking works [32, 33] and P+T method, and
399 then construct a community-based MRS_{α} , because that the microbiome is a complex ecosystem composed
400 of numerous sub-communities, and its influence on the disease development acts at the community instead
401 of the single-microbe level and is disease-dependent. The application in the NYULH COVID-19 cohort

402 demonstrates the superior performance of MRS_{α} in the disease prediction, compared to the standard MRS_S ,
403 which is constructed similarly as PRS, ML-based prediction algorithms, and six alpha diversity measures
404 on the whole microbiome community. The evaluation of MRS_{α} using the GMHI integrated dataset which
405 consists of independent discovery and validation cohorts reveals the notable reproducibility of MRS_{α} in
406 terms of disease prediction.

407 Combining omics datasets that provide biological information from different layers is vital to
408 comprehensively study phenotypes and accurately predict diseases. However, complex data structures, for
409 example, high-dimensionality, sparsity, compositionality, interdependence, and hierarchical tree structures,
410 all make multi-omics data integration challenging. In this paper, the proposed MRS provides a
411 straightforward and flexible way to incorporate multi-omics datasets and explore the microbial interactions
412 with other omics profiles. Integration of the proposed MRS and the risk scores constructed from other omics
413 data increases the ability for disease prediction. Integrations of metagenomic with metatranscriptomic and
414 host transcriptomic datasets from NYULH COVID-19 cohort underline the critical and insightful utility of
415 the constructed risk scores for disease prediction and the promising ability of multi-omics data integration
416 for predictive accuracy improvement. Additionally, the data from TEDDY study illuminates the potential
417 in combining MRS and PRS to explore genetic-microbial interaction and identify the high risk population.

418 Apart from the ANCOM-BC and Shannon index used in the proposed MRS_{α} , there are other differential
419 abundance methods available to identify the signature microbial taxa associated with disease and other
420 alpha diversity indices to characterize the community diversity. Here we investigate how does using two
421 other differential abundance methods (ALDEx2 [56] and Maaslin2 [57] suggested by [32, 33]) and
422 Simpson and observed alpha diversity indices to construct MRS_{α} affect the predictive performance of the
423 MRS framework in terms of AUC value and 95% CI in the discovery and validation cohorts [27] separately.
424 Figure S5 shows that no single MRS_{α} can uniformly perform best for all predictions in the discovery and
425 validation cohorts, as various DA methods have different model assumptions and test hypotheses and
426 various alpha diversities indices have different definitions, while links between microbiome profile with

427 various healthy or disease conditions are different. Specifically, given an alpha diversity index in the second
428 step, DA method has no effect on the prediction performance of MRS_{α} in both discovery and validation
429 cohorts (p -value>0.05) using the Kruskal-Wallis test on the AUC, except for Simpson index in the
430 discovery cohort (p -value=0.03) (Figure S6). MRS_{α} s constructed with ANCOM-BC, ALDEx2, and
431 Maaslin2, which all have been well-recognized [32, 33], have comparable performances. It supports our
432 suggestion that to carry over the evaluation results of the DA tests from an objective benchmark work to
433 guide the selection of DA test in the MRS framework. In terms of comparisons among Shannon, Simpson
434 and Observed indices, Observed index based MRS_{α} has the highest AUCs, followed by Shannon index,
435 while Simpson index has the lowest AUCs, in the discovery cohort (Figure S7). On the other hand, Shannon
436 index consistently has better or comparable AUCs in the validation cohort. Meanwhile, Observed and
437 Simpson indices introduce more variation in the predictive performance of MRS_{α} (Figure S7). Observed
438 index lacks some reproducibility in the validation cohort, compared to its impressive performance in the
439 discovery cohort, probably because it only accounts for species richness. Taken together, Shannon index
440 based MRS_{α} has relatively more robust and consistent prediction performance. With existing discussions
441 [32, 33] and the observations above in this manuscript, we include various DA methods commonly used
442 and recommended in the microbiome association studies and various alpha diversity indices in the MRS R
443 package to let the proposed MRS framework informative and more practically valuable.

444 The findings of this study have some limitations. First, considering microbial profile varies across
445 ethnicities as well as geographies [58-60], it is necessary to evaluate the portability of MRS between
446 populations. More advanced methods will be required to reduce the bias due to ethnical or geographical
447 differences. Second, the microbiome data have versatile characteristics and unique features, such as
448 phylogenetic tree structure, functional structure, hierarchical taxonomy, and dynamic nature, which also
449 play critical roles in analytical accuracy and efficiency [61, 62]. Incorporating such features may improve
450 the accuracy of MRS. Third, derivation and validation of MRS require large scale microbiome studies.
451 However, the high cost of metagenomics sequencing restrict the comprehensive external validation.

452 Despite the above challenges, this paper proposes a practicable way to summarize the microbial profiles
453 and provides promising findings for comprehensive microbiome research to bolster the microbiome's utility
454 as a potential source of novel therapeutic features.

455 **Conclusions**

456 This paper sheds light on the utility of the microbiome data for disease prediction and multi-omics
457 integration by converting the complex microbial profile into a continuous risk score. The proposed MRS
458 tool provides great potential in studying the complex microbial ecosystem, understanding the microbiome's
459 role in disease diagnosis and prognosis, and exploring microbiome's full clinical potential.

460

461 **List of abbreviations**

462 ANCOM-BC: analysis of compositions of microbiomes with bias correction; AUC: area under the receiver
463 operating characteristic curve; CA: colorectal adenoma; CC: colorectal cancer; CD: Crohn's disease; CV:
464 cross-validation; DA: differential abundance; DESeq2: differential expression analysis (v2); GMHI: gut
465 microbiome health index; GWASs: genome-wide association studies; HR: hazard ratio; MetaPhlAn:
466 metagenomic phylogenetic analysis; ML: machine learning; MRS: microbial risk score; MWASs:
467 microbiome-wide association studies; NYULH: NYU Langone Health; PRS: polygenic risk score; PD:
468 phylogenetic diversity; P+T: Pruning and thresholding; QIIME: quantitative insights into microbial ecology;
469 RA: rheumatoid arthritis; ROC: receiver operating characteristic; StrainPhlAn: metagenomic strain-level
470 phylogenetic analysis; TEDDY: the environmental determinants of diabetes in the young; T1D: type 1
471 diabetes; T2D: type 2 diabetes (T2D).

472

473 **Declarations**

474 **Ethics approval and consent to participate**

475 All utilized microbiome datasets are publicly available. No ethics approval or consent to participate was
476 required for this study.

477 **Consent for publication**

478 Not applicable: All utilized microbiome datasets are publicly available. No consent for publication was
479 required for this study.

480 **Availability of data and materials**

481 For the NYULH COVID-19 cohort, all sequencing data used for this analysis are available in the NCBI
482 Sequence Read Archive under project numbers PRJNA688510 and PRJNA687506 (RNA and DNA
483 sequencing, respectively).

484 For the TEDDY study, TEDDY microbiome 16S rRNA gene sequencing data are publicly available in the
485 NCBI database of Genotypes and Phenotypes (dbGaP) with the primary accession code phs001443. v1.p1,
486 in accordance with the dbGaP controlled-access authorization process. Clinical metadata analysed during
487 the current study will be made available in the NIDDK Central Repository at
488 <https://repository.niddk.nih.gov/studies/teddy/?query=teddy>.

489 MRS R package used for the analyses is available at <https://sites.google.com/site/huilinli09/software> and
490 <https://github.com/chanw0/MRS>, together with its manual. We also included the GMHI data and provided
491 the code in the example section to reproduce the results in this manuscript.

492 **Competing interests**

493 The authors declare that they have no competing interests.

494 **Funding**

495 The study was supported in part by NIH grants P20CA252728 and R37 CA244775.

496 **Authors' contributions**

497 CW developed the microbial risk score framework, performed data analyses, and wrote the manuscript.
498 LS performed data analyses in the NYULH COVID-19 cohort and contributed to manuscript writing. JH
499 performed data analyses in the TEDDY cohort and contributed to manuscript writing. BZ, RH, and JA
500 contributed to the biological insights and interpretation, and to manuscript writing. HL contributed to the
501 methodological ideas for the proposed framework, simulations, real data analyses, and manuscript
502 writing. All authors read and approved the final manuscript.

503 **Acknowledgements**

504 Not applicable.

505

506 **References**

- 507 1. Hu J, Koh H, He L, Liu M, Blaser MJ, Li H: **A two-stage microbial association mapping**
508 **framework with advanced FDR control.** *Microbiome* 2018, **6**(1):1-16.
- 509 2. Gilbert JA, Quinn RA, Debelius J, Xu ZZ, Morton J, Garg N, Jansson JK, Dorrestein PC, Knight
510 R: **Microbiome-wide association studies link dynamic microbial consortia to disease.** *Nature* 2016,
511 **535**(7610):94-103.
- 512 3. Koh H, Livanos AE, Blaser MJ, Li H: **A highly adaptive microbiome-based association test**
513 **for survival traits.** *BMC genomics* 2018, **19**(1):1-13.
- 514 4. Koh H, Blaser MJ, Li H: **A powerful microbiome-based association test and a microbial taxa**
515 **discovery framework for comprehensive association mapping.** *Microbiome* 2017, **5**(1):1-15.
- 516 5. Ahn J, Sinha R, Pei Z, Dominianni C, Wu J, Shi J, Goedert JJ, Hayes RB, Yang L: **Human gut**
517 **microbiome and risk for colorectal cancer.** *Journal of the National Cancer Institute* 2013,
518 **105**(24):1907-1911.

- 519 6. Kostic AD, Xavier RJ, Gevers D: **The microbiome in inflammatory bowel disease: current**
520 **status and the future ahead.** *Gastroenterology* 2014, **146**(6):1489-1499.
- 521 7. Hoffmann AR, Proctor L, Surette M, Suchodolski J: **The microbiome: the trillions of**
522 **microorganisms that maintain health and cause disease in humans and companion animals.**
523 *Veterinary Pathology* 2016, **53**(1):10-21.
- 524 8. Kelly TN, Bazzano LA, Ajami NJ, He H, Zhao J, Petrosino JF, Correa A, He J: **Gut microbiome**
525 **associates with lifetime cardiovascular disease risk profile among bogalusa heart study**
526 **participants.** *Circulation research* 2016, **119**(8):956-964.
- 527 9. Gilbert JA, Blaser MJ, Caporaso JG, Jansson JK, Lynch SV, Knight R: **Current understanding**
528 **of the human microbiome.** *Nature medicine* 2018, **24**(4):392-400.
- 529 10. Fattorusso A, Di Genova L, Dell'Isola GB, Mencaroni E, Esposito S: **Autism spectrum**
530 **disorders and the gut microbiota.** *Nutrients* 2019, **11**(3):521.
- 531 11. Integrative H, Proctor LM, Creasy HH, Fettweis JM, Lloyd-Price J, Mahurkar A, Zhou W, Buck
532 GA, Snyder MP, Strauss III JF: **The integrative human microbiome project.** *Nature* 2019,
533 **569**(7758):641-648.
- 534 12. Wang C, Hu J, Blaser MJ, Li H: **Estimating and testing the microbial causal mediation effect**
535 **with high-dimensional and compositional microbiome data.** *Bioinformatics (Oxford, England)* 2020,
536 **36**(2):347-355.
- 537 13. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H, Alm
538 EJ, Arumugam M, Asnicar F *et al*: **Reproducible, interactive, scalable and extensible microbiome**
539 **data science using QIIME 2.** *Nat Biotechnol* 2019, **37**(8):852-857.
- 540 14. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower C,
541 Segata N: **MetaPhlan2 for enhanced metagenomic taxonomic profiling.** *Nature methods* 2015,
542 **12**(10):902-903.
- 543 15. Truong DT, Tett A, Pasolli E, Huttenhower C, Segata N: **Microbial strain-level population**
544 **structure and genetic diversity from metagenomes.** *Genome research* 2017, **27**(4):626-638.

- 545 16. Choi SW, Mak TS-H, O'Reilly PF: **Tutorial: a guide to performing polygenic risk score**
546 **analyses.** *Nature Protocols* 2020, **15**(9):2759-2772.
- 547 17. Wand H, Lambert SA, Tamburro C, Iacocca MA, O'Sullivan JW, Sillari C, Kullo IJ, Rowley R,
548 Dron JS, Brockman D: **Improving reporting standards for polygenic scores in risk prediction studies.**
549 *Nature* 2021, **591**(7849):211-219.
- 550 18. Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI: **The human**
551 **microbiome project.** *Nature* 2007, **449**(7164):804-810.
- 552 19. McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, Aksenov AA,
553 Behsaz B, Brennan C, Chen Y: **American gut: an open platform for citizen science microbiome**
554 **research.** *Msystems* 2018, **3**(3):e00031-00018.
- 555 20. Xavier JB, Young VB, Skufca J, Ginty F, Testerman T, Pearson AT, Macklin P, Mitchell A,
556 Shmulevich I, Xie L: **The cancer microbiome: distinguishing direct and indirect effects requires a**
557 **systemic view.** *Trends in cancer* 2020, **6**(3):192-204.
- 558 21. de Cárcer DA: **A conceptual framework for the phylogenetically constrained assembly of**
559 **microbial communities.** *Microbiome* 2019, **7**(1):1-11.
- 560 22. Coyte KZ, Rao C, Rakoff-Nahoum S, Foster KR: **Ecological rules for the assembly of**
561 **microbiome communities.** *PLoS biology* 2021, **19**(2):e3001116.
- 562 23. Cho I, Blaser MJ: **The human microbiome: at the interface of health and disease.** *Nature*
563 *Reviews Genetics* 2012, **13**(4):260-270.
- 564 24. Thukral AK: **A review on measurement of Alpha diversity in biology.** *Agric Res J* 2017,
565 **54**(1):1-10.
- 566 25. Whittaker RH: **Evolution and measurement of species diversity.** *Taxon* 1972, **21**(2-3):213-251.
- 567 26. Sulaiman I, Chung M, Angel L, Tsay J-CJ, Wu BG, Yeung ST, Krolikowski K, Li Y, Duerr R,
568 Schluger R *et al*: **Microbial signatures in the lower airways of mechanically ventilated COVID-19**
569 **patients associated with poor clinical outcome.** *Nature Microbiology* 2021, **6**(10):1245-1258.

- 570 27. Gupta VK, Kim M, Bakshi U, Cunningham KY, Davis JM, Lazaridis KN, Nelson H, Chia N,
571 Sung J: **A predictive index for health status using species-level gut microbiome profiling.** *Nature*
572 *communications* 2020, **11**(1):1-16.
- 573 28. Lee HS, Burkhardt BR, McLeod W, Smith S, Eberhard C, Lynch K, Hadley D, Rewers M, Simell
574 O, She JX: **Biomarker discovery study design for type 1 diabetes in The Environmental**
575 **Determinants of Diabetes in the Young (TEDDY) study.** *Diabetes/metabolism research and reviews*
576 2014, **30**(5):424-434.
- 577 29. Rewers M, Hyöty H, Lernmark Å, Hagopian W, She J-X, Schatz D, Ziegler A-G, Toppari J,
578 Akolkar B, Krischer J: **The Environmental Determinants of Diabetes in the Young (TEDDY) study:**
579 **2018 update.** *Current diabetes reports* 2018, **18**(12):1-14.
- 580 30. Zheng P, Li Z, Zhou Z: **Gut microbiome in type 1 diabetes: A comprehensive review.**
581 *Diabetes/metabolism research and reviews* 2018, **34**(7):e3043.
- 582 31. Lin H, Peddada SD: **Analysis of compositions of microbiomes with bias correction.** *Nature*
583 *communications* 2020, **11**(1):1-11.
- 584 32. Nearing JT, Douglas GM, Hayes MG, MacDonald J, Desai DK, Allward N, Jones C, Wright RJ,
585 Dhanani AS, Comeau AM: **Microbiome differential abundance methods produce different results**
586 **across 38 datasets.** *Nature Communications* 2022, **13**(1):1-16.
- 587 33. Lin H, Peddada SD: **Analysis of microbial compositions: a review of normalization and**
588 **differential abundance analysis.** *NPJ biofilms and microbiomes* 2020, **6**(1):1-13.
- 589 34. Wilcoxon F: **Individual comparisons by ranking methods.** In: *Breakthroughs in statistics.*
590 Springer; 1992: 196-202.
- 591 35. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, Huttenhower C:
592 **Metagenomic biomarker discovery and explanation.** *Genome biology* 2011, **12**(6):1-18.
- 593 36. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for RNA-**
594 **seq data with DESeq2.** *Genome biology* 2014, **15**(12):1-21.

- 595 37. Robinson MD, McCarthy DJ, Smyth GK: **edgeR: a Bioconductor package for differential**
596 **expression analysis of digital gene expression data.** *Bioinformatics* 2010, **26**(1):139-140.
- 597 38. Mandal S, Van Treuren W, White RA, Eggesbø M, Knight R, Peddada SD: **Analysis of**
598 **composition of microbiomes: a novel method for studying microbial composition.** *Microbial ecology*
599 *in health and disease* 2015, **26**(1):27663.
- 600 39. Kaul A, Mandal S, Davidov O, Peddada SD: **Analysis of microbiome data in the presence of**
601 **excess zeros.** *Frontiers in microbiology* 2017, **8**:2114.
- 602 40. Marcos-Zambrano LJ, Karaduzovic-Hadziabdic K, Loncar Turukalo T, Przymus P, Trajkovic V,
603 Aasmets O, Berland M, Gruca A, Hasic J, Hron K: **Applications of machine learning in human**
604 **microbiome studies: a review on feature selection, biomarker identification, disease prediction and**
605 **treatment.** *Frontiers in microbiology* 2021, **12**:313.
- 606 41. Gou W, Ling C-w, He Y, Jiang Z, Fu Y, Xu F, Miao Z, Sun T-y, Lin J-s, Zhu H-l: **Interpretable**
607 **machine learning framework reveals robust gut microbiome features associated with type 2**
608 **diabetes.** *Diabetes Care* 2021, **44**(2):358-366.
- 609 42. Ke G, Meng Q, Finley T, Wang T, Chen W, Ma W, Ye Q, Liu T-Y: **Lightgbm: A highly**
610 **efficient gradient boosting decision tree.** *Advances in neural information processing systems* 2017,
611 **30**:3146-3154.
- 612 43. Vabalas A, Gowen E, Poliakoff E, Casson AJ: **Machine learning algorithm validation with a**
613 **limited sample size.** *PloS one* 2019, **14**(11):e0224365.
- 614 44. Lamri A, Mao S, Desai D, Gupta M, Paré G, Anand SS: **Fine-tuning of Genome-Wide**
615 **Polygenic Risk Scores and Prediction of Gestational Diabetes in South Asian Women.** *Scientific*
616 *reports* 2020, **10**(1):1-9.
- 617 45. Jost L: **Entropy and diversity.** *Oikos* 2006, **113**(2):363-375.
- 618 46. Gauthier J, Derome N: **Evenness-Richness Scatter Plots: a Visual and Insightful**
619 **Representation of Shannon Entropy Measurements for Ecological Community Analysis.** *Msphere*
620 2021, **6**(2):e01019-01020.

- 621 47. Blaser MJ: **Missing microbes: how the overuse of antibiotics is fueling our modern plagues:**
622 Macmillan; 2014.
- 623 48. Padilla-Martínez F, Collin F, Kwasniewski M, Kretowski A: **Systematic review of polygenic**
624 **risk scores for type 1 and type 2 diabetes.** *International journal of molecular sciences* 2020,
625 **21(5):1703.**
- 626 49. Perry DJ, Wasserfall CH, Oram RA, Williams MD, Posgai A, Muir AB, Haller MJ, Schatz DA,
627 Wallet MA, Mathews CE: **Application of a genetic risk score to racially diverse type 1 diabetes**
628 **populations demonstrates the need for diversity in risk-modeling.** *Scientific reports* 2018, **8(1):1-13.**
- 629 50. Udler MS, McCarthy MI, Florez JC, Mahajan A: **Genetic risk scores for diabetes diagnosis and**
630 **precision medicine.** *Endocrine reviews* 2019, **40(6):1500-1520.**
- 631 51. Harrell FE: **Cox proportional hazards regression model.** In: *Regression modeling strategies.*
632 Springer; 2015: 475-519.
- 633 52. Chatterjee N, Shi J, García-Closas M: **Developing and evaluating polygenic risk prediction**
634 **models for stratified disease prevention.** *Nat Rev Genet* 2016, **17(7):392-406.**
- 635 53. Vatanen T, Franzosa EA, Schwager R, Tripathi S, Arthur TD, Vehik K, Lernmark Å, Hagopian
636 WA, Rewers MJ, She J-X: **The human gut microbiome in early-onset type 1 diabetes from the**
637 **TEDDY study.** *Nature* 2018, **562(7728):589-594.**
- 638 54. Stewart CJ, Ajami NJ, O'Brien JL, Hutchinson DS, Smith DP, Wong MC, Ross MC, Lloyd RE,
639 Doddapaneni H, Metcalf GA: **Temporal development of the gut microbiome in early childhood from**
640 **the TEDDY study.** *Nature* 2018, **562(7728):583-588.**
- 641 55. Kuhn M: **Building predictive models in R using the caret package.** *Journal of statistical*
642 *software* 2008, **28(1):1-26.**
- 643 56. Gloor G: **ALDEx2: ANOVA-Like Differential Expression tool for compositional data.**
644 *ALDEX manual modular* 2015, **20:1-11.**

- 645 57. Mallick H, Rahnavard A, McIver LJ, Ma S, Zhang Y, Nguyen LH, Tickle TL, Weingart G, Ren
646 B, Schwager EH: **Multivariable association discovery in population-scale meta-omics studies.** *PLoS*
647 *computational biology* 2021, **17**(11):e1009442.
- 648 58. Gaulke CA, Sharpton TJ: **The influence of ethnicity and geography on human gut**
649 **microbiome composition.** *Nature medicine* 2018, **24**(10):1495-1496.
- 650 59. Deschasaux M, Bouter KE, Prodan A, Levin E, Groen AK, Herrema H, Tremaroli V, Bakker GJ,
651 Attaye I, Pinto-Sietsma S-J: **Depicting the composition of gut microbiota in a population with varied**
652 **ethnic origins but shared geography.** *Nature medicine* 2018, **24**(10):1526-1531.
- 653 60. He Y, Wu W, Zheng H-M, Li P, McDonald D, Sheng H-F, Chen M-X, Chen Z-H, Ji G-Y, Zheng
654 Z-D-X: **Regional variation limits applications of healthy gut microbiome reference ranges and**
655 **disease models.** *Nature medicine* 2018, **24**(10):1532-1535.
- 656 61. Lozupone C, Knight R: **UniFrac: a new phylogenetic method for comparing microbial**
657 **communities.** *Appl Environ Microbiol* 2005, **71**(12):8228-8235.
- 658 62. Chen J, Bushman FD, Lewis JD, Wu GD, Li H: **Structure-constrained sparse canonical**
659 **correlation analysis with an application to microbiome data analysis.** *Biostatistics* 2013, **14**(2):244-
660 258.

661 **Figure 1.** The workflow of the microbial risk score (MRS) framework. Data Input: a phyloseq-class object
662 is needed, which consists of a feature table (observed count table), a sample metadata, a taxonomy table
663 (optional), and a phylogenetic tree (optional). MRS Algorithm has two steps: Step 1 is to identify a sub-
664 community consisting of the signature microbial taxa with the P+T method and AUC evaluation in the
665 discovery cohort. The black ROC curve which has the largest AUC determines the optimal p -value cutoff.
666 Step 2 is to integrate the identified microbial taxa into a continuous score, i.e., calculate the MRS value for
667 each sample by calculating the diversity of the identified sub-community with the Shannon index. In
668 addition, the constructed MRS is independently validated in the validation cohort. Application: In this
669 manuscript, we perform multi-omics data integration for disease prediction by jointly modeling the
670 proposed MRS and other risk scores constructed from other omics data in two real data cohorts.

671

672 **Figure 2.** The optimal p -value thresholds by P+T method for including taxa in MRS_{α} , MRS_{wS} , and
673 MRS_{unwS} , separately, using the metagenomic data in the NYULH COVID-19 cohort. Specifically, given a
674 cut-off, the taxa with p -values less than the cut-off were selected and defined as a sub-community. The p -
675 values were obtained by ANCOM-BC method. The leave-one-out CV was used for the predictions. MRS_{α} :
676 the negative alpha diversity (Shannon index) was calculated for each sample on the selected sub-community;
677 MRS_{wS} : the weighted sum of relative abundances of the selected taxa with the weights being the coefficients

678 estimated from the ANCOM-BC log-linear model; MRS_{unws} : the sum of relative abundances of the
679 selected taxa.

680

681 **Figure 3.** Evaluation of MRS in the discovery and validation cohorts [27]. A: The AUC values and 95%
682 confidence intervals (CIs) for MRS_{α} s to predict healthy and different disease conditions in discovery and
683 validation cohorts, respectively. B: Venn diagrams of taxa identified in pairwise comparisons of Healthy
684 versus CA, CC, CD, and RA. CA: colorectal adenoma, CC: colorectal cancer, CD: Crohn's disease, and
685 RA: rheumatoid arthritis.

686

687 **Figure 4.** The ROC curves and AUC values for the various risk scores to predict alive and deceased status
688 in the NYULH COVID-19 cohort. A. Predication performance for the individual risk scores constructed
689 based on metagenome (DNA_MRS_{α}), metatranscriptome (RNA_MRS_{α}), and host transcriptome (Host),
690 separately. B. Predication performance based on multiple risk scores using additive model.

691

692 **Figure 5.** Box plots of the score comparisons between alive and deceased group. All risk scores are
693 standardized among all samples, respectively. The statistical significance on group comparison is
694 evaluated by Wilcoxon signed-rank test.

695

696 **Figure 6.** Scatterplots of risk scores based on metagenome, metatranscriptome, and host transcriptome data.
697 A-C: Scatterplots of DNA_MRS_{α} vs RNA_MRS_{α} , DNA_MRS_{α} vs Host, and RNA_MRS_{α} vs Host,
698 respectively. Dotted line denotes the mean of the corresponding risk score across all subjects. D: 3D
699 scatterplot of DNA_MRS_{α} vs RNA_MRS_{α} vs Host.

700

701 **Figure 7.** Heatmaps of Spearman's rank correlations between the top 50 taxa from metagenome and the
702 top 50 genes from host transcriptome, in alive and deceased groups, separately. The top 50 features were
703 selected based on the proportion of selection in all CV iterations.

704

705 **Figure 8.** Results for T1D prediction in the TEDDY study. A. ROC curves and AUC values for
706 predicting T1D status using various risk score. PRS_{hla} is constructed from the HLA alleles alone, and
707 PRS is constructed from all SNPs found in the TEDDY cohort based on the existing PRS algorithm [49].
708 MRS_{α} is the negative alpha diversity (Shannon index) calculated on the selected sub-community, which is
709 selected by ANCOM-BC method and P+T method. B–D. Kaplan–Meier plots for the groups of subjects at
710 high and low risk of developing T1D, based on PRS, MRS_{α} , and the combination of PRS and MRS_{α} ,
711 respectively. Subjects whose risk scores are above the third quartile are defined as high risk, others as
712 low risk, others as low risk.

713

714 **Table 1.** Classification evaluation for subjects having extreme risk categories (labeled as either “High risk”
715 or “Low risk” by both or all three risk scores) in the NYULH COVID-19 cohort.

Combination of the risk scores	Sensitivity	Specificity	Accuracy	F1
DNA_MRS $_{\alpha}$ + RNA_MRS $_{\alpha}$	0.67	0.78	0.71	0.75
DNA_MRS $_{\alpha}$ + Host	0.78	0.65	0.74	0.80
RNA_MRS $_{\alpha}$ + Host	0.48	0.92	0.58	0.63
DNA_MRS $_{\alpha}$ + RNA_MRS $_{\alpha}$ + Host	0.86	0.91	0.88	0.89

716

717 **Table 2.** Association results between the risk scores and the time to death based on the Cox proportional-
718 hazards model in the NYULH COVID-19 cohort.

Risk score	Hazard ratio		<i>p</i> -value
	Estimate	95% confidence interval	
DNA_MRS $_{\alpha}$	1.80	1.36-2.38	3.56E-05
RNA_MRS $_{\alpha}$	1.87	1.11-3.14	0.0179
Host	1.43	1.16-1.76	0.000855
DNA_MRS $_{\alpha}$ +Host	1.54	1.28-1.84	2.52E-06
DNA_MRS $_{\alpha}$ + RNA_MRS $_{\alpha}$	2.57	1.78-3.71	4.46E-07
RNA_MRS $_{\alpha}$ +Host	2.00	1.51-2.64	1.39E-06
DNA_MRS $_{\alpha}$ + RNA_MRS $_{\alpha}$ + Host	1.97	1.58-2.45	1.60E-09

719

720 Additional material

721

722 **Additional file 1: Figure S1.** The ROC curves and AUC values for various ML algorithms to predict the
723 alive or deceased status in the NYULH COVID-19 cohort. A. Predication performance for elastic-net
724 logistic regression (glmnet), penalized discriminant analysis (pda2), regularized random forest (RRF), and
725 neural networks with feature extraction (pcaNNet) methods. B. Predication performance for naive Bayes
726 (naïve_bayes), neural network (nnet), stochastic gradient boosting (gbm), and support vector machines with
727 polynomial kernel (svmPoly) methods.

728 **Figure S2.** The AUC values and 95% CIs for MRS $_{\alpha}$ s to classify healthy and nonhealthy and two disease
729 conditions in the discovery and validation GMHI cohorts [27], respectively. CA: colorectal adenoma, CC:
730 colorectal cancer, CD: Crohn’s disease, and RA: rheumatoid arthritis.

731 **Figure S3.** Heatmaps of Spearman’s rank correlations between the top 50 taxa from metagenome and the
732 top 50 taxa from metatranscriptome, in the alive and deceased groups, separately. The top 50 features were
733 selected based on the proportion of selection in all CV iterations.

734 **Figure S4.** Heatmaps of Spearman’s rank correlations between the top 50 taxa from metatranscriptome and
735 the top 50 genes from host transcriptome, in the alive and deceased groups, separately. The top 50 features
736 were selected based on the proportion of selectin in all CV iterations.

737 **Figure S5.** Comparisons among various MRSs in terms of AUC value and 95% CI in the discovery and
738 validation cohorts [27]. Here candidate taxa are identified by ANCOM-BC [31], ALDEx2 [56], and
739 Maaslin2 [57], and the MRS_{α} s are constructed by Shannon, Simpson, and Observed indices, respectively.
740 DA: differential abundance, CA: colorectal adenoma, CC: colorectal cancer, CD: Crohn's disease, and
741 RA: rheumatoid arthritis.

742 **Figure S6.** The mean and standard derivation of the ranks of MRS_{α} 's AUCs with ANCOM-BC,
743 ALDEx2, and Maaslin2, respectively. For each alpha diversity index in each comparison of two diseases
744 or healthy conditions, the AUCs of MRS_{α} with three DA methods were ranked 1-3. A higher rank
745 represents a higher AUC. For each alpha diversity index, the Kruskal-Wallis test was performed to check
746 difference among three DA methods. All: all samples were used for test. Statistical significance: ns: p -
747 value > 0.05; *: p -value \leq 0.05.

748 **Figure S7.** The mean and standard derivation of the ranks of MRS_{α} 's AUCs with Shannon, Simpson, and
749 Observed indices, respectively. For each DA method in each comparison of two diseases or healthy
750 conditions, the AUCs of MRS_{α} with three indices were ranked 1-3. A higher rank represents a higher
751 AUC. For each DA method, the Kruskal-Wallis test was performed to check difference among three alpha
752 diversity indices. All: all samples were used for test. Statistical significance: ns: p -value > 0.05; *: p -value
753 \leq 0.05; **: p -value \leq 0.01; ***: p -value \leq 0.001; ****: p -value \leq 0.0001.

754

755

756 **Additional file 2: Table S1.** Number of discovery and validation samples used for MRS evaluation and
757 validation from the GMHI multi-study cohort.

758 **Table S2.** AUC values for six common alpha diversity indices on the whole community to predict alive
759 and deceased status in the NYULH COVID-19 cohort.

760 **Table S3** The identified species for MRS_{α} construction in terms of comparisons among healthy, CA, CC,
761 CD, RA, and nonhealthy based on the discovery samples in the GMHI multi-study cohort.

762 **Table S4.** Average and standard deviation of relative abundances of the identified species in Healthy, CA,
763 CC, CD, and RA discovery samples from the GMHI multi-study cohort. The identified species are used for
764 MRS_{α} construction in terms of pairwise comparisons of Healthy versus CA, CC, CD, and RA, respectively.

765 **Table S5.** Average and standard deviation of relative abundances of the identified species in Healthy, CA,
766 CC, CD, and RA validation samples from the GMHI multi-study cohort. The identified species are used for
767 MRS_{α} construction in terms of pairwise comparisons of Healthy versus CA, CC, CD, and RA, respectively.

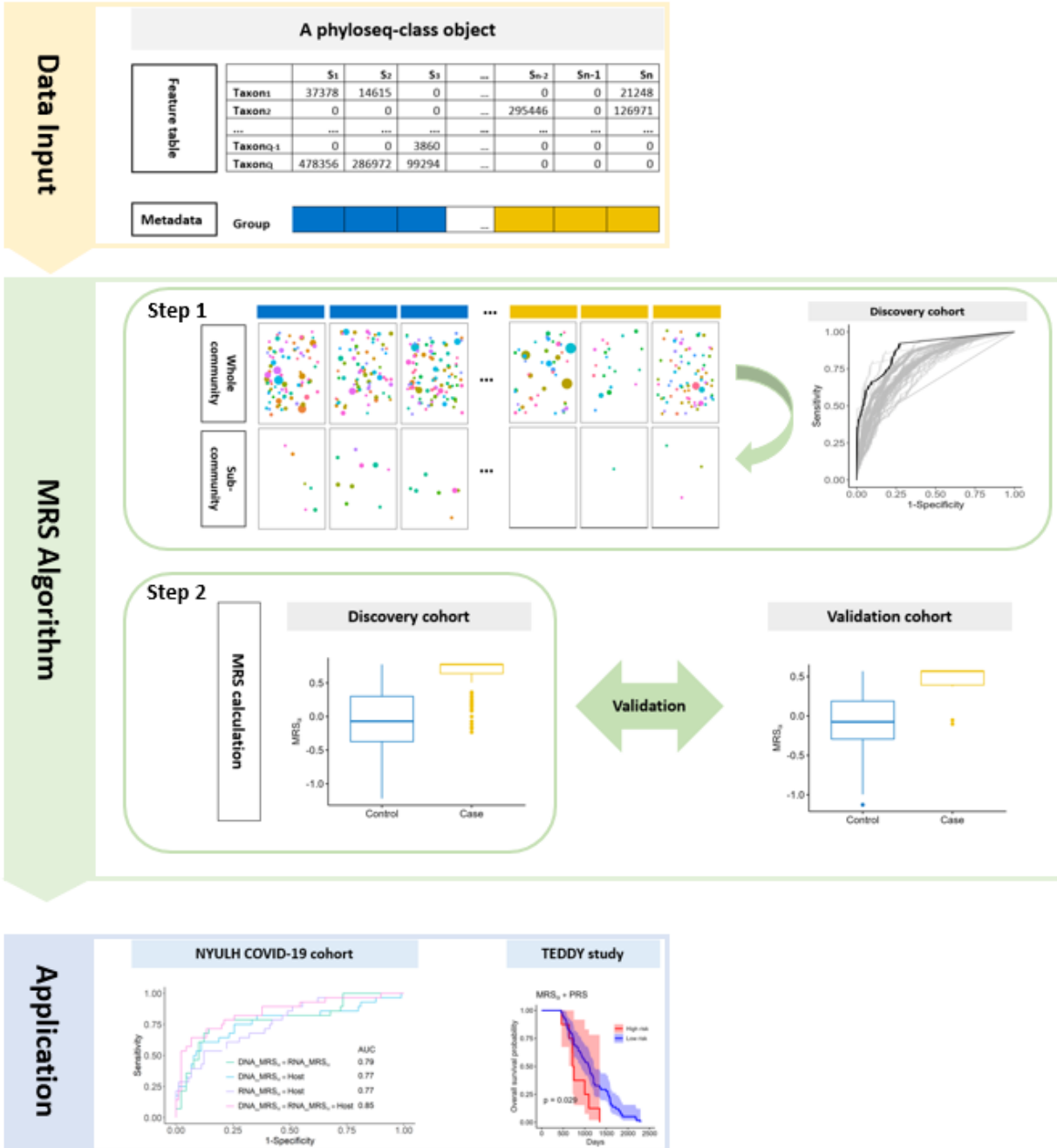
768 **Table S6.** Factors used for metagenomic, metatranscriptomic and host transcriptomic risk scores in the
769 NYULH COVID-19 cohort.

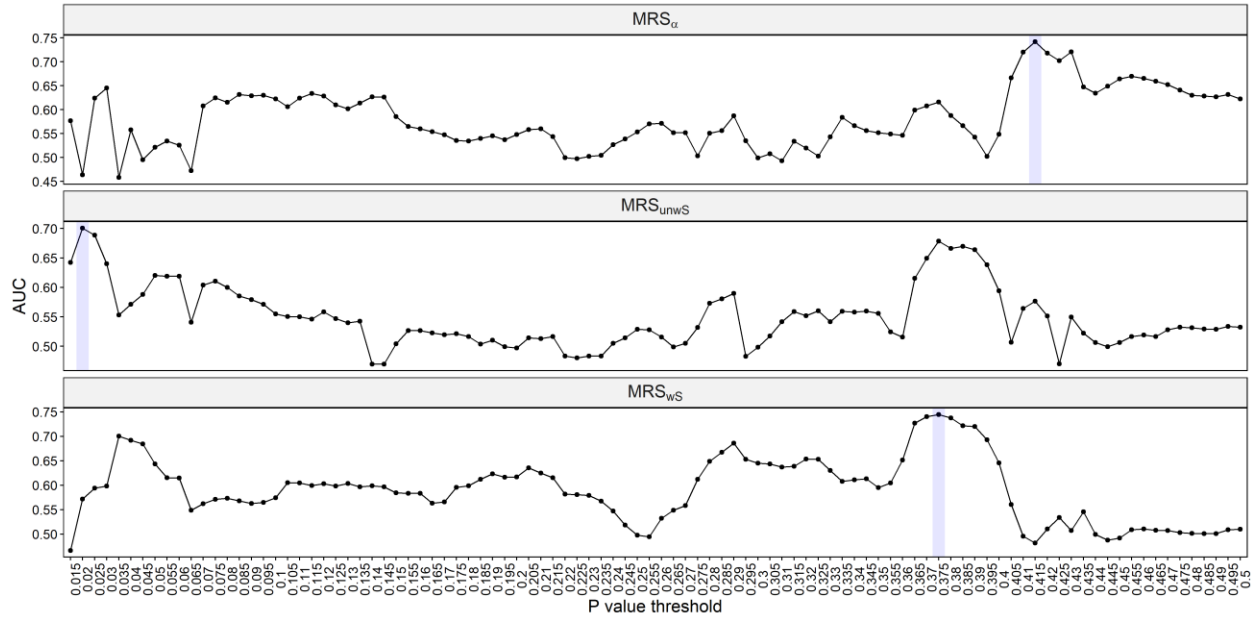
770

771 **Additional file 3: Section S1** Computational details for risk scores

772

773



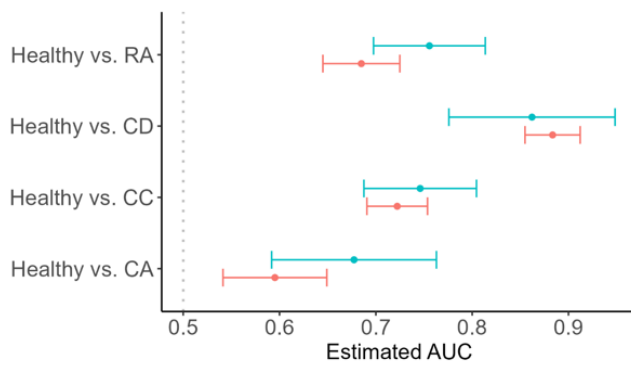


775

776

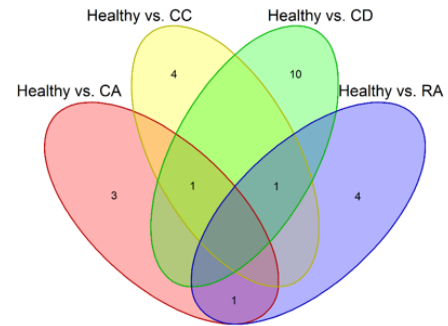
777

A



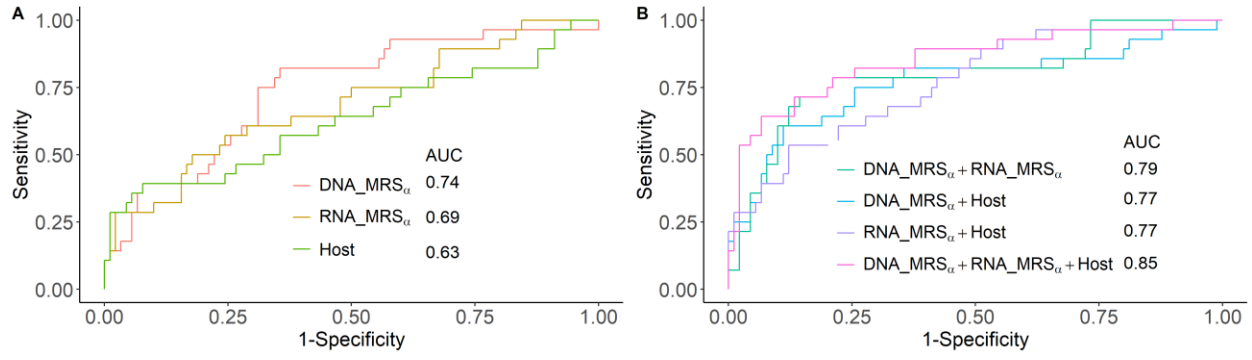
—●— Discovery —●— Validation

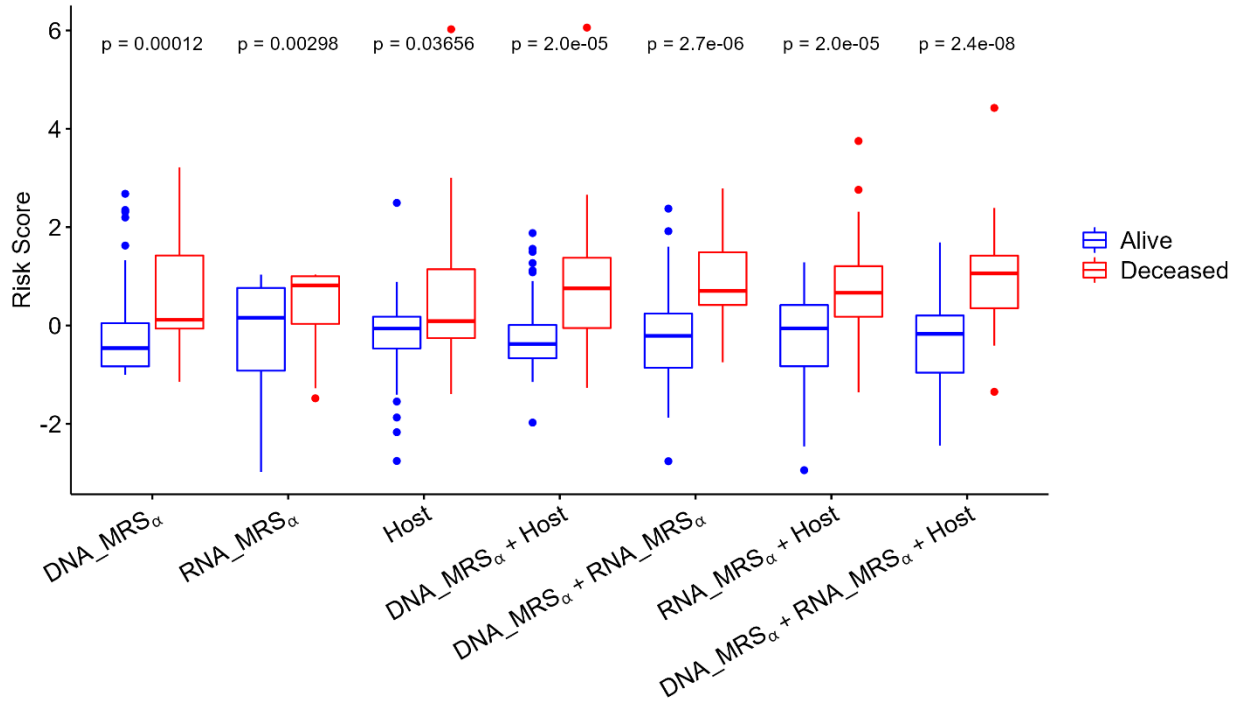
B

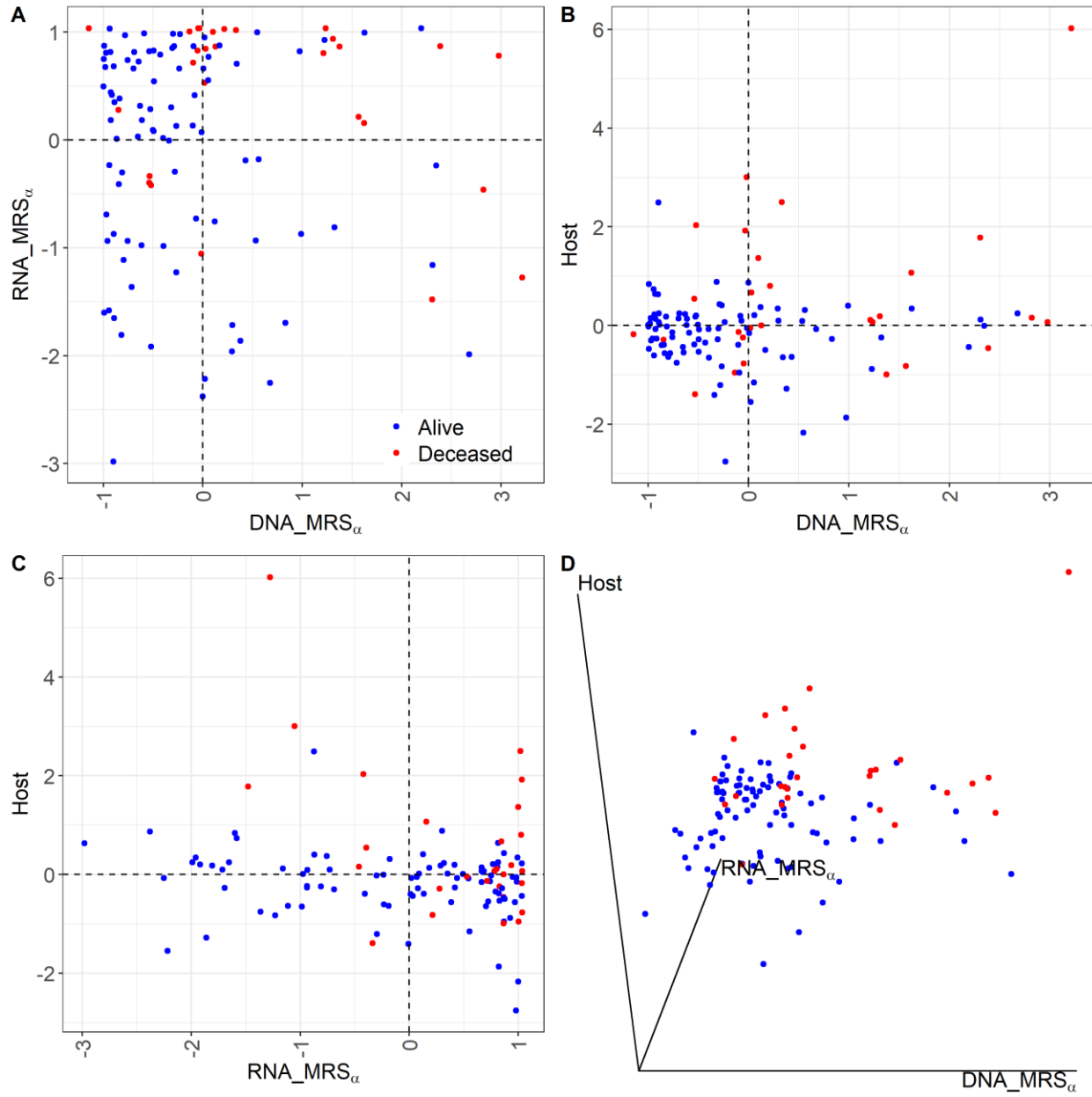


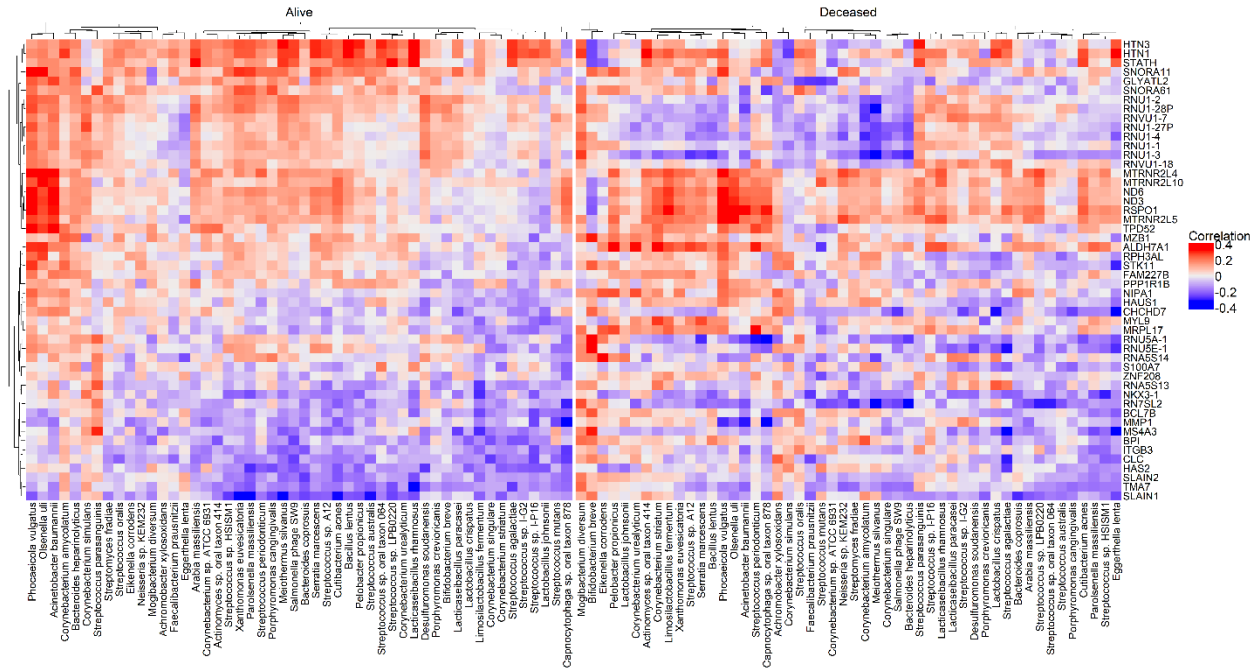
778

779





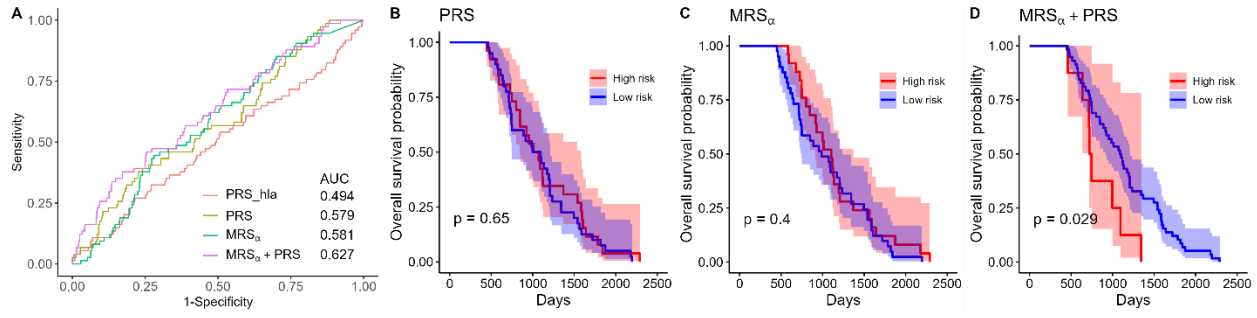




783

784

785



786