

1 **GAN-GMHI: Generative Adversarial Network for high**
2 **discrimination power in microbiome-based disease**
3 **prediction**

4

5 Yuxue Li[#], Gang Xie[#], Yuguo Zha[#], Kang Ning^{*}

6

7 *MOE Key Laboratory of Molecular Biophysics, Hubei Key Laboratory of Bioinformatics and*
8 *Molecular-imaging, Center of Artificial Intelligence Biology, Department of Bioinformatics and*
9 *Systems Biology, College of Life Science and Technology, Huazhong University of Science and*
10 *Technology, Wuhan 430074, China*

11

12 [#]Equal contribution.

13 ^{*}Corresponding author.

14 E-mail: ningkang@hust.edu.cn (Ning K)

15

16 **Running title:** *Li Y / GAN-GMHI for Microbiome-based Disease Prediction*

17

18 **Abstract**

19 Gut microbiome-based health index (GMHI) has been applied with success, while the
20 discrimination powers of GMHI varied for different diseases, limiting its utility on a
21 broad-spectrum of diseases. In this work, a generative adversarial network (GAN)
22 model is proposed to improve the discrimination power of GMHI. Built based on the
23 batch corrected data through GAN, GAN-GMHI has largely reduced the batch effects,
24 and profoundly improved the performance for distinguishing healthy individuals and
25 different diseases. GAN-GMHI has provided a solution to unravel the strong
26 association of gut microbiome and diseases, and indicated a more accurate venue
27 toward microbiome-based disease monitoring. The code for GAN-GMHI is available
28 at <https://github.com/HUST-NingKang-Lab/GAN-GMHI>.

29

30 **Keywords:** Gut microbiome-based health index; Generative adversarial network;
31 Batch effect; Microbiome-based disease

32

33 **Importance**

34 The association of gut microbiome and diseases has been proven for many diseases,
35 while the transformation of such association to a robust and universal disease
36 prediction model has remained illusive, largely due to the batch effects presents in
37 multiple microbiome cohorts. Our analyses have indicated a plausible venue, which is
38 based on GAN technique, towards batch effect removal for microbiome datasets.
39 GAN-GMHI is a novel method built based on the batch corrected data through GAN,
40 as well as GMHI for prediction of a broad-spectrum of diseases.

41

42 **Introduction**

43 There are important links between many complex chronic diseases and the human gut
44 microbiome [1]. Specific sets of gut microbes could directly or indirectly influence
45 complex chronic diseases, such as the microbiome dysbiosis in the development of
46 rheumatoid arthritis [2], thus it is natural that the gut microbiome could be utilized for

47 disease prediction [1,3]. However, a general microbiome-based index for the
48 prediction of a broad-spectrum of diseases is lacking.

49 Previous work has reported the gut microbiome health index (GMHI) [1], a
50 robust index for assessing health status, based on the species-level taxonomic profile
51 of stool metagenomic sequencing samples. GMHI values can be used to classify
52 samples as healthy ($GMHI > 0$), non-healthy ($GMHI < 0$), or neither ($GMHI = 0$), and
53 results in the previous work have shown strong reproducibility on the validation
54 dataset. However, GMHI has limited power to distinguish samples from different
55 diseases: as the stool metagenomes in that study were collected from over 40
56 published studies, it is nearly impossible to exclude experimental and technical
57 inter-study batch effects [1]. The batch effect refers to the technical difference caused
58 by the processing and measurement of samples in different batches that are not related
59 to any biological variation recorded during the test [4].

60 A number of traditional tools for batch effect removal have been developed, such
61 as ComBat [5]. However, traditional batch effect removal tools have shown
62 limitations when faced with heterogeneous data of human gut microbiome from
63 hundreds of studies and thousands of individuals. To address this major hurdle, here
64 we introduced GAN-GMHI, based on the generative adversarial network (GAN), for
65 improved discrimination power of GMHI. GAN was applied to reduce the batch
66 effects on a large collection of gut microbiome samples from multiple cohorts
67 containing both healthy and non-healthy individuals. Then GMHI could be applied to
68 the batch corrected data for prediction. Compared with the original GMHI,
69 GAN-GMHI makes the distribution of GMHI values within the group more
70 concentrated and the distinction between healthy and non-healthy samples more
71 clearly. The effectiveness of GAN for cross-cohort batch correction has been
72 demonstrated: the prediction accuracy of GAN-GMHI has been improved to 88.70%
73 for distinguishing healthy individuals and non-healthy individuals, compared to the
74 accuracy of 70.95%, 72.00% achieved by GMHI and ComBat-GMHI. In summary,
75 batch effect does exist in datasets from different sources, and GMHI can better predict
76 the status of health based on GAN corrected datasets.

77

78 **Method**

79 The GAN-GMHI framework consists of three stages. First, a dataset containing
80 phenotype and batch information for all samples is constructed. Second, GAN is
81 applied to guide the batch effect correction of raw data. Third, the corrected dataset is
82 output as the training dataset for GMHI prediction (Figure S1).

83 The batch effect removal method of iMAP [6], a GAN method previously applied
84 on single-cell RNA-Seq data, was adapted for batch effect removal in this study. It is
85 worth noting that the datasets to be batch-corrected by GAN must be classified based
86 on the phenotype first, and the sub datasets of each phenotype are regrouped
87 according to the batch. Such processing ensures that the unwanted technical variations
88 among different datasets are eliminated, and the biological differences between
89 different phenotypes are retained. We should emphasize that GAN-GMHI applied the
90 same core functions as iMAP. However, we optimized the structure of the model to
91 better fit the microbiome abundance data (see the script at
92 [https://github.com/HUST-NingKang-Lab/GAN-GMHI/blob/main/scripts/DNN\(GAN\)](https://github.com/HUST-NingKang-Lab/GAN-GMHI/blob/main/scripts/DNN(GAN).ipynb)
93 [.ipynb](#)). The model contains a generator and a discriminator. The generator is an
94 autoencoder structure which consists of seven layers and millions of parameters. The
95 discriminator is a dense neural network which consists of three layers. Additionally,
96 we compared GAN with other three batch effect removal tools (ComBat [5], Seurat3
97 [7], and Harmony [8]) on the same dataset.

98

99 **Results**

100 We have performed a comprehensive analysis on the integrated dataset of 2636
101 healthy and 1711 non-healthy (including 12 disease phenotypes) individuals' stool
102 metagenomes from 34 published studies [1]. All of these samples are used as a
103 discovery dataset. Additionally, we have used 679 samples (118 healthy and 561
104 non-healthy) as a validation dataset [1]. The discovery and validation datasets
105 configuration is the same as in Gupta *et al.* [1].

106 We first assessed and compared prediction results based on the discovery cohort
107 (training data). By comparison of the species-level GMHI before batch effect
108 correction (referred to as RAW), after batch effect correction by GAN, and after batch
109 effect correction by ComBat (**Figure 1A**) for distinguishing samples from healthy and
110 non-healthy individuals, we observed that the prediction accuracy is largely improved
111 after batch effect correction: GAN-GMHI achieved an overall accuracy of 88.70%,
112 which is higher than GMHI's accuracy of 70.95%. Moreover, batch effect correction
113 by GAN is better than that by other three batch effect correction methods (Tables S1
114 and S2): GAN-GMHI (accuracy of 88.70%) outperformed ComBat-GMHI (accuracy
115 of 72.00%), Seurat3-GMHI (accuracy of 70.36%), and Harmony-GMHI (accuracy of
116 44.65%). In summary, we emphasize that GAN reduced batch effect with high fidelity,
117 and augmented the gut microbiome-based health index by profoundly improved
118 discrimination power.

119 It has been reported that there is a significant change in the alpha diversity (*i.e.*,
120 Shannon index) of gut microbiome in non-healthy individuals. Therefore, we
121 compared the abilities of GAN-GMHI and Shannon diversity indicators to
122 differentiate the gut microbiome of healthy and non-healthy individuals. The results
123 demonstrated that GAN-GMHI could yield clearer separation compared with Shannon
124 diversity in differentiating healthy and non-healthy individuals (Figure 1B).

125 Results on more than ten non-healthy phenotypes have also shown the advantage
126 of GAN-GMHI over GMHI as regard to differentiating these diverse groups of
127 phenotypes. When GMHI was applied, the GMHI values were dispersed over a wide
128 range, and GMHI values for healthy samples were slightly higher than those for
129 non-healthy samples except for symptomatic arteriosclerosis. On the other hand, when
130 GAN-GMHI was applied, the GMHI values were concentrated for each group, and
131 the healthy group was significantly higher than the 12 non-healthy phenotypes ($P <$
132 0.001 for all non-healthy groups), and the third quartile of GMHI was lower than 0
133 for all non-healthy phenotypes (Figure 1C). Moreover, it is easier for clinical
134 interpretation based on the results of GMHI. For example, on type 2 diabetes (T2D),

135 GAN-GMHI has captured *Lactobacillus* as biomarkers, which are well founded by
136 published works [9].

137 Additionally, we compared GAN-GMHI and GMHI on the validation dataset.
138 Cross-cohort batch correction by GAN improved the performance for distinguishing
139 healthy and non-healthy individuals. The prediction accuracy of GAN-GMHI on the
140 validation dataset is 73.05%, compared to the 72.61% of GMHI (Table S1).

141 Furthermore, GAN is not only applicable for the GMHI disease prediction model,
142 but could also be easily adapted to other models, such as random forest (RF). It has
143 been observed that GMHI and RF exhibit similar performance on the validation
144 dataset, while results of GMHI are easier to interpret clinically (Table S1). We
145 emphasize that although the results of GAN-GMHI and GAN-RF also have similar
146 accuracies on the validation dataset, GAN-GMHI has inherited the interpretability of
147 the GMHI method, and thus is more suitable for clinical interpretation. For example,
148 GAN-GMHI has captured *Lactobacillus* as biomarkers on T2D, which are well
149 founded by published works [9].

150 Finally, the computational expense of GAN-GMHI is feasible even on a regular
151 laptop. For example, in an experiment including 298 samples from four batches of the
152 colorectal cancer phenotype, the total running time, including training and predicting,
153 is no more than one hour and the maximum usage of memory is less than 2 gigabytes.
154

155 **Conclusion**

156 The association of gut microbiome and diseases has been proven for many diseases,
157 while the transformation of such association to a robust and universal disease
158 prediction model has remained illusive, largely due to the batch effects presents in
159 multiple microbiome cohorts. Our analyses have indicated a plausible venue, which is
160 based on GAN technique, towards batch effect removal for microbiome datasets.
161 GAN-GMHI is a novel method built based on the batch corrected data through GAN,
162 as well as GMHI for prediction of a broad-spectrum of diseases. Our study showed
163 that GAN-GMHI could largely reduce the batch effect, and profoundly improved the

164 performance for distinguishing healthy and non-healthy individuals. Batch effect
165 correction by GAN was also better than that by ComBat. In summary, GAN
166 augmented the gut microbiome-based health index, and GAN-GMHI has indicated a
167 more accurate venue towards microbiome-based disease monitoring.

168

169 **Code availability**

170 The source code is available at the GitHub
171 (<https://github.com/HUST-NingKang-Lab/GAN-GMHI>), and BioCode
172 (<https://ngdc.cncb.ac.cn/biocode/tools/7275>).

173

174 **CRedit author statement**

175 **Yuxue Li:** Investigation, Visualization, Writing - Original Draft. **Gang Xie:**
176 Investigation, Writing - Original Draft. **Yuguo Zha:** Writing - Review & Editing.
177 **Kang Ning:** Writing - Review & Editing, Conceptualization, Supervision, Funding
178 acquisition. All authors read and approved the final manuscript.

179

180 **Competing interests**

181 The authors declare that they have no competing interests.

182

183 **Acknowledgments**

184 We are grateful to Mingyue Cheng and Hui Chong for their insightful discussions.
185 This work was partially supported by National Natural Science Foundation of China
186 (Grant Nos. 32071465, 31871334 and 31671374), and the National Key R&D
187 Program (Grant No. 2018YFC0910502).

188

189 **ORCID**

190 0000-0001-7450-3178 (Yuxue Li)

191 0000-0001-5672-732X (Gang Xie)

192 0000-0003-3702-9416 (Yuguo Zha)

193 0000-0003-3325-5387 (Kang Ning)

194

195 **References**

196 [1] Gupta VK, Kim M, Bakshi U, Cunningham KY, Davis JM, Lazaridis KN, et al. A
197 predictive index for health status using species-level gut microbiome profiling. *Nat*
198 *Commun* 2020;11:4635.

199 [2] Bergot A-S, Giri R, Thomas R. The microbiome and rheumatoid arthritis. *Best*
200 *Pract Res Clin Rheumatol* 2019;33:101497.

201 [3] Shreiner AB, Kao JY, Young VB. The gut microbiome in health and in disease.
202 *Curr Opin Gastroenterol* 2015;31:69–75.

203 [4] Hornung R, Boulesteix A-L, Causeur D. Combining location-and-scale batch
204 effect adjustment with data cleaning by latent factor adjustment. *BMC Bioinformatics*
205 2016;17:27.

206 [5] Stein CK, Qu P, Epstein J, Buros A, Rosenthal A, Crowley J, et al. Removing
207 batch effects from purified plasma cell gene expression microarrays with modified
208 ComBat. *BMC Bioinformatics* 2015;16:63.

209 [6] Wang D, Hou S, Zhang L, Wang X, Liu B, Zhang Z. iMAP: integration of multiple
210 single-cell datasets by adversarial paired transfer networks. *Genome Biol* 2021;22:63.

211 [7] Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, et al.
212 Comprehensive integration of single-cell data. *Cell* 2019;177:1888–902.

213 [8] Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, et al. Fast,
214 sensitive and accurate integration of single-cell data with Harmony. *Nat Methods*
215 2019;16:1289–96.

216 [9] Wang Y, Ouyang M, Gao X, Wang S, Fu C, Zeng J, et al. Phoea,
217 Pseudoflavonifactor and *Lactobacillus intestinalis*: Three potential biomarkers of gut
218 microbiota that affect progression and complications of obesity-induced type 2
219 diabetes mellitus. *Diabetes Metab Syndr Obes* 2020;13:835–850.

220

221 **Figure legends**

222 **Figure 1 Comparison of GAN-GMHI with other methods under different**
223 **settings**

224 **A.** Violin plots of GMHI for the healthy and non-healthy groups before (left, referred
225 to as RAW), after batch correction by GAN (middle), and after batch correction by
226 ComBat (right). *200, $P < 1E-200$; *300, $P < 1E-300$; RAW-GMHI, original GMHI;
227 GAN-GMHI, GMHI with GAN enhancement; ComBat-GMHI, GMHI with ComBat
228 enhancement. **B.** the distribution of the RAW (top) and GAN corrected (bottom)
229 GMHI and Shannon diversity. RAW-Shannon, original Shannon index;
230 GAN-Shannon, Shannon index with GAN enhancement. **C.** Violin plots of GMHI
231 index for the healthy and 12 non-healthy phenotypes before (left, referred to as RAW),
232 after batch correction by GAN (middle), and after batch correction by ComBat (right).
233 The red box, an example to show the different results by GMHI, GAN-GMHI, and
234 ComBat-GMHI. ***, $P < 0.001$; ****, $P < 0.0001$; ns, not significant; ACVD,
235 atherosclerotic cardiovascular disease; CA, colorectal adenoma; CC, colorectal cancer;
236 CD, Crohn's disease; IGT, impaired glucose tolerance; OB, obesity; OW, overweight;
237 RA, rheumatoid arthritis; SA, symptomatic arteriosclerosis; T2D, type 2 diabetes; UC,
238 ulcerative colitis; UW, underweight.

239

240 **Supplementary materials**

241 **Figure S1 Schematic diagram of the technological framework**

242 **Table S1 The accuracy of the two classifiers for healthy and non-healthy**

243 **classification**

244 **Table S2 The overall prediction accuracy of different methods before and after**

245 **batch correction**

246 **File S1 Details about the GAN-GMHI method**

247

