

1 **Generative Adversarial Network augmented the gut**
2 **microbiome-based health index by profoundly improved**
3 **discrimination power**

4
5 *Yuxue Li^{1,§}, Gang Xie^{1,§}, Yuguo Zha^{1,§}, Kang Ning^{1,*}*

6 *1 Key Laboratory of Molecular Biophysics of the Ministry of Education, Hubei Key Laboratory of*
7 *Bioinformatics and Molecular-imaging, Center of AI Biology, Department of Bioinformatics and*
8 *Systems Biology, College of Life Science and Technology, Huazhong University of Science and*
9 *Technology, Wuhan 430074, Hubei, China*

10 [§] These authors contributed equally to this work.

11 * Corresponding author. Tel: +86 27 87793041, e-mail: ningkang@hust.edu.cn (Kang Ning).

12

13 **Abstract**

14 **Summary:** Gut microbiome-based health index (GMHI) has been applied with success, while the
15 discrimination powers of GMHI varied for different diseases, limiting its utility on a broad
16 spectrum of diseases. In this work, a Generative Adversarial Network (GAN) model is proposed to
17 improve the discrimination power of GMHI. Built based on the batch corrected data through GAN
18 (<https://github.com/HUST-NingKang-Lab/GAN-GMHI>), GAN-GMHI has largely reduced the
19 batch effects, and profoundly improved the performance for distinguishing healthy individuals and
20 different diseases. GAN-GMHI has provided results to support the strong association of gut
21 microbiome and diseases, and indicated a more accurate venue towards microbiome-based disease
22 monitoring.

23 **Availability and implementation:** GAN-GMHI is publicly available on GitHub:
24 <https://github.com/HUST-NingKang-Lab/GAN-GMHI>.

25 **Contact:** ningkang@hust.edu.cn

26 **Supplementary information:** Supplementary data are available at Bioinformatics online.

27

28 **Introduction**

29 There are important links between many complex chronic diseases and the human gut microbiome
30 (Gupta, et al., 2020). Specific sets of gut microbes could directly or indirectly influence the
31 complex chronic diseases, such as the microbiome dysbiosis in the development of rheumatoid
32 arthritis (Bergot, et al., 2019), thus it is nature that gut microbiome could be utilized for disease

33 prediction (Gupta, et al., 2020; Shreiner, et al., 2015). However, a general microbiome-based
34 index for prediction of a broad spectrum of diseases is lacking.

35

36 A previous work has reported the Gut Microbiome Health Index (GMHI) (Gupta, et al., 2020), a
37 robust index for assessing health status, based on the species-level taxonomic profile of stool
38 metagenomic sequencing samples. GMHI values can be used to classify samples as healthy
39 ($GMHI > 0$), non-healthy ($GMHI < 0$), or neither ($GMHI = 0$), and its results have shown strong
40 reproducibility on the validation datasets. However, GMHI has limited power to distinguish
41 samples from different diseases, largely due to the existence of batch effects: as the stool
42 metagenomes in that study were collected from over 40 published studies, it is nearly impossible
43 to exclude experimental and technical inter-study batch effects (Gupta, et al., 2020). Thus, the
44 overall prediction accuracy of GMHI is still far from perfect: 70.72% for distinguishing healthy
45 individuals and non-healthy individuals.

46

47 Therefore, we introduced GAN-GMHI, based on the Generative Adversarial Network (GAN), for
48 improved discrimination power of GMHI. GAN was applied to reduce the batch effects on a large
49 collection of gut microbiome samples from multiple cohorts containing both health and disease
50 individuals. Then GMHI could be applied on the batch corrected data for prediction. Compared
51 with original GMHI, GAN-GMHI makes the distribution of GMHI values within the group more
52 concentrated and the distinction between healthy and non-healthy samples more clearly. The
53 effectiveness of GAN for cross-cohort batch correction has been demonstrated: the prediction
54 accuracy of GAN-GMHI has been improved to 88.70% for distinguishing healthy individuals and
55 non-healthy individuals, compared to the accuracy of 70.95% achieved by GMHI. In summary,
56 batch effect does exist in data sets from different sources, and GMHI can better predict the status
57 of health based on GAN corrected data sets.

58

59 **Methods**

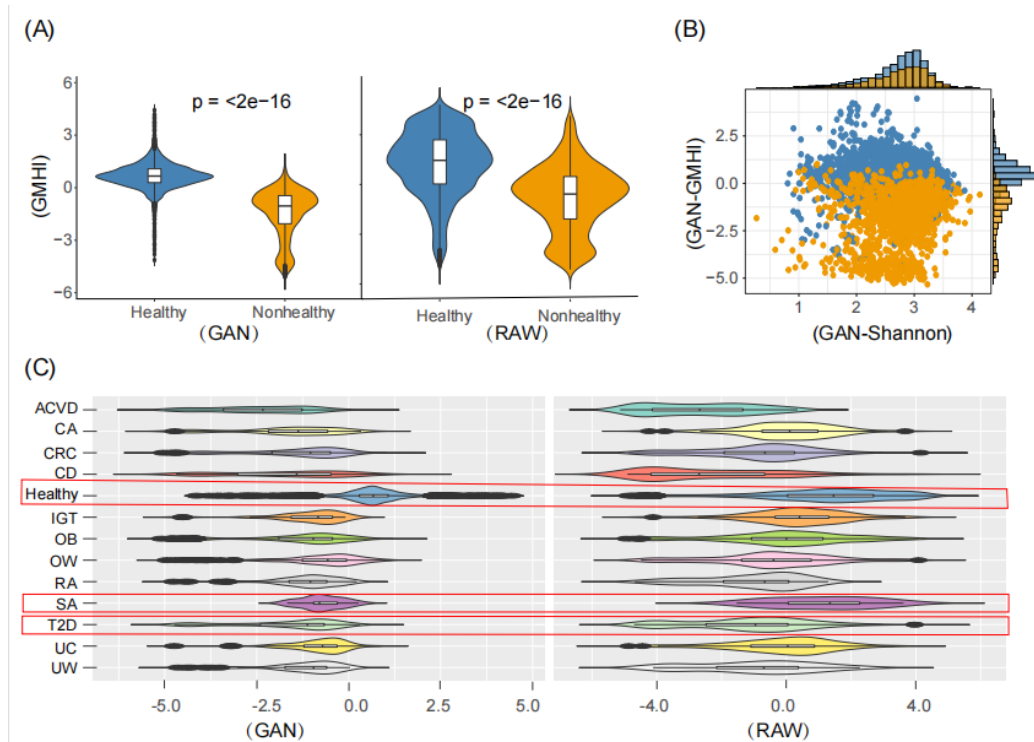
60 Our GAN-GMHI framework consists of three stages, constructing a dataset containing phenotype
61 and batch information for all samples, and then GAN guiding the batch effect correction of raw
62 data, the corrected datasets are output as the training data set for GMHI prediction
63 (**Supplementary Figure 1**). The batch effect removal method of iMAP (Wang, et al., 2021), a
64 GAN method previously applied on single-cell RNA-Seq data, was adapted for batch effect
65 removal in this study. It is worth noting that the datasets to be batch-corrected by GAN must be
66 classified based on the phenotype first, and the sub-data sets of each phenotype are regrouped
67 according to the batch. To ensure that the unwanted technical variations among different datasets
68 are eliminated, but the biological differences between different phenotypes are not diminished.

69

70 **Results**

71 We have performed a comprehensive analysis on the integrated dataset of 2,636 healthy and 1,711
72 diseased (including 12 disease phenotypes) individuals' stool metagenomes from 34 published
73 studies (Gupta, et al., 2020). All of these samples are used as analytical datasets. Additionally, we

74 have used 679 samples (118 healthy and 561 diseased) as validation datasets (Gupta, et al., 2020).
75 The analytical and validation datasets configuration is the same as in (Gupta, et al., 2020).
76
77 We have first assessed and compared prediction results based on the discovery cohort (training
78 data). By comparison of the species-level GMHI before (RAW) and after (GAN) batch correction
79 (**Figure 1 (A)**) for distinguishing samples from healthy and non-healthy individuals, we observed
80 that the accuracy for prediction of the healthy and diseased groups after correction is 87.03% and
81 91.29% (with overall accuracy of 88.70%), respectively, compared with 75.61% and 63.76%
82 before correction (overall accuracy of 70.95%), which has proven the advantage of GAN-GMHI
83 over GMHI (**Supplementary Table 1**). Additionally, we compared the abilities of GAN-GMHI
84 and Shannon diversity indicators to differentiate the gut microbiome of healthy and non-healthy
85 individuals. The results demonstrated that GAN-GMHI could yield clearer separation compared
86 with Shannon's diversity in differentiating healthy and non-healthy individuals (**Figure 1 (B)**).
87 Furthermore, results on comparison among the healthy group and the 12 non-healthy phenotypes
88 have showed that: when GMHI was applied, the GMHI values were dispersed over a wide range,
89 and GMHI values for healthy samples were slightly higher than those for non-healthy samples
90 except for SA. On the other hand, when GAN-GMHI was applied, the GMHI values were
91 concentrated for each group, and the healthy group was significantly higher than the 12 disease
92 phenotypes (p -value < 0.05 for all disease groups), and the third quartile of GMHI was lower than
93 0 for all disease phenotypes (**Figure 1 (C)**). Moreover, GMHI's results are easier for clinical
94 interpretation. For example, on Type 2 diabetes (T2D), GAN-GMHI has captured *Lactobacillus* as
95 biomarkers, which are well founded by published works (Wang, et al.).



96

97 **Figure 1. Comparison of GAN-GMHI with other methods under different settings.** (A) Violin
98 plots of GMHI for the healthy and non-healthy groups before (left) and after (right) batch

99 correction by GAN. (B) the distribution of corrected GMHI and Shannon diversity. (C) Violin
100 plots of GMHI index for the healthy and 12 disease phenotypes before (right) and after (left) batch
101 correction by GAN. ACVD: Arteriosclerosis Cardiovascular Disease, CA: colorectal adenoma, CC:
102 colorectal cancer, CD: Crohn's disease, IGT: Impaired glucose tolerance, OB: obesity, OW:
103 overweight, RA: rheumatoid arthritis, SA: Symptomatic arteriosclerosis, T2D: Type 2 Diabetes,
104 UC: ulcerative colitis, UW: underweight.

105

106 Additionally, we have compared GAN-GMHI and GMHI on the validation datasets. Cross-cohort
107 batch correction by GAN profoundly improved the performance for distinguishing healthy
108 individuals and different diseases. The prediction accuracy of GAN-GMHI has been significantly
109 improved to 88.70% for distinguishing healthy individuals and non-healthy individuals, compared
110 to the accuracy of 70.95% achieved by GMHI, and GAN-GMHI still outperforms GMHI on the
111 independent validation cohort (**Supplementary Table 1**).

112

113 Moreover, GAN is not only applicable for GMHI disease prediction model, but could also be
114 easily adapted to other models, such as Random Forest (RF). It has already been observed that
115 GMHI and RF have similar performances on the validation datasets, while GMHI's results are
116 easier for clinical interpretation (**Supplementary Table 1**). We emphasize that although the results
117 of GAN-GMHI and GAN-RF also have similar accuracies on the validation datasets, GAN-GMHI
118 has inherited the interpretability of the GMHI method, and thus is more suitable for clinical
119 interpretation. For example, on Type 2 diabetes (T2D), GAN-GMHI has captured *Lactobacillus* as
120 biomarkers, which are well founded by published works (Wang, et al.).

121

122 **Conclusion**

123 The association of gut microbiome and diseases has been proven for many diseases, while
124 transformation of such association to a robust and universal disease prediction model has
125 remained illusive, largely due to the batch effects presents in multiple microbiome cohorts.
126 GAN-GMHI is a novel method built based on the batch corrected data through GAN, as well as
127 GMHI for prediction of a broad spectrum of diseases. Results have shown that it has largely
128 reduced the batch effects, and profoundly improved the performance for distinguishing disease
129 and healthy individuals. In summary, Generative Adversarial Network augmented the gut
130 microbiome-based health index, and GAN-GMHI has indicated a more accurate venue towards
131 microbiome-based disease monitoring.

132

133 **Acknowledgments**

134 The authors are grateful to Mingyue Cheng and Hui Chong for insightful discussions.

135

136 **Funding**

137 This work was partially supported by National Science Foundation of China grant 81774008,
138 81573702, 31871334, and the Ministry of Science and Technology's national key research and
139 development program grant (No. 2018YFC0910502).

140

141 *Conflict of Interest: none declared.*

142

143 **References**

144 Bergot, A.-S., Giri, R. and Thomas, R. The microbiome and rheumatoid arthritis. *Best Practice &*
145 *Research Clinical Rheumatology* 2019;33(6):101497.

146 Gupta, V.K., *et al.* A predictive index for health status using species-level gut microbiome
147 profiling. *Nature Communications* 2020;11(1):4635.

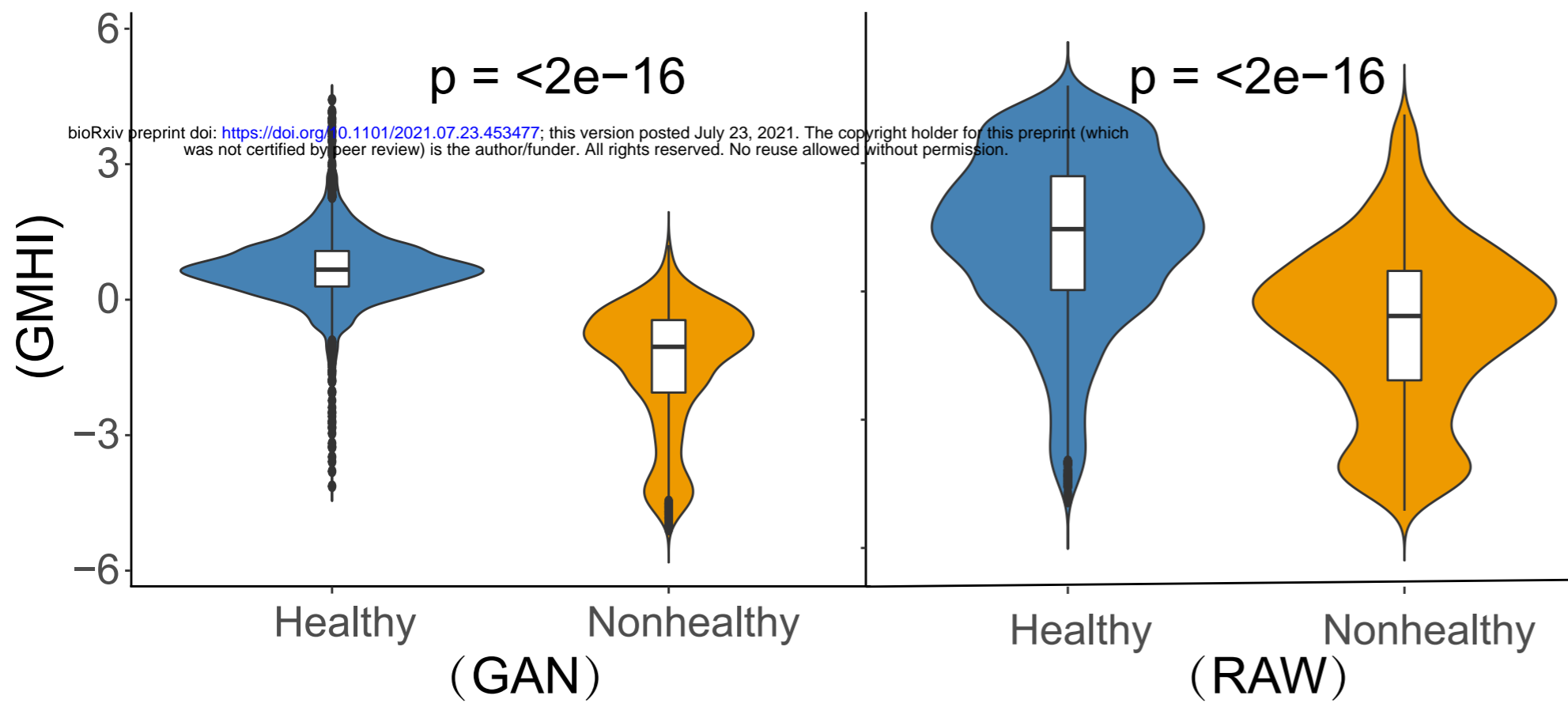
148 Shreiner, A.B., Kao, J.Y. and Young, V.B. The gut microbiome in health and in disease. *Current*
149 *Opinion in Gastroenterology* 2015;31(1).

150 Wang, D., *et al.* iMAP: integration of multiple single-cell datasets by adversarial paired transfer
151 networks. *Genome Biology* 2021;22(1):63.

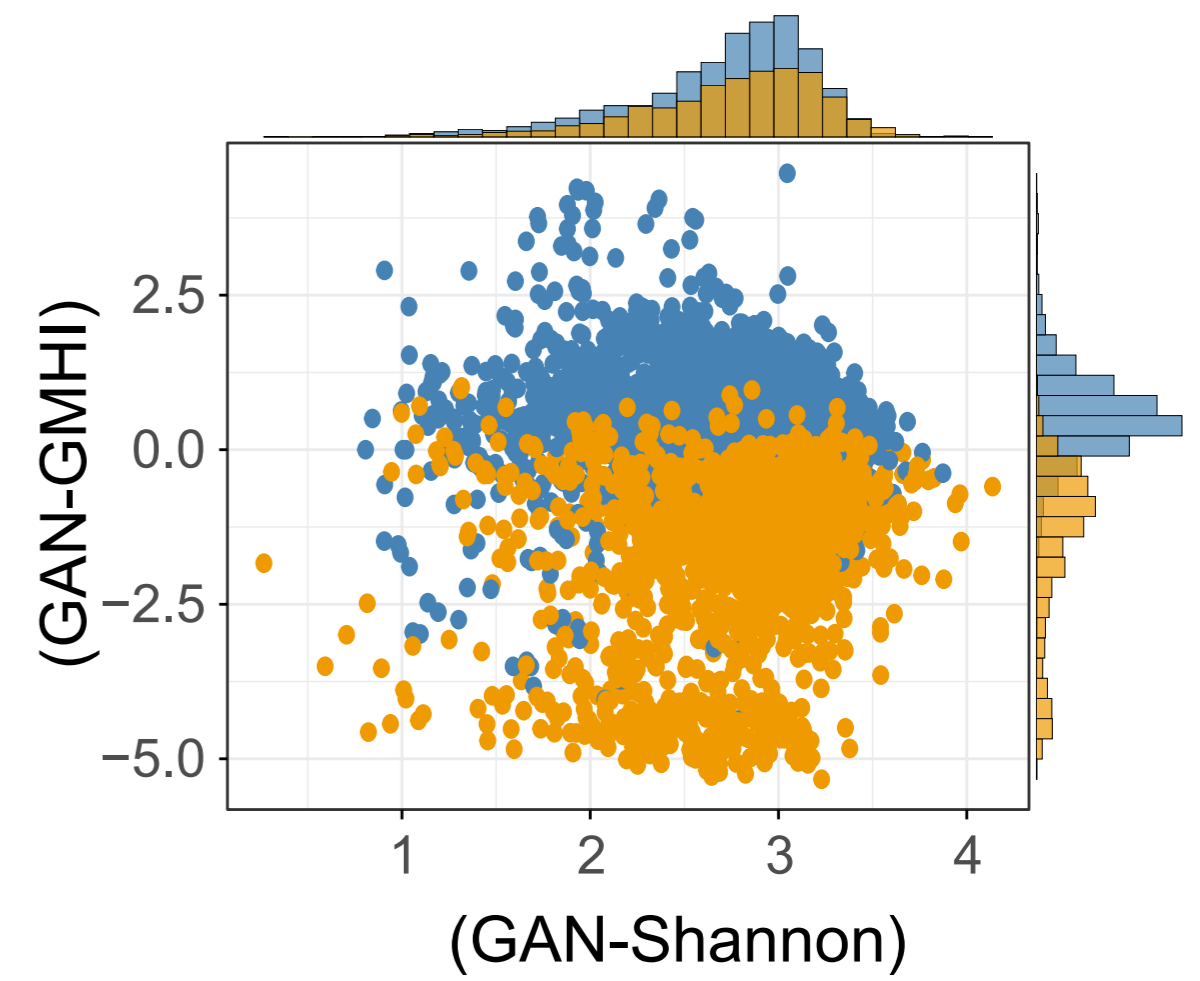
152 Wang, Y.A.-O.X., *et al.* Phocea, Pseudoflavonifractor and Lactobacillus intestinalis: Three
153 Potential Biomarkers of Gut Microbiota That Affect Progression and Complications of
154 Obesity-Induced Type 2 Diabetes Mellitus. (1178-7007 (Print)).

155

(A)



(B)



(C)

