

## Machine Learning Strategy for Gut Microbiome-Based Diagnostic Screening of Cardiovascular Disease

Sachin Aryal<sup>1</sup>, Ahmad Alimadadi<sup>1</sup>, Ishan Manandhar<sup>1</sup>, Bina Joe<sup>1</sup>, Xi Cheng<sup>1</sup>

**Abstract**—Cardiovascular disease (CVD) is the number one leading cause for human mortality. Besides genetics and environmental factors, in recent years, gut microbiota has emerged as a new factor influencing CVD. Although cause-effect relationships are not clearly established, the reported associations between alterations in gut microbiota and CVD are prominent. Therefore, we hypothesized that machine learning (ML) could be used for gut microbiome-based diagnostic screening of CVD. To test our hypothesis, fecal 16S ribosomal RNA sequencing data of 478 CVD and 473 non-CVD human subjects collected through the American Gut Project were analyzed using 5 supervised ML algorithms including random forest, support vector machine, decision tree, elastic net, and neural networks. Thirty-nine differential bacterial taxa were identified between the CVD and non-CVD groups. ML modeling using these taxonomic features achieved a testing area under the receiver operating characteristic curve (0.0, perfect antidiscrimination; 0.5, random guessing; 1.0, perfect discrimination) of  $\approx 0.58$  (random forest and neural networks). Next, the ML models were trained with the top 500 high-variance features of operational taxonomic units, instead of bacterial taxa, and an improved testing area under the receiver operating characteristic curves of  $\approx 0.65$  (random forest) was achieved. Further, by limiting the selection to only the top 25 highly contributing operational taxonomic unit features, the area under the receiver operating characteristic curves was further significantly enhanced to  $\approx 0.70$ . Overall, our study is the first to identify dysbiosis of gut microbiota in CVD patients as a group and apply this knowledge to develop a gut microbiome-based ML approach for diagnostic screening of CVD. (*Hypertension*. 2020;76:1555-1562. DOI: 10.1161/HYPERTENSIONAHA.120.15885.) • [Data Supplement](#)

**Key Words:** artificial intelligence ■ cardiovascular disease ■ diagnosis ■ gut microbiome ■ machine learning ■ metagenomic sequencing

Cardiovascular disease (CVD) refers to a number of morbid conditions such as heart failure,<sup>1</sup> hypertension,<sup>2</sup> and atherosclerosis,<sup>3</sup> which could develop simultaneously or may lead to each other.<sup>4,5</sup> Worldwide, by 2030, CVD death toll is estimated to surpass 23.6 million.<sup>1</sup> Multiple clinical tests, including ECG,<sup>6</sup> chest x-ray,<sup>7</sup> and echocardiogram,<sup>8</sup> are routinely required for a comprehensive evaluation of cardiovascular health. Therefore, a convenient screening test for an overall evaluation of cardiovascular health could save diagnostic time and facilitate a timely therapeutic intervention.<sup>9</sup>

Machine learning (ML), a major branch of artificial intelligence, has been successfully used for diagnostic testing and prediction of a variety of diseases such as cancer,<sup>10</sup> diabetes mellitus,<sup>11</sup> and inflammatory bowel disease.<sup>12</sup> For example, ML models have been trained with gut microbiota features to classify healthy and inflammatory bowel disease subjects.<sup>13</sup> Since dysregulated gut microbiota is observed in several types of CVD, such as hypertension,<sup>14–20</sup> heart

failure,<sup>21</sup> and atherosclerosis,<sup>22</sup> we hypothesized that supervised ML models could be trained with gut microbiota data for diagnostic screening of CVD. To test this hypothesis, we evaluated the capacity of different supervised ML models to detect and differentiate gut microbiome signatures from fecal 16S metagenomics data obtained from 478 CVD and 473 non-CVD subjects through the American Gut Project. To our knowledge, our study is the first to demonstrate the promising potential of artificial intelligence via ML modeling for a convenient diagnostic screening of CVD based on fecal microbiota composition.

### Methods

The authors declare that all supporting data are available within the article and in the [Data Supplement](#).

### Data Collection and Processing

The workflow of the whole study is summarized in Figure 1A. Human 16S ribosomal RNA sequencing data were collected through

Received July 3, 2020; first decision July 23, 2020; revision accepted August 20, 2020.

From the Bioinformatics and Artificial Intelligence Laboratory, Center for Hypertension and Precision Medicine, Program in Physiological Genomics, Department of Physiology and Pharmacology, University of Toledo College of Medicine and Life Sciences, Toledo, OH.

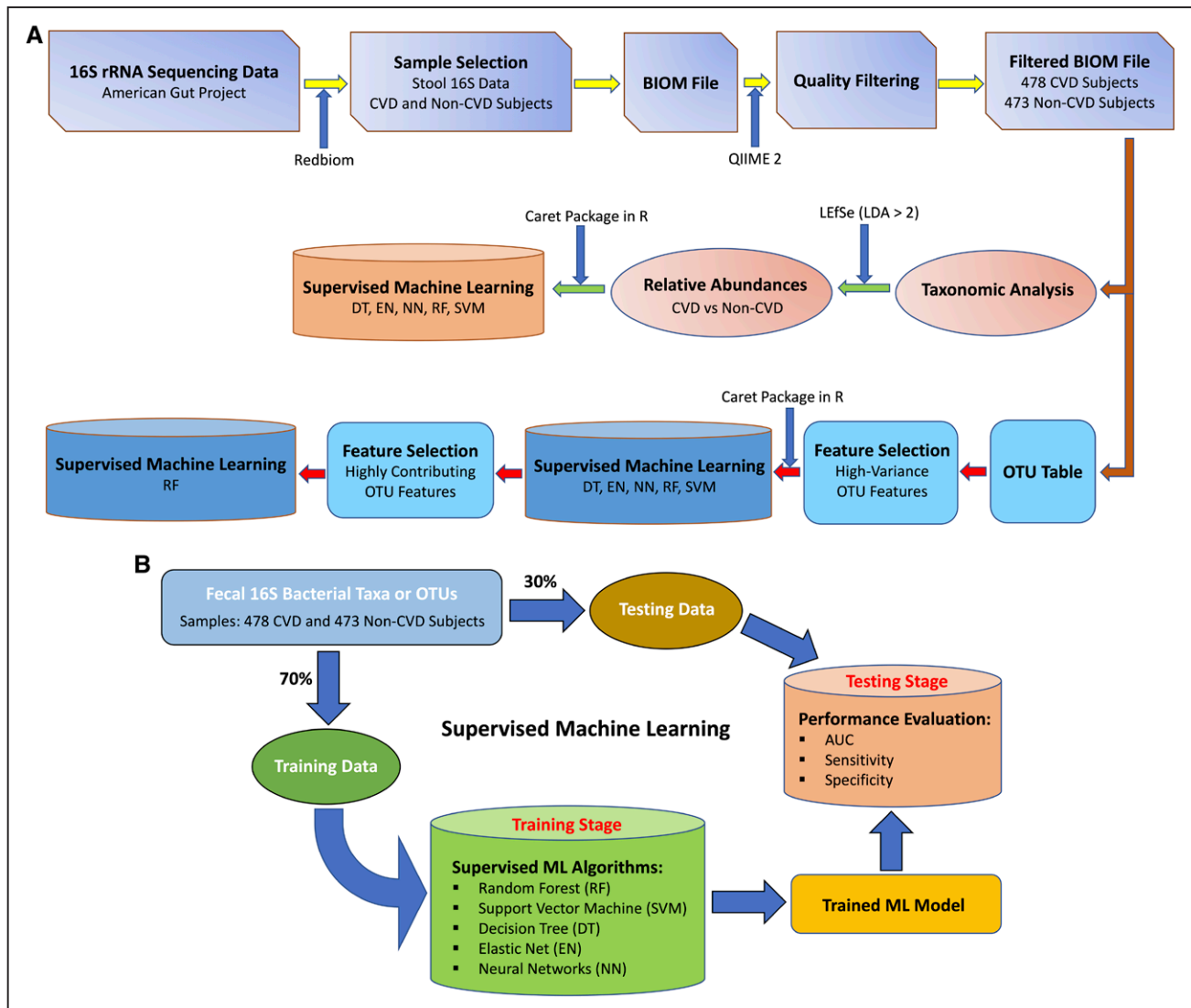
The Data Supplement is available with this article at <https://www.ahajournals.org/doi/suppl/10.1161/HYPERTENSIONAHA.120.15885>.

Correspondence to Xi Cheng, Bioinformatics and Artificial Intelligence Laboratory, Center for Hypertension and Precision Medicine, Program in Physiological Genomics, Department of Physiology and Pharmacology, University of Toledo College of Medicine and Life Sciences, Block Health Science Bldg, Rm 320, 3000 Arlington Ave, Toledo, OH 43614, Email [xi.cheng@utoledo.edu](mailto:xi.cheng@utoledo.edu) or Bina Joe, Bioinformatics and Artificial Intelligence Laboratory, Center for Hypertension and Precision Medicine, Program in Physiological Genomics, Department of Physiology and Pharmacology, University of Toledo College of Medicine and Life Sciences, Block Health Science Bldg, Rm 237, 3000 Arlington Ave, Toledo, OH 43614, Email [bina.joe@utoledo.edu](mailto:bina.joe@utoledo.edu)

© 2020 American Heart Association, Inc.

*Hypertension* is available at <https://www.ahajournals.org/journal/hyp>

DOI: 10.1161/HYPERTENSIONAHA.120.15885



**Figure 1.** The study workflow. **A**, Overall analysis. **B**, Supervised machine learning (ML). AUC indicates area under the receiver operating characteristic curve; CVD, cardiovascular disease; DT, decision tree; EN, elastic net; LDA, linear discriminant analysis; NN, neural networks; OTU, operational taxonomic unit; RF, random forest; and SVM, support vector machine.

the American Gut Project<sup>23</sup> using Redbiom.<sup>24</sup> Of a total of 16998 stool samples (as of February 11, 2020) under Qiita study ID 10317, 613 CVD samples were collected from the participants diagnosed (by a medical professional) with CVD, and 16385 non-CVD samples were collected from the participants with no CVD. Of 16385 non-CVD samples, 602 samples were randomly selected to match the final sample size of the CVD group after quality filtering. Metadata and BIOM files of the samples were downloaded using the redbiom fetch function with the context Deblur-Illumina-16S-V4-150nt-780653. The BIOM file was further processed using QIIME 2 (version 2019.10) for quality filtering to discard the samples with a total frequency <10000. The table of operational taxonomic units (OTUs) was generated using the filtered BIOM file with the BIOM format tool.<sup>25</sup> The stool 16S data collected from 478 CVD and 473 non-CVD subjects were obtained for subsequent analyses.

### Taxonomic Analysis

Taxonomic assignment was performed using QIIME 2 with a pre-trained Naive Bayes classifier on the Greengenes (version 13.8) 99% OTUs.<sup>26</sup> Linear discriminant analysis effect size<sup>27</sup> via Galaxy/Hutlab (<https://huttenhower.sph.harvard.edu/galaxy/>) was used to identify differentially abundant taxonomic features. Taxonomical features with a linear discriminant analysis score >2.0 were plotted with the linear discriminant analysis effect size bar graph and cladogram.

### Supervised ML Modeling

The process of supervised ML is summarized in Figure 1B. Five different supervised ML algorithms were trained with the features of bacterial taxa or OTUs using the caret R package<sup>28</sup>: decision tree, elastic net, neural networks, random forest (RF), and support vector machine with radial kernel. Kernlab,<sup>29</sup> randomForest,<sup>30</sup> rpart,<sup>31</sup> and glmnet<sup>32</sup> were used as the helper R packages. Data were assigned into training (70%) and testing (30%) datasets after the whole dataset was shuffled. To reduce the computational complexity and the dimensionality of the feature space, OTU-wise variance was calculated for each OTU as a preliminary task for the selection of OTU features and the top 500 OTUs with the highest variance across all the samples were selected for training the ML models. Training performance of the different ML models was evaluated by 10-fold cross-validation, and the process was repeated ten times. Hyperparameter tuning was automatically executed by caret testing 10 different values for each hyperparameter. In the testing phase, prediction performance of each ML model was evaluated by the performance parameters including area under the receiver operating characteristic curves (AUC), sensitivity, and specificity. The entire process, representing a Monte Carlo procedure,<sup>33</sup> comprising of data shuffling, data splitting, training, and testing, was independently performed for 50 iterations. The box plot representations of the values of AUC, sensitivity, and specificity were generated using the ggplot2 package<sup>34</sup> in R.

### Identification of Highly Contributing OTU Features

Highly contributing OTU features (HCOFs) were selected on the basis of variable importance scores (ranged from 0 to 100; 0, no contribution to the model; 100, contributing most to the model) calculated using the `varImp` function<sup>28</sup> from the `caret` R package. Importance scores of top OTU features were plotted using the `ggplot2` package in R. To evaluate how the selected HCOFs were able to classify the CVD and non-CVD groups, only selected HCOFs were used for ML modeling as described above.

### Statistical Analysis

Linear discriminant analysis effect size<sup>27</sup> was used to perform the Kruskal-Wallis test for differential analysis of bacterial taxa among different groups, and the linear discriminant analysis score  $>2.0$  was defined as the threshold for selecting the discriminative features. The values of mean and SD of AUC, sensitivity, and specificity were computed from the 50 independent iterations of ML modeling.

## Results

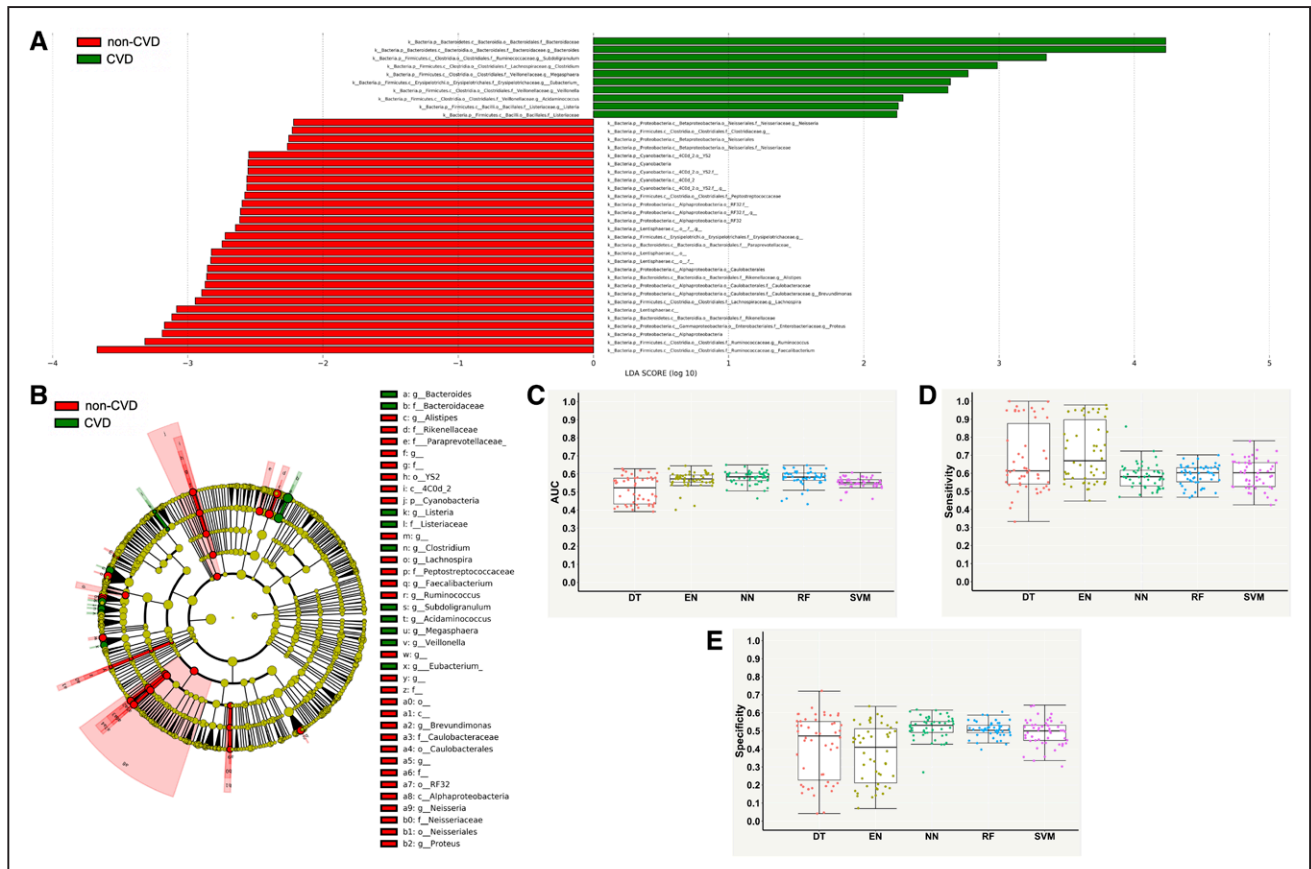
### Differential Bacterial Taxa Between the CVD and Non-CVD Groups

Significant differences in gut microbiota were observed between the CVD and non-CVD subjects (Figure 2A and 2B). A total of 39 taxonomic features (linear discriminant analysis,  $>2$ ) were found to be enriched in either CVD or non-CVD group (Figure 2A; Table S1 in the Data Supplement). For example, at the bacterial genus level, *Bacteroides*,

*Subdoligranulum*, *Clostridium*, *Megasphaera*, *Eubacterium*, *Veillonella*, *Acidaminococcus*, and *Listeria* were more abundant in the CVD group (Figure 2A). In contrast, *Faecalibacterium*, *Ruminococcus*, *Proteus*, *Lachnospira*, *Brevundimonas*, *Alistipes*, and *Neisseria* were more abundant in the non-CVD group (Figure 2A). Differential enrichments in several major bacterial taxa in the CVD and non-CVD groups and their phylogenetic relationships are presented using the cladogram (Figure 2B).

### Supervised ML Models Trained With Enriched Taxonomic Features

Supervised ML models were trained with the 39 differential taxonomic features for predictive classification and diagnostics of the CVD and non-CVD subjects. Table 1 and Figure 2C through 2E present performance measures of the 5 different ML algorithms evaluated on the testing dataset for the CVD versus non-CVD classification. RF and neural networks performed better than other models, but they only achieved an AUC of  $\approx 0.58$ , followed by elastic net ( $\approx 0.57$  AUC), support vector machine ( $\approx 0.55$  AUC), and decision tree ( $\approx 0.51$  AUC; Table 1; Figure 2C). RF and neural networks had lower sensitivity but higher specificity than elastic net, decision tree, and support vector machine (Table 1; Figure 2D and 2E).



**Figure 2.** Differential bacterial taxa between the groups of cardiovascular disease (CVD) and non-CVD and performance measures of supervised machine learning models for classifying the CVD and non-CVD subjects using differential taxonomic features. **A**, Linear discriminant analysis effect size bar graph showing differential bacterial taxa. **B**, Cladogram showing phylogenetic relationships of differential bacterial taxa. **C**, Area under the receiver operating characteristic curve (AUC). **D**, Sensitivity. **E**, Specificity. Each point in the box plot represents the corresponding performance measure in one iteration (total 50 iterations). DT indicates decision tree; EN, elastic net; LDA, linear discriminant analysis; NN, neural network; RF, random forest; and SVM, support vector machine.

**Table 1. Performance Measures of Supervised Machine Learning Models for Classifying the Cardiovascular Disease and Non-Cardiovascular Disease Subjects Using Differential Taxonomic Features and Top 500 High-Variance OTU Features**

Features	Algorithms	AUC	Sensitivity	Specificity
Bacterial Taxa	DT	0.51±0.07	0.68±0.18	0.41±0.18
	EN	0.57±0.04	0.71±0.17	0.37±0.16
	NN	0.58±0.04	0.59±0.07	0.52±0.06
	RF	0.58±0.04	0.59±0.06	0.51±0.04
	SVM	0.55±0.03	0.60±0.08	0.49±0.07
High-variance OTUs	DT	0.52±0.08	0.57±0.10	0.53±0.11
	EN	0.56±0.05	0.56±0.09	0.55±0.09
	NN	0.48±0.04	0.59±0.30	0.46±0.28
	RF	0.65±0.03	0.70±0.05	0.50±0.04
	SVM	0.57±0.04	0.60±0.07	0.52±0.09

Values are presented as mean±SD (calculated from 50 iterations). In each iteration, the entire process of data shuffling, data splitting, training, and testing was independently performed to compute for all the performance parameters. AUC indicates area under the receiver operating characteristic curve; DT, decision tree; EN, elastic net; NN, neural networks; OTU, operational taxonomic unit; RF, random forest; and SVM, support vector machine.

### Supervised ML Models Trained With High-Variance OTUs

Next, supervised ML models were trained with the top 500 high-variance OTU features, instead of taxonomic features, to test whether the diagnostic classification could be further improved. Interestingly, the testing AUC of RF was significantly improved to  $\approx 0.65$ , and its sensitivity was also significantly increased to  $\approx 0.70$  despite no significant improvement of specificity (Table 1; Figure 3). However, the AUC and specificity of neural networks significantly decreased to  $\approx 0.48$  and  $\approx 0.46$ , respectively (Table 1; Figure 3). No significant improvements in the performance measures of elastic net, decision tree, and support vector machine were observed (Table 1; Figure 3).

### Supervised ML Models Trained With HCOFs

To further improve the diagnostic classification of the RF model and also reduce the dimensionality of the OTU feature space, HCOFs were further selected from the top 500 high-variance OTU features. Variable importance scores (ranged from 0 to 100) of OTUs were calculated, and the top 100 HCOFs with the highest scores were selected for training the RF model (Figure 4A; Table S2). The RF algorithm was then reimplemented using the top 20, 25, 50, 75, and 100 HCOFs, respectively. The RF models trained with the top 20 and top 25 HCOFs not only reduced the dimensionality of the feature space but also achieved a further improvement of testing AUC ( $\approx 0.70$ ) and performed slightly better than other 3 RF models trained with  $\geq 50$  HCOFs (Table 2; Figure 4B). The RF model trained with the top 25 HCOFs had slightly higher sensitivity and specificity than the model trained with the top 20 HCOFs (Table 2; Figure 4C and 4D). Therefore, we concluded that the RF model trained with only top 25 HCOFs features could achieve a good diagnostic classification power of predicting and identifying the subjects with CVD.

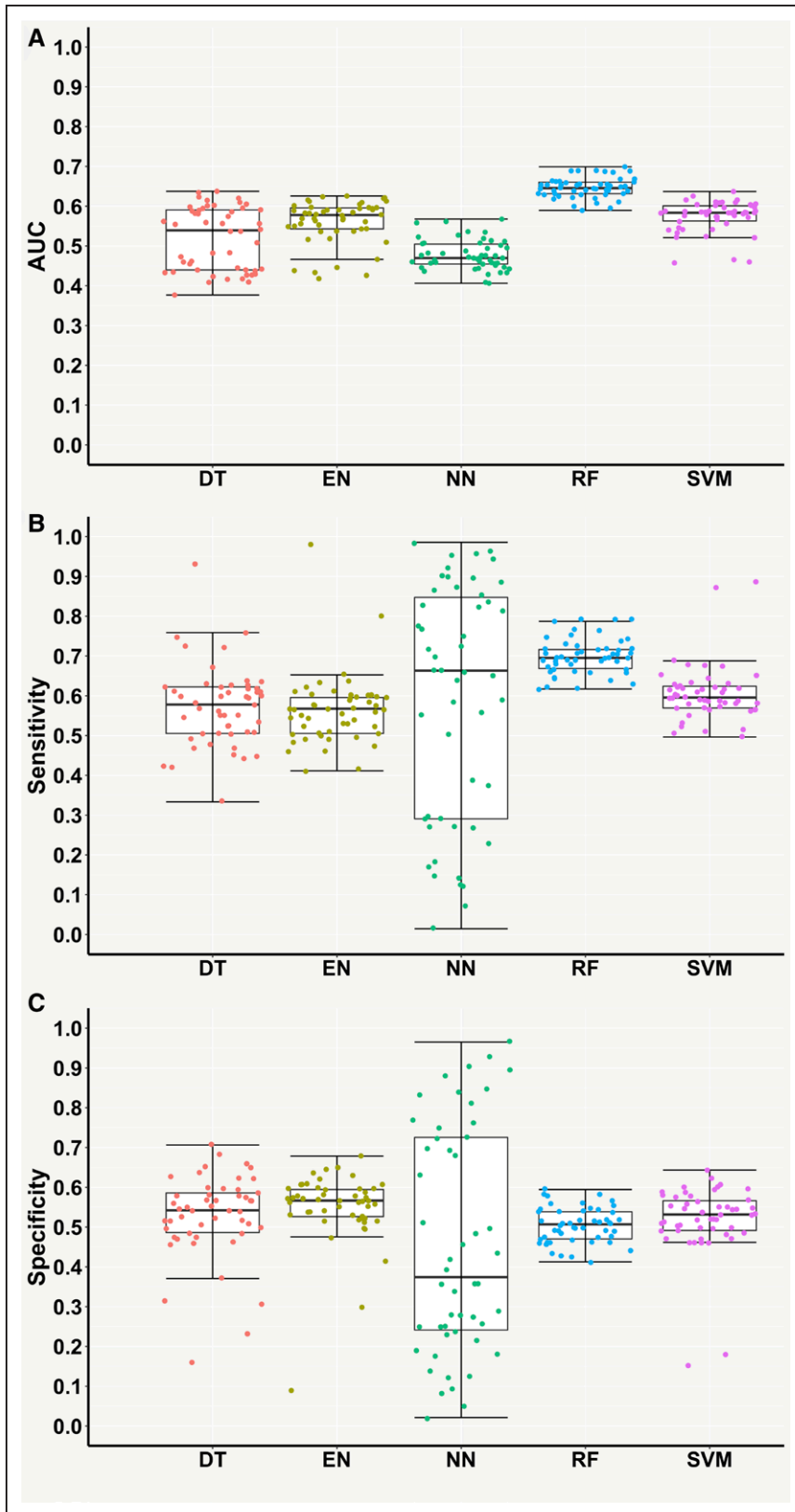
### Discussion

Mounting evidence points to a strong link between cardiovascular health and gut microbiota.<sup>35–37</sup> Albeit being highly variable between individuals, gut microbiota has been successfully

used as a feature to differentiate between health and disease in a variety of illnesses.<sup>13,38,39</sup> Therefore, in this study, we asked whether gut microbiome data can be used to diagnose CVD in humans. CVD is a broad term including a range of morbid conditions from hypertension and atherosclerosis to heart failure. As such, the host molecular mechanisms underlying a broad group of subjects classified as having CVD vary widely. Even so, we asked whether there are any early warning signs that are trackable across all of the clinical conditions, which belong under the broad class called as CVD. To this end, given the recent literature on a strong association between gut microbial communities and a variety of CVD,<sup>17,21,22,40</sup> we examined whether an alteration in gut microbial composition could serve as a common differentiator between subjects with any form of CVD and those with normal cardiovascular health. Remarkably, not only were we able to detect distinct microbial signatures (Figure 2A and 2B) but we were also successful in applying gut microbiome data as training modules for supervised ML modeling to differentiate between these 2 groups with a promising predictive diagnostics potential.

The approach of utilizing 16S metagenomics data for disease prediction using supervised ML is not new<sup>38,41–43</sup>; however, its application in CVD is novel. One of the strengths of our study is that it was conducted with a large sample size consisting of 478 CVD and 473 non-CVD human subjects. While larger sample sizes are better under a controlled setting of restricting them by a single feature such as age, for example, the cohort we used here was not limited by any features. The entire cohort was well represented by a dynamic range of various features such as ages, sexes, dietary habits, and lifestyles.<sup>23</sup> Thereby, the experimental design was more permissive to contribute to a high degree of within-group variability for a rigorous examination of the capacity of ML models using gut microbiota as the sole feature for diagnostic classification of non-CVD versus CVD. However, we have to point out the limitation that gut microbiome can be influenced by other features such as diet and medication, but those data are not fully available in the American Gut Project for a comprehensive evaluation of their impact in our current ML analysis. Moreover,

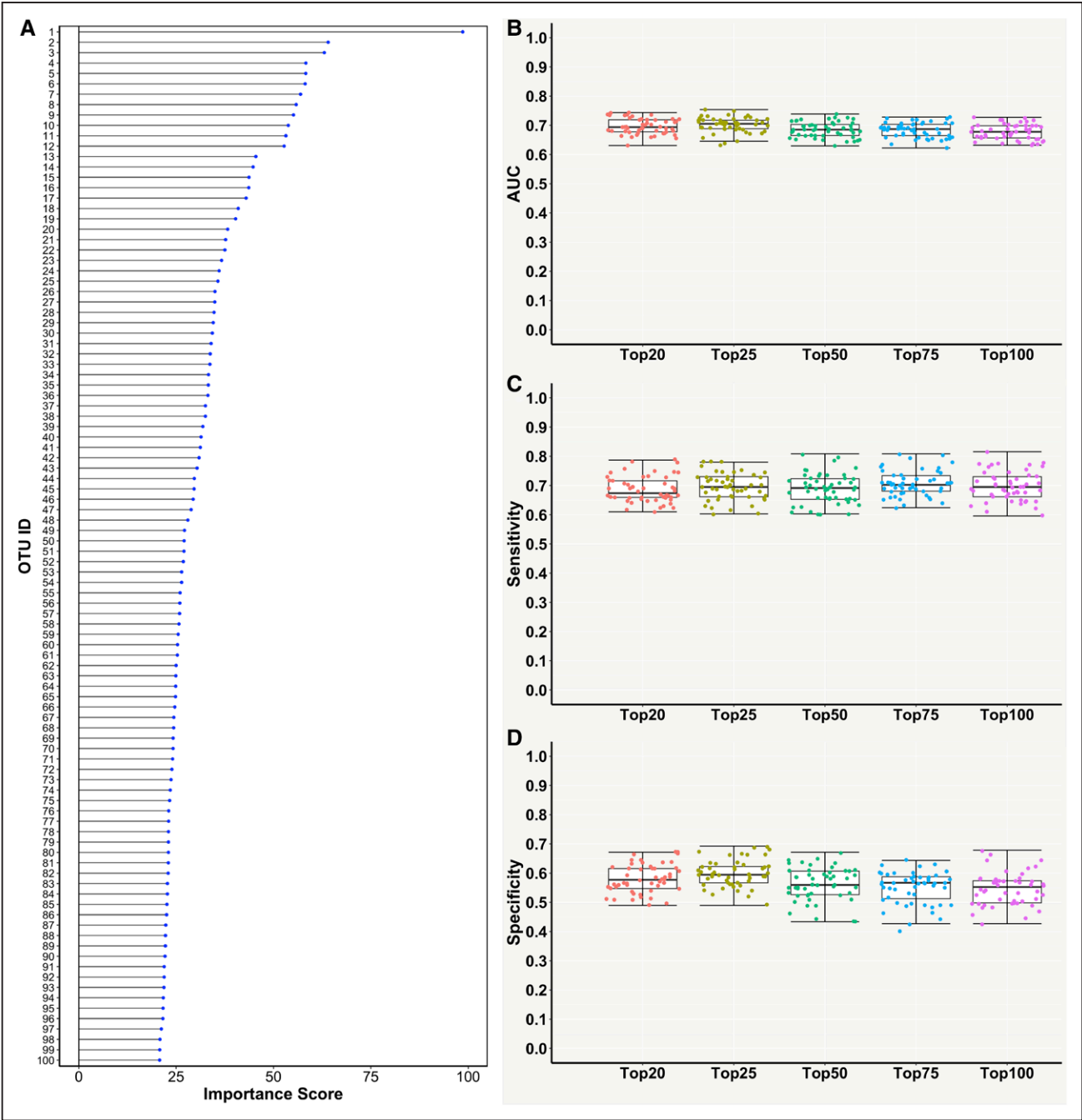




**Figure 3.** Performance measures of supervised machine learning models for classifying the cardiovascular disease (CVD) and non-CVD subjects using the top 500 high-variance operational taxonomic unit features. **A**, Area under the receiver operating characteristic curve (AUC). **B**, Sensitivity. **C**, Specificity. Each point in the box plot represents the corresponding performance measure in one iteration (total 50 iterations). DT indicates decision tree; EN, elastic net; NN, neural networks; RF, random forest; and SVM, support vector machine.

even though we only used the fecal 16S data collected from the CVD participants indicated by “diagnosed by a medical professional (doctor, physician assistant)” and the non-CVD participants indicated by “I do not have this condition” in the database of the American Gut Project, we could not rule out the

possibility of misreported or undiagnosed CVD cases. Despite this, remarkably, differential gut microbiol signatures were detectable between the CVD and non-CVD groups. These data point to a core set of altered gut microbiota as a common denominator for a variety of clinical presentations of CVD.



**Figure 4.** Performance measures of the random forest (RF) model for classifying the cardiovascular disease (CVD) and non-CVD subjects using the top highly contributing operational taxonomic unit features (HCOFs). **A**, Variable importance scores (ranged from 0 to 100) of the top 100 HCOFs. **B**, Area under the receiver operating characteristic curve (AUC). **C**, Sensitivity. **D**, Specificity. Each point in the box plot represents the corresponding performance measure in one iteration (total 50 iterations). ID indicates identifier; and OTU, operational taxonomic unit.

Initial ML modeling using these differential taxonomic features was not satisfactory and only achieved  $\approx 0.58$  AUC (Table 1; Figure 2C), indicating that the identified differential bacterial taxa were not sufficient as reliable features in the ML-based decision-making process. As OTUs differentiate bacteria based on DNA sequence similarity and represents a more informative feature than taxonomic assignment, we further tested whether OTU features could be used to train ML models to improve their prediction power. It should be noted that our study did not normalize OTU data across all the samples as we aimed to test the capacity and adaptability of ML

models trained with raw OTU data to classify and predict new unknown samples without the need for repeated processing of all the previous samples with the new samples in future. Top 500 high-variance OTU features, representing those most variable OTUs within all the CVD and non-CVD samples to provide rich feature information, were used for ML modeling, and an improved testing AUC,  $\approx 0.65$ , was achieved by the RF model (Table 1; Figure 3A). Since OTUs performed better than known taxa, it is also likely that a vast majority of the microbes that are common to all the forms of CVD are perhaps yet unknown for their taxonomic assignments.

**Table 2. Performance Measures of the Random Forest Model for Classifying the CVD Subjects and Non-CVD Subjects Using the Highly Contributing Operational Taxonomic Unit Features**

Top Features	AUC	Sensitivity	Specificity
Top 20	0.70±0.03	0.69±0.04	0.58±0.05
Top 25	0.70±0.03	0.70±0.05	0.60±0.05
Top 50	0.69±0.03	0.69±0.05	0.56±0.06
Top 75	0.68±0.03	0.71±0.04	0.55±0.06
Top 100	0.68±0.03	0.70±0.05	0.55±0.06

Values are presented as mean±SD (calculated from 50 iterations). In each iteration, the entire process of data shuffling, data splitting, training, and testing was independently performed to compute for all the performance parameters. AUC indicates area under the receiver operating characteristic curve; and CVD, cardiovascular disease.

To reduce the dimensionality of the feature space and further improve the predictive diagnostics performance, we calculated the variable importance scores of the top 500 high-variance OTUs and selected the top 100 OTUs with the highest scores as the most highly contributing features for retraining the RF model. A final testing AUC of ≈0.70 was achieved with only 25 OTU features, which were used to train the RF model (Table 2; Figure 4B). Importantly, these **high-contributing OTUs (Figure 4A; Table S2) for ML modeling could be considered as new biomarkers for future mechanistic research and clinical application.**

The current ML study differs from the prior reported ML approaches in that we used microbial composition data of stool sample, whereas almost all reported prior studies are based on health records.<sup>44–49</sup> One of those reported accuracies is through supervised ML modeling trained with multiple clinical factors, including age, sex, smoking habit, systolic blood pressure, total cholesterol, HDL (high-density lipoprotein) cholesterol, blood pressure treatment, and diabetes mellitus, to predict CVD risks, wherein an AUC of ≈0.76 was achieved.<sup>49</sup> By comparison, our study has achieved a promising AUC of ≈0.70 with a single parameter of stool gut microbiome data. While this demonstrates the promising potential of applying microbiome-based ML for predicting CVD, in future, it will be of interest to further calibrate and improve predictive capability of ML modeling by including more samples from different sources or stratifying specific types of CVD incorporated with combinatorial features such as health records, in addition to gut microbiome data.

## Perspectives

To our knowledge, our study is the first to demonstrate the promising potential of artificial intelligence via ML modeling for a convenient diagnostic screening of CVD based on fecal microbiota composition. As multiple clinical tests, such as ECG, chest X-ray, and blood work, are usually required for a comprehensive evaluation of cardiovascular health, our gut microbiome-based supervised ML approach is promising for initial routine cardiovascular health monitoring before proceeding with those various clinical tests for proper diagnosis of specific kinds of CVD. Moreover, the **ML-based feature selection approach that we described by identifying highly contributing OTUs further expands the biomarker toolkit for CVD.** Our feature selection results show that a small number of highly informative OTUs not only reduce computational

complexity of ML modeling but also further improve their diagnostic classification performances. These highly contributing OTUs could be further investigated for their pathophysiological and mechanistic implications in cardiovascular health.

## Acknowledgments

X. Cheng acknowledges the funding support from the P30 Core Center Pilot Grant from NIDA Center of Excellence in Omics, Systems Genetics, and the Addictome. B. Joe acknowledges grant support from the National Heart, Lung, and Blood Institute (HL143082).

## Sources of Funding

This work was supported by the Dean's Postdoctoral to Faculty Fellowship from University of Toledo College of Medicine and Life Sciences to X. Cheng.

## Disclosures

None.

## References

1. Bonnefont-Rousselot D. Resveratrol and cardiovascular diseases. *Nutrients*. 2016;8:250. doi: 10.3390/nu8050250
2. Cheriyan J, O'Shaughnessy KM, Brown MJ. Primary prevention of CVD: treating hypertension. *BMJ Clin Evid*. 2010;2010:0214.
3. Frostegård J. Immunity, atherosclerosis and cardiovascular disease. *BMC Med*. 2013;11:117. doi: 10.1186/1741-7015-11-117
4. Agmon Y, Khandheria BK, Meissner I, Schwartz GL, Petterson TM, O'Fallon WM, Gentile F, Whisnant JP, Wiebers DO, Seward JB. Independent association of high blood pressure and aortic atherosclerosis: a population-based study. *Circulation*. 2000;102:2087–2093. doi: 10.1161/01.cir.102.17.2087
5. Guglin M, Khan H. Pulmonary hypertension in heart failure. *J Card Fail*. 2010;16:461–474. doi: 10.1016/j.cardfail.2010.01.003
6. Hadjem M, Salem O, Naït-Abdesselam F. An ECG monitoring system for prediction of cardiac anomalies using WBAN. In: 2014 IEEE 16th International Conference on E-Health Networking, Applications and Services (Healthcom). IEEE; 2014:441–446.
7. Iijima K, Hashimoto H, Hashimoto M, Son BK, Ota H, Ogawa S, Eto M, Akishita M, Ouchi Y. Aortic arch calcification detectable on chest X-ray is a strong independent predictor of cardiovascular events beyond traditional risk factors. *Atherosclerosis*. 2010;210:137–144. doi: 10.1016/j.atherosclerosis.2009.11.012
8. Chanthong P, Lapphra K, Saihongthong S, Sricharoenchai S, Wittawatmongkol O, Phongsamart W, Rungmaitree S, Kongstan N, Chokephaibulkit K. Echocardiography and carotid intima-media thickness among asymptomatic HIV-infected adolescents in Thailand. *AIDS*. 2014;28:2071–2079. doi: 10.1097/QAD.0000000000000376
9. Bakirhan NK, Ozcelikay G, Ozkan SA. Recent progress on the sensitive detection of cardiovascular disease markers by electrochemical-based biosensors. *J Pharm Biomed Anal*. 2018;159:406–424. doi: 10.1016/j.jpba.2018.07.021
10. Ramos-Pollán R, Guevara-López MA, Suárez-Ortega C, Díaz-Herrero G, Franco-Valiente JM, Rubio-Del-Solar M, González-de-Posada N, Vaz MA, Loureiro J, Ramos I. Discovering mammography-based machine learning classifiers for breast cancer diagnosis. *J Med Syst*. 2012;36:2259–2269. doi: 10.1007/s10916-011-9693-2
11. Tapak L, Mahjub H, Hamidi O, Poorolajal J. Real-data comparison of data mining methods in prediction of diabetes in Iran. *Health Inform Res*. 2013;19:177–185. doi: 10.4258/hir.2013.19.3.177
12. Mossotto E, Ashton JJ, Coelho T, Beattie RM, MacArthur BD, Ennis S. Classification of paediatric inflammatory bowel disease using machine learning. *Sci Rep*. 2017;7:2427. doi: 10.1038/s41598-017-02606-2
13. Hacilar H, Nalbantoğlu OU, Bakir-Güngör B. Machine learning analysis of inflammatory bowel disease-associated metagenomics dataset. In: 2018 3rd International Conference on Computer Science and Engineering (UBMK). IEEE; 2018:434–438.
14. Mell B, Jala VR, Mathew AV, Byun J, Waghulde H, Zhang Y, Haribabu B, Vijay-Kumar M, Pennathur S, Joe B. Evidence for a link between gut microbiota and hypertension in the Dahl rat. *Physiol Genomics*. 2015;47:187–197. doi: 10.1152/physiolgenomics.00136.2014
15. Jose PA, Raj D. Gut microbiota in hypertension. *Curr Opin Nephrol Hypertens*. 2015;24:403–409. doi: 10.1097/MNH.0000000000000149

16. Li J, Zhao F, Wang Y, Chen J, Tao J, Tian G, Wu S, Liu W, Cui Q, Geng B, et al. Gut microbiota dysbiosis contributes to the development of hypertension. *Microbiome*. 2017;5:14. doi: 10.1186/s40168-016-0222-x
17. Sun S, Lulla A, Sioda M, Winglee K, Wu MC, Jacobs DR Jr, Shikany JM, Lloyd-Jones DM, Launer LJ, Fodor AA, et al. Gut microbiota composition and blood pressure. *Hypertension*. 2019;73:998–1006. doi: 10.1161/HYPERTENSIONAHA.118.12109
18. Yang T, Santisteban MM, Rodriguez V, Li E, Ahmari N, Carvajal JM, Zadeh M, Gong M, Qi Y, Zubcevic J, et al. Gut dysbiosis is linked to hypertension. *Hypertension*. 2015;65:1331–1340. doi: 10.1161/HYPERTENSIONAHA.115.05315
19. Yan Q, Gu Y, Li X, Yang W, Jia L, Chen C, Han X, Huang Y, Zhao L, Li P, et al. Alterations of the gut microbiome in hypertension. *Front Cell Infect Microbiol*. 2017;7:381. doi: 10.3389/fcimb.2017.00381
20. Chakraborty S, Mandal J, Cheng X, Galla S, Hindupur A, Saha P, Yeoh BS, Mell B, Yeo JY, Vijay-Kumar M, et al. Diurnal timing dependent alterations in gut microbial composition are synchronously linked to salt-sensitive hypertension and renal damage. *Hypertension*. 2020;76:59–72. doi: 10.1161/HYPERTENSIONAHA.120.14830
21. Cui X, Ye L, Li J, Jin L, Wang W, Li S, Bao M, Wu S, Li L, Geng B, et al. Metagenomic and metabolomic analyses unveil dysbiosis of gut microbiota in chronic heart failure patients. *Sci Rep*. 2018;8:635. doi: 10.1038/s41598-017-18756-2
22. Karlsson FH, Fåk F, Nookaew I, Tremaroli V, Fagerberg B, Petranovic D, Bäckhed F, Nielsen J. Symptomatic atherosclerosis is associated with an altered gut metagenome. *Nat Commun*. 2012;3:1245. doi: 10.1038/ncomms2266
23. McDonald D, Hyde E, Debelius JW, Morton JT, Gonzalez A, Ackermann G, Aksenov AA, Behsaz B, Brennan C, Chen Y, et al. American gut: an open platform for citizen science microbiome research. *mSystems*. 2018;3:e00031–e00118. doi: 10.1128/mSystems.00031-18
24. McDonald D, Kaehler B, Gonzalez A, DeReus J, Ackermann G, Marotz C, Huttley G, Knight R. redbiom: a rapid sample discovery and feature characterization system. *mSystems*. 2019;4:e00215–e00219. doi: 10.1128/mSystems.00215-19
25. McDonald D, Clemente JC, Kuczynski J, Rideout JR, Stombaugh J, Wendel D, Wilke A, Huse S, Hufnagle J, Meyer F, et al. The Biological Observation Matrix (BIOM) format or: how I learned to stop worrying and love the ome-ome. *Gigascience*. 2012;1:7. doi: 10.1186/2047-217X-1-7
26. Bokulich NA, Kaehler BD, Rideout JR, Dillon M, Bolyen E, Knight R, Huttley GA, Gregory Caporaso J. Optimizing taxonomic classification of marker-gene amplicon sequences with QIIME 2's q2-feature-classifier plugin. *Microbiome*. 2018;6:90. doi: 10.1186/s40168-018-0470-z
27. Segata N, Izard J, Waldron L, Gevers D, Miropolsky L, Garrett WS, Huttenhower C. Metagenomic biomarker discovery and explanation. *Genome Biol*. 2011;12:R60. doi: 10.1186/gb-2011-12-6-r60
28. Kuhn M. Building predictive models in R using the caret package. *J Stat Softw*. 2008;28:1–26.
29. Karatzoglou A, Smola A, Hornik K, Zeileis A. kernlab - an S4 package for kernel methods in R. *J Stat Softw*. 2004;1:3–17.
30. Liaw A, Wiener M. Classification and regression by randomforest. *Forest*. 2001;23:18–22.
31. Therneau T, Atkinson B, Ripley B, Ripley MB. Package 'rpart'. 2015. cran ma ic ac uk/web/packages/rpart/rpart.pdf. Accessed April 20, 2016.
32. Jurka TP, Collingwood L, Boydston AE, Grossman E, van Atteveldt W. RTextTools: a supervised learning package for text classification. *R J*. 2013;5:6–11.
33. Andrieu C, De Freitas N, Doucet A, Jordan MI. An introduction to MCMC for machine learning. *Mach Learn*. 2003;50:5–43.
34. Wickham H. *Ggplot2: Elegant Graphics for Data Analysis*. Springer; 2016.
35. Stock J. Gut microbiota: an environmental risk factor for cardiovascular disease. *Atherosclerosis*. 2013;229:440–442. doi: 10.1016/j.atherosclerosis.2013.05.019
36. Li XS, Obeid S, Klingenberg R, Gencer B, Mach F, Räber L, Windecker S, Rodondi N, Nanchen D, Muller O, et al. Gut microbiota-dependent trimethylamine N-oxide in acute coronary syndromes: a prognostic marker for incident cardiovascular events beyond traditional risk factors. *Eur Heart J*. 2017;38:814–824. doi: 10.1093/eurheartj/ehw582
37. Heianza Y, Ma W, Manson JE, Rexrode KM, Qi L. Gut microbiota metabolites and risk of major adverse cardiovascular disease events and death: a systematic review and meta-analysis of prospective studies. *J Am Heart Assoc*. 2017;6:e004947. doi: 10.1161/JAHA.116.004947
38. Zeller G, Tap J, Voigt AY, Sunagawa S, Kultima JR, Costea PI, Amiot A, Böhm J, Brunetti F, Habermann N, et al. Potential of fecal microbiota for early-stage detection of colorectal cancer. *Mol Syst Biol*. 2014;10:766. doi: 10.15252/msb.20145645
39. Frank DN, St Amand AL, Feldman RA, Boedeker EC, Harpaz N, Pace NR. Molecular-phylogenetic characterization of microbial community imbalances in human inflammatory bowel diseases. *Proc Natl Acad Sci U S A*. 2007;104:13780–13785. doi: 10.1073/pnas.0706625104
40. Wang Z, Klipfell E, Bennett BJ, Koeth R, Levison BS, Dugar B, Feldstein AE, Britt EB, Fu X, Chung YM, et al. Gut flora metabolism of phosphatidylcholine promotes cardiovascular disease. *Nature*. 2011;472:57–63. doi: 10.1038/nature09922
41. Douglas GM, Hansen R, Jones CMA, Dunn KA, Comeau AM, Bielawski JP, Tayler R, El-Omar EM, Russell RK, Hold GL, et al. Multi-omics differentially classify disease state and treatment outcome in pediatric Crohn's disease. *Microbiome*. 2018;6:13. doi: 10.1186/s40168-018-0398-3
42. Wingfield B, Coleman S, McGinnity TM, Bjourson AJ. A metagenomic hybrid classifier for paediatric inflammatory bowel disease. In: 2016 International Joint Conference on Neural Networks (IJCNN). *IEEE*; 2016:1083–1089.
43. Chen W, Cheng YM, Zhang SW, Pan Q. Supervised method for periodontitis phenotypes prediction based on microbial composition using 16S rRNA sequences. *Int J Comput Biol Drug Des*. 2014;7:214–224. doi: 10.1504/IJCBD.2014.061647
44. Elsayad AM, Fakhr M. Diagnosis of cardiovascular diseases with bayesian classifiers. *JCS*. 2015;11:274–282.
45. Papaloukas C, Fotiadis DI, Likas A, Michalis LK. An ischemia detection method based on artificial neural networks. *Artif Intell Med*. 2002;24:167–178. doi: 10.1016/s0933-3657(01)00100-2
46. Khalaf AF, Owis MI, Yassine IA. A novel technique for cardiac arrhythmia classification using spectral correlation and support vector machines. *Expert Syst Appl*. 2015;42:8361–8368.
47. Masetic Z, Subasi A. Congestive heart failure detection using random forest classifier. *Comput Methods Programs Biomed*. 2016;130:54–64. doi: 10.1016/j.cmpb.2016.03.020
48. Tsoi KKF, Chan NB, Yiu KKL, Poon SKS, Lin B, Ho K. Machine learning clustering for blood pressure variability applied to Systolic Blood Pressure Intervention Trial (SPRINT) and the Hong Kong Community Cohort. *Hypertension*. 2020;76:569–576. doi: 10.1161/HYPERTENSIONAHA.119.14213
49. Weng SF, Reys J, Kai J, Garibaldi JM, Qureshi N. Can machine-learning improve cardiovascular risk prediction using routine clinical data? *PLoS One*. 2017;12:e0174944. doi: 10.1371/journal.pone.0174944

## Novelty and Significance

### What Is New?

- Our study analyzed the large-scale gut microbiota data collected from a significant number of human cardiovascular disease (CVD) and non-CVD subjects and reported distinct gut microbiome features associated with cardiovascular health and disease, without any further subclassification into the various types of CVD.
- Further, this is the first study that demonstrates the successful application of artificial intelligence via gut microbiome-based machine learning modeling for potential diagnostic screening of CVD.

### What Is Relevant?

- Hypertension is one of the most significant risk factors for developing almost all kinds of CVD, and thus our gut microbiome-based supervised machine learning approach can be potentially used for routine monitoring and evaluation of hypertension-involved cardiovascular deterioration.

### Summary

Differential composition of gut microbiota was identified in human subjects diagnosed with and without CVD. Gut microbiome-based supervised machine learning modeling has been demonstrated as a promising novel approach for diagnostic screening of CVD.