

1 **DeepMicro: deep representation learning for disease prediction based on microbiome data**

2 Min Oh¹ and Liqing Zhang^{1,*}

3

4 ¹Department of Computer Science, Virginia Tech, Blacksburg, VA, USA

5 *Correspondence: lqzhang@cs.vt.edu

6

7

8 **Abstract**

9 Human microbiota plays a key role in human health and growing evidence supports the potential use of
10 microbiome as a predictor of various diseases. However, the high-dimensionality of microbiome data,
11 often in the order of hundreds of thousands, yet low sample sizes, poses great challenge for machine
12 learning-based prediction algorithms. This imbalance induces the data to be highly sparse, preventing
13 from learning a better prediction model. Also, there has been little work on deep learning applications to
14 microbiome data with a rigorous evaluation scheme. To address these challenges, we propose DeepMicro,
15 a deep representation learning framework allowing for an effective representation of microbiome profiles.
16 DeepMicro successfully transforms high-dimensional microbiome data into a robust low-dimensional
17 representation using various autoencoders and applies machine learning classification algorithms on the
18 learned representation. In disease prediction, DeepMicro outperforms the current best approaches based
19 on the strain-level marker profile in five different datasets. In addition, by significantly reducing the
20 dimensionality of the marker profile, DeepMicro accelerates the model training and hyperparameter
21 optimization procedure with 8X-30X speedup over the basic approach. DeepMicro is freely available at
22 <https://github.com/minoh0201/DeepMicro>.

23

24

25 **Introduction**

26 As our knowledge of microbiota grows, it becomes increasingly clear that the human microbiota plays a
27 key role in human health and diseases¹. The microbial community, composed of trillions of microbes, is a
28 complex and diverse ecosystem living on and inside a human. These commensal microorganisms benefit
29 humans by allowing them to harvest inaccessible nutrients and maintain the integrity of mucosal barriers
30 and homeostasis. Especially, the human microbiota contributes to the host immune system development,
31 affecting multiple cellular processes such as metabolism and immune-related functions^{1,2}. They have
32 been shown to be responsible for carcinogenesis of certain cancers and substantially affect therapeutic
33 response³. All these emerging evidences substantiate the potential use of microbiota as a predictor for
34 various diseases⁴.

35 The development of high-throughput sequencing technologies has enabled researchers to capture a
36 comprehensive snapshot of the microbial community of interest. The most common components of the
37 human microbiome can be profiled with 16S rRNA gene sequencing technology in a cost-effective way⁵.
38 Comparatively, shotgun metagenomic sequencing technology can provide a deeper resolution profile of
39 the microbial community at the strain level^{6,7}. As the cost of shotgun metagenomic sequencing keeps
40 decreasing and the resolution increasing, it is likely that a growing role of the microbiome in human health
41 will be uncovered from the mounting metagenomic datasets.

42 Although novel technologies have dramatically increased our ability to characterize human microbiome
43 and there is evidence suggesting the potential use of the human microbiome for predicting disease state,
44 how to effectively utilize the human microbiome data faces several key challenges. Firstly, effective
45 dimensionality reduction that preserves the intrinsic structure of the microbiome data is required to
46 handle the high dimensional data with low sample sizes, especially the microbiome data with strain-level
47 information that often contain hundreds of thousands of gene markers but for only some hundred or
48 fewer samples. With a low number of samples, large number of features can cause the curse of
49 dimensionality, usually inducing sparsity of the data in the feature space. Along with traditional
50 dimensionality reduction algorithms, autoencoder that learns a low-dimensional representation by
51 reconstructing the input⁸ can be applied to exploit microbiome data. Secondly, given the fast amounting
52 metagenomic data, there is an inadequate effort in adapting machine learning algorithms for predicting
53 disease state based on microbiome data. In particular, deep learning is a class of machine learning
54 algorithms that builds on large multi-layer neural networks, and that can potentially make effective use

55 of metagenomic data. With the rapidly growing attention from both academia and industry, deep learning
56 has produced unprecedented performance in various fields, including not only image and speech
57 recognition, natural language processing, and language translation but also biological and healthcare
58 research⁹. A few studies have applied deep learning approaches to abundance profiles of the human gut
59 microbiome for disease prediction^{10,11}. However, there has been no research utilizing strain-level profiles
60 for the purpose. Comparatively, strain level profiles, often containing hundreds of thousands of gene
61 markers' information, should be more informative for accurately classifying the samples into patient and
62 healthy control groups across different types of diseases than abundance profiles that usually contain only
63 a few hundred bacteria's abundance information¹². Lastly, to evaluate and compare the performance of
64 machine learning models, it is necessary to introduce a rigorous validation framework to estimate their
65 performance over unseen data. Pasolli et al., a study that built classification models based on microbiome
66 data, utilized a 10-fold cross-validation scheme that tunes the hyper-parameters on the test set without
67 using a validation set¹². This approach may overestimate model performance as it exposes the test set to
68 the model in the training procedure^{13,14}.

69 To address these issues, we propose DeepMicro, a deep representation learning framework that deploys
70 various autoencoders to learn robust low-dimensional representations from high-dimensional
71 microbiome profiles and trains classification models based on the learned representation. We applied a
72 thorough validation scheme that excludes the test set from hyper-parameter optimization to ensure
73 fairness of model comparison. Our model surpasses the current best methods in terms of disease state
74 prediction of inflammatory bowel disease, type 2 diabetes in the Chinese cohort as well as European
75 women cohort, liver cirrhosis, and obesity. DeepMicro is open-sourced and publicly available software to
76 benefit future research, allowing researchers to obtain a robust low-dimensional representation of
77 microbiome profiles with user-defined deep architecture and hyper-parameters.

78

79 **Methods**

80 *Dataset and Extracting Microbiome Profiles*

81 We considered publicly available human gut metagenomic samples of six different disease cohorts:
82 inflammatory bowel disease (IBD), type 2 diabetes in European women (EW-T2D), type 2 diabetes in
83 Chinese (C-T2D) cohort, obesity (Obesity), liver cirrhosis (Cirrhosis), and colorectal cancer (Colorectal). All

84 these samples were derived from whole-genome shotgun metagenomic studies that used Illumina paired-
 85 end sequencing technology. Each cohort consists of healthy control and patient samples as shown in Table
 86 1. IBD cohort has 25 individuals with inflammatory bowel disease and 85 healthy controls ¹⁵. EW-T2D
 87 cohort has 53 European women with type 2 diabetes and 43 healthy European women ¹⁶. C-T2D cohort
 88 has 170 Chinese individuals with type 2 diabetes and 174 healthy Chinese controls ¹⁷. Obesity cohort has
 89 164 obese patients and 89 non-obese controls ¹⁸. Cirrhosis cohort has 118 patients with liver cirrhosis and
 90 114 healthy controls ¹⁹. Colorectal cohort has 48 colorectal cancer patients and 73 healthy controls ²⁰. In
 91 total, 1,156 human gut metagenomic samples, obtained from MetAML repository ¹², were used in our
 92 experiments.

93

94 Table 1. Human gut microbiome datasets used for disease state prediction

Disease	Dataset name	# total samples	# of healthy controls	# of patient samples	Data source references
Inflammatory Bowel Disease	IBD	110	85	25	¹⁵
Type 2 Diabetes	EW-T2D	96	43	53	¹⁶
	C-T2D	344	174	170	¹⁷
Obesity	Obesity	253	89	164	¹⁸
Liver Cirrhosis	Cirrhosis	232	114	118	¹⁹
Colorectal Cancer	Colorectal	121	73	48	²⁰

95

96 Two types of microbiome profiles were extracted from the metagenomic samples: 1) strain-level marker
 97 profile and 2) species-level relative abundance profile. MetaPhlan2 was utilized to extract these profiles
 98 with default parameters ⁷. We utilized MetAML to preprocess the abundance profile by selecting species-
 99 level features and excluding sub-species-level features ¹². The strain-level marker profile consists of binary
 100 values indicating the presence (1) or absence (0) of a certain strain. The species-level relative abundance
 101 profile consists of real values in [0,1] indicating the percentages of the species in the total observed
 102 species. The abundance profile has a few hundred dimensions, whereas the marker profile has a much
 103 larger number of dimensions, up to over a hundred thousand in the current data (Table 2).

104 Table 2. The number of dimensions of the preprocessed microbiome profiles

Profile type	IBD	EW-T2D	C-T2D	Obesity	Cirrhosis	Colorectal
marker profile	91,756	83,456	119,792	99,568	120,553	108,034
abundance profile	443	381	572	465	542	503

105

106

107 *Deep Representation Learning*

108 An autoencoder is a neural network reconstructing its input x . Internally, its general form consists of an
109 encoder function $f_{\phi}(\cdot)$ and a decoder function $f'_{\theta}(\cdot)$ where ϕ and θ are parameters of encoder and
110 decoder functions, respectively. An autoencoder is trained to minimize the difference between an input
111 x and a reconstructed input x' , the reconstruction loss (e.g., squared error) that can be written as follows:

$$112 \quad L(x, x') = \|x - x'\|^2 = \left\| x - f'_{\theta} \left(f_{\phi}(x) \right) \right\|^2.$$

113 After training an autoencoder, we are interested in obtaining a latent representation $z = f_{\phi}(x)$ of the
114 input using the trained encoder. The latent representation, usually in a much lower-dimensional space
115 than the original input, contains sufficient information for reconstructing the original input as close as
116 possible. We utilized this representation to train classifiers for disease prediction.

117 For the DeepMicro framework, we incorporated various deep representation learning techniques,
118 including shallow autoencoder (SAE), deep autoencoder (DAE), variational autoencoder (VAE), and
119 convolutional autoencoder (CAE), to learn a low-dimensional embedding for microbiome profiles. Note
120 that the diverse combinations of hyper-parameters defining the structure of autoencoders (e.g., the
121 number of units and layers) have been explored in a grid fashion as described below, however, users are
122 not limited to the tested hyper-parameters and can use their own hyper-parameter grid fitted to their
123 data.

124 Firstly, we utilized SAE, the simplest autoencoder structure composed of the encoder part where the input
125 layer is fully connected with the latent layer, and the decoder part where the output layer produces
126 reconstructed input x' by taking weighted sums of outputs of the latent layer. We introduced a linear
127 activation function for the latent and output layer. Other options for the loss and activation functions are
128 available for users (such as binary cross-entropy and sigmoid function). Initial values of the weights and
129 bias were initialized with Glorot uniform initializer²¹. We examined five different sizes of dimensions for
130 the latent representation (32, 64, 128, 256, and 512).

131 In addition to the SAE model, we implemented the DAE model by introducing hidden layers between the
132 input and latent layers as well as between the latent and output layers. All of the additional hidden layers

133 were equipped with Rectified Linear Unit (ReLU) activation function and Glorot uniform initializer. The
134 same number of hidden layers (one layer or two layers) were inserted into both encoder and decoder
135 parts. Also, we gradually increased the number of hidden units. The number of hidden units in the added
136 layers was set to the double of the successive layer in the encoder part and to the double of the preceding
137 layer in the decoder part. With this setting, model complexity is controlled by both the number of hidden
138 units and the number of hidden layers, maintaining structural symmetry of the model. For example, if the
139 latent layer has 512 hidden units and if two layers are inserted to the encoder and decoder parts, then
140 the resulting autoencoder has 5 hidden layers with 2048, 1024, 512, 1024, and 2048 hidden units,
141 respectively. Similar to SAE, we varied the number of hidden units in the latent layer as follows: 32, 64,
142 128, 256, 512, thus, in total, we tested 10 different DAE architectures (Table S2).

143 A variational autoencoder (VAE) learns probabilistic representations z given input x and then use these
144 representations to reconstruct input x' ²². Using variational inference, the true posterior distribution of
145 latent embeddings (i.e., $p(z|x)$) can be approximated by the introduced posterior $q_\phi(z|x)$ where ϕ are
146 parameters of an encoder network. Unlike the previous autoencoders learning an unconstrained
147 representation, VAE learns a generalized latent representation under the assumption that the posterior
148 approximation follows Gaussian distribution. The encoder network encodes the means and variances of
149 the multivariate Gaussian distribution. The latent representation z can be sampled from the learned
150 posterior distribution $q_\phi(z|x) \sim N(\mu, \Sigma)$. Then the sampled latent representation is passed into the
151 decoder network to generate the reconstructed input $x' \sim g_\theta(x|z)$ where θ are the parameters of the
152 decoder.

153 To approximate the true posterior, we need to minimize the Kullback-Leibler (KL) divergence between the
154 introduced posterior and the true posterior,

$$155 \quad KL(q_\phi(z|x)||p(z|x)) = -ELBO(\phi, \theta; x) + \log(p(x)),$$

156 rewritten as

$$157 \quad \log(p(x)) = ELBO(\phi, \theta; x) + KL(q_\phi(z|x)||p(z|x)),$$

158 where $ELBO(\phi, \theta; x)$ is an evidence lower bound on the log probability of the data because the KL term
159 must be greater than or equal to zero. It is intractable to compute the KL term directly but minimizing the
160 KL divergence is equivalent to maximizing the lower bound, decomposed as follows:

161 $ELBO(\phi, \theta; x) = \mathbb{E}_{q_{\phi}(z|x)}[\log(g_{\theta}(x|z))] - KL(q_{\phi}(z|x)||p(z)).$

162 The final objective function can be induced by converting the maximization problem to the minimization
163 problem.

164
$$L(\phi, \theta; x) = -\mathbb{E}_{q_{\phi}(z|x)}[\log(g_{\theta}(x|z))] + KL(q_{\phi}(z|x)||p(z))$$

165 The first term can be viewed as a reconstruction term as it forces the inferred latent representation to
166 recover its corresponding input and the second KL term can be considered as a regularization term to
167 modulate the posterior of the learned representation to be Gaussian distribution. We used ReLU
168 activation and Glorot uniform initializer for intermediate hidden layers in encoder and decoder. One
169 intermediate hidden layer was used and the number of hidden units in it varied from 32, 64, 128, 256, to
170 512. The latent layer was set to 4, 8, or 16 units. Thus, altogether we tested 15 different model structures.

171 Instead of fully connected layers, a convolutional autoencoder (CAE) is equipped with convolutional layers
172 in which each unit is connected to only local regions of the previous layer²³. A convolutional layer consists
173 of multiple filters (kernels) and each filter has a set of weights used to perform convolution operation that
174 computes dot products between a filter and a local region²⁴. We used ReLU activation and Glorot uniform
175 initializer for convolutional layers. We did not use any pooling layer as it may generalize too much to
176 reconstruct an input. The n -dimensional input vector was reshaped like a squared image with a size of
177 $d \times d \times 1$ where $d = \lfloor \sqrt{n} \rfloor + 1$. As $d^2 \geq n$, we padded the rest part of the reshaped input with zeros. To
178 be flexible to an input size, the filter size of the first convolutional layer was set to 10% of the input width
179 and height, respectively (i.e. $[0.1d] \times [0.1d]$). For the first convolutional layer, we used 25% of the filter
180 size as the size of stride which configures how much we slide the filter. For the following convolutional
181 layers in the encoder part, we used 10% of the output size of the preceding layer as the filter size and 50%
182 of this filter size as the stride size. All units in the last convolutional layer of the encoder part have been
183 flattened in the following flatten layer which is designated as a latent layer. We utilized convolutional
184 transpose layers (deconvolutional layers) to make the decoder symmetry to the encoder. In our
185 experiment, the number of filters in a convolutional layer was set to half of that of the preceding layer for
186 the encoder part. For example, if the first convolutional layer has 64 filters and there are three
187 convolutional layers in the encoder, then the following two convolutional layers have 32 and 16 filters,
188 respectively. We varied the number of convolutional layers from 2 to 3 and tried five different numbers

189 of filters in the first convolutional layer (4, 8, 16, 32, and 64). In total, we tested 10 different CAE model
190 structures.

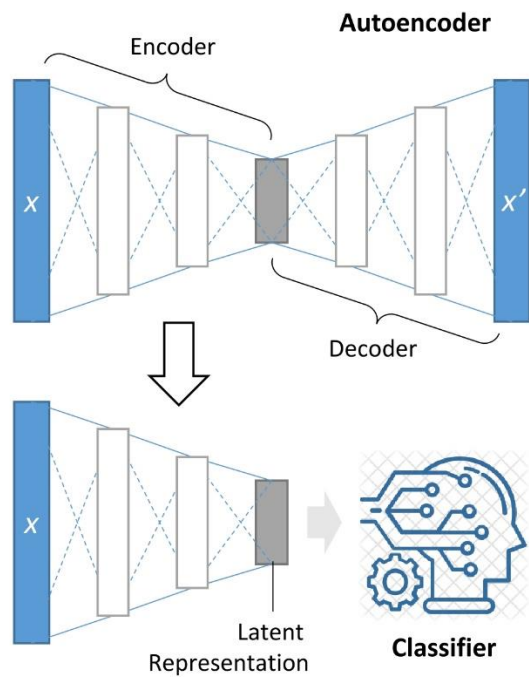
191 To train deep representation models, we split each dataset into a training set, a validation set, and a test
192 set (64% training set, 16% validation set, and 20% test set; Figure S1). Note that the test set was withheld
193 from training the model. We used the early-stopping strategy, that is, trained the models on the training
194 set, computed the reconstruction loss for the validation set after each epoch, stopped the training if there
195 was no improvement in validation loss during 20 epochs, and then selected the model with the least
196 validation loss as the best model. We used mean squared error for reconstruction loss and applied
197 adaptive moment estimation (Adam) optimizer for gradient descent with default parameters (learning
198 rate: 0.001, epsilon: 1e-07) as provided in the original paper²⁵. We utilized the encoder part of the best
199 model to produce a low-dimensional representation of the microbiome data for downstream disease
200 prediction.

201

202 *Prediction of disease states based on the learned representation*

203 We built classification models based on the encoded low-dimensional representations of microbiome
204 profiles (Figure 1). Three machine learning algorithms, support vector machine (SVM), random forest (RF),
205 and Multi-Layer Perceptron (MLP), were used. We explored hyper-parameter space with grid search. SVM
206 maximizes the margin between the supporting hyperplanes to optimize a decision boundary separating
207 data points of different classes²⁶. In this study, we utilized both radial basis function (RBF) kernel and a
208 linear kernel function to compute decision margins in the transformed space to which the original data
209 was mapped. We varied penalty parameter C (2^{-5} , 2^{-3} , ..., 2^5) for both kernels as well as kernel coefficient
210 γ (2^{-15} , 2^{-13} , ..., 2^3) for RBF kernel. In total, 60 different combinations of hyper-parameters were
211 examined to optimize SVM (Table S2).

212



213

214 Figure 1. DeepMicro framework. An autoencoder is trained to map the input X to the low-dimensional
215 latent space with the encoder and to reconstruct X with the decoder. The encoder part is reused to
216 produce a latent representation of any new input X that is in turn fed into a classification algorithm to
217 determine whether the input is the positive or negative class.

218

219 RF builds multiple decision trees based on various sub-samples of the training data and merges them to
220 improve the prediction accuracy. The size of sub-samples is the same as that of training data but the
221 samples are drawn randomly with replacement from the training data. For the hyper-parameter grid of
222 RF classifier, the number of trees (estimators) was set to 100, 300, 500, 700, and 900, and the minimum
223 number of samples in a leaf node was altered from 1 to 5. Also, we tested two criteria, Gini impurity and
224 information gain, for selecting features to split a node in a decision tree. For the maximum number of
225 features considered to find the best split at each split, we used a square root of n and a logarithm to base
226 2 of n (n is the sample size). In total, we tested 100 combinations of hyper-parameters of RF.

227 MLP is an artificial neural network classifier that consists of an input layer, hidden layers, and an output
228 layer. All of the layers are fully connected to their successive layer. We used ReLu activations for all hidden
229 layers and sigmoid activation for the output layer that has a single unit. The number of units in the hidden

230 layers was set to half of that of the preceding layer except the first hidden layer. We varied the number
231 of hidden layers (1, 2, and 3), the number of epochs (30, 50, 100, 200, and 300), the number of units in
232 the first hidden layer (10, 30, 50, 100), and dropout rate (0.1 and 0.3). In total, 120 hyper-parameter
233 combinations were tested in our experiment.

234 We implemented DeepMicro in Python 3.5.2 using machine learning and data analytics libraries, including
235 Numpy 1.16.2, Pandas 0.24.2, Scipy 1.2.1, Scikit-learn 0.20.3, Keras 2.2.4, and Tensorflow 1.13.1. Source
236 code is publicly available at the git repository (<https://github.com/minoh0201/DeepMicro>).

237

238 *Performance Evaluation*

239 To avoid an overestimation of prediction performance, we designed a thorough performance evaluation
240 scheme (Figure S1). For a given dataset (e.g. Cirrhosis), we split it into training and test set in the ratio of
241 8:2 with a given random partition seed, keeping a ratio between classes in both training and test set to be
242 the same as that of the given dataset. Using only the training set, a representation learning model was
243 trained. Then, the learned representation model was applied to the training set and test set to obtain
244 dimensionality-reduced training and test set. After the dimensionality has been reduced, we conducted
245 5-fold cross-validation on the training set by varying hyper-parameters of classifiers. The best hyper-
246 parameter combination for each classifier was selected by averaging an accuracy metric of the five
247 different results. The area under the receiver operating characteristics curve (AUC) was used for
248 performance evaluation. We trained a final classification model using the whole training set with the best
249 combination of hyper-parameters and tested it on the test set. This procedure was repeated five times by
250 changing the random partition seed at the beginning of the procedure. The resulting AUC scores were
251 averaged and the average was used to compare model performance.

252

253

254 **Results**

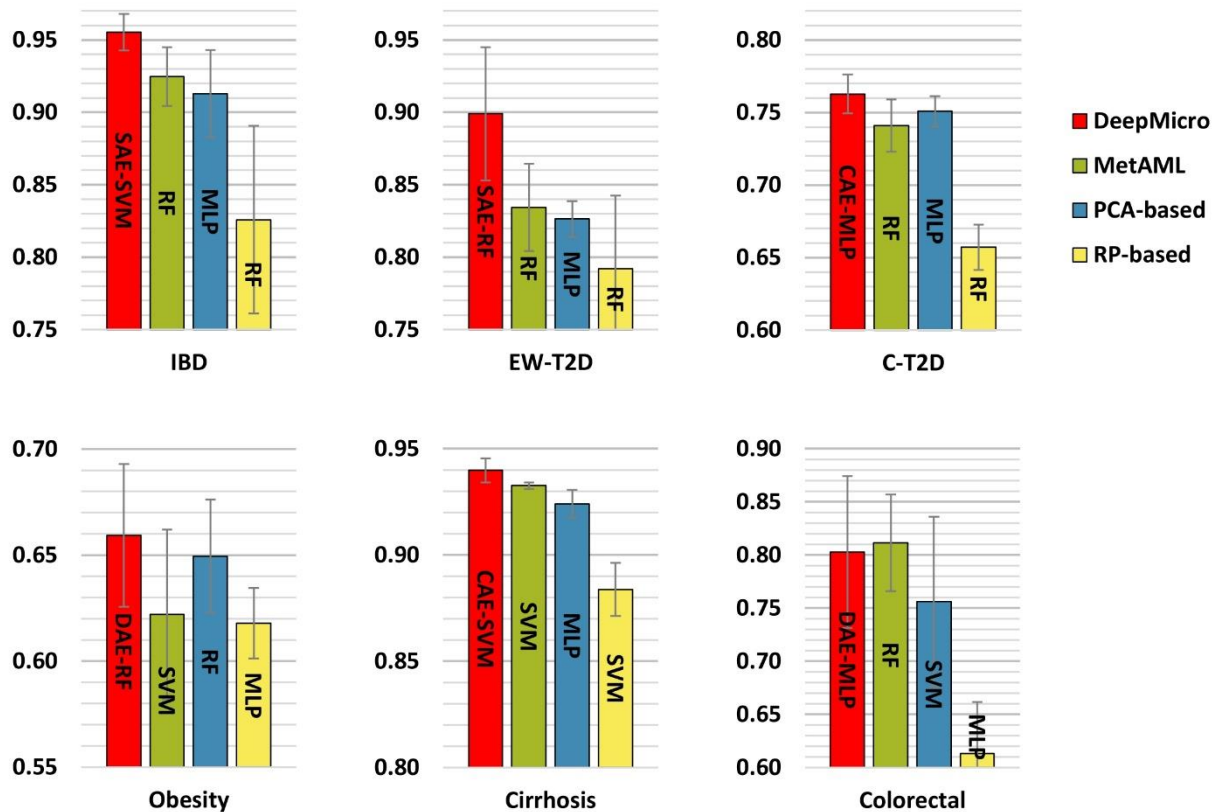
255 We developed DeepMicro, a deep representation learning framework for predicting individual phenotype
256 based on microbiome profiles. Various autoencoders (SAE, DAE, VAE, and CAE) have been utilized to learn
257 a low-dimensional representation of the microbiome profiles. Then three classification models including

258 SVM, RF, and MLP were trained on the learned representation to discriminate between disease and
259 control sample groups. We tested our framework on six disease datasets (Table 1), including inflammatory
260 bowel disease (IBD), type 2 diabetes in European women (EW-T2D), type 2 diabetes in Chinese (C-T2D),
261 obesity (Obesity), liver cirrhosis (Cirrhosis), and colorectal cancer (Colorectal). For all the datasets, two
262 types of microbiome profiles, strain-level marker profile and species-level relative abundance profile, have
263 been extracted and tested (Table 2). Also, we devised a thorough performance evaluation scheme that
264 isolates the test set from the training and validation sets in the hyper-parameter optimization phase to
265 compare various models (See Methods and Figure S1).

266 We compared our method to the current best approach (MetAML) that directly trained classifiers, such
267 as SVM and RF, on the original microbiome profile¹². We utilized the same hyper-parameters grid used in
268 MetAML for each classification algorithm. In addition, we tested Principal Component Analysis (PCA) and
269 Gaussian Random Projection (RP), using them as the replacement of the representation learning to
270 observe how traditional dimensionality reduction algorithms behave. For PCA, we selected the principal
271 components explaining 99% of the variance in the data²⁷. For RP, we set the number of components to
272 be automatically adjusted according to Johnson-Lindenstrauss lemma (eps parameter was set to 0.5)²⁸⁻³⁰.

273 We picked the best model for each approach in terms of prediction performance and compared the
274 approaches across the datasets. Figure 2 shows the results of DeepMicro and the other approaches for
275 the strain-level marker profile. DeepMicro outperforms the other approaches for five datasets, including
276 IBD (AUC = 0.955), EW-T2D (AUC = 0.899), C-T2D (AUC = 0.763), Obesity (AUC = 0.659), and Cirrhosis (AUC
277 = 0.940). For Colorectal dataset, DeepMicro has slightly lower performance than the best approach
278 (DeepMicro's AUC = 0.803 vs. MetAML's AUC = 0.811). The marker profile-based models generally
279 perform better than the abundance profile-based models (Figure S8 and S2). The only exception is Obesity
280 dataset for which the abundance-based DeepMicro model shows better performance (AUC = 0.674). Note
281 that as AUC could be misleading in an imbalanced classification scenario³¹, we also evaluated the area
282 under the precision-recall curve (AUPRC) for the imbalanced data set IBD and observed the same trend
283 between AUC and AUPRC (Table S3).

284

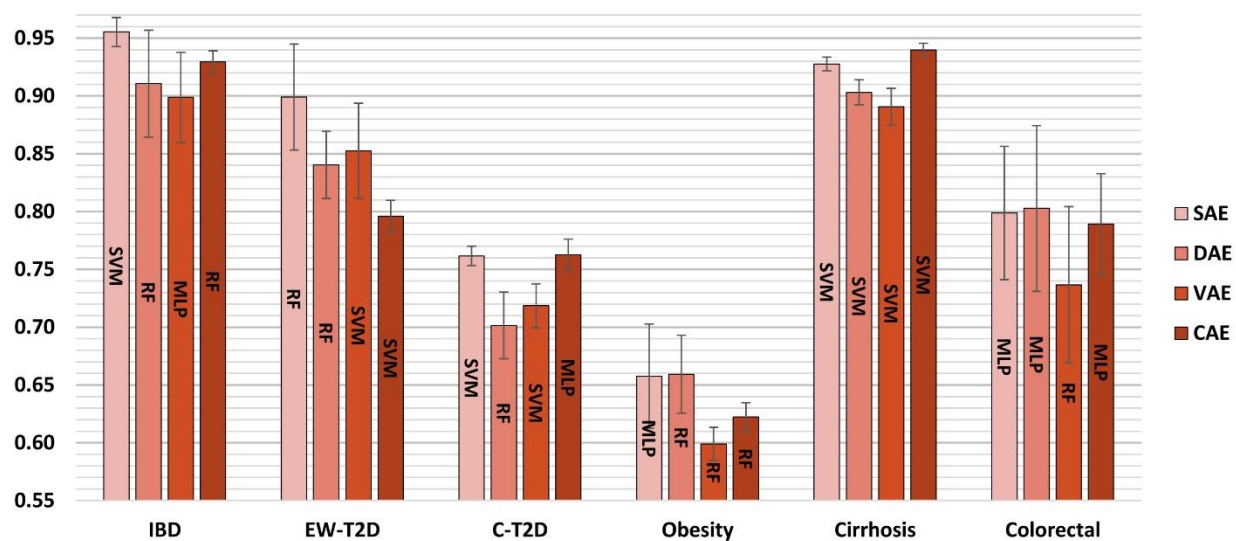


285

286 Figure 2. Disease prediction performance for marker profile-based models. Prediction performance of
287 various methods built on marker profile has been assessed with AUC. MetAML utilizes support vector
288 machine (SVM) and random forest (RF), and the superior model is presented (green). Principal component
289 analysis (PCA; blue) and gaussian random projection (RP; yellow) have been applied to reduce dimensions
290 of datasets before classification. DeepMicro (red) applies shallow autoencoder (SAE), deep autoencoder
291 (DAE), variational autoencoder (VAE), and convolutional autoencoder (CAE) for dimensionality reduction.
292 Then SVM, RF, and multi-layer perceptron (MLP) classification algorithms have been used.

293

294 For marker profile, none of the autoencoders dominate across the datasets in terms of getting the best
295 representation for classification. Also, the best classification algorithm varied according to the learned
296 representation and to the dataset (Figure 3). For abundance profile, CAE dominates over the other
297 autoencoders with RF classifier across all the datasets (Figure S3).



298

299 Figure 3. Disease prediction performance for different autoencoders based on marker profile (assessed
300 with AUC). Classifiers used: support vector machine (SVM), random forest (RF), and multi-layer perceptron
301 (MLP); Autoencoders used: shallow autoencoder (SAE), deep autoencoder (DAE), variational autoencoder
302 (VAE), and convolutional autoencoder (CAE)

303

304 We also directly trained MLP on the dataset without representation learning and compared the prediction
305 performance with that of the traditional approach (the best between SVM and RF). It is shown that MLP
306 performs better than MetAML in three datasets, EW-T2D, C-T2D, and Obesity, when marker profile is used
307 (Figure S4). However, when abundance profile is used, the performance of MLP was worse than that of
308 the traditional approach across all the datasets (Figures S5).

309 Furthermore, we compared running time of DeepMicro on marker profiles with a basic approach not using
310 representation learning. For comparison, we tracked both training time and representation learning time.
311 For each dataset, we tested the best performing representation learning model producing the highest
312 AUC score (i.e. SAE for IBD and EW-T2D, DAE for Obesity and Colorectal, and CAE for C-T2D and Cirrhosis;
313 Table S1). We fixed the seed for random partitioning of the data, and applied the formerly used
314 performance evaluation procedure where 5-fold cross-validation is conducted on the training set to
315 obtain the best hyper-parameter with which the best model is trained on the whole training set and is
316 evaluated on the test set (See Methods). The computing machine we used for timestamping is running on
317 Ubuntu 18.04 and equipped with an Intel Core i9-9820X CPU (10 cores), 64 GB Memory, and a GPU of

318 NVIDIA GTX 1080 Ti. We note that our implementation utilizes GPU when it learns representations and
 319 switches to CPU mode to exhaustively use multiple cores in a parallel way to find best hyper-parameters
 320 of the classifiers. Table 3 shows the benchmarking result on marker profile. It is worth noting that
 321 DeepMicro is 8X to 30X times faster than the basic approach (17X times faster on average). Even if MLP is
 322 excluded from the benchmarking because it requires heavy computation, DeepMicro is up to 5X times
 323 faster than the basic (2X times faster on average).

324

325 Table 3. Time benchmark for DeepMicro and basic approaches without representation learning (in sec)

Method		IBD	EW-T2D	C-T2D	Obesity	Cirrhosis	Colorectal
Basic approach	SVM*	126	85	1705	711	777	187
	RF	42	41	99	79	72	50
	MLP	3,776	2,449	12,057	8,186	8,593	4,508
Total elapsed		3,943	2,575	13,861	8,976	9,442	4,745
DeepMicro	RL	74	194	554	113	521	215
	SVM	2	2	8	8	17	2
	RF	28	28	47	33	40	30
	MLP	103	93	188	137	287	105
Total elapsed		207	317	798	291	864	352

326 *RL: Representation Learning; SVM: Support Vector Machine; RF: Random Forest; MLP: Multi-layer Perceptron

327

328 Discussion

329 We developed a deep learning framework transforming a high-dimensional microbiome profile into a low-
 330 dimensional representation and building classification models based on the learned representation. At
 331 the beginning of this study, the main goal was to reduce dimensions as strain-level marker profile has too
 332 many dimensions to handle, expecting that noisy and unnecessary information fades out and the refined
 333 representation becomes tractable for downstream prediction. Firstly, we tested PCA on marker profile
 334 and it showed a slight improvement in prediction performance for C-T2D and Obesity but not for the
 335 others. The preliminary result indicates that either some of the meaningful information was dropped or
 336 noisy information still remains. To learn meaningful feature representations, we trained various
 337 autoencoders on microbiome profiles. Our intuition behind the autoencoders was that the learned
 338 representation should keep essential information in a condensed way because autoencoders are forced
 339 to prioritize which properties of the input should be encoded during the learning process. We found that
 340 although the most appropriate autoencoder usually allows for better representation that in turn results

341 in better prediction performance, what kind of autoencoder is appropriate highly depends on problem
342 complexity and intrinsic properties of the data.

343 In the previous study, it has been shown that adding healthy controls of the other datasets could improve
344 prediction performance assessed by AUC ¹². To check if this finding can be reproduced, for each dataset,
345 we added control samples of the other datasets only into the training set and kept the test set the same
346 as before. Figure S6 shows the difference between the best performing models built with and without
347 additional controls. In general, prediction performance dropped (on average by 0.037) once negative
348 (control) samples are introduced to the training set across the datasets in almost all approaches except
349 only a few cases (Figure S6). In contrast to the previous study, the result indicates that the insertion of
350 only negative samples into the training set may not help to improve the classification models, and a
351 possible explanation might be that changes in the models rarely contribute to improving the classification
352 of positive samples ³². Interestingly, if we added negative samples into the whole data set before split it
353 into training and test set, we usually observed improvements in prediction performance. However, we
354 found that these improvements are trivial because introducing negative samples into the test set easily
355 reduces false positive rate (as the denominator of false positive rate formula is increased), resulting in
356 higher AUC scores.

357 Even though adding negative samples might not be helpful for a better model, it does not mean that
358 additional samples are meaningless. We argue that more samples can improve prediction performance,
359 especially when a well-balanced set of samples is augmented. To test this argument, we gradually
360 increased the proportion of the training set and observed how prediction performance changed over the
361 training sets of different sizes. Generally, improved prediction performance has been observed as more
362 data of both positive and negative samples are included (Figure S7). With the continued availability of
363 large samples of microbiome data, the deep representation learning framework is expected to become
364 increasingly effective for both condensed representation of the original data and also downstream
365 prediction based on the deep representation.

366 DeepMicro is publicly available software which offers cutting-edge deep learning techniques for learning
367 meaningful representations from the given data. Researchers can apply DeepMicro to their high-
368 dimensional microbiome data to obtain a robust low-dimensional representation for the subsequent
369 supervised or unsupervised learning. For predictive problems increasingly studied with microbiome data
370 such as drug response prediction, forensic human identification, and food allergy prediction, deep

371 representation learning might be useful in terms of boosting the model performance. Moreover, it might
372 be worthwhile to use the learned representation for clustering analysis. The distance between data points
373 in the latent space can be a basis for clustering microbiome samples and it could help capture the shared
374 characteristics within a group which are difficult to be identified in the original data space. DeepMicro has
375 been used to deal with microbiome data but it is not limited to a specific type of data and its application
376 can be extended to various omics data, such as genome and proteome data.

377

378

379 **Data availability**

380 All data and codes are available at <https://github.com/minoh0201/DeepMicro>.

381 **Reference**

- 382 1 Cho, I. & Blaser, M. J. The human microbiome: at the interface of health and disease. *Nature*
383 *Reviews Genetics* **13**, 260 (2012).
- 384 2 Huttenhower, C. *et al.* Structure, function and diversity of the healthy human microbiome.
385 *nature* **486**, 207 (2012).
- 386 3 McQuade, J. L., Daniel, C. R., Helmink, B. A. & Wargo, J. A. Modulating the microbiome to
387 improve therapeutic response in cancer. *The Lancet Oncology* **20**, e77-e91 (2019).
- 388 4 Eloe-Fadrosh, E. A. & Rasko, D. A. The human microbiome: from symbiosis to pathogenesis.
389 *Annual review of medicine* **64**, 145-163 (2013).
- 390 5 Hamady, M. & Knight, R. Microbial community profiling for human microbiome projects: tools,
391 techniques, and challenges. *Genome research* **19**, 1141-1152 (2009).
- 392 6 Scholz, M. *et al.* Strain-level microbial epidemiology and population genomics from shotgun
393 metagenomics. *Nature methods* **13**, 435 (2016).
- 394 7 Truong, D. T. *et al.* MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nature*
395 *methods* **12**, 902 (2015).
- 396 8 Kramer, M. A. Nonlinear principal component analysis using autoassociative neural networks.
397 *AIChE journal* **37**, 233-243 (1991).
- 398 9 Min, S., Lee, B. & Yoon, S. Deep learning in bioinformatics. *Briefings in bioinformatics* **18**, 851-
399 869 (2017).
- 400 10 Nguyen, T. H., Chevalleyre, Y., Prifti, E., Sokolovska, N. & Zucker, J.-D. Deep learning for
401 metagenomic data: using 2d embeddings and convolutional neural networks. *arXiv preprint*
402 *arXiv:1712.00244* (2017).
- 403 11 Nguyen, T. H., Prifti, E., Chevalleyre, Y., Sokolovska, N. & Zucker, J.-D. Disease classification in
404 metagenomics with 2d embeddings and deep learning. *arXiv preprint arXiv:1806.09046* (2018).
- 405 12 Pasolli, E., Truong, D. T., Malik, F., Waldron, L. & Segata, N. Machine learning meta-analysis of
406 large metagenomic datasets: tools and biological insights. *PLoS computational biology* **12**,
407 e1004977 (2016).
- 408 13 Cawley, G. C. & Talbot, N. L. On over-fitting in model selection and subsequent selection bias in
409 performance evaluation. *Journal of Machine Learning Research* **11**, 2079-2107 (2010).
- 410 14 Varma, S. & Simon, R. Bias in error estimation when using cross-validation for model selection.
411 *BMC bioinformatics* **7**, 91 (2006).
- 412 15 Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic sequencing.
413 *nature* **464**, 59 (2010).
- 414 16 Karlsson, F. H. *et al.* Gut metagenome in European women with normal, impaired and diabetic
415 glucose control. *Nature* **498**, 99 (2013).
- 416 17 Qin, J. *et al.* A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*
417 **490**, 55 (2012).
- 418 18 Le Chatelier, E. *et al.* Richness of human gut microbiome correlates with metabolic markers.
419 *Nature* **500**, 541 (2013).
- 420 19 Qin, N. *et al.* Alterations of the human gut microbiome in liver cirrhosis. *Nature* **513**, 59 (2014).
- 421 20 Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal cancer.
422 *Molecular systems biology* **10**, 766 (2014).
- 423 21 Glorot, X. & Bengio, Y. in *Proceedings of the thirteenth international conference on artificial*
424 *intelligence and statistics*. 249-256.
- 425 22 Kingma, D. P. & Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*
426 (2013).

- 427 23 Li, F., Qiao, H. & Zhang, B. Discriminatively boosted image clustering with fully convolutional
428 auto-encoders. *Pattern Recognition* **83**, 161-173 (2018).
- 429 24 Krizhevsky, A., Sutskever, I. & Hinton, G. E. in *Advances in neural information processing*
430 *systems*. 1097-1105.
- 431 25 Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint*
432 *arXiv:1412.6980* (2014).
- 433 26 Cortes, C. & Vapnik, V. Support-vector networks. *Machine learning* **20**, 273-297 (1995).
- 434 27 Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. *The London,*
435 *Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **2**, 559-572 (1901).
- 436 28 Bingham, E. & Mannila, H. in *Proceedings of the seventh ACM SIGKDD international conference*
437 *on Knowledge discovery and data mining*. 245-250 (ACM).
- 438 29 Dasgupta, S. Experiments with random projection. *arXiv preprint arXiv:1301.3849* (2013).
- 439 30 Dasgupta, S. & Gupta, A. An elementary proof of the Johnson-Lindenstrauss lemma.
440 *International Computer Science Institute, Technical Report* **22**, 1-5 (1999).
- 441 31 Saito, T. & Rehmsmeier, M. The precision-recall plot is more informative than the ROC plot
442 when evaluating binary classifiers on imbalanced datasets. *PloS one* **10**, e0118432 (2015).
- 443 32 Mazurowski, M. A. *et al.* Training neural network classifiers for medical decision making: The
444 effects of imbalanced datasets on classification performance. *Neural networks* **21**, 427-436
445 (2008).

446

447

448 **List of abbreviations**

- 449 - **IBD**: inflammatory bowel disease
- 450 - **EW-T2D**: type 2 diabetes in European women
- 451 - **C-T2D**: type 2 diabetes in Chinese
- 452 - **Obesity**: obesity
- 453 - **Cirrhosis**: liver cirrhosis
- 454 - **Colorectal**: colorectal cancer
- 455 - **SAE**: shallow autoencoder
- 456 - **DAE**: deep autoencoder
- 457 - **VAE**: variational autoencoder
- 458 - **CAE**: convolutional autoencoder
- 459 - **ReLU**: rectified linear unit
- 460 - **KL**: Kullback-Leibler
- 461 - **SVM**: support vector machine
- 462 - **RF**: random forest
- 463 - **MLP**: multi-layer perceptron
- 464 - **RBF**: radial basis function
- 465 - **AUC**: area under the receiver operating characteristics curve
- 466 - **PCA**: Principal Component Analysis
- 467 - **RP**: Gaussian Random Projection

468

469

470

471 **Acknowledgments**

472 This work is partially supported by the funding from Data and Decisions Destination Area at Virginia Tech.

473 Also, we thank Dr. Bert Huang for the constructive discussion.

474

475 **Author contributions**

476 MO designed the study, collected data, implemented the software, and performed experiments. MO and

477 LZ interpreted the results and wrote the manuscript. All authors read and approved the final manuscript.

478

479 **Competing interests**

480 The authors declare no competing interests.

481