

Deep learning based multi-batch calibration for classification in various omics

Qian Wang (✉ wang.qian@sjtu.edu.cn)

Shanghai Jiao Tong University

Jingyang Niu

Shanghai Jiao Tong University

Wei Xu

Shanghai Jiao Tong University

Dongming Wei

Shanghai JiaoTong University

Kun Qian

Shanghai Jiao Tong University <https://orcid.org/0000-0003-1666-1965>

Article

Keywords:

Posted Date: December 2nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1017151/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Deep learning based multi-batch calibration for classification in**
2 **various omics**

3 JingYang Niu^{1,#}, Wei Xu, DongMing Wei, Kun Qian, Qian Wang^{1,*}

4

5 **Abstract**

6 **Background:** The amount of available biological data has exploded since the
7 emergence of high-throughput technologies, which is not only revolting the way
8 we recognize molecules and diseases but also bringing novel analytical challenges
9 to bioinformatics analysis. In the last decade, deep learning has become a
10 dominant technique in data science. However, classification accuracy is plagued
11 with domain discrepancy. Notably, in the presence of multiple batches, domain
12 discrepancy typically happens between individual batches. The recently proposed
13 pair-wise adaptation approach may be suboptimal as it fails to eliminate the
14 external factors across multiple batches and takes the classification task into
15 account simultaneously.

16 **Results:** We propose a joint deep learning framework for integrating batch effect
17 removal and classification upon various omics data. To this end, we validate it on
18 two private metabolomics (MALDI MS) datasets and one public transcriptomics
19 (scRNA-seq) dataset. Especially for the former, we have achieved the highest
20 diagnostic accuracy (ACC), with notable ~10% improvement than over state-of-
21 the-art methods. Overall, these results indicate that our approach removes batch
22 effect more effectively than conventional methods and yields more accurate

23 classification results for smart diagnosis.

24

25 **Introduction**

26 Computational analysis of high-throughput omics data (*i.e.*, genomics,
27 transcriptomics, proteomics, metabolomics and radiomics) has become popular
28 over recent decades¹. Taking metabolomics as an example, it hunts for
29 quantitative descriptions of complex biological samples (usually urine or blood),
30 and associates clinical observations of diseases with temporal fluctuations of
31 metabolites². By measuring and modeling metabolism alternations in biological
32 samples, metabolomics offers relevance to fresh biological insight into diseases
33 and therapies^{2,3}. Similar landscapes also appear in the field of transcriptomics.
34 Single-cell RNA sequencing (scRNA-seq) could identify cell lineages and
35 disentangle cellular heterogeneity in complex tissues by characterizing high-
36 throughput gene expression profiles for cell types and states⁴⁻⁶, thereby further
37 revealing unexplored biological diversity with valuable pathological information.
38 Many technologies used in biology — including high throughput ones such as
39 microarrays, mass spectrometers and next generation sequencing — depend on
40 a complicated set of reagents and hardware, along with highly trained personnel,
41 to produce accurate measurements⁷. However, batch effects may occur when
42 different technicians are responsible for different subsets of the experiments or
43 diverse reagents, chips or instruments are used. Exactly, batch effects are sub-
44 groups of measurements with qualitatively different behavior across the

45 conditions unrelated to a study's biological or scientific variables⁷. In matrix
46 assisted laser desorption/ionization mass spectrometry (MALDI MS)⁸, for instance,
47 batch effects (if not adequately dealt with) could subsequently lead to serious
48 concerns about the validity of the biological conclusions^{7,9} if the serum samples
49 for a patient were repeatedly processed in different plates. Therefore, it is
50 necessary to identify and remove the batch effects before proceeding to the
51 downstream analysis.

52 It is still a challenge to machine learning¹⁰ to perform computational analysis,
53 considering many measurements (typically corresponding to the feature
54 dimensionality) and usually limited number of samples (or sample size) in high-
55 throughput omics data. The recent leap of deep learning, has outperformed many
56 conventional machine learning techniques at revealing prognostic subtypes in
57 high-risk neuroblastoma¹¹, analyzing lung adenocarcinoma prognostication¹², and
58 revealing hidden high-resolution cellular subtypes⁴. Research on batch effects is
59 no exception, and some methods based on deep learning have emerged recently.
60 The first application of deep learning in batch effect removal was the ResNet
61 method¹³. While deep network has powerful capability of approximating highly
62 nonlinear mapping, the solution of ResNet is unsupervised (i.e., without knowing
63 the class labels of samples). BERMUDA⁴ requires a clustering similarity matrix
64 before training the deep transfer learning network, thus the choice of clustering
65 method and the quality of similarity matrix owns a significant impact on the final
66 results. With the evolution of deep network architecture, there are more and more

67 breakthroughs in GANs¹⁴ in the past three years. A typical method applied to this
68 field is the NormAE¹⁵ developed by Rong et al. in 2020. Its basic idea lies in
69 constructing an adversarial training procedure between a nonlinear AE to remove
70 batch effects and a discriminator to distinguish the source of domain based on
71 the latent space. However, our original intention is to classify real biological
72 categories, which is essential for diagnostics, prognostics, and identification of
73 metabolic biomarkers. Afterwards, the DESC algorithm invented by Li et al.⁵ also
74 relies too much on the auxiliary role of clustering assignment probability.
75 Unfortunately, the classification efficacy would not be improved only by
76 decreasing the mismatching across different batches.

77 In fact, the rattling bottleneck of batch effect has drawn extensive prior research
78 in the early years, which could split into a variety of different perspectives. First of
79 all, according to the design principle, there are two conventional ways to suppress
80 batch effect, namely location-scale (LS) and matrix-factorization (MF) methods¹⁶.
81 For instance, ComBat¹⁷ is a popular LS-based approach, which employs a Bayesian
82 framework to model the data by parameterizing location and scale for each batch
83 and feature independently¹⁸. However, the assumption of normal data distribution
84 for each batch in the LS methods may be over-simplified to treat complicated
85 batch effects as additive and multiplicative components. As an alternative to the
86 LS category, surrogate variable analysis (SVA) provides an MF way to remove
87 batch effect¹⁸. The MF approaches assumes that the data variation induced by
88 batch effect is independent of the biological interpretation of interest^{16,18}. However,

89 such assumptions may not be valid in practice.

90 Apart from the above, traditional algorithms could also be divided by application
91 scenarios. For instance, in non-targeted metabolomics diagnosis, there is a need
92 to construct a discriminative model using existing source batches and apply them
93 to predict the labels of future target batches. Ratio_G¹⁹ and fSVA²⁰ adjust data for
94 enhancement of prediction performance in a predictive model, while most of the
95 conventional studies such as ComBat¹⁷ establish statistical differences at the
96 population level, which are ignorant of the subsequent classification task when
97 modeling batch effect.

98 Most existing algorithms are based on pairwise analysis in which samples from
99 two batches are considered at a time. Recently, a variety of new methods have
100 emerged in the field of single-cell RNA sequencing (scRNA-seq), including mutual
101 nearest neighbors (MNN)²¹, canonical correlation analysis (CCA)²², and Seurat 3.0²³.
102 The common point of these methods lies that they all exploit the idea of nearest
103 neighbors to identify the similar clusters between single cells across two batches,
104 thereby integrating them into a shared space. Meanwhile, we also have developed
105 an algorithm recently that specifically designed for improving classification
106 accuracy while removing batch effects. However, it is a pairwise mapping
107 approach in nature. For data with more than two batches, as typically encountered
108 in real-world settings, these approaches could only calibrate pair by pair.
109 Consequently, it is desirable to develop an algorithm that could simultaneously
110 accommodate samples from all batches at once.

111 To make the subsequent learning-based classification more convincing, we
112 propose a novel deep learning framework to integrate multi-batch calibration and
113 sample classification. We apply the proposed method to two metabolomics
114 diagnosis applications and one transcriptomics classification scenario, respectively,
115 and demonstrate its superior performance in both batch effect removal and
116 classification capability. The contributions of our paper are summarized as follows:

- 117 • We borrow the idea of traditional GAN that the discriminator and
118 reconstructor(s) are adversarial to the calibrator in different epoch steps,
119 namely “walk in two steps”, and we also make a breakthrough that the
120 discriminator directly distinguishes accurate biological labels rather than
121 domain information.
- 122 • Compared with the framework we designed before, it breaks through the
123 limitation of two batches and could calibrate multiple domains synchronously.
124 Not only does it achieve higher classification accuracy than pairwise ones, but
125 it also more convenient and time-saving.
- 126 • Our approach is not only suitable for the binary classification in biological
127 diagnosis of non-targeted metabolomics, but also extends to the multi-
128 classification of cell subtypes under scRNA-seq of transcriptomics.

129

130 **Results**

131 **1. Overview of the method**

132 In this section, we describe our approach that supports batch effect removal for

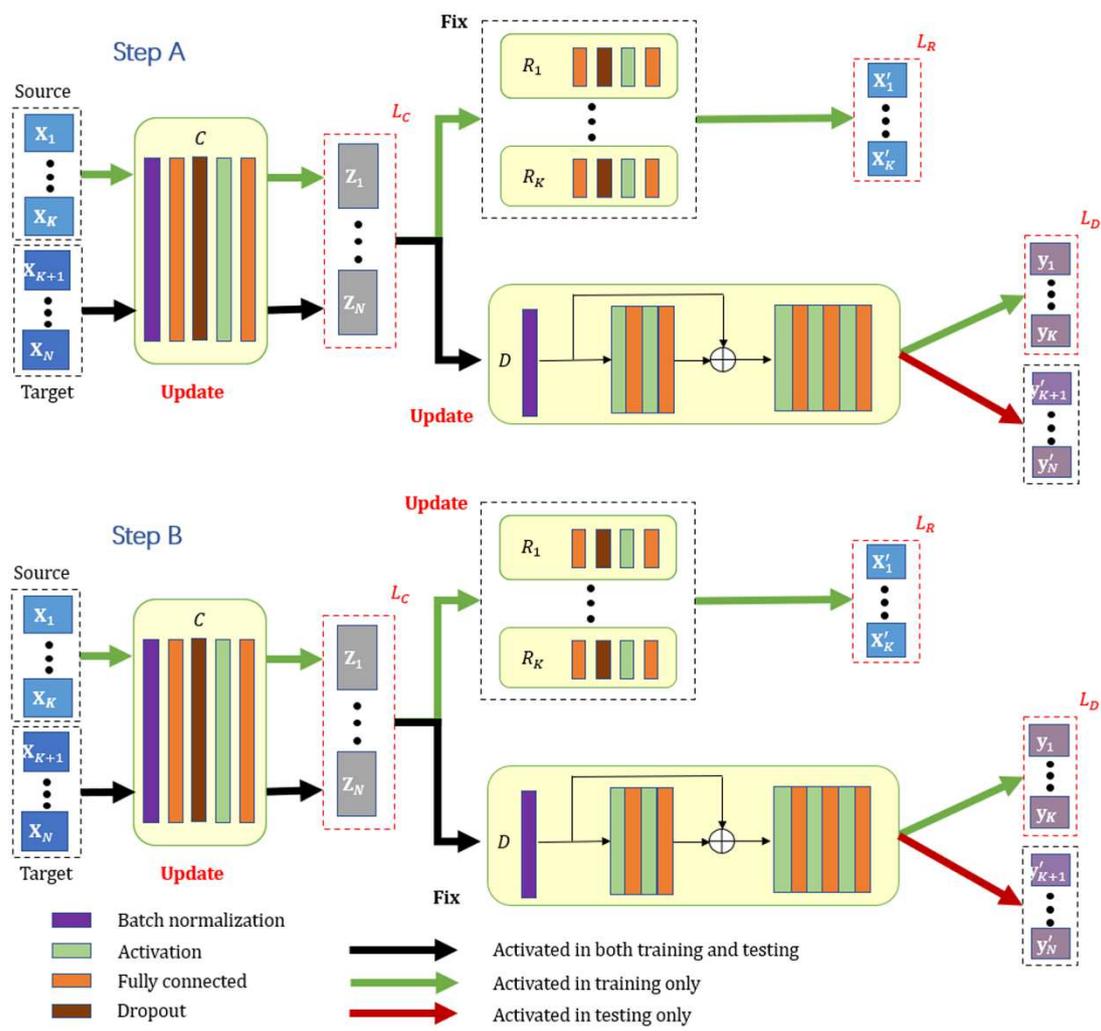
133 multiple batches simultaneously. Here we portray the framework, statistical
134 validation, and theoretical deduction of an automated batch adjustment tool
135 designed to minimize batch effects and well suited for better generalization for
136 relevant classification application.

137 Many works on batch effect removal are applicable to single-source-single-target
138 scenario only. That is, given a source batch where subjects are diagnosed already,
139 one may train a model and then apply it to the target batch. If the batch effect is
140 properly calibrated, the subjects in the target batch can be correctly classified for
141 diagnosis. However, in many real-world settings, one may seek to apply the
142 trained model to multiple (future) target batches, or boost the classification
143 performance by using multiple source batches for training.

144 Our framework is naturally suitable for such a multi-source-multi-target
145 circumstance. Suppose there are N batches before calibration,
146 $\{\mathbf{X}_1, \dots, \mathbf{X}_K, \mathbf{X}_{K+1}, \dots, \mathbf{X}_N\}$, in which $\{\mathbf{X}_1, \dots, \mathbf{X}_K\}$ represent K source batches of
147 known labels while $\{\mathbf{X}_{K+1}, \dots, \mathbf{X}_N\}$ are $(N - K)$ unlabeled target batches. The
148 application scenario of our method is to train a model from known labels of source
149 batches, such that the samples in target batches can infer their labels by
150 classification.

151 For this objective, our first goal is to find a calibrator C to accept all batches of
152 raw data and generate $\{\mathbf{Z}_1, \dots, \mathbf{Z}_K, \mathbf{Z}_{K+1}, \dots, \mathbf{Z}_N\}$, where $\mathbf{Z}_i = C(\mathbf{X}_i)$. The calibrated
153 data $\{\mathbf{Z}_1, \dots, \mathbf{Z}_K, \mathbf{Z}_{K+1}, \dots, \mathbf{Z}_N\}$ then pass through the discriminator D , which yields
154 the class label per sample. Notice that D is trained only on source batches with its

155 corresponding class labels $\mathbf{y}_1, \dots, \mathbf{y}_K$ as supervision, and then we leverage the
 156 trained D to classify the corrected samples in target batches and infer predicted
 157 labels $\mathbf{y}_{K+1}', \dots, \mathbf{y}_N'$, where $\mathbf{y}_i' = D(\mathbf{Z}_i)$. Meanwhile, to make sure that the encodes
 158 of latent space powerful and well-functioning representation, we enforce all
 159 sample data of source batches $\mathbf{Z}_1, \dots, \mathbf{Z}_K$ to be fully reconstructed from the
 160 calibrated latent space (the reason not reconstruct target batches will be
 161 confirmed in the experimental part), e.g., by passing them through K individual
 162 reconstructors R_1, \dots, R_K to obtain the reconstruction results $\mathbf{X}'_1, \dots, \mathbf{X}'_K$, namely
 163 $\mathbf{X}'_i = R_i(\mathbf{Z}_i)$. [Figure 1](#) illustrates the overall framework of our proposed method.
 164



166 **Figure 1. The architecture of our proposed framework and its adversarial training steps.**

167 The source batches $\mathbf{X}_1, \dots, \mathbf{X}_K$ and the target batches $\mathbf{X}_{K+1}, \dots, \mathbf{X}_N$ are processed through the
168 same calibrator \mathcal{C} , to ensure all batches are tightly distributed in the latent space. The source
169 batches supervise the training of discriminator D in step A, which hereafter predicts the class
170 labels for target batches in testing phase. Reconstructors (R_1, \dots, R_K) are used to recover the
171 input source batches from latent encoding in step B, which guarantees the representative
172 latent features.

173

174 In order to make the entire model to converge effectively, two adversarial training
175 steps are involved. In step A, we only renew the model weights of calibrator \mathcal{C} and
176 discriminator D without updating the model weights in reconstructor(s) R . Next,
177 we only train \mathcal{C} and R with a fixed D in step B. The iterative epochs of the two
178 steps are adjusted according to their respective convergence conditions. With the
179 help of this strategy, the resulted model allows simultaneous alignment multiple
180 batches, from sources to targets.

181

182 **2. Evaluation metrics**

183 This study proposes a joint deep learning framework to perform multi-batch
184 calibration in multi-category applications. In order to prove in turn that the
185 framework is suitable for multi-source single-target, multi-source multi-target
186 and other omics multi-classification problems, we report experimental results
187 using three datasets from two high-throughput technologies, *i.e.*, two private
188 MALDI MS datasets and one public scRNA-seq dataset. We verify our framework
189 by means of comparing it to several most representative algorithms in the

190 literature. Detailed assessment will be reported below.

191 Our evaluation principally focuses on removing batch effect and classification
192 performance. Particularly, for batch effect removal, we adopt MMD as a
193 quantitative metric and present their results in the form of boxplot. We also utilize
194 t-SNE^{24,25} to visualize the distribution of the high-dimensional data before and
195 after calibration by each method. Furthermore, we apply four metrics to assess
196 the classification performance, consisting of Accuracy (ACC), F-score, Area Under
197 Curve (AUC), and Matthews correlation coefficient (MCC)^{19,26}. The ACC, F-score,
198 MCC are defined below:

$$ACC = \frac{TP + TN}{TP + FN + FP + TN} \quad (5)$$

$$F\text{-score} = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (6)$$

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}} \quad (7)$$

199 where TP, FP denote the true/false positives and TN, FN denote the true/false
200 negatives.

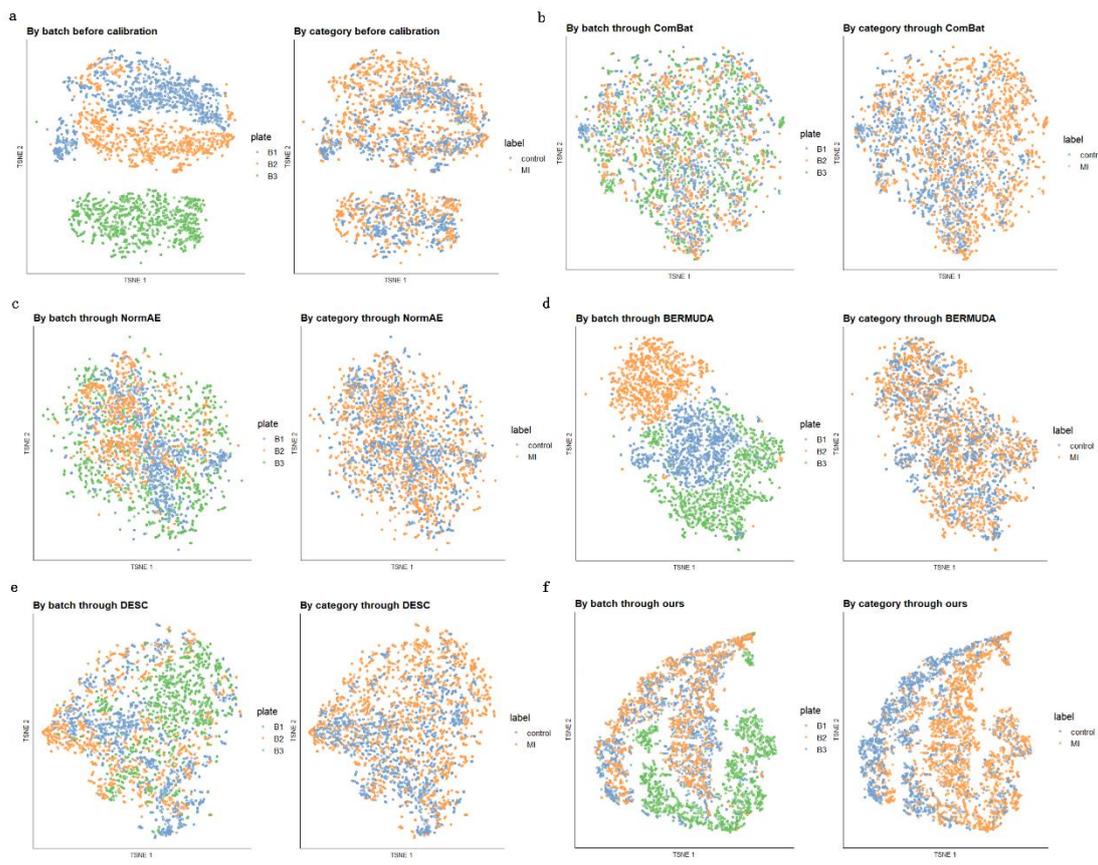
201

202 **3. multi-source-single-target situation**

203 **Dataset.** In order to prove that our approach has obvious utility on multi-source-
204 single-target data, a total of 796 individuals were recruited, including 322 controls
205 and 474 acute myocardial infarction (MI) patients from Shanghai Chest Hospital
206 Affiliated to Shanghai Jiao Tong University. All MI patients had specific
207 histopathological diagnoses (99th cardiac troponin I (cTnI) levels and
208 electrocardiogram (ECG)). The 322 individuals in the control group had no clinical

209 evidence of cardiovascular diseases such as coronary heart disease (CHD) but may
210 suffer from other diseases. There is no significant difference in age and sex
211 between controls and MI patients. All of the investigation protocols in this study
212 were approved by the institutional ethics committees of the Shanghai Chest
213 Hospital and School of Biomedical Engineering, Shanghai Jiao Tong University.
214 According to the Helsinki Declaration, written informed consent to participate in
215 the study were provided from all individuals, and the use of their biological
216 samples was approved for analysis.

217 **Prove the existence of batch effect.** In order to prove that these data are indeed
218 disturbed by batch effect, we first explore the number of differential features in
219 each batch, under which there are significant differences in the expression levels
220 of the case and control groups (e.g., patients and healthy people). We quest the
221 differential peaks for each batch's case and control group according to BPSC²⁷, a
222 differential expression gene probe algorithm based on Beta-Poisson model.
223 Figure S1a shows that the three batches have only 82 common differential
224 features (about 200 difference peaks are extracted in each batch and take the
225 intersection of three batches) under the condition that the p-value is set to 0.05.
226 It could be concluded that this dataset exists serious batch effects because the
227 differential peaks that separate the case and control group in three batches are
228 distinctive.
229



230

231

Figure 2. t-SNE Visualization of three batches on the private MI data of MALDI MS. a raw data, **b** ComBat, **c** NormAE, **d** BERMUDA, **e** DESC, **f** our calibration. In the left half of each section, different colors highlight the three batches. In the right half, different colors identify the actual labels of samples (Class 1: MI; Class 0: control).

234

235

236

We further visualize the data distribution before and after batch effect calibration

237

by each algorithm. [Figure 2](#) shows the t-SNE plot of the raw inputs and every

238

calibrated data for 3980 MI samples. Instead of distributing with respect to the

239

batch factor, the same classification label samples should be close to each other

240

in the feature space ideally. However, as in the left part of [Figure 2a](#), one may

241

observe that the three batches (plates, in different colors) are distributed in various

242

patterns, suggesting an apparent batch effect that separates them and may hinder

243

downstream classification (c.f. the right half, colored in accordance to class labels).

244

On the contrary, after being calibrated as in the left part of [Figure 2b~2f](#) by each

245 method, the three batches share thoroughly entangled distributions except
246 BERMUDA ([Figure 2d](#)), implying most of their in-between mismatch due to batch
247 effect has been removed. Although ours does not seem to be as uniform as
248 ComBat ([Figure 2b](#)), NormAE ([Figure 2c](#)) and DESC ([Figure 2e](#)) in terms of batch
249 mixing, however, as in the right half of [Figure 2f](#), the samples are naturally more
250 compact in accordance to their labels than other approaches. This is partially due
251 to the discriminator in the two-step strategy, which helps span the feature space
252 to remove batch effect and facilitate the classification of labels.

253 **Quantitative comparison.** The effectiveness of our framework can further be
254 verified by the classification performance quantitatively. By testing multiple
255 aspects of performance, rather than relying on a single measure, the cross-batch
256 predictions contain $\binom{3}{1}$ classification tasks (choosing 2 out of 3 batches as the
257 training set, and the other batch as the test one) and each batch takes turns once,
258 to test the robustness of our approach. As shown in [Table 1a](#), the quantitative
259 results show that our algorithm can improve in every group, which has increased
260 by as much as 10 percentage points. Especially, the batch effects are strong when
261 batch (1,2) as source and batch 3 as target, the improvement of ACC, F-score,
262 AUC, MCC is as high as 17.2%, 8.2%, 22.1%, 40.2% after applying our method for
263 both endpoints, respectively. Moreover, we also make comparisons with the prior
264 pairwise calibration algorithm developed by ourselves, as shown in the middle
265 column of [Table 1a](#), For each plate, we first use the other two plates to calibrate
266 respectively and then get average. It can be concluded that calibrating multiple

267 batches simultaneously yields the dual advantages of accuracy and calculation
 268 time.

269

270 **Table 1. Classification results on the MI data of MALDI MS**

a

	Before Calibration			Pairwise Calibration			Joint Calibration		
	1	2	3	1	2	3	1	2	3
ACC	0.721	0.686	0.599	0.826	0.777	0.754	0.876	0.807	0.771
F-score	0.708	0.693	0.743	0.856	0.828	0.806	0.884	0.853	0.825
AUC	0.745	0.725	0.519	0.814	0.749	0.729	0.886	0.776	0.740
MCC	0.498	0.446	0.128	0.642	0.515	0.487	0.763	0.579	0.530

b

Source	Target	Baseline	ComBat	NormAE	BERMUDA	DESC	Recon_T	Ours	CrossValid
1, 2	3	0.664	0.779	0.756	0.686	0.775	0.746	0.774	0.937
1, 3	2	0.726	0.774	0.708	0.716	0.804	0.751	0.807	0.919
2, 3	1	0.723	0.795	0.724	0.738	0.756	0.879	0.867	0.948
Average		0.704	0.783	0.729	0.713	0.778	0.792	0.816	0.935

271

272 **a** The result of four indicators evaluating the cross-batch prediction. The middle column is
 273 the result of averaging pairwise calibration, and the right column is the multi-batch calibration
 274 developed in this text. **b** Comparison of classification accuracy with multiple source batches
 275 for training and only one target batch for testing. Note that “Recon_T” denotes an ablation
 276 experiment that reconstruct all target batches.

277

278 Afterwards, we select several latest and most representative tools including
 279 ComBat¹⁷, NormAE¹⁵, BERMUDA⁴, and DESC⁵ for further comparison. Accuracy of
 280 cross-batch prediction in the sample level is utilized to assess the effectiveness of
 281 each method. The comparing results are reported in [Table 1b](#). Our performance
 282 is superior over all other approaches, accompanied by an improvement of
 283 3.3~10.3% than others on average. Compared with NormAE, for instance, a novel
 284 algorithm tailored for specific metabolomics data types, our method has also

285 achieved 8.7% improvement (72.9% vs. 81.6% overall) taking advantage of label
286 supervision from the source data. From the results in the third-to-last column, we
287 conclude that the accuracy of reconstructing all batches (including source and
288 target batches) is slightly lower than that for only source batches.

289 The classification performance could even be comparable to the situation when
290 batch effect is theoretically ruled out. Notably, we compute the in-batch
291 classification accuracy by conducting 10-fold cross-validation within every batch.
292 These results are regarded as a reference to the classification performance without
293 being interfered by batch effect, which are listed in the last column of [Table 1b](#)
294 that represent corresponding results within the target batch. We observe that our
295 method produces the results that get much closer to the ceilings where batch
296 effect is completely ruled out.

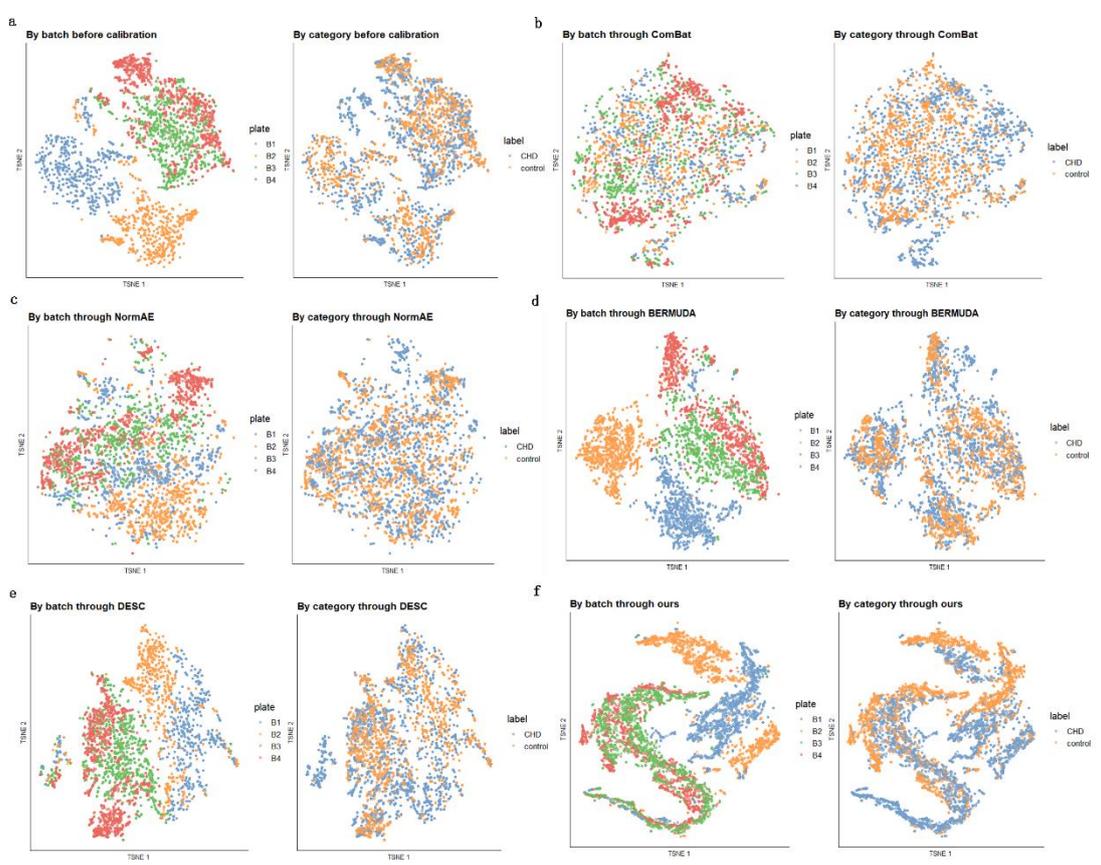
297

298 **4. multi-source multi-target situation**

299 **Dataset.** For further certifying that our algorithm is equally applicable to multi-
300 source and multi-target data, 1203 individuals were recruited in the same way as
301 previous, consisting of 562 healthy controls and 641 coronary heart disease (CHD)
302 patients from the Shanghai Chest Hospital Affiliated to Shanghai Jiao Tong
303 University. For controls, serum samples from 562 healthy volunteers who had no
304 clinical evidence of cardiovascular diseases and other major disease were
305 collected as controls. The same preprocessing pipeline as the MI's was utilized in
306 this experiment. For four batches of this dataset totally, we collected 1559, 1478,

307 1523, 1449 samples for 312, 296, 305, 290 subjects, respectively.

308 **Prove the existence of batch effect.** In this experiment, we prove the existence
309 of batch effect by computing the MMD values between source and target batches
310 before and after calibration by each algorithm. We firstly compute from a subset
311 of 500 samples randomly drawn from all samples available in one batch and then
312 take over 30 permutation runs shown by the form of boxplot. Since it could only
313 be calculated in pairwise, we randomly select batch 3 and batch 4 as an example
314 and show in Figure S1b. Apparently, our calibration decreases the MMD value
315 between the two batches (i.e., 0.152 ± 0.007 after being processed by our method,
316 which is lower than raw data and the results calibrated by other methods). The
317 above results prove that our algorithm can suppress batch effect effectively.
318



320 **Figure 3. t-SNE Visualization of four batches on the private CHD data of MALDI MS. a**
321 raw data, **b** ComBat, **c** NormAE, **d** BERMUDA, **e** DESC, **f** our calibration. The left part of each
322 method is colored by batch indices. In the right half, the samples are colored by disease labels.

323

324 We also visualize the distribution by projecting the raw and calibrated data by
325 each method to the same space by t-SNE, as shown in [Figure 3](#). The four batches
326 are mixed together (in the left part of [Figure 3f](#)) compared with the case before
327 calibration (the left part in [Figure 3a](#)). In this regard, ours is better than BERMUDA
328 ([Figure 3d](#)) and comparable to DESC ([Figure 3e](#)), but inferior to ComBat ([Figure](#)
329 [3b](#)) and NormAE ([Figure 3c](#)). On the other hand, if colored by the class labels (CHD
330 vs. control), one may observe that the samples are almost inseparable before
331 calibration (in the right half of [Figure 3a](#)), similar scenarios also occur after
332 calibration of other methods ([Figure 3b~3e](#)), yet much more separable through
333 our calibration (the right half in [Figure 3f](#)). This is exactly what we expect. Such
334 results present different tunes but are rendered with equal skill as the experiments
335 on previous MI data.

336 **Quantitative comparison.** We then evaluate the classification performance
337 quantitatively. Given two source batches for training and another two target
338 batches for test, we take the average of different training sets under three cases
339 when each batch as the test set to facilitate intuitive comparison, as shown in
340 [Table 2a](#). Specifically, we also display the comparisons of pairwise calibration in
341 the middle column, which prove that our framework of joint multi-batch
342 calibration is superior to previous pairwise ones in most cases. The original
343 classification results before averaging can be found in Table S1. For example,

344 when batch (1,2) are used as the training set and batch (3,4) for test, the metrics
 345 of ACC, F-score, AUC, and MCC have increased by (13.0%, 19.4%), (3.7%, 6.1%),
 346 (16.0%, 19.3%) and (25.6%, 29.9%), respectively.

347

348 **Table 2. Classification results on the CHD data of MALDI MS**

a

Target	Before Calibration				Pairwise Calibration				Joint Calibration			
	1	2	3	4	1	2	3	4	1	2	3	4
ACC	0.732	0.756	0.715	0.639	0.783	0.721	0.715	0.731	0.825	0.815	0.775	0.746
F-score	0.679	0.809	0.768	0.720	0.779	0.755	0.722	0.756	0.820	0.838	0.799	0.752
AUC	0.744	0.735	0.702	0.640	0.787	0.710	0.716	0.732	0.829	0.808	0.770	0.746
MCC	0.527	0.526	0.428	0.334	0.581	0.467	0.439	0.474	0.664	0.631	0.552	0.495

b

Source	Target	Baseline		ComBat		NormAE		BERMUDA		DESC		Recon_T		Ours	
1, 2	3, 4	0.609	0.559	0.678	0.645	0.521	0.511	0.638	0.665	0.716	0.696	0.753	0.734	0.751	0.726
1, 3	2, 4	0.775	0.679	0.782	0.759	0.662	0.661	0.695	0.691	0.730	0.675	0.683	0.714	0.748	0.718
1, 4	2, 3	0.715	0.794	0.781	0.787	0.667	0.763	0.743	0.658	0.695	0.792	0.802	0.762	0.828	0.804
2, 3	1, 4	0.709	0.659	0.818	0.793	0.739	0.665	0.702	0.700	0.696	0.744	0.719	0.728	0.769	0.758
2, 4	1, 3	0.778	0.742	0.771	0.789	0.654	0.687	0.761	0.651	0.764	0.801	0.822	0.799	0.865	0.812
3, 4	1, 2	0.649	0.763	0.588	0.578	0.577	0.558	0.599	0.571	0.663	0.752	0.776	0.778	0.822	0.831
Average		0.703		0.731		0.639		0.673		0.727		0.756		0.786	

349

350 **a** The result of four indicators evaluating the cross-batch prediction about before calibration,
 351 averaging pairwise calibration, and after multi-batch calibration. **b** Comparison of
 352 classification accuracy with multiple source batches for training and multiple target batches
 353 for testing. The implication of "Recon_T" is the same as [Table 1b](#).

354

355 Next, we still evaluate the methods involved in the previous section. Since the
 356 source and target batch are paired, there exists $\binom{4}{2}$ cross-batch prediction tasks.

357 A complete list of the approaches analyzed in this section is provided in [Table 2b](#).

358 The performance of our method is optimal on all experimental groups. Particularly,
 359 the last combination, namely batch (3,4) as source and batch (1,2) as target, yields
 360 the most improvement for our method (15.9%, 7.9%) compared to the second-
 361 ranking approach (DESC). Furthermore, from an ablation experiment in the

362 penultimate column, the results of reconstructing all batches are inferior to that
363 only reconstructing source batches. These results in overall indicate that our
364 approach not only be suitable for multi-source single-target situations, but also
365 performs well on multi-source multi-target circumstances.

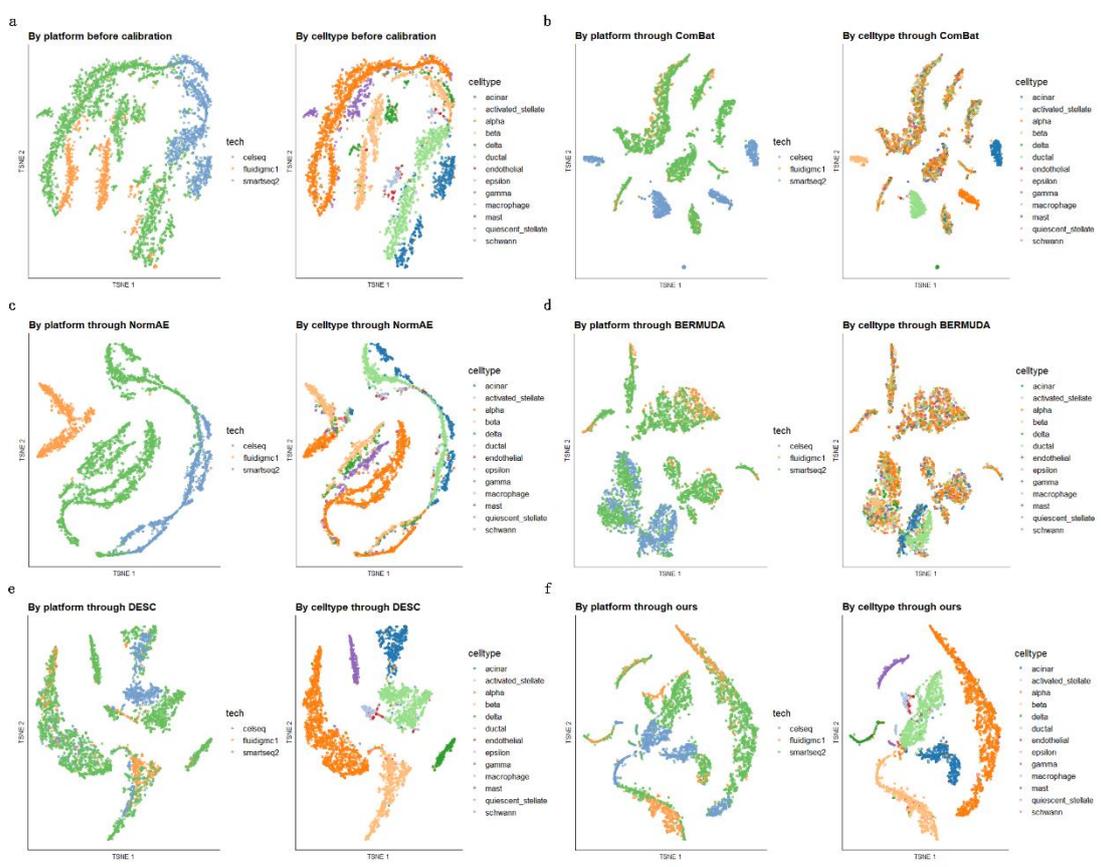
366

367 **5. scRNA-seq:**

368 **Dataset.** In order to prove that our method could also achieve better performance
369 in the multi-classification application of other omics besides metabolomics, we
370 focused on the human pancreatic data through several scRNA-seq protocols. We
371 have combined three publicly available datasets generated using CelSeq from
372 Gene Expression Omnibus (GSE81076)²⁸, Fluidigm C1 (GSE86469)²⁹ and SMART-
373 Seq2 from the European Bioinformatics Institute (E-MTAB-5061)³⁰ and each
374 platform represents one batch. The pancreas is a highly heterogeneous tissue with
375 several well-defined cell types. 13 cell type labels have been taken either from the
376 provided metadata or been derived according to the methodology described in
377 the original publication, annotated as “acinar”, “activated stellate”, “alpha”, “beta”,
378 “delta”, “ductal”, “endothelial”, “epsilon”, “gamma”, “macrophage”, “mast”,
379 “quiescent stellate”, “schwann” (further details of data preprocessing as follows).

380 **Prove the existence of batch effect.** We also visually compare this data set and
381 the t-SNE results are shown in [Figure 4](#). Notice that in the real-world data, the
382 multiple training and test instances often do not come from the same underlying
383 distribution, but we have projected them onto the same subspace. In the left part

384 of [Figure 4a](#), the underlying cause of such clustering is the sequencing platform
 385 and this is clear because we plot the points marking each point according to
 386 protocol numbers. The left part of [Figure 4b~4f](#) presents results after calibration
 387 by each method, where we can observe that the differences are much smaller
 388 after our calibration than ComBat ([Figure 4b](#)), NormAE ([Figure 4c](#)) and
 389 comparable to BERMUDA ([Figure 4d](#)), DESC ([Figure 4e](#)). In the right half of [Figure](#)
 390 [4f](#), the points cluster belonging to the same cell-types appears more compact
 391 than before calibration and by most algorithms (the right half in [Figure 4a~d](#)), and
 392 equal to DESC on this dataset. It implies that the calibration process makes the
 393 actual category marker structure significantly more similar.
 394



395
 396 **Figure 4. t-SNE visualization of the public pancreas data of scRNA-seq. a** raw data, **b**
 397 ComBat, **c** NormAE, **d** BERMUDA, **e** DESC, **f** our calibration. In the left part of each section,

398 different colors highlight the batches (sequencing platforms). In the right half, different colors
 399 identify 13 actual labels of cell-types.

400

401 **Quantitative comparison.** In this section, the accuracy of cross-batch prediction
 402 is still used to evaluate the effectiveness of those algorithms mentioned above.
 403 There are $\binom{3}{1}$ classification tasks and the comparing results are reported in [Table](#)
 404 [3](#). Due to a multi-classification task (13 categories), it is much more difficult to
 405 improve accuracy than pair-wise classification. Nevertheless, our approach still
 406 wins out of the others. Take the second group (batch (1,3) as source and batch 2
 407 as target) as an example, ours outperforms the second-ranking algorithm
 408 (NormAE) 24.3 percentage points. In addition, we also make comparisons with the
 409 previous pairwise calibration algorithm developed by ourselves, as shown in the
 410 "Pairwise" column of [Table 3](#). For each platform, we first use the other two
 411 platforms to calibrate separately and then average. These results overall indicate
 412 that our approach not only be more suitable for binary classification diagnosis of
 413 metabolomics than other methods, but also suitable for multi-classification
 414 scenarios of other omics.

415

416 **Table 3. Classification results on the pancreas data of scRNA-seq**

Source	Target	Baseline	ComBat	NormAE	BERMUDA	DESC	Pairwise	Recon_T	Ours
1、2	3	0.162	0.147	0.432	0.187	0.135	0.233	0.321	0.428
1、3	2	0.226	0.544	0.215	0.128	0.208	0.183	0.370	0.458
2、3	1	0.001	0.161	0.221	0.118	0.163	0.165	0.050	0.227
Average		0.130	0.284	0.289	0.144	0.169	0.194	0.247	0.371

417

418 Comparison of classification accuracy with multiple categories on the scRNA-seq of
419 transcriptomics. Note that "Pairwise" represents the two-batch calibration method designed
420 by ourselves.

421

422 **Discussion**

423 In recent decades, high throughput omics analysis technology has become much
424 more mature. However, batch effects are ubiquitous in high-throughput
425 experiments (e.g., RNA sequencing, metabolomics), the source of which is far-
426 ranging, including different platforms, different periods of the same platform,
427 different reagents and times of the same sample, etc. To this end, we introduce
428 an end-to-end deep learning framework and demonstrate that our proposed
429 method can effectively remove batch effect by the experiments on two
430 metabolomics and one transcriptomics datasets. Our framework outperforms all
431 compared methods in terms of classification accuracy through the experiments.
432 This is partially attributed to multiple modules that are interacted with each other
433 in the architecture. In addition, the "two-step" adversarial strategy we adopted
434 during training also dramatically facilitates the compromise between the
435 classification and reconstruction modules, something other algorithms fails to
436 account for. The mutual restraint of these modules is substantially improving the
437 overall performance of our network.

438 A considerable number of computational methods in genomics and
439 transcriptomics have been developed to remove batch effects. However, they
440 might be less effective in improving the accuracy of classification based on
441 different omics, because it is not easy to generalize their findings to the sample

442 level for the sake of individualized biological diagnosis and treatment³¹. Although
443 NormAE is based on metabolomics, it is not superior to others in MALDI MS
444 diagnosis. A possible reason lies that it doesn't take advantage of the supervisory
445 information provided by the source labels. DESC takes subsequent classification
446 tasks into account, so it can yield better performance than NormAE or BERMUDA.
447 Although DESC is an adversarial model and the application scenario also involves
448 classification, and BERMUDA utilizes the deep transfer learning network, both
449 depend too much on the clustering similarity matrix. Overall, those methods are
450 intrinsically driven by data size/scale that is susceptible to bias introduced by
451 batch effects and cannot effectively address these deviations especially
452 concerning the need for classification.

453 There are some deficiencies in this work that cannot be underestimated. In the
454 cases where exist multiple cell types in transcriptome, the improvement of some
455 batches is not apparent. Generally, assumptions made by learning algorithms that
456 identical cell types in different batches are often violated, resulting in degradation
457 of the algorithms' performance during inference of test data. We might as well
458 calculate the distribution distance of test set with those labels-known source data
459 and choose the closest ones as training set under this circumstance from now on.
460 Furthermore, conclusions reached in this study are based on the application of
461 batch effect removal in the context of cross-batch prediction. Beyond the fact that
462 this approach is no longer unsupervised and requires domain knowledges, the
463 amount of labeled data that might be needed to achieve reasonable performance

464 could be large. These are still outstanding questions and need further
465 investigation.

466

467 **Conclusion**

468 The rise of omics techniques has resulted in an explosion of high throughput data
469 in modern biomedical research. However, these analytical barriers are further
470 compounded, since the bias introduced by the non-biological nature of the batch
471 effect can be strong enough to mask, or confound actual biological differences.
472 Therefore, it is necessary to develop a novel tool, which will be time efficient,
473 incorporate flexibility for data types, investigator-driven batch adjustment
474 approach choices, and the ability to evaluate such adjustment approaches, with
475 great application potential for the analysis of large samples in clinical studies.

476 We have introduced a novel end-to-end learning framework for simultaneous
477 multi-batch calibration and classification, and we conduct adversarial training by
478 “walk-in two steps” strategy. Upon the private MALDI MS and public scRNA-seq
479 datasets, we confirm that our framework can effectively suppress batch effect and
480 accomplish classification, outperforming the second-ranking algorithm by a
481 substantial margin on many groups. We have witnessed the applications to
482 transcriptomics and metabolomics datasets here and released publicly available
483 codes on the GitHub. Furthermore, in the same way that general deep learning
484 techniques, operating on raw data outperform traditional algorithms tailored for
485 specific data types, involving domain knowledge, or massive pre-processing, we

486 demonstrate that our proposed algorithm and experimental results really
487 promising and may open new horizon for removing batch effects in biological
488 datasets.

489

490 **Methods**

491 **1. Loss Function**

492 **Calibration loss MMD.** Our calibrator \mathcal{C} is responsible for lessening the
493 discrepancy between source and target batches by matching them into a
494 common space. It comprises a batch normalization (BN), a fully connected (FC), a
495 dropout layer, a Tanhshrink activation and again an FC layers. The number of
496 nodes in each hidden layer remains the same as the feature dimensionality of
497 inputs. In addition, it is critical to define a measurement of divergence among
498 source and target batches' distributions. Therefore, we train the calibrator \mathcal{C} by
499 minimizing the maximum mean discrepancy (MMD) in the latent space:

$$500 \quad \mathcal{L}_{\mathcal{C}} = \binom{N}{2}^{-1} \sum_{i=1}^{N-1} \sum_{j=i+1}^N \|C(\mathbf{X}_i) - C(\mathbf{X}_j)\|_1, \quad (1)$$

501 where \mathbf{X}_i and \mathbf{X}_j indicate the i -th source batch and the j -th target batch
502 respectively, and $\|\cdot\|_1$ is the L1-norm operator. MMD value approaches zero
503 when the underlying distributions of the observed samples are highly similar.

504 **Reconstruction loss MSE.** In order to prevent losing the intrinsic biological states
505 encoded in the latent code \mathbf{Z} , we require $R(\mathbf{Z})$ to transform the encodes \mathbf{Z} back
506 to the reconstructed \mathbf{X}' that reflects the characteristics of the original data. All of
507 reconstructors have the same structure that consists of first and last two FC layers,

508 a dropout layer and a Tanhshrink activation. The reconstruction loss was
 509 calculated as the residual in L2-norm between the input prior to calibration and
 510 the reconstructed output of the encoder-decoder backbone for self-learning:

$$511 \quad \mathcal{L}_R = \frac{1}{K} \sum_{i=1}^K \|\mathbf{X}_i - R_i(C(\mathbf{X}_i))\|_2^2. \quad (2)$$

512 Note that K denotes the number of source batches and each reconstructor
 513 corresponds to a certain batch.

514 **Classification loss CE.** In addition, we introduce D to be a task-specific class label
 515 classifier. It should be noted that we distinguish the category labels for the source
 516 batches in the training stage. The class label information of the target batches is
 517 only used when evaluating the prediction performance.

518 The internal structure of discriminator D has similar components but different
 519 arrangement with C or R . Specifically, it involves three partitions: the first one only
 520 contains a BN layer; the second partition is a residual structure, including two
 521 alternating Leaky ReLU activation layers and FC whose number of nodes switching
 522 between input scale and 100; the last one contains three alternating Leaky ReLU
 523 activation layers and FC, in which the node reduced from input scale to 50, 16,
 524 and finally to the number of categories. When the application scenario is only
 525 binary classification considered, the last layer utilizes the Sigmoid activation
 526 function. We calculate the loss as cross-entropy (CE) between its output and
 527 sample biological labels:

$$528 \quad \mathcal{L}_D = -\frac{1}{K} \sum_{i=1}^K \sum_j \sum_{c=1}^M \mathbf{y}_{jc} \log D(C(\mathbf{x}_{jc})). \quad (3)$$

529 where K and M represent the number of batches or categories respectively, and

530 y_{jc} can be assigned 1 when the true category of sample j is equal to c , otherwise
531 0.

532 **2. Training Process**

533 In order to control the iterative epochs in training the discriminator and
534 reconstructor(s) more conveniently, we adopt a "two-step" strategy in the training
535 process. We set two hyperparameters λ_1, λ_2 to control the step A and step B
536 respectively, as shown in formula (4). Note that the parameters of calibrator C are
537 always kept updating. The benefits of doing so lies that when the reconstruction
538 converges slowly, we increase the number of rounds of step A, and on the
539 contrary, increase the number of iterations of step B. These two steps are repeated
540 back and forth until reaching convergence of the model. The total loss of our
541 framework is then calculated by considering the calibrator, reconstructor, and
542 discriminator as a whole, i.e., to minimize

$$543 \quad \mathcal{L} = \lambda_1 \cdot (\alpha_1 \cdot \mathcal{L}_C + \beta \cdot \mathcal{L}_D) + \lambda_2 \cdot (\alpha_2 \cdot \mathcal{L}_C + \gamma \cdot \mathcal{L}_R), \quad (4)$$

544 where λ_1, λ_2 denote the number of reincarnations aiming at step A and step B,
545 and $\alpha_1, \beta, \alpha_2, \gamma$ are scalar weights for each component network in two steps,
546 respectively.

547 We implement the proposed scheme with PyTorch (version 1.8.1+cu102) and
548 Sklearn (version 0.21.3) framework. The downstream analysis has been carried out
549 using Python (version 3.6.8), and R (version 4.0.4) for visualization. For details, we
550 use ADAM for training with default settings (i.e., the exponential decay rate of the
551 first/second moment estimation). To achieve fair comparison, the types and

552 numbers of classification network layers remain the same before and after
553 calibration. The gradient update rule of mini-batch in deep learning is used to
554 train our model where corresponding losses are calculated from a sampled “mini-
555 batch” during each iteration of the corresponding step. All the experiments are
556 run on the same host with 16GB memory and an Nvidia RTX 2080Ti GPU.

557 **3. Information for Other Methods**

558 Corresponding open source code could be found about those algorithms
559 involved in comparative experiments. The ComBat has been implemented by
560 ComBat() function into R software package sva
561 (<http://bioconductor.org/packages/3.5/bioc/html/sva.html>). The source codes of
562 NormAE algorithm are publicly available at <https://github.com/luyiyun/NormAE>.
563 Since our data based MALDI MS or scRNA-seq instead of LC MS in NormAE, which
564 not exist so-called injection order, therefore, it is eliminated in training and testing
565 process. In addition, the mass quality control was conducted using standard
566 molecules on the stage of serum plates, so it doesn't appear at the preprocessing
567 matrix. The source codes of BERMUDA are publicly available at
568 <https://github.com/txWang/BERMUDA>. In principle, we compute the similarity
569 matrix based on the MetaNeighbor algorithm by following recommended
570 protocols in its reports. An open-source implementation of the DESC can be
571 downloaded from <https://eleozzr.github.io/desc/>. We generally yield a .h5ad file
572 through preparing an AnnData object based on the processing pipeline for
573 following analysis.

574 **4. MI Dataset**

575 The collection process for the MALDI MS data was introduced by Huang, et al.³².
576 According to the protocol, ~2 mL of blood was collected by venipuncture and
577 centrifuged at 3000 rpm for 10 min. Then, the serum was transferred to a
578 microtube and stored at -80°C. For each subject, we repeated detection five times
579 to enhance the reproducibility and stability. The samples with a relative standard
580 deviation (RSD) less than 5% were discarded in subsequent quality control.
581 Consequently, each subject would eventually yield 1-5 (mostly 5) samples. While
582 each sample could have a predicted label in test, the diagnosis result should be
583 ensembled to the subject level as the median of all samples in one subject. In total,
584 for the three batches, we collected 1330, 1305, 1345 samples for 266, 261, 269
585 subjects, respectively. Note that the following results were all reported at the
586 sample level. All data was processed through smoothing filter, baseline correction,
587 peak extraction and peak alignment³². As for each sample, the m/z range was set
588 from 100 to 300³³, and about 200 features were extracted through data
589 preprocessing.

590 In the experiment of classification before calibration, the learning rate, number of
591 epochs, size of mini-batch and learning step are set to 10^{-4} , 100, 128 and 10^4
592 respectively. In order to prevent overfitting the network, the L2 weight decay is
593 set to 5×10^{-5} during training. For training our network after calibration, the
594 learning rate is set to 10^{-3} and the number of epochs is set to 100. We set the
595 mini-batch size to 128 and the coefficients of losses are $\alpha_1 = 10$, $\beta = 10$ in step

596 A for $\lambda_1 = 1$ and $\alpha_2 = 10$, $\gamma = 0.01$ in step B for $\lambda_2 = 1$ by grid-search. The
597 learning step, L2 weight decay are set to 10^4 , 5×10^{-5} during training. In order
598 to compare all the results on the same benchmark, the classification network of
599 in-batch 10-fold cross-validation shares the same framework as cross-batch
600 experiments. The learning rate, number of epochs, size of mini-batch, learning
601 step and L2 weight decay are set to 10^{-3} , 10, 128, and $100, 5 \times 10^{-5}$, respectively.
602 For most of other methods, including ComBat, BERMUDA and DESC, we evaluate
603 their performance based on the default parameters and by following
604 recommended pipeline in their tutorials. As for NormAE, in order to ensure the
605 convergence of the model, except for the (lr_rec, lr_disc_b), epoch and batch_size
606 which are set to (0.0002, 0.0001), (200, 100, 150) and 200, other parameters are
607 defaulted.

608 **5. CHD Dataset**

609 In this set of experiment, the structure of classification network is also consistent
610 with the previous rule, namely keeping the same before and after calibration. The
611 learning rate, number of epochs, size of mini-batch, learning step and L2 weight
612 decay are set to 10^{-4} , 50, 128, 10^4 and 5×10^{-5} , respectively. For parameters
613 after calibration, the learning rate, number of epochs, size of mini-batch and
614 learning step are set to 10^{-4} , 50, 128 and 10^4 , respectively. In addition, the hyper-
615 parameters of the coefficients of losses are set to $\alpha_1 = 10$, $\beta = 10$ in step A for
616 $\lambda_1 = 4$ and $\alpha_2 = 10$, $\gamma = 0.01$ in step B for $\lambda_2 = 1$ by grid-search. The L2 weight
617 decay is set to 5×10^{-5} during training, the same with first experiment. We

618 evaluate most of other methods except NormAE based on the default parameters
619 and tutorials as prior. As for NormAE, the (lr_rec, lr_disc_b), epoch and batch_size
620 are set to (0.0002, 0.0001), (100, 50, 100), 200, and other parameters are defaulted.

621 **6. scRNA-seq Dataset**

622 The data batches were introduced in the form of Seurat²³ R objects featuring
623 standardized annotations, which relies on anchor cells between pairs of datasets.
624 However, misidentification of anchors from different batches might have led to
625 reduced classification accuracy. Therefore, we need segment objects according to
626 different protocols and transform them into Single Cell Experiment (SCE) objects.
627 Metadata and counts were extracted from the SCE R objects and used for
628 performing standard preprocessing. Next, we explored common genes by “scran”
629 R package and only genes that were detected in all three experiments were kept.
630 Through a series of procedures such as filtering low-quality cells, standardizing,
631 and selecting the most informative genes by calculating the degree of variation
632 for each gene, the resulting dataset consists of three batches for a total of 4036
633 cells with 325 genes each.

634 For classification before calibration, the learning rate, number of epochs, size of
635 mini-batch, learning step and L2 weight decay are set to 10^{-5} , 50, 128, 10^4 and
636 5×10^{-5} , respectively. For framework after calibration, the learning rate, number
637 of epochs, size of mini-batch, learning step and L2 weight decay are set to
638 5×10^{-6} , 55, 128, 10^4 and 5×10^{-5} , respectively. Moreover, the hyper-
639 parameters of the coefficients of losses are set to $\alpha_1 = 10$, $\beta = 10$ in step A for

640 $\lambda_1 = 1$ and $\alpha_2 = 10$, $\gamma = 0.1$ in step B for $\lambda_2 = 1$. The ComBat, BERMUDA and
641 DESC still utilize the default parameters and tutorials. In NormAE, the (lr_rec,
642 lr_disc_b), epoch and batch_size are set to (0.0002, 0.0001), (100, 50, 100), 200,
643 while other parameters are yet defaulted.

644

645 **Abbreviations**

646 scRNA-seq , single-cell RNA sequencing

647 MS, mass spectrometry

648 MALDI, matrix assisted laser desorption/ionization

649 SVA, surrogate variable analysis

650 MMD, maximum mean discrepancy

651 MSE, mean square error

652 AE, Autoencoder

653 GAN, Generative Adversarial Network

654 MCC, Matthews correlation coefficient

655 MI, myocardial infarction

656 CHD, coronary heart disease

657

658 **References**

659 1 Md. Mohaiminul Islam, Y. W. a. P. H. Deep Learning Models for Predicting Phenotypic
660 Traits and Diseases from Omics Data. *Artificial Intelligence - Emerging Trends and*
661 *Applications*, doi:10.5772/intechopen.75311 (2018).

662 2 Nicholson, J. K. & Lindon, J. C. Metabonomics. *Nature* **455**, 1054-1056,
663 doi:10.1038/4551054a (2008).

664 3 Zenobi, R. Single-Cell Metabolomics: Analytical and Biological Perspectives. *Science* **342**,

665 1243259, doi:10.1126/science.1243259 (2013).

666 4 Wang, T. *et al.* BERMUDA: a novel deep transfer learning method for single-cell RNA
667 sequencing batch correction reveals hidden high-resolution cellular subtypes. *Genome*
668 *Biol* **20**, 165, doi:10.1186/s13059-019-1764-6 (2019).

669 5 Li, X. *et al.* Deep learning enables accurate clustering with batch effect removal in single-
670 cell RNA-seq analysis. *Nature Communications* **11**, 2338, doi:10.1038/s41467-020-
671 15851-3 (2020).

672 6 Stegle, O., Teichmann, S. A. & Marioni, J. C. Computational and analytical challenges in
673 single-cell transcriptomics. *Nat Rev Genet* **16**, 133-145, doi:10.1038/nrg3833 (2015).

674 7 Leek, J. T. *et al.* Tackling the widespread and critical impact of batch effects in high-
675 throughput data. *Nat Rev Genet* **11**, 733-739, doi:10.1038/nrg2825 (2010).

676 8 Cohen, S. L. & Chait, B. T. Influence of Matrix Solution Conditions on the MALDI-MS
677 Analysis of Peptides and Proteins. *Analytical Chemistry* **68**, 31-37, doi:10.1021/ac9507956
678 (1996).

679 9 Akey, J. M., Biswas, S., Leek, J. T. & Storey, J. D. On the design and analysis of gene
680 expression studies in human populations. *Nature Genetics* **39**, 807-808,
681 doi:10.1038/ng0707-807 (2007).

682 10 Zhang, Z. *et al.* Deep learning in omics: a survey and guideline. *Brief Funct Genomics* **18**,
683 41-57, doi:10.1093/bfpg/ely030 (2019).

684 11 Zhang, L. *et al.* Deep Learning-Based Multi-Omics Data Integration Reveals Two
685 Prognostic Subtypes in High-Risk Neuroblastoma. *Frontiers in Genetics* **9**, 477,
686 doi:10.3389/fgene.2018.00477 (2018).

687 12 Lee, T.-Y., Huang, K.-Y., Chuang, C.-H., Lee, C.-Y. & Chang, T.-H. Incorporating deep
688 learning and multi-omics autoencoding for analysis of lung adenocarcinoma
689 prognostication. *Computational Biology and Chemistry* **87**, 107277,
690 doi:<https://doi.org/10.1016/j.compbiolchem.2020.107277> (2020).

691 13 Shaham, U. *et al.* Removal of batch effects using distribution-matching residual networks.
692 *Bioinformatics* **33**, 2539-2546, doi:10.1093/bioinformatics/btx196 (2017).

693 14 Goodfellow, I. J. *et al.* in *Proceedings of the 27th International Conference on Neural*
694 *Information Processing Systems - Volume 2* 2672-2680 (MIT Press, Montreal, Canada,
695 2014).

696 15 Rong, Z. *et al.* NormAE: Deep Adversarial Learning Model to Remove Batch Effects in
697 Liquid Chromatography Mass Spectrometry-Based Metabolomics Data. *Anal Chem* **92**,
698 5082-5090, doi:10.1021/acs.analchem.9b05460 (2020).

699 16 Lazar, C. *et al.* Batch effect removal methods for microarray gene expression data
700 integration: a survey. *Brief Bioinform* **14**, 469-490, doi:10.1093/bib/bbs037 (2013).

701 17 Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data
702 using empirical Bayes methods. *Biostatistics* **8**, 118-127, doi:10.1093/biostatistics/kxj037
703 (2007).

704 18 Leek, J. T. & Storey, J. D. Capturing Heterogeneity in Gene Expression Studies by Surrogate
705 Variable Analysis. *PLoS Genetics* **3**, e161, doi:10.1371/journal.pgen.0030161 (2007).

706 19 Luo, J. *et al.* A comparison of batch effect removal methods for enhancement of prediction
707 performance using MAQC-II microarray gene expression data. *Pharmacogenomics J* **10**,
708 278-291, doi:10.1038/tpj.2010.57 (2010).

- 709 20 Parker, H. S., Corrada Bravo, H. & Leek, J. T. Removing batch effects for prediction
710 problems with frozen surrogate variable analysis. *PeerJ* **2**, e561, doi:10.7717/peerj.561
711 (2014).
- 712 21 Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell RNA-
713 sequencing data are corrected by matching mutual nearest neighbors. *Nat Biotechnol* **36**,
714 421-427, doi:10.1038/nbt.4091 (2018).
- 715 22 Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
716 transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*
717 **36**, 411-420, doi:10.1038/nbt.4096 (2018).
- 718 23 Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e1821,
719 doi:10.1016/j.cell.2019.05.031 (2019).
- 720 24 Laurens van der Maaten, G. H. Visualizing Data using t-SNE. *Journal of Machine Learning*
721 *Research* **9**, 2579-2605 (2008).
- 722 25 Maaten, L. V. D. Accelerating t-SNE using tree-based algorithms. *J. Mach. Learn. Res.* **15**,
723 3221-3245 (2014).
- 724 26 Shi, L. *et al.* The MicroArray Quality Control (MAQC)-II study of common practices for the
725 development and validation of microarray-based predictive models. *Nat Biotechnol* **28**,
726 827-838, doi:10.1038/nbt.1665 (2010).
- 727 27 Vu, T. N. *et al.* Beta-Poisson model for single-cell RNA-seq data analyses. *Bioinformatics*
728 **32**, 2128-2135, doi:10.1093/bioinformatics/btw202 (2016).
- 729 28 Muraro, Mauro J. *et al.* A Single-Cell Transcriptome Atlas of the Human Pancreas. *Cell*
730 *Systems* **3**, 385-394.e383, doi:<https://doi.org/10.1016/j.cels.2016.09.002> (2016).
- 731 29 Lawlor, N. *et al.* Single-cell transcriptomes identify human islet cell signatures and reveal
732 cell-type-specific expression changes in type 2 diabetes. *Genome Res* **27**, 208-222,
733 doi:10.1101/gr.212720.116 (2017).
- 734 30 Segerstolpe, Å. *et al.* Single-Cell Transcriptome Profiling of Human Pancreatic Islets in
735 Health and Type 2 Diabetes. *Cell Metabolism* **24**, 593-607,
736 doi:<https://doi.org/10.1016/j.cmet.2016.08.020> (2016).
- 737 31 Yamada, R., Okada, D., Wang, J., Basak, T. & Koyama, S. Interpretation of omics data
738 analyses. *Journal of Human Genetics*, doi:10.1038/s10038-020-0763-5 (2020).
- 739 32 Huang, L. *et al.* Machine learning of serum metabolic patterns encodes early-stage lung
740 adenocarcinoma. *Nat Commun* **11**, 3556, doi:10.1038/s41467-020-17347-6 (2020).
- 741 33 Qiu, Y. *et al.* Serum Metabolite Profiling of Human Colorectal Cancer Using GC-TOFMS
742 and UPLC - QTOFMS. *Journal of Proteome Research* **8**, 4844-4850,
743 doi:10.1021/pr9004162 (2009).
- 744

745 **Authors' contributions**

746 Conceptualization: J.N., Q.W. Methodology: J.N., Q.W. Acquisition of data and
747 materials: W.X., K.Q. Formal analysis: J.N., Q.W. Writing: J.N. Review and editing,

748 D.W. Supervision: D.W., Q.W. Funding acquisition: Q.W. All authors approved the
749 final version.

750

751 **Funding**

752 Not applicable

753

754 **Availability of data and materials**

755 Complete code and part of data are publicly available at:

756 <https://github.com/n778509775/NWCQ.git>

757

758 **Ethics approval and consent to participate**

759 Not applicable

760

761 **Competing interests**

762 Not applicable

763

764 **Author details**

765 ¹School of Biomedical Engineering, Shanghai Jiao Tong University, Shanghai,

766 200030, China

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementaryinformation.pdf](#)