

# 1 **BatchEval Pipeline: Batch Effect Evaluation Workflow for** 2 **Multiple Datasets Joint Analysis**

3

4 Chao Zhang<sup>1</sup>, Qiang Kang<sup>1</sup>, Mei Li<sup>1</sup>, Hongqing Xie<sup>1</sup>, Shuangfang Fang<sup>1,2</sup>, Xun Xu<sup>3,\*</sup>

5

6 <sup>1</sup> BGI Research, Shenzhen, 518103, China

7 <sup>2</sup> BGI Research, Beijing, 102601, China

8 <sup>3</sup> BGI Research, Wuhan, 430074, China

9 \* Corresponding: [xuxun@genomics.cn](mailto:xuxun@genomics.cn)

10

## 11 **ABSTRACT**

12 As genomic sequencing technology continues to advance, it becomes increasingly important to  
13 perform multiple dataset joint analysis of transcriptomics to understand complex biological  
14 systems. However, batch effect presents challenges for dataset integration, such as sequencing  
15 measured by different platforms and datasets collected at different times. Here, we develop a  
16 BatchEval Pipeline, which is used to evaluate batch effect of dataset integration and output a  
17 comprehensive report. This report consists of a series of HTML pages for the assessment findings,  
18 including a main page, a raw dataset evaluation page and several built-in methods evaluation  
19 pages. The main page exhibits basic information of integrated datasets, comprehensive score of  
20 batch effect and the most recommended method for batch effect removal to current datasets. The  
21 residual pages exhibit the evaluation details of raw dataset and evaluation results of many built-in  
22 batch effect removal methods after removing batch effect. This comprehensive report enables  
23 researchers to accurately identify and remove batch effects, resulting in more reliable and  
24 meaningful biological insights from integrated datasets. In summary, BatchEval Pipeline  
25 represents a significant advancement in batch effect evaluation and is a valuable tool to improve  
26 the accuracy and reliability of the experimental results.

27 **Availability & Implementation:** The source code of BatchEval is available at  
28 <https://github.com/STOmics/BatchEval>.

29

## 30 INTRODUCTION

31 Advancements in gene sequencing technology have facilitated the integrated analysis of  
32 multiple batches of gene transcription data, resulting in more reliable information extracted from  
33 datasets. However, batch effects caused by various sequencing platforms, experimental designs,  
34 experimentalists, laboratory circumstances, and experimental reagent batches are often neglected  
35 during joint analysis. The batch effect introduces technical biases into sequencing to reduce the  
36 dependability of downstream analysis, which must be addressed. Several efficient approaches  
37 have been proposed to minimize or lessen technical biases and batch effect in integrated datasets,  
38 while retaining the most significant biological variance. These approaches include non-linear  
39 models, such as Seurat's classic correlation analysis [1], linear regression models, such as Combat  
40 [2], and implementations based on matching mutual nearest neighbors, such as MNNs [3]. In  
41 some differential gene expression analysis models, such as MSAT [4], DESeq [5] and Limma [6],  
42 batch effect is integrated into the model as a significant component to eliminate it from the dataset.  
43 More recently, works have been published to improve the integration of spatially resolved  
44 transcriptomics by exploiting spatial information, such as spatiAlign [7] and PRECAST [8].

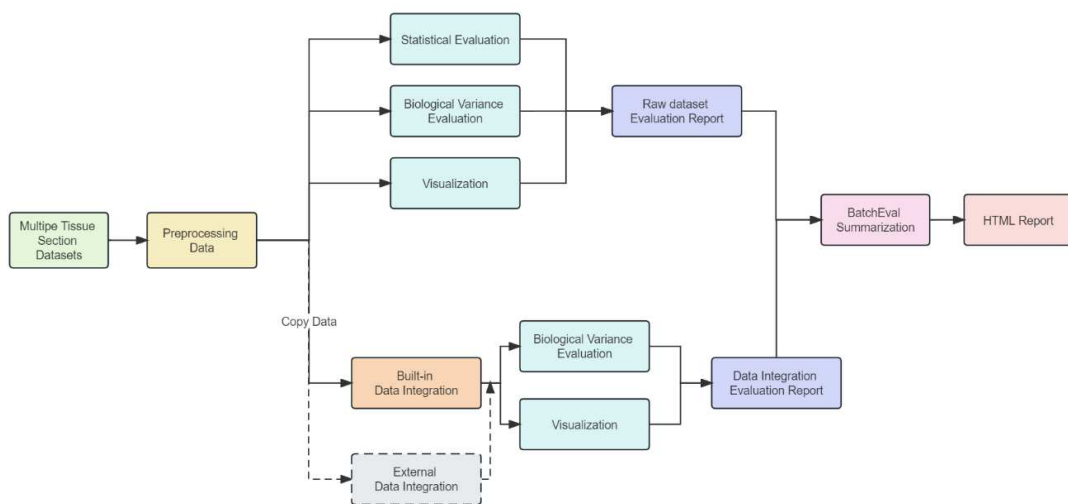
45 Visualization tools are useful for initially assessing dataset integration and quickly comparing  
46 results after removing batch effect. Density distribution visualizations show differences in gene  
47 expression counts between different tissue sections, and overlays of probability density functions  
48 compare variations across sections. Quantitative statistics, such as the chi-square test [9] and the  
49 local inverse Simpson index (*LISI*) [10], precisely quantifies dataset integration. When integrating  
50 multiple different tissue sections, it may not be clear whether batch effect exists or which removal  
51 method should be used, thus requiring detailed analysis. Even though user-friendly software like  
52 BatchQC has been published [11], there is still a shortage of tools suitable for analyzing large-  
53 scale dataset integration.

54 To address these issues, we develop the BatchEval Pipeline<sup>1</sup>, a comprehensive evaluation  
55 batch effect workflow for large-scale dataset integration and output a summary evaluation report.  
56 It evaluates dataset integration from different perspectives related to data mixing and biological

---

<sup>1</sup> <https://github.com/STOmics/BatchEval.git>

57 variance preservation after batch effect removal. The report consists of a series of HTML pages,  
58 including a main page, a raw integration dataset evaluation page and a series of pages of  
59 evaluation results for each built-in batch effect removal method. The main page exhibits the  
60 details of integration datasets, comprehensive score of evaluation batch effect for raw datasets, and  
61 summaries batch effect removal results using different built-in methods (including single-cell-  
62 based methods, such as Harmony [10] and BBKNN [12], and spatially resolved transcriptomics  
63 method, such as spatiAlign [7]), and the most recommended method for testing datasets to remove  
64 batch effect. Users can conveniently access the detailed report of raw integrated datasets and  
65 removal batch effect method through super link in the main page. The raw dataset page consists of  
66 statistical evaluation, biological variance preservation metric score and visualization. And the  
67 batch effect removal method pages also consist of biological variance preservation metric score  
68 and visualization. All of the detailed pages can conveniently go back to the main page through the  
69 top left button. BatchEval Pipeline workflow is shown as in Figure 1. The main contributions of  
70 this study are as follows:



71  
72

Figure 1. BatchEval Pipeline workflow

73 1. The BatchEval Pipeline carefully analyzes batch effect from integrated datasets,  
74 considering various factors related to data mixing and biological variance preservation between  
75 dataset integration before and after.

76 2. The BatchEval Pipeline can provide extensible remove batch effect benchmarking and  
77 output an evaluation report, which includes single-cell-based and spatially resolved  
78 transcriptomics methods. BatchEval Pipeline can also recommend suitable batch effect removal

79 methods according to different datasets.

80 3. The evaluation results of the BatchEval Pipeline are crucial in determining whether and  
81 how to correct batch effects in datasets.

## 82 **RESULTS**

### 83 **Dataset**

84 We collected two groups of spatially resolved transcriptomics datasets to evaluate BatchEval  
85 Pipeline. (1) Two mouse brain olfactory bulb (OB) datasets, which measured by Stereo-seq [13]  
86 and 10x Genomics Visium [14], respectively. Both OB datasets contained five identical cell type  
87 annotations. The Stereo-seq dataset included 1123 spots, each of which had an average of 6028  
88 genes and an average gene expression of 16681. In contrast, the 10x Genomics Visium dataset  
89 included 1184 spots, each of which had an average of 4580 genes and an average gene expression  
90 of 16916. (2) Five time-series mouse embryonic brain datasets, which were measured by Stereo-  
91 seq. These brain sections were collected at different embryonic tissues from E9.5 to E15.5, which  
92 included a total of 81181 spots/cells and 22864 genes in the merged dataset.

### 93 **Statistical analysis with BatchEval Pipeline**

94 Batch effect and sequencing outcomes are typically complex and influenced by various  
95 factors, such as experimental conditions, operators, reagents, and timing, which are not directly  
96 related to the experiment's goal. However, the batch effect can impact the accuracy of downstream  
97 results. The dataset includes a variety of batch effects, as well as associations between linear  
98 and non-linear processes, which can affect the distribution of gene expression probability density  
99 for each spot across subsequent data batches. Additionally, different sequencing techniques can  
100 result in varying sequencing depths across data batches. To address these issues, BatchEval  
101 Pipeline performs Min-Max normalization and logarithmic mapping preprocessing on each  
102 spot/cell gene expression levels and integrates multiple batches of gene expression data into low-  
103 dimensional representations.

104 BatchEval Pipeline employs the Kruskal-Wallis H test [15] to evaluate the variation in the  
105 average level of gene expression across different tissue sections and performs variance analysis on  
106 gene expression total counts for each tissue section. However, when comparing data from different  
107 platforms, such as Stereo-seq and 10x Genomics Visium, it is important to note that they may have

108 considerable differences and not satisfy the homogeneity of variance assumption. Therefore, it  
109 may not be clear how gene expression consistency differs between the data from two platforms. To  
110 determine if gene expression data from several batches originated from the same distribution, a  
111 Kolmogorov-Smirnov Test [16] is performed. By observation, the data from Stereo-seq and 10x  
112 Genomics Visium expression do not fit together properly (Table 1).

113 Table 1. Statistical evaluation results of mouse olfactory bulb dataset before data integration

Variation Analysis	$n_{batch}$	$n_{sample}$	$F$	$p$ value	$F$ ref (2, 3116)
	3	3119	252.6876	0.	2.9986
K-S Test	$n_{sample}$		$k-s$ stat	$p$ value	
batch0-1	1935		0.6924	0.	
batch0-2	1996		0.6115	0.	
batch1-2	2307		0.1579	0.	
Cramer's Test	Pearson Correlation Coefficient		Cramer's V Coefficient		
	0.9270		0.8190		

114 The BatchEval Pipeline utilizes a contingency table to analyze the correlation between  
115 experimental conditions and dataset batches. For example, Cramer's V correlation coefficient [17]  
116 in Table 1 is calculated based on it. By assessing the correlation between different experimental  
117 conditions for each batch, the pipeline determines that the experimental conditions in the Stereo-  
118 seq and 10x Genomics Visium dataset batches are closely related, with a Cramer's V correlation  
119 coefficient of approximately 0.8190 (the statistical evaluation results of mouse embryonic brain  
120 datasets are shown Supplementary Table S1).

### 121 **Biological variance preservation evaluation with BatchEval Pipeline**

122 BatchEval Pipeline adopts a non-linear neural network classifier to estimate data mixing  
123 across multiple tissue sections. The model takes cells/spots gene expression matrices as inputs and  
124 predicts which tissue sections for each cell/spot come from. If the prediction accuracy is low, it  
125 means that the integrated dataset is mixed well, otherwise, the integrated dataset is mixed poor.  
126 As shown in Table 2, the batch/domain estimate score accept rate of mouse embryonic brain  
127 datasets is around 2%, indicating that the classifier can successfully differentiate each spot/cell  
128 come from which tissue sections. Therefore, it is concluded that there is a definite batch effect  
129 between each time-series mouse embryonic brain dataset.

130 Exploratory transcriptome downstream analysis, especially low-dimensional embedding  
131 techniques such as principal component analysis (PCA) and clustering, are commonly used in

132 batch effect evaluation methods. These methods are used to assess the efficiency of dataset  
 133 integration, which can be divided into two parts: retaining biological variances and testing the  
 134 validity of dataset integration/data mixing. BatchEval Pipeline employs k-BET test [9] to evaluate  
 135 the mixing level between different tissue sections in the neighborhood surrounding each spot/cell.  
 136 This method first projects integrated datasets into a low-dimensional embedding by PCA or  
 137 Uniform Manifold Approximation and Projection (UMAP) [18], then applies Pearson’s Chi-square  
 138 test. If the neighborhood of each spot/cell is from various tissue sections, the relevant data mixing  
 139 level is more acceptable, indicating that the data mixing well. In the low-dimensional embedding,  
 140 the dataset can be clearly distinguished, and the batch effect is evident. As shown in Table 2, the  
 141 distributions of time-series mouse embryonic brain datasets differ significantly, and the k-BET  
 142 score accept rate is approximately 0.0085 ( $p = 0.0027 < 0.05$ ).

143 Table 2. Biological variance evaluation results of mouse embryonic brain datasets before data integration

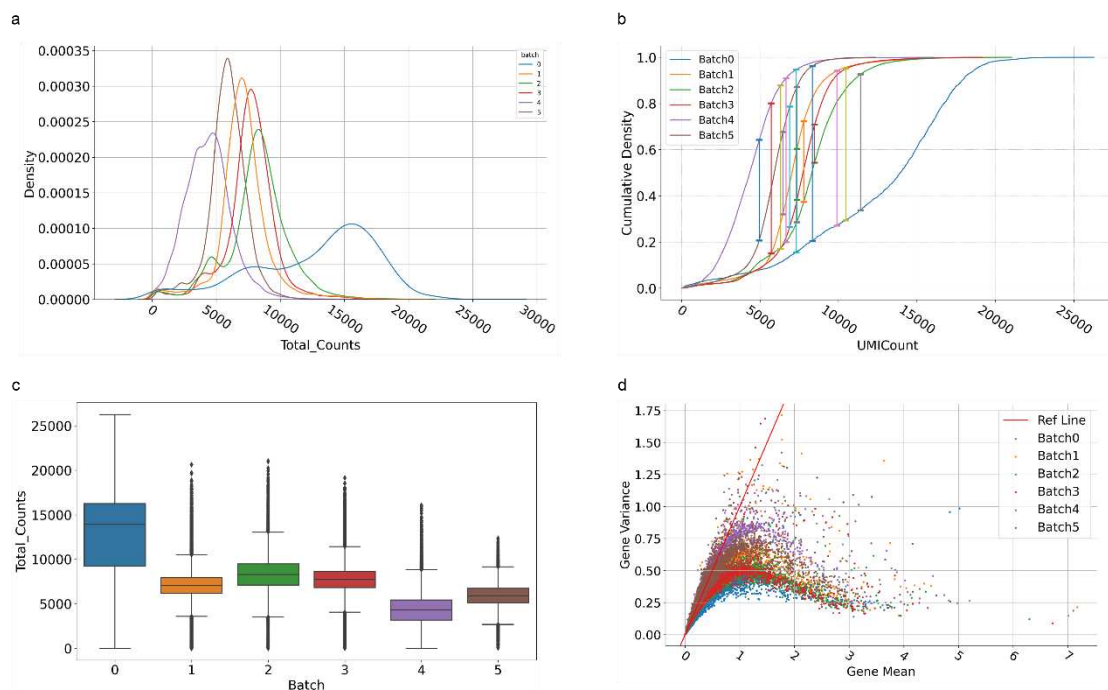
Batch/Domain Estimate score	$n\_batch$	$n\_sample$	train size	accept rate
	5	81181	56826	0.0194
k-BET score	chi mean	95% $p$ value	accept rate	reject rate
	58.0453	0.0027	0.0085	0.9915
Local inverse Simpson’s index		$iLISI$	$cLISI$	F1 score
		0.0154	0.0386	0.0303
Silhouette score		$iSS$	$cSS$	F1 score
		0.6675	0.6155	0.4318

144 The BatchEval Pipeline utilizes the  $LISI$  [10] and silhouette coefficient ( $SS$ ) to estimate the  
 145 data mixing and biological variance preservation while remove batch effect, presuming that cell  
 146 type labels are accessible for the data. After integrating and projecting the data into a low-  
 147 dimensional embedding, the Pipeline computes the  $LISI$  and  $SS$  using two different groupings,  
 148 respectively: (1) grouping using different tissue sections as the batch  $LISI$  ( $iLISI$ )/batch  $SS$  ( $iSS$ )  
 149 score, and (2) grouping known cell types or clustering types as the cell/domain-type  $LISI$  score  
 150 ( $cLISI$ ) or cell/domain-type  $SS$  score ( $cSS$ ). And then using F1 score to summarize the  $LISI$  and  $SS$   
 151 score, respectively. A larger F1 score of  $LISI/SS$  suggests better dataset integration that preserves  
 152 the biological variations between domain cluster types while removing batch effect across  
 153 multiple tissues. As shown in Table 2, the F1 scores of  $LISI$  and  $SS$  are not well, indicating that  
 154 the raw datasets of time-series mouse embryonic brain datasets are mixed worse, and it requires  
 155 further batch effect removal processing (the statistical evaluation results of mouse olfactory bulb

156 dataset are shown Supplementary Table S2).

## 157 Visualization of BatchEval Pipeline

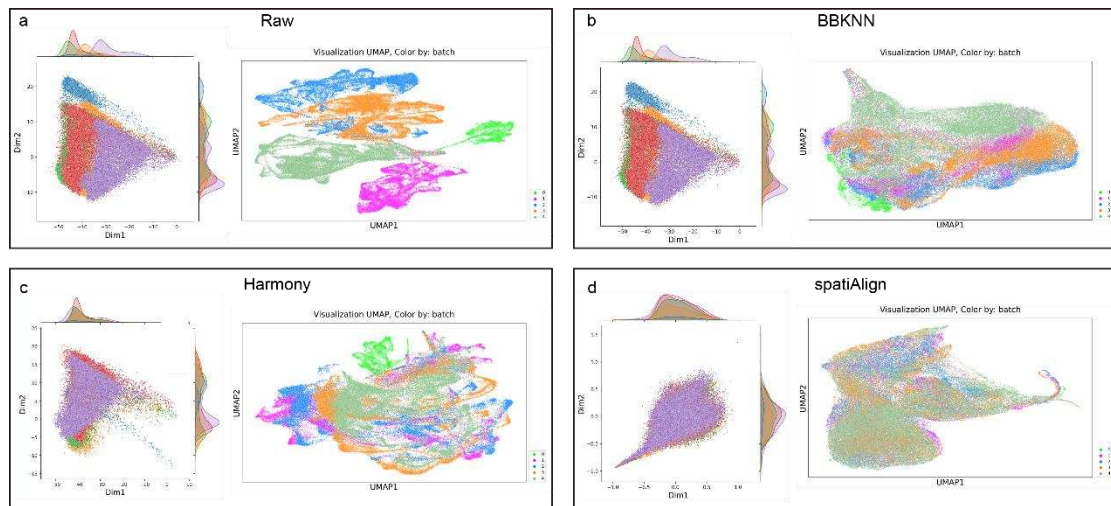
158 In Figure 2, we exhibit the statistical analysis results of mouse embryonic brain datasets  
159 before data integration. The kernel distribution curve of gene expression total counts (Figure. 2a)  
160 and cumulative density curve (CDF) (Figure. 2b) are significant differences. And if the maximum  
161 difference value between CDF curve is lagrer, the more significant of batch effect. The total  
162 counts of gene expression for each spot/cell also have differences (Figure. 2c). Moreover, we  
163 evaluate the mean and variance of each gene (Figure. 2d), and we observed that “batch 4” genes  
164 are more discrete than in other tissues.



165  
166 Figure 2. Visualization of statistical analysis for time-series mouse embryonic brain datasets. (a) Kernel  
167 distribution curve of total gene expression counts for each tissue section, different color represents different  
168 tissues. (b) Cumulative density curve of total gene expression counts for each tissue section. (c) Boxplot of total  
169 gene expression counts. (d) The scatter plot of mean gene and gene variance.

170 Furthermore, BatchEval Pipeline employs PCA to reduce the dimension of merged dataset.  
171 The scatter plots and edge probability density (PDF) curves are used to illustrate the variations in  
172 distribution across multiple data batches in PCA dimensions. As shown in Figure 3a, the first two  
173 PCA dimensions are used as horizontal and vertical coordinates, and a PDF curve is fitted to  
174 capture the variations in the distribution of gene expression in different batches. The BatchEval

175 Pipeline can quickly display the level of data mixing across multiple batches using UMAP and  
176 various color displays based on batch metadata. We utilized the integrated dataset in UMAP space  
177 and colored by “batch”. As shown in Figures 3b, c and d, we observed that spatiAlign can  
178 efficiently mix time-series mouse embryonic brain datasets, however, other methods cannot mix  
179 datasets as well as spatiAlign.



180  
181 Figure 3. Visualization of integration for raw datasets and built-in batch effect removal methods. (a) Visualization  
182 of raw dataset. (b) Visualization of dataset integrated by BBKNN. (c) Visualization of dataset integrated by  
183 Harmony. (d) Visualization of dataset integrated by spatiAlign. All sub-figures left, joint plot of PCA, which using  
184 first two PC component to plot, and right, UMAP plot, colored by different data tissue sections.

185 Since the distribution of gene expression levels in different batches is usually unknown, we  
186 also employed quantile-quantile plots (Q-Q plot) to evaluate data distribution, which can be used  
187 to assess whether data from different batches follow the same distribution [19]. Data from  
188 multiple batches were coupled with each other to generate Q-Q maps for gene expression. When  
189 the gene expression of the matched data was less different, the cumulative density distribution of  
190 the corresponding gene expression was more tightly spaced.

### 191 **BatchEval score evaluation**

192 We employ the biological variance preservation score between batch effect removal before  
193 and after and data integration mixing score to summaries the batch effect of dataset integration.  
194 For example, as shown in Table 3, we utilize mouse embryonic brain datasets for test, and the  
195 conclusion is that this dataset should further perform batch effect removal processing and the most  
196 recommended method for users is “spatiAlign”. This information is essential for researchers to

197 understand the reliability of downstream analyses and to make informed decisions about the  
198 suitability of the dataset for their research goals (the summarizes evaluation results of mouse  
199 olfactory bulb are shown Supplementary Table S3).

200 Table 3. The summary evaluation of time-series mouse embryonic brain datasets

	Raw	spatiAlign	Harmony	BBKNN
k-BET score	0.0085	0.7419	0.1055	0.1598
95% <i>p</i> value	0.0027	0.3183	0.0266	0.0478
<i>BatchEval score</i>	0.1605	0.6315	0.5064	0.3086
Conclusion	This dataset has batch effect and requires further processing and recommend “spatiAlign”. More details of “spatiAlign” can be found in ' <a href="https://github.com/STOmics/Spatialign.git">https://github.com/STOmics/Spatialign.git</a> '.			

## 201 METHOD

202 The gene expression profiles are stored in a matrix with rows representing cells/spots and  
203 columns representing genes in the dataset. To ensure consistent and reliable results in subsequent  
204 analyses, the datasets are preprocessed using Min-Max normalization and log mapping, which  
205 standardizes the values to the same range and transforms them to a logarithmic scale. The  
206 preprocessing can be calculated using the follow parameterization:

$$207 \hat{x} = \log 1p \left( \frac{x - x_{\min}}{x_{\max} - x_{\min}} \right) \quad (1)$$

208 where  $x$  is the original expression of each spot/cell,  $x_{\min}$  and  $x_{\max}$  are the minimum and  
209 maximum expressions captured in dataset, respectively.

### 210 Statistical test

211 The Kruskal-Wallis H test [15] mixes the gene expression of each spot in different batches,  
212 sorts them from the smallest to largest, and records the ordinal number (rank). If the sorted spot  
213 gene expression values are the same, the corresponding rank is the same. Then, rank sum of each  
214 value is calculated with the Student’s T test. The value of  $H$  indicates the distribution of rank in  $k$   
215 batches. The larger the value of  $H$ , the greater the difference in rank is. The calculation formula is  
216 shown as follows:

$$217 H = \frac{12}{N(N+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(N+1) \quad (2)$$

218 where  $N$  denotes the number of all spots in different batches,  $k$  denotes the number of data

219 batches,  $n_i$  denotes the number of spots in batch  $i$ , and  $R_i^2$  denotes the square of the rank sum of  
220 each spot in batch  $i$ .

221 Kolmogorov-Smirnov Test [16] can be used to test whether the cumulative density  
222 distributions of two datasets are different:

$$223 \quad \text{stat} = \sup_x |F_1(x) - F_2(x)| \quad (3)$$

224 where  $x$  denotes the gene expression of spot and  $F(\cdot)$  denotes cumulative density function.

225 Cramer's V correlation coefficient [17] can be used to determine correlation between several  
226 experimental conditions in each dataset. Generating a cross-tabulation between each experimental  
227 condition and batch is the first step. By default, the BatchEval Pipeline generates the cross-  
228 tabulation using several experimental conditions for every batch. The calculation formula is as  
229 follows:

$$230 \quad \phi_c = \sqrt{\frac{\chi^2}{N(k-1)}} \quad (4)$$

231 where  $\chi^2$  denotes Pearson chi-square statistic,  $N$  is the total number of samples tested in the  
232 cross-tabulation table, and  $k$  number of categories of smaller variables in cross-tabulation table.

233 Cramer's V correlation coefficient reflects the magnitude of correlation between data from  
234 different batches under different experimental conditions.

235 Neural network models can fit complex relationships between datasets. BatchEval Pipeline  
236 constructs domain consistency assessment neural network models for estimating the magnitude of  
237 variation in different batches of datasets:

$$238 \quad \hat{x} = \text{relu}\left(\sum_k w_k x_k + b\right) \quad (5)$$

239 where  $\text{relu} = \max(0, x)$ ,  $x$  denotes the gene vector of each spot,  $w$  denotes the model weight,  
240 and  $b$  denotes the model offset. The number of spots included in each batch will be inconsistent,  
241 and BatchEval Pipeline employs Focal Loss [20] to measure the difference between model  
242 prediction results and labels, and the stochastic gradient descent method to update model  
243 parameters. The Focal Loss is calculated as follows:

$$244 \quad \text{Focal\_Loss} = -\alpha (1 - P_i)^\gamma \log P_i \quad (6)$$

245 where,  $\alpha$ ,  $\gamma$  are hyperparameters and  $P_t$  is predicted probability value of model output. The  
246 more accurate results predicted by the model, the higher level of differentiability of the data  
247 domain in different batches and more obvious the batch effect is.

#### 248 **Metric score**

249 The k-BET test [9] is based on the Pearson chi-square test to test the level of mixing of  
250 different batches of data within a local area and is calculated as follows:

$$251 \quad \kappa_j^k = \sum_{i=1}^l \frac{n_{ji}^k - f_i \times k}{f_i \times k} \sim \chi_{l-1}^2 \quad (7)$$

252 where  $\chi_{l-1}^2$  denotes the  $\chi^2$  distribution with degree of freedom  $l-1$  and  $n_{ji}^k$  denotes the number  
253 of cells in subset  $j$  of size  $k$  in batch  $i$ .  $f$  denotes the probability distribution of batch  $i$   
254 samples in the overall population.

255 To simultaneously evaluate the separation of each cell/batch cluster and mixing of multiple  
256 datasets, we calculate the average local inverse Simpson's index (*LISI*) [10] score of the datasets  
257 using two different groupings: (1) grouping using different datasets as the batch *LISI* score (*iLISI*),  
258 and (2) grouping known cell types or clustering types as the cell/domain-type *LISI* score (*cLISI*).  
259 In dataset integration, a larger value of *iLISI* indicates better mixing of datasets, and smaller *cLISI*  
260 indicates better preservation of the biological variance between cell/domain-types. These two  
261 metrics can be summarized using the F1 score as follows:

$$262 \quad \text{F1 score}_{LISI} = \frac{2 \times (1 - cLISI) \times iLISI}{(1 - cLISI) + iLISI} \quad (8)$$

263 where  $cLISI = \text{MinMaxNorm}(1 / \sum_{b=1}^B p^2(b))$ ,  
264  $iLISI = \text{MinMaxNorm}(1 / \sum_{c=1}^C p^2(c))$ ,  $B$  and  $C$  are the cell/domain-type or batch marker,  
265  $p(b)$  and  $p(c)$  are the probability value that the cell/domain-type or batch in the local area,  
266 respectively.

267 Furthermore, we also implement silhouette coefficient (*SS*) to evaluate the separation of each  
268 domain cluster and mixing multiple datasets. Similar to F1 score of *LISI*, we calculate the average  
269 *SS* of the datasets using two grouping label, domain cluster type (*cSS*), different dataset type (*iSS*),  
270 respectively. These two metrics also can be summarized using F1 score as follows,

271 
$$\text{F1 score}_{SS} = \frac{2 \times (1 - iSS') \times cSS'}{(1 - iSS') + cSS'} \quad (9)$$

272 where,  $iSS' = 1 + iSS / 2$  and  $cSS' = 1 + cSS / 2$ . A larger F1 score of  $SS$  suggests better  
273 dataset integration that preserves the biological variations between domain cluster types while  
274 removing batch effect across multiple tissues.

### 275 **BatchEval score**

276 Batch effect is confounded in the data, and BatchEval Pipeline provides a comprehensive  
277 assessment of the batch effect of aggregated data from different dimensions. BatchEval Pipeline is  
278 designed to calculate the final batch effect score by weighting the mean value with the formula as  
279 follows.

280 
$$\text{BatchEval score} = \text{Mean}(\text{F1 score}_{LSI}, \text{F1 score}_{SS}, (1 - \text{DomainAcc})) \quad (10)$$

281 where *DomainAcc* is the accuracy of neural network domain classifier.

## 282 **DISCUSSION**

283 Although many approaches have been developed to remove batch effect, there is still a lack  
284 of effective methods to evaluate batch effect for large-scale dataset integration, especially spatially  
285 resolved transcriptomics. The sources and effects of batch effect can vary greatly from experiment  
286 to another, and it is essential to analyze the most common potential sources of batch effect to  
287 improve the effectiveness of their removal and to facilitate the integration of data from different  
288 batches. Additionally, there is a need for a more efficient way to accurately determine the extent to  
289 which data is affected by batch effect, as well as to remove batch effect more accurately.

290 The BatchEval Pipeline simplifies batch pre-testing by offering a comprehensive evaluation  
291 report for data integration, making it particularly beneficial for multi-omics studies with multiple  
292 datasets or samples collected at different times. It provides statistical testing, batch effect metrics  
293 evaluation, and visualization, allowing researchers to efficiently explore and correct for batch  
294 effect in their data.

295 Although substantial progress has been made in identifying and evaluating batch effect, there  
296 is still much for improvement to enhance the accuracy and effectiveness of batch effect removal.  
297 BatchEval Pipeline is a powerful tool for evaluation of integrated large-scale gene expression  
298 datasets. It provides a quantitative measure of biological variance preservation and data

299 integration mixing, and the conclusion indicates whether the batch effect is significant or not.  
300 Using the BatchEval Pipeline, users can objectively evaluate the presence and severity of batch  
301 effects in their integrated datasets. This feature makes the tool particularly valuable for  
302 researchers, who need to analyze large datasets, as it provides an easy and reliable way to assess  
303 data quality and ensures that downstream analyses are robust and reliable.

## 304 **AVAILABILITY OF SOURCE CODE AND REQUIREMENTS**

305 Project name: BatchEval

306 Project home page: <https://github.com/STOmics/BatchEval>

307 Operating system(s): Platform independent

308 Programming language: Python

309 Tutorials: <https://batcheval.readthedocs.io/en/latest/index.html>

310 License: MIT License

## 311 **DATA AVAILABILITY**

312 The mouse olfactory bulb and embryonic brain datasets measured by Stereo-seq can be download  
313 from: <https://db.cngb.org/stomics/mosta>, and the 10x Genomics Visium datasets can be download  
314 from: <https://www.10xgenomics.com/resources/datasets/adult-mouse-olfactory-bulb-1-standard>.

## 315 **DECLARATIONS**

### 316 **Ethics approval and consent to participate**

317 Not applicable.

### 318 **Competing interests**

319 The authors declare they have no competing interests.

### 320 **Funding**

321 This work was supported by the National Key R&D Program of China (2022YFC3400400).

### 322 **Authors' contributions**

323 Conceptualization: CZ; Project administration and supervision: XX; Software: CZ; Data  
324 collection, processing, and application: ZC and HX; Project coordination: SF; Manuscript writing  
325 and figure generation: ZC and QK; Manuscript review: QK and ML.

326 **Acknowledgements**

327 We thank China National GeneBank for providing data support for this study.

328 **SUPPLEMENTARY MATERIAL**

329 Supplementary Table S1. Statistical evaluation results of mouse embryonic brain datasets before data integration

Variation Analysis	$n\_batch$	$n\_sample$	$F$	$p$ value	$F$ ref (2, 3116)
	5	81181	20433.7507	0.	2.3720
K-S Test	$n\_sample$		$k-s$ stat	$p$ value	
batch0-1	16673		0.6706	0.	
batch0-2	20986		0.5892	0.	
batch0-3	24791		0.6569	0.	
batch0-4	28064		0.7920	0.	
batch1-2	31437		0.3494	0.	
batch1-3	35242		0.2209	0.	
batch1-4	38515		0.6507	0.	
batch2-3	39555		0.1641	0.	
batch2-4	42828		0.7101	0.	
batch3-4	46633		0.7073	0.	
Cramer's Test	Pearson Correlation Coefficient		Cramer's V Coefficient		
	0.9775		0.9005		

330 Supplementary Table S2. Biological variance evaluation results of mouse olfactory bulb dataset before data

331 integration

Batch/Domain Estimate	$n\_batch$	$n\_sample$	train size	accept rate
score	3	3119	2183	0
k-BET score	chi mean	95% $p$ value	accept rate	reject rate
	29.8851	0	0	1
Local inverse Simpson's index	$iLISI$	$cLISI$	F1 score	
	0.0034	0.0888	0.0067	
Silhouette score	$iSS$	$cSS$	F1 score	
	0.7842	0.4911	0.2999	

332 Supplementary Table S3. The summary evaluation of mouse olfactory dataset

	Raw	spatiAlign	Harmony	BBKNN
k-BET score	0	0.6916	0.0901	0.1828
95% $p$ value	0	0.2870	0.0228	0.0546
<i>BatchEval</i> score	0.1022	0.7035	0.5256	0.3305
Conclusion	This dataset has batch effect and requires further processing and recommend "spatiAlign". More details of 'spatiAlign' can be found in ' <a href="https://github.com/STOmics/Spatialign.git">https://github.com/STOmics/Spatialign.git</a> '.			

333 **REFERENCES**

334 1. Hao Y, Hao S, Andersen-Nissen E *et al*: **Integrated analysis of multimodal single-cell**

- 335 **data.** *Cell* 2021, **184**(13):3573-3587. e3529.
- 336 2. Johnson WE, Li C, Rabinovic A: **Adjusting batch effects in microarray expression**  
337 **data using empirical Bayes methods.** *Biostatistics* 2007, **8**(1):118-127.
- 338 3. Haghverdi L, Lun AT, Morgan MD *et al*: **Batch effects in single-cell RNA-sequencing**  
339 **data are corrected by matching mutual nearest neighbors.** *Nature biotechnology*  
340 2018, **36**(5):421-427.
- 341 4. Finak G, McDavid A, Yajima M *et al*: **MAST: a flexible statistical framework for**  
342 **assessing transcriptional changes and characterizing heterogeneity in single-cell**  
343 **RNA sequencing data.** *Genome biology* 2015, **16**(1):1-13.
- 344 5. Love MI, Huber W, Anders S: **Moderated estimation of fold change and dispersion for**  
345 **RNA-seq data with DESeq2.** *Genome biology* 2014, **15**(12):1-21.
- 346 6. Ritchie ME, Phipson B, Wu D *et al*: **limma powers differential expression analyses for**  
347 **RNA-sequencing and microarray studies.** *Nucleic acids research* 2015, **43**(7):e47-e47.
- 348 7. Zhang C, Liu L, Zhang Y *et al*: **spatiAlign: An Unsupervised Contrastive Learning**  
349 **Model for Data Integration of Spatially Resolved Transcriptomics.** *bioRxiv* 2023.
- 350 8. Liu W, Liao X, Luo Z *et al*: **Probabilistic embedding, clustering, and alignment for**  
351 **integrating spatial transcriptomics data with PRECAST.** *Nature Communications*  
352 2023, **14**(1):296.
- 353 9. Büttner M, Miao Z, Wolf FA *et al*: **A test metric for assessing single-cell RNA-seq**  
354 **batch correction.** *Nature methods* 2019, **16**(1):43-49.
- 355 10. Korsunsky I, Millard N, Fan J *et al*: **Fast, sensitive and accurate integration of single-**  
356 **cell data with Harmony.** *Nature methods* 2019, **16**(12):1289-1296.
- 357 11. Manimaran S, Selby HM, Okrah K *et al*: **BatchQC: interactive software for evaluating**  
358 **sample and batch effects in genomic data.** *Bioinformatics* 2016, **32**(24):3836-3838.
- 359 12. Polański K, Young MD, Miao Z *et al*: **BBKNN: fast batch alignment of single cell**  
360 **transcriptomes.** *Bioinformatics* 2020, **36**(3):964-965.
- 361 13. Chen A, Liao S, Cheng M *et al*: **Spatiotemporal transcriptomic atlas of mouse**  
362 **organogenesis using DNA nanoball-patterned arrays.** *Cell* 2022, **185**(10):1777-1792  
363 e1721.
- 364 14. Stahl PL, Salmen F, Vickovic S *et al*: **Visualization and analysis of gene expression in**

- 365 **tissue sections by spatial transcriptomics.** *Science* 2016, **353**(6294):78-82.
- 366 15. MacFarland TW, Yates JM, MacFarland TW *et al*: **Kruskal–Wallis H-test for oneway**  
367 **analysis of variance (ANOVA) by ranks.** *Introduction to nonparametric statistics for*  
368 *the biological sciences using R* 2016:177-211.
- 369 16. Massey Jr FJ: **The Kolmogorov-Smirnov test for goodness of fit.** *Journal of the*  
370 *American statistical Association* 1951, **46**(253):68-78.
- 371 17. Cramér H: **Mathematical methods of statistics**, vol. 26: Princeton university press;  
372 1999.
- 373 18. McInnes L, Healy J, Melville J: **Umap: Uniform manifold approximation and**  
374 **projection for dimension reduction.** *arXiv preprint arXiv:180203426* 2018.
- 375 19. Tsai D-M, Yang C-H: **A quantile–quantile plot based pattern matching for defect**  
376 **detection.** *Pattern Recognition Letters* 2005, **26**(13):1948-1962.
- 377 20. Lin T-Y, Goyal P, Girshick R *et al*: **Focal loss for dense object detection.** In:  
378 *Proceedings of the IEEE international conference on computer vision: 2017*; 2017: 2980-  
379 2988.
- 380