1    Defining effective strategies to integrate multi-sample single-nucleus ATAC-seq datasets via a

2    multimodal-guided approach

3    Kathryn Weinand[1,2,3,4,5], Erica M. Langan[1,6], Michelle Curtis[1,2,3,4], Soumya Raychaudhuri[1,2,3,4,5,*]

4

5    [1] Broad Institute of MIT and Harvard, Cambridge, MA, USA.

6    [2] Division of Rheumatology, Inflammation, and Immunity, Department of Medicine, Brigham and

7    Women's Hospital and Harvard Medical School, Boston, MA, USA.

8    [3] Division of Genetics, Department of Medicine, Brigham and Women's Hospital and Harvard

9    Medical School, Boston, MA, USA.

10    [4] Center for Data Sciences, Brigham and Women's Hospital and Harvard Medical School,

11    Boston, MA, USA.

12    [5] Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA.

13    [6] Harvard-MIT Program in Health Sciences and Technology, Harvard Medical School, Boston,

14    MA, USA

15

16    *Correspondence to:
17    Soumya Raychaudhuri
18    Harvard New Research Building
19    77 Avenue Louis Pasteur, Suite 250
20    Boston, MA 02115
21    soumya@broadinstitute.org
22    Ph: 617-525-4484 Fax: 617-525-4488
23

24 **Abstract**

25 **Background**

26 Chromatin accessibility, measured via single-nucleus Assay for Transposase-Accessible

27 Chromatin with sequencing (snATAC-seq), can reveal the underpinnings of transcriptional

28 regulation across heterogeneous cell states. As the number and scale of snATAC-seq datasets

29 increases, we need robust computational pipelines to integrate samples within a dataset and

30 datasets across studies. These integration pipelines should correct cell-state-obfuscating

31 technical effects while conserving underlying biological cell states, as has been shown for

32 single-cell RNA-seq (scRNA-seq) pipelines. However, scRNA-seq integration methods have

33 performed inconsistently on snATAC-seq datasets, potentially due to sparsity and genomic

34 feature differences.

35 **Results**

36 Using single-nucleus multimodal datasets profiling ATAC and RNA simultaneously, we can

37 measure snATAC-seq integration method performance by comparison to independently

38 integrated snRNA-seq gold standard embeddings and annotations. Here, we benchmark 58

39 pipelines, incorporating 7 integration methods plus 1 embedding correction method with 5

40 feature sets. Using our command-line tool, we assessed 5 multimodal datasets at 3 different

41 resolutions using 2 novel metrics to determine the best practices for multi-sample snATAC-seq

42 integration. ATAC features outperformed Gene Activity Score (GAS) features, and embedding

43 correction with Harmony was generally useful. SnapATAC2, PeakVI, and ArchR's iterative

44 Latent Semantic Indexing (LSI) performed well.

45 **Conclusions**

46 We recommend SnapATAC2 + Harmony with pre-defined ENCODE candidate *cis*-regulatory

47 element (cCRE) features as a first-pass pipeline given its metric performance, generalizability of

48 features, and method resource-efficiency. This and other high-performing pipelines will guide

49 future comprehensive gene regulation maps.

50    **Keywords**: chromatin accessibility, snATAC-seq, integration, benchmark, bio-conservation,

51    batch correction, multimodal datasets

52

53 **Background**

54      Single-nucleus ATAC-seq data (snATAC-seq) defines the open chromatin landscape of

55  individual cells[1] and has emerged as an important complement to single-cell RNA-seq (scRNA-

56  seq) studies. Open chromatin indicates active regulatory regions such as promoters and

57  enhancers to elucidate the transcriptional regulation of gene expression programs. As scalable

58  snATAC-seq platforms have become affordable and available[2–4], the number of these datasets

59  has increased dramatically. With this expansion, it is crucial to determine the extent of both

60  shared and distinct cell states, transcription factors (TFs), and ultimately developmental

61  mechanisms[4–6] and disease therapies[3,7–9] across datasets and diseases. However, it is

62  challenging to incorporate cells from multiple datasets in different studies[5,6], encompassing

63  diverse experimental protocols and tissues, as these technical differences often obscure

64  important biological signals. Even in individual large studies[3,5,7], combining samples obtained

65  across separate batches may result in technical bias. Therefore, effective data integration

66  across snATAC-seq samples and studies is essential.

67      Technical confounding has also been observed in single-cell and single-nucleus RNA-

68  seq studies[10–14]. By assuming that cellular cell states are shared across subsets of samples,

69  investigators have created highly effective scRNA-seq integration algorithms that remove

70  technical effects while conserving the underlying biology. A recent benchmark[15] suggested that

71  these scRNA-seq methods effectively balanced conserving biological signals (bio-conservation)

72  and performing batch correction.

73      However, these same methods have not demonstrated reliable performance for

74  snATAC-seq data[15], perhaps due to differences in genomic features. While scRNA-seq analysis

75  initially uses a pre-defined set of ~20,000 protein-coding genes, snATAC-seq analysis does not

76  have a universal feature set. Peaks of localized signals called within a dataset are the most

77  common feature type for snATAC-seq data. However, snATAC-seq data typically has fewer

78  reads mapping to a larger feature set (~100,000 peaks)[16] than scRNA-seq data. Therefore,

79    single cell integration methods, largely derived for scRNA-seq data, struggle with the increased

80    sparsity when applied to snATAC-seq data.

81        To assess bio-conservation and batch correction for a snATAC-seq integration pipeline,

82    a benchmarking study requires gold standard datasets with clear ground truths. Previous

83    snATAC-seq integration benchmark studies have used datasets with cell types defined using

84    the tested modality[15], with broad cell type annotations that are relatively easy to distinguish[15],

85    with confounded cell types and batches[17], or avoided batch effects altogether[18]. However, we

86    can strategically utilize multimodal datasets that simultaneously capture gene expression and

87    chromatin accessibility measurements for the same cells. In these instances, we can define a

88    gold standard cell type and embedding by integrating snRNA-seq profiles with state-of-the-art

89    methods shown to be effective in both bio-conservation and batch correction[15,19]. This

90    embedding can be used as an independent high-resolution standard to benchmark snATAC-seq

91    integration pipelines. To assess the quality of snATAC-seq integration, we use metrics

92    assessing integration at the most local level by comparing to the independent snRNA-seq

93    integrated embeddings for the same cells.

94        Here, we benchmarked 58 snATAC-seq data integration pipelines across 7 methods, 5

95    feature sets, and 1 embedding correction[20] to define best practices for the genomics community

96    (**Fig. 1a-b**; **Additional file 1: Table S1**). We assessed each pipeline using 2 novel local metrics

97    (**Fig. 1c-f**) to quantify efficacy in integration while preserving subtle functional differences

98    between cells. We also used 2 established metrics for comparison. Using our command-line

99    tool, we applied these strategies across 5 diverse blood and tissue datasets[21,22] at increasing

100   resolutions that represent increasing integration difficulty (**Fig. 1g-h**). Overall, SnapATAC2[23] +

101   Harmony[20] using a pre-defined set of ENCODE candidate *cis*-regulatory element[19] (cCRE)

102   features had the best performance. ArchR[24] and PeakVI[25] were also among the better-

103   performing methods, though the latter was more resource-intensive. Furthermore, we found that

5

104 ATAC features (peaks, cCREs, tiles) consistently performed better than Gene Activity Score

105 (GAS) features.

106

107 **Results**

108 **Benchmark Design**

109 We benchmarked pipelines consisting of three parts: one of five feature sets, one of

110 seven integration methods, and the presence or absence of a correction method (**Additional**

111 **file 1: Table S1**); this was followed by post-processing and quantification of performance

112 metrics (**Fig. 1a**). The seven integration methods used the feature matrices to create an initial

113 embedding while the correction method, Harmony[20], adjusted those embeddings to account for
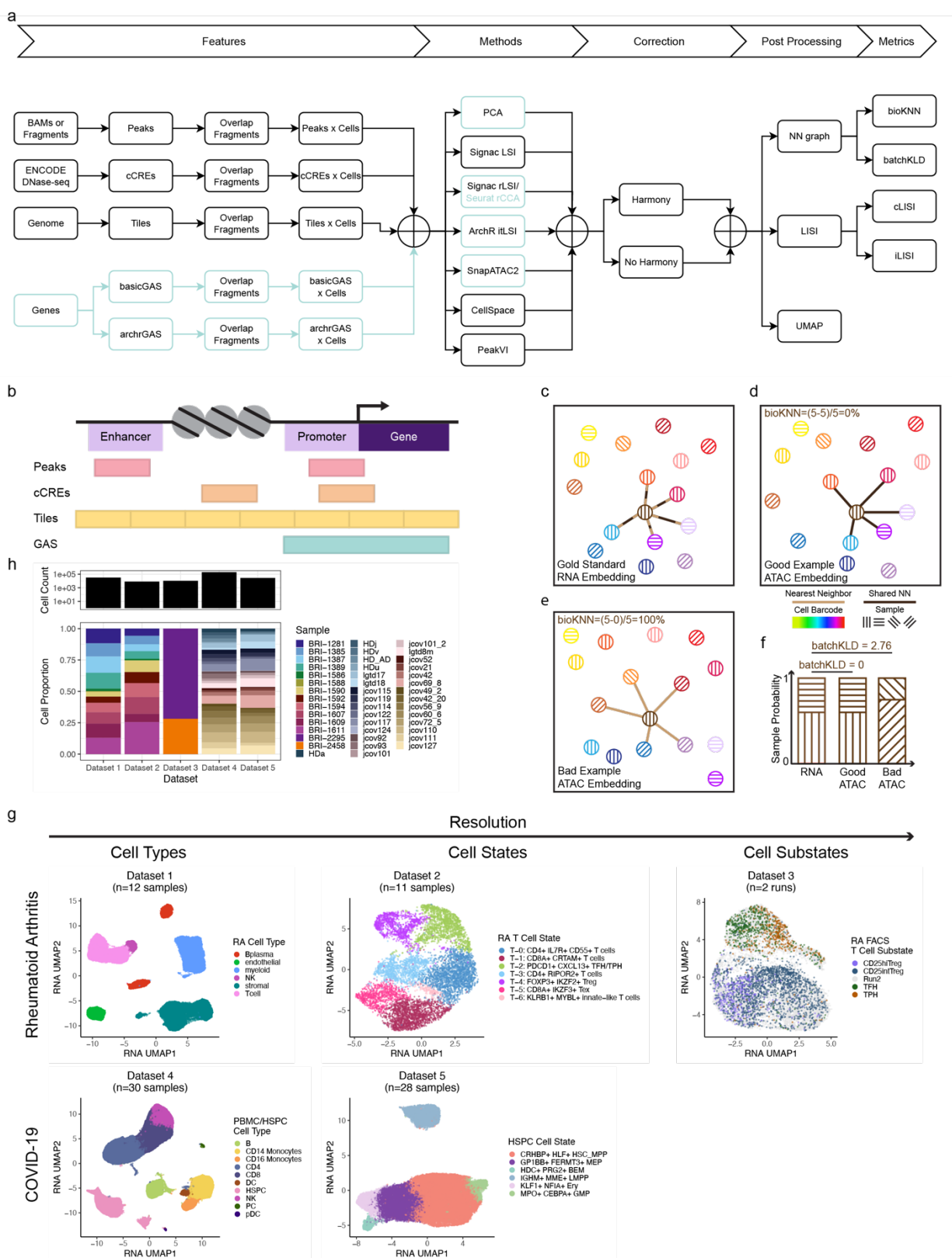
114 batch effects.

**Fig. 1.** Benchmark design.
**a.** Overview of snATAC-seq integration benchmark. GAS features only used methods highlighted in blue and used Seurat rCCA instead of Signac rLSI (**Methods**). **b.** For an example

119　locus, a visual comparison of the four feature types tested: peaks, cCREs, tiles, GASs. **c.-e.**
120　bioKNN metric. Cell identity is denoted by color. **c.** For the brown cell, the nearest neighbors
121　(NN) defined by snRNA-seq. NN edges are striped to correspond with both **d.** and **e.**. **d.** In a
122　good example, the brown cell's NN defined by snATAC-seq; since all top 5 NN are shared in
123　both modalities, denoted by dark brown NN edges, the bioKNN is (5-5)/5=0%, the best possible
124　value. **e.** In a bad example, the brown cell's NN defined by snATAC-seq; since no top 5 NN are
125　shared in both modalities, denoted by tan NN edges, the bioKNN is (5-0)/5=100%, the worst
126　possible value. **f.** Batch Kullback-Leibler Divergence (batchKLD) metric. Sample identity is
127　denoted by pattern. The batchKLD between top NN from **c.** and **d.** was 0, the best possible
128　value. The batchKLD between top NN from **c.** and **e.** was 2.76, this example's worst possible
129　value (**Methods**). **g.** We tested five multimodal datasets from two studies[21,22]: (1) RA
130　inflammatory tissue broad cell types, (2) RA inflammatory tissue T cell states, (3) RA PBMC T
131　cell substates, (4) COVID-19 PBMC broad cell types, and (5) COVID-19 PBMC HSPC states.
132　Dataset 3 cell phenotypes were determined using FACS (**Methods**). Otherwise, the cell
133　phenotypes were determined by multimodal snRNA-seq. UMAPs were defined using snRNA-
134　seq. **h.** The total (log10) cell counts across all 5 Benchmarking Datasets (**top**) as well as the
135　proportion of cells per sample (**bottom**). The samples in Datasets 4/5 are grouped by healthy
136　(blue), non-COVID-19 (slate), mild/mild_late COVID-19 (pinkish gray), early COVID-19 (pink),
137　and late COVID-19 (yellow).
138
139　　　　Of the 5 different feature sets, 3 were ATAC-specific features and 2 were gene-like GAS

140　features (**Fig. 1b**; **Table 1**). Peaks called using each dataset's BAM or fragment files were the

141　most dataset-specific feature set; these features described the intrinsic signal within a dataset

142　but transfer poorly to other datasets. In contrast, tiles, non-overlapping fixed-length sliding

143　windows tiled across the entire genome, were highly transferrable. Tiles, however, were

144　memory-intensive, with a set size an order of magnitude larger than called peaks. Furthermore,

145　arbitrary tile boundaries may disrupt meaningful biological signals. cCREs were regulatory

146　regions called from thousands of ENCODE bulk DNase-seq experiments spanning diverse cell

147　types[19]; they had the generalizability of tiles, but the biological value of peaks, though they may

148　miss some dataset-specific biology. We fixed each of these three snATAC-seq features to be

149　500 bp to avoid length as a confounding factor (**Methods**). For the GAS features, we

150　aggregated signal in and around genes in two different ways. The "basicGAS" score simply

151　counted fragments that overlapped the 2 kb promoter and gene body. Increasing complexity,

152　"archrGAS" score used ArchR's[24] default exponential decay model, which weights tiles

153　overlapping gene windows, with high weights in the 5 kb promoter and gene body, but

154　decreased weights further away; it was also linearly scaled for gene length. As with genes,

8

155 GASs were transferable across datasets. We note that since linking of ATAC reads to target

156 genes remains challenging, GAS scores may not accurately reflect the biological signal of a

157 gene. Indeed, within broad cell types using multimodal snATAC-/snRNA-seq data, we

158 previously[26] found that there was minimal correlation between real gene expression and five

159 different GAS methods, including the two used here.

160 **Table 1.** Feature summary.

| Feature | Peak | cCRE | Tile | GAS |
|---|---|---|---|---|
| Dataset-specific features | Yes | No | No | No |
| Localized signal | Yes | Maybe | No | No |
| Biologically meaningful | Yes | Yes | No | Yes |
| Transferrable | No | Yes | Yes | Yes |
| Standard size | Enforced | Enforced | Yes | No |
| Feature counts | ~100K | ~600K | ~6M | ~20K |
| Overlapping features | Yes | Yes | No | Yes |
| Tn5 bias correction | Yes | No | No | No |

161 cCRE: candidate cis-regulatory element; GAS: gene activity score.

162       We assessed seven integration methods for initial embedding generation: (1) Principal

163 Component Analysis (PCA), (2) Signac[27] using latent semantic indexing (LSI), (3) Signac using

164 reciprocal LSI (rLSI), (4) ArchR[24] using iterative LSI (itLSI), (5) SnapATAC2[23] (SA2) using

165 matrix-free spectral embedding, (6) CellSpace[28] (CS) using neural embedding, and (7) PeakVI[25]

166 (pVI) using a variational autoencoder (**Table 2**). These methods reflected linear (1-4) and non-

167 linear (5-7) options. Only rLSI/rCCA and PeakVI specified a batch variable for explicit batch

168 correction. Since CellSpace, PeakVI, Signac LSI, and Signac rLSI were implemented for use

169 with ATAC features, we did not test GAS features with them. However, to test GAS features in a

170 framework similar to Signac, we used Seurat[29], a method originally designed to analyze scRNA-

171 seq data, with GAS features inputs. We specifically tested Seurat's batch-aware reciprocal

172 canonical correlation analysis (rCCA) method, which is related to the anchor integration

173 approach used by rLSI.

174

9

175 **Table 2.** Methods summary.

| Method | Original Study | Number of Citations | Type | Linear | Batch Variable Specified | Features Tested | Note |
|---|---|---|---|---|---|---|---|
| PCA | Multiple studies (e.g., irlba in R) | NA | Principal Components Analysis | Yes | No | Peak, cCRE, Tile, GAS | |
| Signac | Stuart et al., Nature Methods, 2021 | 625 | Latent Semantic Indexing | Yes | No | Peak*, cCRE, Tile | GAS features processed by Seurat would default to PCA |
| ArchR | Granja, Corces, et al., Nature Genetics, 2021 | 546 | Reciprocal Latent Semantic Indexing | Yes | Yes | Peak*, cCRE, Tile, GAS | GAS features were processed by Seurat with CCAIntegration |
| SnapATAC2 | Zhang et al., Nature Methods, 2024 | 18 | Iterative Latent Semantic Indexing | Yes | No | Peak, cCRE, Tile*, GAS | |
| CellSpace | Tayyebi et al., Nature Methods, 2024 | 2 | Matrix-free spectral embedding | No | No | Peak*, cCRE, Tile, GAS | |
| PeakVI | Ashuach et al., Cell Reports Methods, 2022 | 39 | Neural embedding | No | No | Peak*, cCRE, Tile* | k-mers (8-mers by default) are sampled from inputted features |
| Harmony | Korsunsky et al., Nature Methods, 2019 | 3750 | Variational autoencoder | No | Yes | Peak*, cCRE, Tile | |
| | | | Soft k-means clustering | Yes | Yes | Peak, cCRE, Tile, GAS | Correction method used with embeddings from all other methods |

176 Number of citations on each journal article's website as of 3/18/25. Asterisks in 'Features
177 Tested' column denote the method's preferred feature(s). Note that for Seurat, its citation is Hao

178 et al., Nature Biotechnology, 2024, it has 628 citations, and it was originally built for scRNA-seq
179 datasets.
180

181     To assess the value of additional batch correction, we tested the inclusion of a high-

182 performing[15] embedding correction method, Harmony[20] (**Table 2**). While Harmony was initially

183 developed for scRNA-seq data, it is often used in conjunction with the assessed snATAC-seq

184 methods[7,21,22]. Harmony corrected each previously-defined embedding to integrate samples

185 within a dataset. We denoted each embedding without Harmony as "Method" and its Harmony-

186 corrected embedding as "Method + Harmony".

187     We assessed each snATAC-seq embedding by comparing it to a gold-standard

188 embedding using two per-cell nearest-neighbor (NN) metrics we developed: bio-conservation

189 KNN (bioKNN) and batch Kullback-Leibler divergence (batchKLD). We chose local-based

190 metrics rather than cluster-based metrics since the latter can obscure inaccuracies by

191 aggregation. Our metrics relied on multimodal data, which provided two readouts, ATAC and

192 RNA, of the same underlying biological identity. Assuming modality-specific batch differences

193 were adequately addressed, each cell should be surrounded by similar cells in both modalities.

194 These shared cells can be identified by cell barcodes present in both modalities, thus avoiding

195 the need for cell type annotations. Using well-established scRNA-seq integration methods, we

196 first created a gold standard embedding from the multimodal RNA, upon which we defined the

197 NN for each cell (**Fig. 1c**). Then, for each snATAC-seq integration strategy, we used its

198 embedding to define per-cell NN (**Fig. 1d-e**). Finally, for each cell, we counted how many of its

199 ATAC top K NN overlapped with its RNA top K NN and quantified bioKNN as the percentage of

200 K cells that were not shared across modalities (**Methods**). A lower bioKNN value reflected a

201 more aligned integration where the underlying biology was shared across modalities.

202 Furthermore, since all the shared cells could be from the same sample and hence poorly

203 integrated, we also assessed batch correction using batchKLD (**Fig. 1f**) between the per-cell

204 sample distributions of the top K NN for the gold standard RNA embedding (**Fig. 1c**) and the top

205    K NN for each ATAC embedding (**Fig. 1d-e**). A lower batchKLD value was favorable, meaning

206    the sample distributions between the two modalities were more similar; the upper bound was

207    blunted by a dataset-driven pseudocount (**Methods**).

208        For comparison, we also quantified commonly used local benchmarking metrics[15,18,23,29],

209    introduced in Korsunsky et al.[20], defined by Local Inverse Simpson's Index (LISI) scores: cell

210    type LISI (cLISI) for bio-conservation and integration LISI (iLISI) for batch correction. A low cLISI

211    score was preferable, signifying that cell types were not mixing and biological variation was

212    preserved. A high iLISI score indicated more sample mixing and suggested an effective

213    correction of sample-specific technical variation. While useful, these metrics had some pitfalls

214    addressed by our NN metrics. Unlike the bioKNN metric, cLISI required cell type annotations,

215    which are often challenging to define in an orthogonal manner. iLISI assumed uniform sample

216    mixing for each cell type was desired, which may not be true in all biological contexts: for

217    example, a cell type might be missing in healthy controls but present in disease samples.

218    batchKLD used the integrated snRNA-seq sample distribution to account for such instances. We

219    calculated all four metrics for each cell.

220

221    **Datasets tested**

222        We tested each of the 58 strategies on 5 different 10x multiome datasets (**Fig. 1g-h**;

223    **Table 3**), for a total of 290 snATAC-seq embeddings. Since the resolution changed by dataset,

224    we denoted cell types, states, and substates collectively as "cell phenotypes." As the resolution

225    increased and cell phenotypes became more similar to each other, it should become harder for

226    these pipelines to distinguish between them. Using each dataset's quality-controlled cell set

227    (**Additional file 2: Fig. S1-5a-b**), we reprocessed the snRNA-seq data from feature selection to

228    UMAP for uniformity (**Additional file 2: Fig. S1-5c-d**; **Methods**). We reassigned gold standard

229    biological cell phenotypes with a PCA + Harmony mRNA clustering resolution that most closely

230    resembled the original author-defined annotations (**Methods**; **Additional file 2: Fig. S1-5e-i**).

231     Within that re-processing, we generated NN graphs for use in our NN metrics for primary

232     analysis. To assess sensitivity of our NN metrics across snRNA-seq integration pipelines, we

233     also generated a Seurat rCCA RNA embedding for a secondary analysis (**Methods**). We

234     excluded samples with fewer than 100 cells since they did not perform well with Seurat rCCA

235     default parameter settings. For primary analysis, we generated NN graphs with 200 cells,

236     though we also used 50 and 100 neighbors to confirm trends.

237     **Table 3.** Dataset summary.

| Benchmarking Dataset | Original Study | Source | Cell Phenotype | Cell Count | Sample Count |
|---|---|---|---|---|---|
| 1 | Weinand et al., Nat Commun, 2024 | RA & OA inflammatory tissue | Broad cell types | 31,547 | 12 |
| 2 | Weinand et al., Nat Commun, 2024 | RA & OA inflammatory tissue | T cell states | 8,069 | 11 |
| 3 | Weinand et al., Nat Commun, 2024 | RA PBMC | T cell substates | 10,669 | 2 runs of 4 pooled samples each |
| 4 | Cheong et al., Cell, 2023 | COVID-19 & healthy PBMC | Broad cell types | 197,360 | 30 |
| 5 | Cheong et al., Cell, 2023 | COVID-19 & healthy PBMC | HSPC states | 27,979 | 28 |

238     RA: rheumatoid arthritis; OA: osteoarthritis; PBMC: peripheral blood mononuclear cell; COVID-
239     19: coronavirus disease 2019; HSPC: hematopoietic stem and progenitor cell
240
241         The first three datasets originated from our previous study[21] investigating rheumatoid

242     arthritis synovial tissue and spanned the gamut from broad cell types ("Dataset 1") to T cell

243     states ("Dataset 2") to T cell substates ("Dataset 3") (**Fig. 1g**). Dataset 1 had 12 samples for a

244     total of 31,547 cells while Dataset 2 totaled 8,069 cells after dropping a sample for low cell

245     count (**Fig. 1h**). One sample in each came from an osteoarthritis patient. Dataset 3 was of

246     particular interest since its cell phenotypes were determined via Fluorescence-Activated Cell

247     Sorting (FACS) surface markers: $CD4^+CD127^-CD25^{hi}$ regulatory T cells (Treg), $CD4^+CD127^-$

248     $CD25^{int}$ Treg, $CD4^+CD25^-PD1^+CXCR5^+$ T Follicular Helper cells (TFH), and $CD4^+CD25^-$

13

249    PD1$^+$CXCR5$^-$ T Peripheral Helper cells (TPH). While mRNA clusters aligned well to the protein

250    hashtags and differential genes were found between the two Treg substates and between the

251    TFH/TPH substates in the original study, we were unable to find many differential promoter

252    peaks[21]. Thus, it should be the most challenging to phenotypically characterize in snATAC-seq

253    embeddings. However, it did have the fewest number of batches with 2 runs of 4 pooled

254    samples each, for a total of 10,669 cells (**Fig. 1h**; **Methods**).

255            The last two datasets came from Cheong et al., 2023[22], where the authors applied a

256    newly developed method for peripheral blood mononuclear cell analysis with progenitor input

257    enrichment (PBMC-PIE). They used this method to profile both mature immune cell types

258    ("Dataset 4") and hematopoietic stem and progenitor cells (HSPCs; "Dataset 5") from 30 PBMC

259    samples across 5 patient stratifications: healthy, non-COVID-19 critical illness, mild COVID-19,

260    early COVID-19, and late COVID-19 (**Methods**; **Fig. 1g-h**). Of note, Dataset 4 stress-tested the

261    time and memory limits for our pipelines as it had almost 200K cells. It was generally more

262    challenging to determine cell phenotypes in Dataset 5 for its 27,979 progenitor cells instead of

263    fully differentiated cell states.

264

265    **Benchmarking 58 different pipelines**

266            Using our command-line tool, we ran each of the 58 pipelines on the five datasets and

267    quantified metrics for bio-conservation and batch correction (**Additional file 2: Fig. S6**).

268            To illustrate results, we focused on Dataset 5 (**Fig. 2a**). We observed specific pipeline

269    choices affected bio-conservation and batch correction. Generally, the best-performing pipelines

270    for this dataset used SnapATAC2 or itLSI. GAS features predominately resulted in poor bio-

271    conservation while the inclusion of Harmony usually improved batch integration. For example,

272    SnapATAC2 + Harmony with peaks generally performed very well (**Fig. 2b**), with distinct cell

273    states and diverse samples that corresponded well with the snRNA-seq data (**Fig. 1g**;

274    **Additional file 2: Fig. S5c**). However, replacing peak features with archrGAS features resulted

275 in a similar mean batchKLD metric, but a worse mean bioKNN metric (**Fig. 2a**); this was

276 reflected in the UMAP where cell states were mixed together despite comparable sample

277 diversity (**Fig. 2c**). Conversely, using SnapATAC2 with peaks but removing the Harmony

278 correction step resulted in a worse mean batchKLD metric (**Fig. 2a**); in the UMAP, there were

279 groups of cells originating from a single sample (**Fig. 2d**) that had correspondingly high

280 batchKLD metrics (**Fig. 2e**) since the snRNA-seq data did not contain these singular sample

281 regions (**Additional file 2: Fig. S5c**). Including Harmony largely did not change the bio-

282 conservation metric, with the best bioKNN values in the smaller cell states, BEM and GMP (**Fig.**
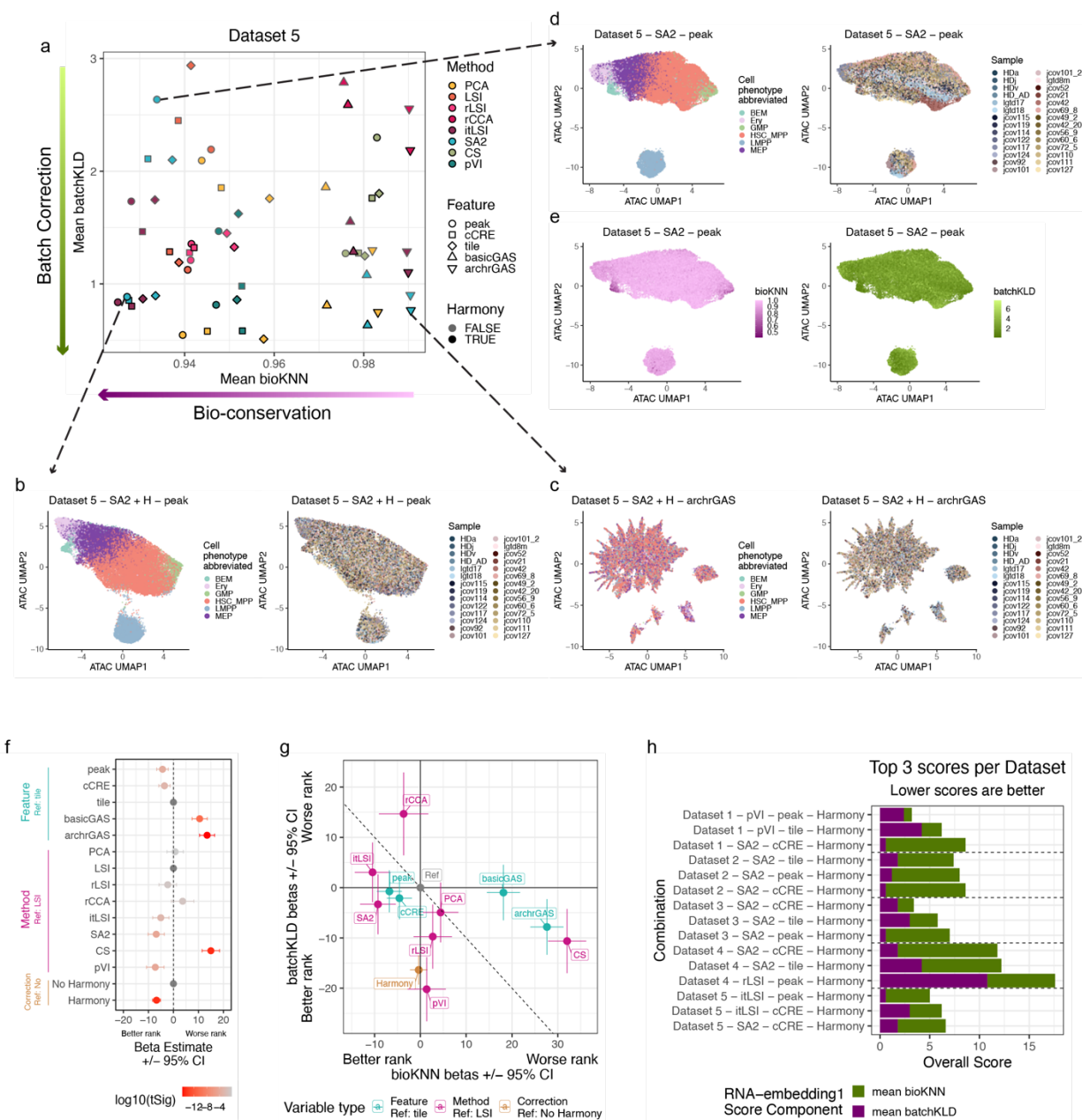
283 **2e**).

**Fig. 2.** Pipeline rankings using NN metrics with RNA-embedding1.
**a.** Mean bio-conservation (**x-axis**; bioKNN) and batch correction (**y-axis**; batchKLD) metrics averaged across Dataset 5 cells. **b.-d.** snATAC-seq UMAPs colored by cell phenotype (**left**) and sample (**right**) for Dataset 5 cells using **b.** SnapATAC2 + Harmony with peak features, **c.** SnapATAC2 + Harmony with archrGAS features, **d.** SnapATAC2 with peak features. **e.** snATAC-seq UMAPs colored by bioKNN (**left**) and batchKLD (**right**) for Dataset 5 cells using SnapATAC2 with peak features. In both cases, lower values (darker colors) are desired. **f.** Combined linear model relating overall score (60% ranked mean bioKNN + 40% ranked mean batchKLD) to dataset, feature, method, and correction combinations (**Methods**). **g.** Separated linear models as in **f.**, but ranked mean bioKNN on x-axis and ranked mean batchKLD on y-axis (**Methods**). **h.** Top 3 overall scores per dataset, colored by ranked mean bioKNN and ranked mean batchKLD components.

16

**General trends across five datasets**

To understand the overall trends, we examined the mean bio-conservation and batch correction metrics across all 58 pipelines and 5 datasets. For each dataset, we ranked performance (1 best; 58 worst) for each of the pipelines and examined a composite rank of bio-conservation and batch correction, weighted by 60% and 40%, respectively[15] (**Methods**); we also ranked the individual metrics. We assessed overall performance using a multivariate linear model, where we determined how features, methods, and embedding correction affected the ranked mean scores across the five datasets[15] (**Fig. 2f**). We denoted our primary analysis focusing on NN metrics utilizing 200 NN within an embedding generated from snRNA-seq datasets using PCA with Harmony correction, as RNA-embedding1. We also did secondary analyses with NN metrics using Seurat gene embeddings and 200 NN ('RNA-embedding2') (**Additional file 2: Fig. S7a**) and LISI metrics (**Additional file 2: Fig. S7b**).

First, we found that feature set ($p<0.00227$) impacted performance in the following order: archrGAS (worst; beta= +13.45 rank), basicGAS, tile (reference rank=0), cCRE, peak (best; beta= -4.36). Second, across datasets and pipelines, the choice of method nominally affected rank performance ($p<0.12$). We observed that rLSI, itLSI, SnapATAC2, and PeakVI improved rank performance when compared to the reference LSI (beta= -2.27, -5.05, -6.89, -7.26, respectively). Globally, SnapATAC2 was hurt by the GAS features that were not included in the PeakVI feature sets tested. In contrast, PCA had little impact on the rank (beta= +0.67) while rCCA had worsened ranks (beta= +3.67). CellSpace had the worst rank performance and the biggest effect size across methods (beta= +14.97). Third, the decision to include Harmony as an additional batch correction step generally improved the overall rank performance ($p=3.36e-15$; beta= -6.77). We observed similar patterns with NN metrics with RNA-embedding2 and LISI metrics (**Additional file 2: Fig. S7a-b**). A notable exception was for Signac rLSI, where it only became significantly different from the LSI reference in the NN metrics using RNA-embedding2 defined from Seurat rCCA ($p=0.00025$) compared to the RNA-embedding1 ($p=0.20$) and LISI

324    (p=0.16) metrics, suggesting that the choice of RNA embedding may have a subtle impact.

325    Seurat rCCA results also changed minimally across metric types, but their ranks were already

326    penalized by the worst-preforming GAS features.

327         To understand the effect of pipeline choices on bio-conservation and batch integration

328    separately, we applied the same approach, but testing a linear model for each metric (**Methods**;

329    **Fig. 2g**; **Additional file 2: Fig. S7c-d**). We observed that the choice of feature set had the

330    greatest effect on bio-conservation metrics (p<0.0011). SnapATAC2 was the only method to

331    achieve better rankings for both bio-conservation and batch correction for both NN metrics.

332    Unsurprisingly, Harmony primarily improved batch correction metrics (p=1.09e-23), as described

333    in the Dataset 5 analysis above (**Fig. 2b,d**).

334

335    **ATAC-specific features surpassed Gene Activity Scores**

336         One of the most striking results of this study was that ATAC-specific features, rather

337    than GAS scores, generally improved performance across pipelines and metrics (**Fig. 2f**;

338    **Additional file 2: Fig. S7a-b**). None of the top 3 combinations per dataset included the GAS

339    features (**Fig. 2h**; **Additional file 2: Fig. S7e-f**). Indeed, using the RNA-embedding1 NN

340    metrics, the best pipeline containing GAS features was Dataset 2 using itLSI + Harmony with

341    basicGAS features (score=25.2); it did not surpass the first quartile of scores for that dataset

342    (score=20.45; **Additional file 1: Table S2**). Consistent with this observation, the corresponding

343    UMAPs illustrated inappropriate mixing between the two primarily CD8+ T cell populations (T-1

344    and T-5) and between the TFH/TPH and Treg populations (T-2 and T-4) while showing pockets

345    of cells segregated from two samples (**Fig. 3a, left**).

346         Amongst ATAC features, we generally observed that peaks and cCREs yielded the best

347    overall performance, followed closely tiles (**Fig. 2f**). This pattern was generally consistent

348    across methods (**Fig. 3b**; **Additional file 2: Fig. S8a-b**). For the above Dataset 2 itLSI +

349    Harmony pipeline, simply switching the basicGAS features for peak features both improved the

18

350    score (12.4) and resulted in more visually distinct T cell states (**Fig. 3a, right**). We note that

351    peaks are the most commonly used snATAC-seq feature, and most of these methods were

352    likely developed assuming peak features. That being said, the majority of dataset-specific peaks

353    (mean=70%) overlapped at least one cCRE (**Additional file 2: Fig. S8c**), perhaps accounting

354    for the generally small differences in metric rankings between them. While tiles overlapped with

355    all other features since tiles accounted for the whole genome, they were not centered at the

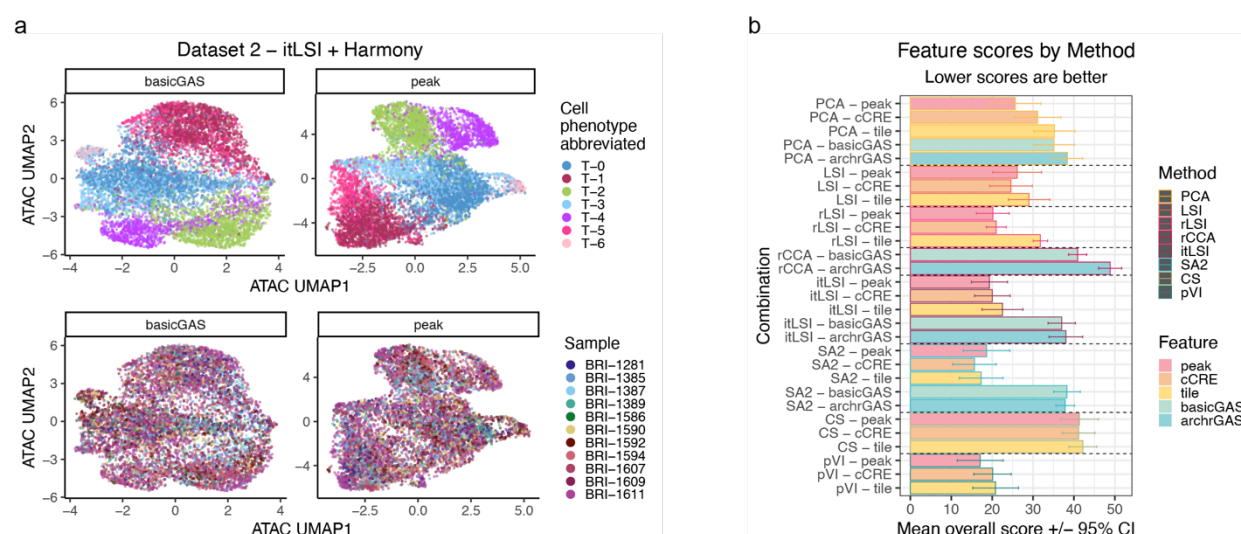356    summit of accessibility unlike the peaks (**Methods**).



357
358    **Fig. 3.** ATAC features (peaks, cCRE, tiles) outperformed GAS features (basicGAS, archrGAS).
359    **a.** Best performing GAS feature pipeline (**left**): Dataset 2 using basicGAS features and the itLSI
360    method with Harmony correction. Substituting peak features within that pipeline (**right**).
361    Corresponding UMAPs colored by cell phenotype (**top**) and sample (**bottom**). **b.** The mean
362    overall scores (60% bio-conservation bioKNN rank + 40% batch correction batchKLD rank)
363    across datasets and correction choice for each method and feature combination. Lower scores
364    are better. Error bars are the 95% confidence interval (CI).
365

366    **SnapATAC2 performed best using NN metrics**

367         Using the RNA-embedding1 NN metrics, SnapATAC2 had the best-ranked combination

368    for 3/5 datasets, and was among the top 3 for the other 2 datasets (**Fig. 2h**). Its success was

369    most apparent in the bio-conservation metrics (**Fig. 2g**). While all ATAC features performed well

370    for SnapATAC2, it did best with cCREs (**Fig. 3b**, **4**; **Additional file 1: Table S2**). SnapATAC2

371    could handle GAS features with limited modification (**Methods**), but as with all other methods,

19

372 the GAS features hurt performance (**Additional file 2: Fig. S6a**; **Additional file 1: Table S2**).

373 SnapATAC2 + Harmony with ATAC features generally performed well when using the RNA-

374 embedding2 NN metrics as well (**Additional file 2: Fig. S6b**).
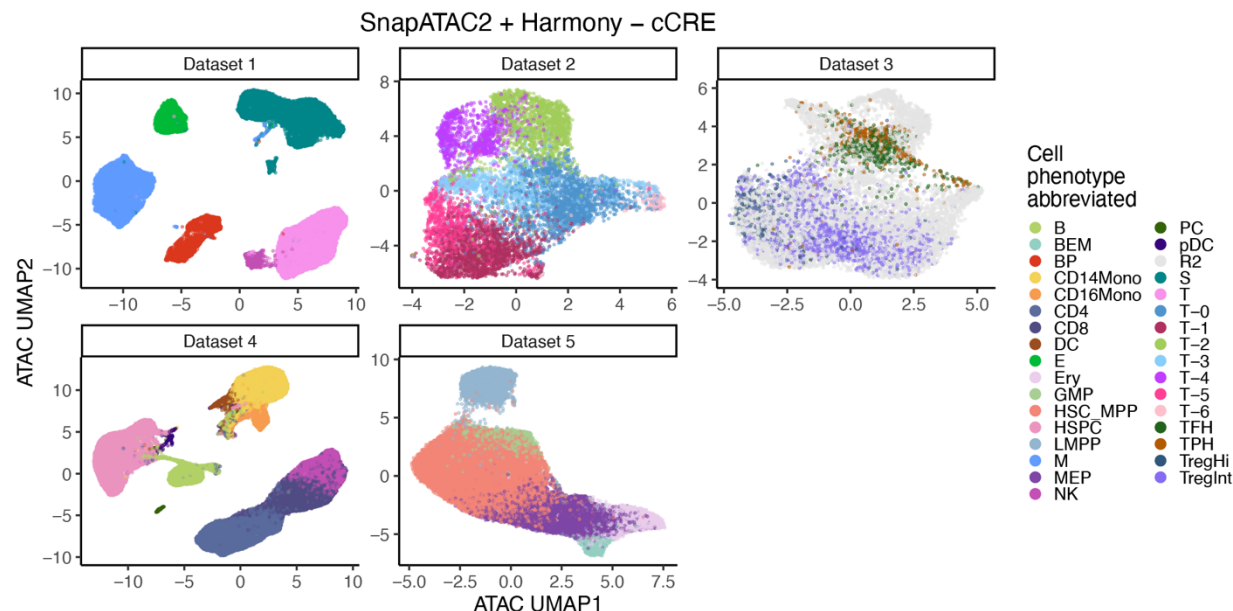


375
376 **Fig. 4.** SnapATAC2 preformed best for NN metrics with RNA-embedding1.
377 snATAC-seq UMAPs colored by cell phenotype for all Benchmarking Datasets using
378 SnapATAC2 + Harmony with cCRE features.
379

380 **PeakVI performed well in the rheumatoid arthritis datasets, but biased by fragment**

381 **counts**

382　　　　In every dataset, we observed that PeakVI + Harmony used with peaks outperformed all

383 other PeakVI pipelines via the RNA-embedding1 NN metrics (**Additional file 1: Table S2**). In

384 fact, this combination was the best-performing pipeline for Dataset 1 (**Fig. 2h**). In the

385 corresponding UMAP, NK cells separated from T cells and mural cells segregated from stromal

386 cells in the lower tail (**Fig. 5a**). It also performed well for Datasets 2 and 3, primarily owing to

387 effective batch correction (**Fig. 2g**; **Additional file 1: Table S2**). However, PeakVI's

388 performance decreased in the COVID-19 datasets, due to worse bio-conservation rankings

389 (**Additional file 1: Table S2**). PeakVI struggled to distinguish pDCs in Dataset 4 and BEM in

20

390    Dataset 5 (**Fig. 5a**). It also displayed some mixed HSPC states for the late COVID-19 samples

391    in Dataset 4 (**Fig. 5b**).

392         One potential reason that PeakVI performance was suboptimal was confounding by

393    fragment count per cell. We noticed that PeakVI often grouped high fragment count cells

394    together in the same region of the UMAP (**Fig. 5c**). Quantifying this observation, the high

395    fragment counts correlated more with worse bioKNN metrics, corresponding to a worse bio-

396    conservation ranking (**Additional file 2: Fig. S9**). Indeed, in Dataset 3, we saw a group of T

397    cells with higher fragment counts and mixed substates, illustrating poor bio-conservation (**Fig.**

398    **5a,c**). This phenomenon was also seen in the pipelines using PeakVI with peak features and

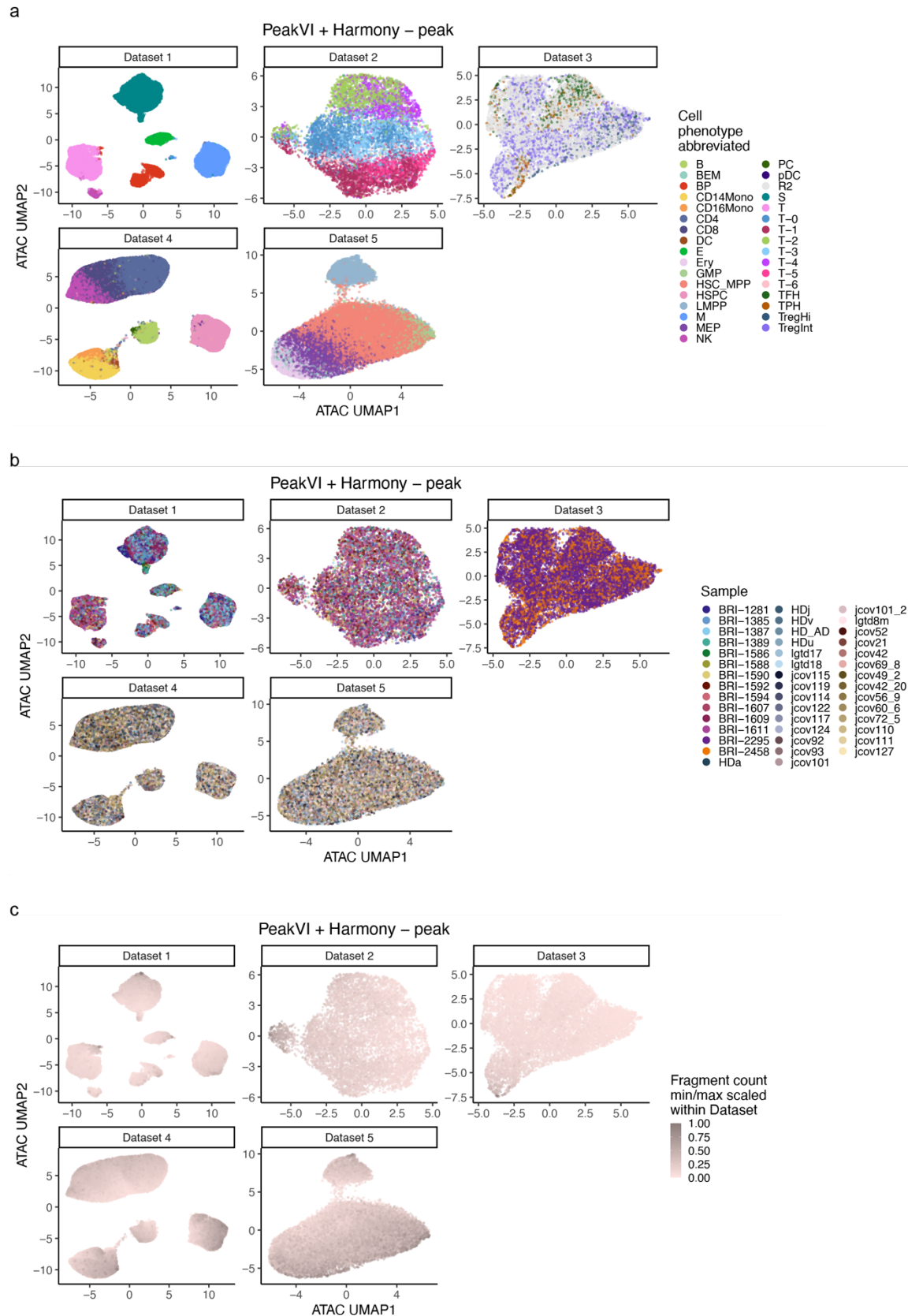399    PeakVI + Harmony with cCRE features (**Additional file 2: Fig. S10a-b**).

**Fig. 5.** PeakVI biased by fragment count.

402  snATAC-seq UMAPs colored by **a.** cell phenotype, **b.** sample, and **c.** fragment count for all
403  Benchmarking Datasets using PeakVI + Harmony with peak features.
404

405  **ArchR itLSI favored cell states over cell types**

406      Using either peaks or cCRE, ArchR's iterative LSI + Harmony ranked best for Dataset 5

407  in the RNA-embedding1 NN metrics (score=5, 6.2, respectively; **Fig. 2h**), and demonstrated

408  good performance with the other metrics as well (**Additional file 2: Fig. S7e-f**). They both

409  showed very good separation between the LMPP cluster and the rest on the UMAPs with

410  generally good batch correction (**Additional file 2: Fig. S11a**). itLSI's best scores were well

411  within the top quartile for the other sub-cell-type resolution datasets: 12.6 for Dataset 2 (peaks

412  with Harmony) and 12.4 for Dataset 3 (tiles with Harmony), out of a theoretical range between 1

413  and 58 within datasets (**Additional file 2: Fig. S11b**).

414      However, itLSI performed worse on Dataset 1 and Dataset 4, which had easier to

415  distinguish broad cell types. itLSI + Harmony with cCRE for Dataset 1 was slightly better than

416  the top quartile of scores (score=18.6 vs 19.25) while Dataset 4 using itLSI + Harmony with

417  peaks did not exceed the top quartile (score=24.6 vs 23.85**; Additional file 2: Fig. S11b**). In

418  these datasets, itLSI's performance was hurt by poor batch effect correction (**Additional file 1:**

419  **Table S2**). Indeed, Dataset 4's top pipeline had areas of its UMAPs dominated by single

420  samples, most notably jcov93 within CD4 T & B cells and lgtd18 within T/NK & HSPC cells

421  (**Additional file 2: Fig. S11c**, **left**). Interestingly, the un-Harmonized version scored similarly

422  (score=25.2) with both pipelines having inappropriate mixing of the myeloid, HSPC, B, and pDC

423  cell types in the late COVID-19 samples (**Additional file 2: Fig. S11c**).

424

425  **CellSpace performed poorly across datasets, features, correction, and metrics**

426      CellSpace was the worst-ranked integration method across all metrics (**Fig. 2f**;

427  **Additional file 2: Fig. S7a-b**). The best RNA-embedding1 overall ranking CellSpace achieved

428  was for Dataset 4, peaks, without Harmony (score=32.2; **Additional file 1: Table S2**). However,

429    even with distinct cell types in Dataset 4, CellSpace remained as a singular entity in the UMAP,

430    with intermixed domains for CD4 T + CD8 T + NK cells, B + Plasma cells, CD14 Monocytes +

431    CD16 Monocytes + Dendritic cells, HSPC + plasmacytoid dendritic cells (**Additional file 2: Fig.**

432    **S12a**). However, it was very well integrated at the sample level even without Harmony

433    correction (**Additional file 2: Fig. S12b**), suggesting an over-integration to the point of unclear

434    cell types.

435

436    **Harmony improves batch correction for most methods**

437        As Harmony was created to correct batch structure while maintaining biological identity,

438    we expected that it would improve batch correction rankings while minimally affecting bio-

439    conservation rankings (**Fig. 2g**). Indeed, we observed that using Harmony generally improved

440    most methods (**Fig. 6a; Additional file 2: Fig. S13**). The methods that did not inherently

441    address batch effects, PCA, LSI, SnapATAC2, generally benefitted the most. rLSI, rCCA, and

442    PeakVI, which all explicitly corrected for batch, had a more modest benefit with Harmony. The

443    two methods with some implicit batch correction, CellSpace and itLSI, were the most likely to

444    have adverse effects post-Harmony. CellSpace's poor interaction with Harmony was seen

445    mostly clearly in cell state Datasets 2 and 5, across all ATAC features, with peaks as the most

446    extreme (**Fig. 6a**). However, we saw batch effects for CellSpace both with and without Harmony

447    for Dataset 2 peaks (**Fig. 6b**). Iterative LSI was primarily negatively affected by Harmony in the

448    cell type datasets for GAS features (**Fig. 6a**). This was typified by Dataset 1 when using itLSI +

449    Harmony with basicGAS features, where stronger batch effects were seen in the post-Harmony

450    UMAPs (**Fig. 6c**). However, when itLSI was used with the non-GAS features, including

451    Harmony resulted in similar or better batch correction (**Fig. 6a**).
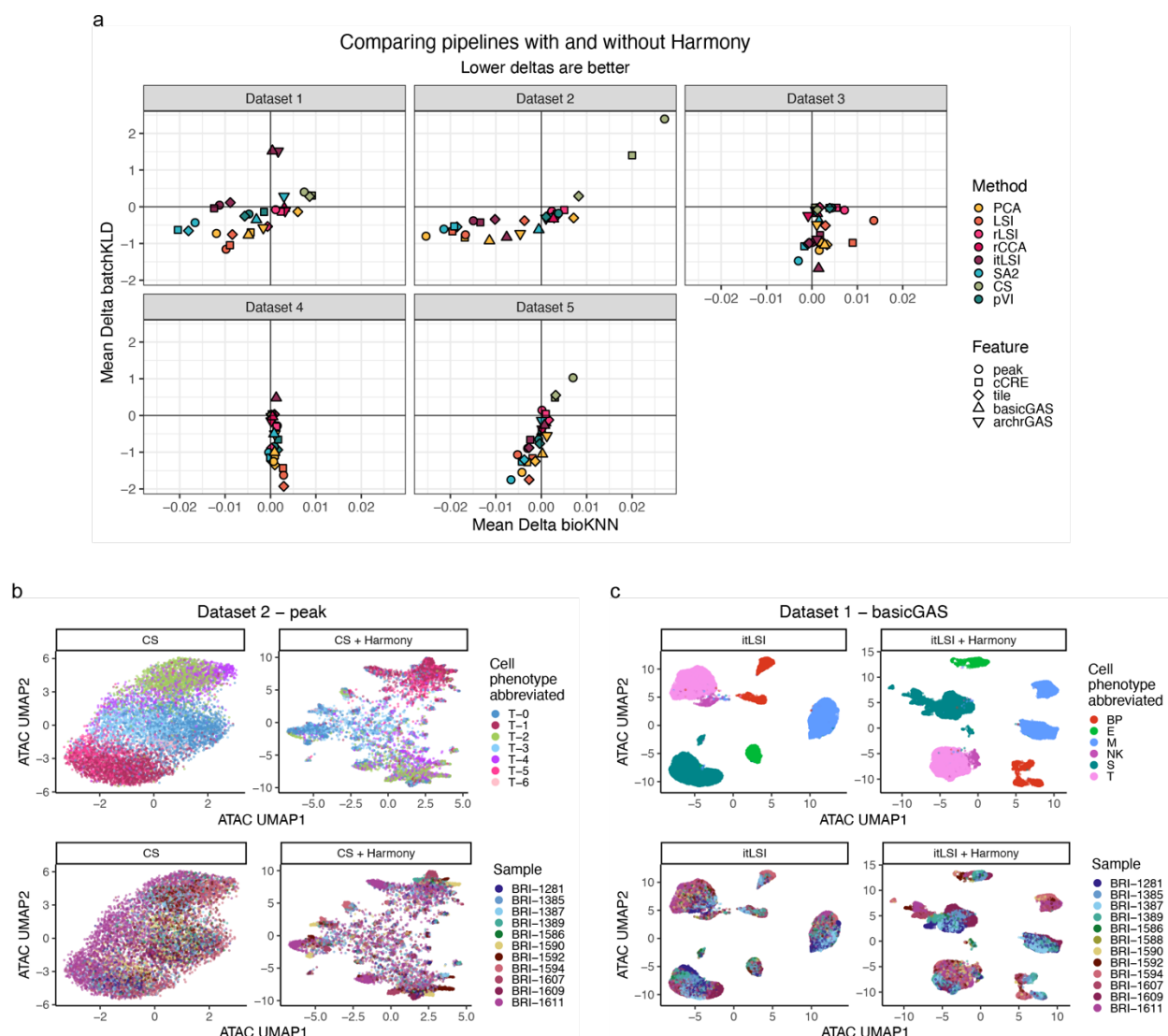
452
453 **Fig. 6.** Harmony generally improves batch correction.
454 **a.** Mean bio-conservation (**x-axis**; bioKNN) and batch correction (**y-axis**; batchKLD) Harmony
455 deltas averaged across method/feature combinations in all Benchmarking Datasets. Harmony
456 deltas were calculated as the Harmony metrics subtracted by the no Harmony metrics for the
457 same method/feature combination. Lower deltas were better. **b.** snATAC-seq UMAPs colored
458 by cell phenotype (**top**) and sample (**bottom**) for Dataset 2 with peak features using CellSpace
459 (**left**) or CellSpace + Harmony (**right**). **c.** snATAC-seq UMAPs colored by cell phenotype (**top**)
460 and sample (**bottom**) for Dataset 1 with basicGAS features using ArchR itLSI (**left**) or ArchR
461 itLSI + Harmony (**right**).
462

463 **Job requirements varied greatly by method**

464        We assessed time and memory requirements for both the overall pipeline job as well as

465 the individual steps spanning pre-processing, embedding generation, Harmony correction, post-

466 processing, and metric/visualization calculation (**Methods**; **Fig. 7a**; **Additional file 2: Fig. S14-**

467     **15**). For both time and memory, we found more variation across methods than by input feature.

468     One of the most resource intensive pipelines was rLSI with tiles, requiring 121.6 CPU-hours and

469     545.7 GB for the ~200K cells of Dataset 4. PeakVI and CellSpace, two nonlinear methods, also

470     required more time to complete than most. In contrast, the third nonlinear method tested,

471     SnapATAC2, was among the most time and memory efficient. LSI, rCCA, and itLSI were also

472     among the fastest. All methods, excluding rLSI, had roughly similar memory requirements. The

473     use of tiles increased memory usage, likely due to the order of magnitude more features than

474     peaks, cCREs, or GASs. The feature selection step we implemented with CellSpace and PCA

475     (**Methods**) took more memory and often more time than generating the embedding, so a more

476     straightforward feature count cutoff as those used in some of the other methods would decrease

477     the overall job requirements (**Additional file 2: Fig. S14-15**). Also, the metric/visualization steps

478     for Dataset 4 were of high memory and time burden (**Additional file 2: Fig. S14-15**), with the

479     former not necessary to most end-users. Of note, since ArchR required a project to be created

480     before our pipeline could be implemented within it, it was not included in the aggregated plot. It

481     maxed out at 617.2 CPU-min and 82.8 GB for Dataset 4 (**Fig. 7b**), thus adding a small, but not

482     negligible, addition to itLSI's requirements. In all cases, Harmony correction was extremely

483     efficient with a minimal impact on both time and memory requirements.
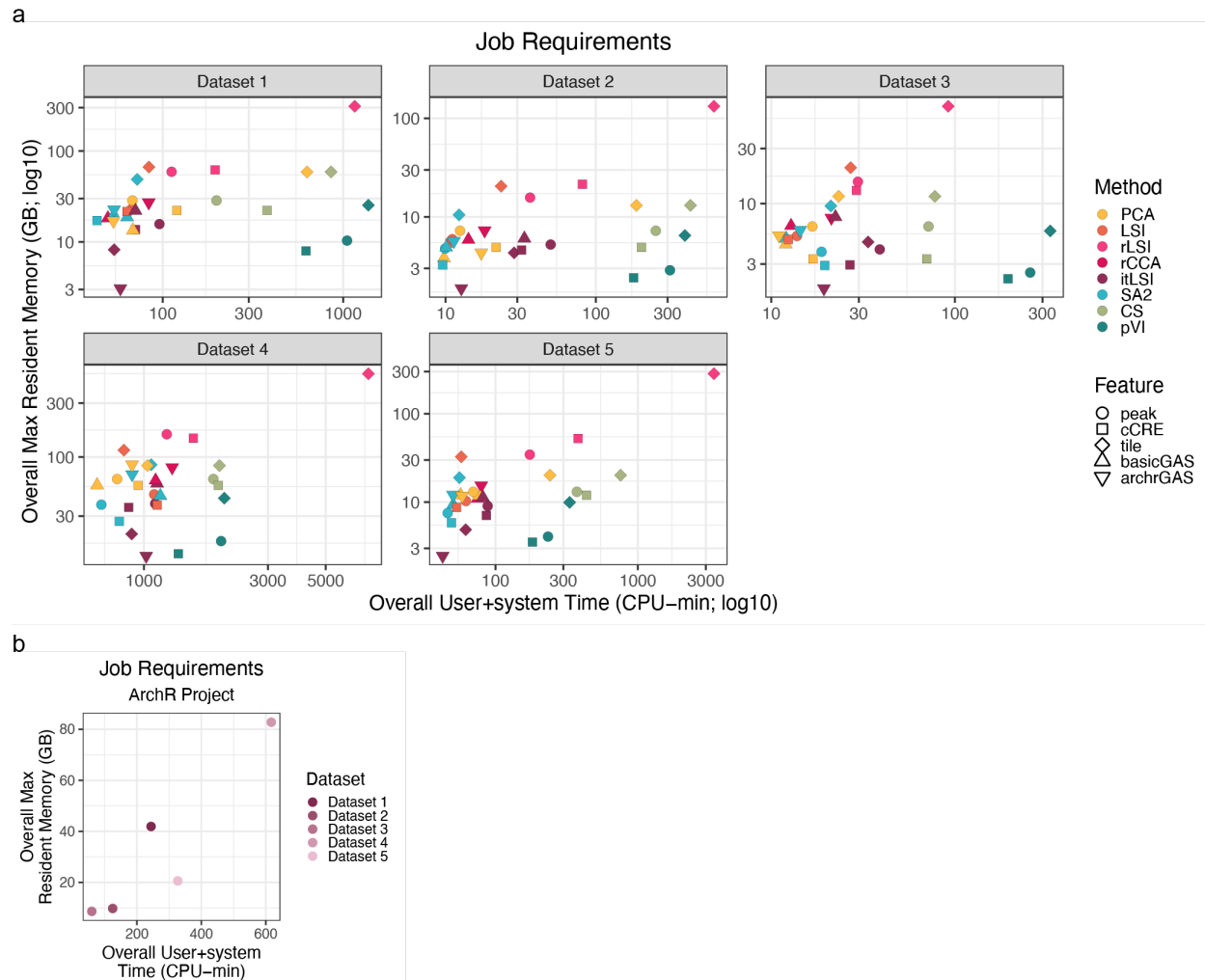
**Fig. 7.** Time and memory requirements vary by method.
Time (**x-axis**; user+system time) and memory (**y-axis**; max resident memory) requirements for
**a.** entire pipelines and **b.** ArchR project creation for all Benchmarking Datasets. Values
calculated with /usr/bin/time (**Methods**).

**NN metrics better for rarer cell phenotypes dominated by fewer samples**

The two RNA embeddings used with the NN metrics produced very concordant rankings,

with a mean bioKNN correlation of 0.97 and a mean batchKLD correlation of 0.91 (**Additional**

**file 2: Fig. S16a-b**). The only notable slight deviation was in Dataset 3 batchKLD metrics, where

RNA-embedding2, generated using the rCCA method, favored rCCA/rLSI snATAC-seq

embeddings. This likely explains the rLSI discrepancy found with the combined linear model

(**Fig. 2f**; **Additional file 2: Fig. S7a-b**).

27

497    The RNA-embedding1 NN metrics and LISI metrics were generally concordant and

498    prioritized pipelines in a similar fashion (**Additional file 2: Fig. S17a-b**). Mean bio-conversation

499    metric correlation across datasets was 0.92; batch correction metrics were slightly less

500    correlated at a mean of 0.75.

501    However, there were a few discrepancies. The most dissimilar pipeline in the ranked

502    batch correction metrics was rCCA with GAS features, with or without Harmony (**Additional file**

503    **2: Fig. S17b**). These pipelines performed well using iLISI, but demonstrated worse performance

504    when assessed with batchKLD. We note that these pipelines were also among the worse-

505    performers in bio-conservation metrics (**Additional file 2: Fig. S17a**), as exemplified by Dataset

506    5 rCCA with archrGAS features with indistinguishable HSPC states in addition to the very well-

507    mixed samples (**Additional file 2: Fig. S17c-d**). This suggested that these batch-aware rCCA

508    pipelines were likely over-integrating towards iLISI's uniform sample mixing beyond the local

509    sample distribution expected by the gold standard snRNA-seq embedding, measured by

510    batchKLD.

511    Furthermore, we saw some interesting examples of cells with good batchKLD values

512    and poor iLISI values, as in the SnapATAC2 + Harmony with cCREs pipeline (**Additional file 2:**

513    **Fig. S18a**). The most notable disparity was in Dataset 4 (**Additional file 2: Fig. S18b**), where

514    the plasma cells were primarily derived from a single sample (**Additional file 2: Fig. S18c**),

515    seen in the UMAPs for both ATAC (**Fig. 4a**; **Additional file 2: Fig. S18d**) and RNA (**Fig. 1g**;

516    **Additional file 2: Fig. S4c**). In this case of confounded batch and biology, a good result should

517    not have uniform batch mixing, meaning iLISI was an imperfect measure. However, since the

518    gold standard RNA embedding also grouped plasma cells despite belonging to primarily one

519    sample, the sample distributions of the plasma cells' RNA NN were similar to those of the ATAC

520    NN, as signified by a good batchKLD metric. We saw a similar effect with erythroid progenitors

521    in Dataset 5 (**Additional file 2: Fig. S18e-f**). Therefore, we concluded that batchKLD served as

522     a better metric in cases with rarer cell types that are dominated by fewer samples, which can be

523     where the most interesting biology resides.

524

525     **Discussion**

526     In this study, we extensively benchmarked 58 snATAC-seq integration pipelines across

527     5 datasets, 5 features, 7 methods, 1 embedding correction, and 4 metrics. The performance

528     analysis of these pipelines gave insights into the most effective strategies for single cell

529     chromatin accessibility integration. Based on this, we recommend using ATAC features, such as

530     peaks and cCREs, as opposed to GAS features, which force chromatin accessibility data to look

531     like genes. Furthermore, using Harmony a batch correction step was usually very helpful. We

532     found SnapATAC2 to be the best-performing method in general, though other methods came

533     close in performance, such as PeakVI and ArchR's itLSI. PCA and LSI were standard methods

534     that performed reasonably well in our benchmark. Our overall best-performing pipeline was

535     SnapATAC2 + Harmony with cCRE, which we recommend for most purposes. In addition to

536     high cross-sample data integration, it had limited time and memory investments, and cCREs

537     render this pipeline to be easily generalizable.

538     Our study offers a powerful and effective benchmarking strategy. By utilizing multimodal

539     datasets, we could use snRNA-seq embeddings as a standard by which to benchmark snATAC-

540     seq integration. To this end, we introduced two new metrics. We assessed bio-conservation by

541     comparing barcodes for the same cells across modalities in our bioKNN metric; this approach

542     negated the need for potentially imprecise per-cell phenotypic annotations usually defined by

543     functional clusters with hard borders. We measured batch correction using batchKLD, which

544     compared per-cell local sample distributions between modalities and could account for

545     instances where biological states were not uniformly mixed across samples. We note that while

546     RNA embeddings are generally considered to represent key biological features, they can

547     themselves be variable. Thus, we tested two independent embeddings built with two popular

29

548     approaches within our NN metrics. Our conclusions were generally consistent across RNA

549     embedding choices (**Additional file 2: Fig. S16**). Additionally, we acknowledge that RNA

550     embeddings may be imperfect, and that if there is shared technical variation across modalities,

551     our strategy would not account for it. However, our metrics were also largely consistent with

552     more standard local metrics, cLISI and iLISI (**Additional file 2: Fig. S17a-b**). In certain

553     instances, those metrics may lead to erroneous conclusions where biological states and

554     technical batches co-occur together. For example, the plasma cells in Dataset 4 largely

555     originated from one sample, which iLISI penalized for insufficient mixing, while lsKLD did not as

556     it was validated in the snRNA-seq data (**Additional file 2: Fig. S18b-d**).

557         Computational analysis of chromatin accessibility data has a fundamental challenge:

558     there is no consistent, reproducible feature set used within and across datasets. Given a new

559     dataset, users often make many specific methodological and parameter choices to call peaks.

560     Reads that do not overlap the final peak set are often discarded, making it possible that the

561     downstream biological conclusions are sensitive to these choices. In terms of data integration,

562     peaks called in one dataset are not easily transferrable across datasets; hence, peaks often

563     need to be recalled and datasets reprocessed. A key finding of this study was the discovery that

564     peaks and ENCODE cCREs achieved similar performance (**Fig. 2f**, **3b**). This suggested that

565     dataset-specific peaks offer little integration advantage over the dataset-agnostic cCREs.

566     Notably, even though there were no ENCODE DNase-seq datasets assayed with RA synovial

567     tissue or COVID-19 blood, the cCREs retained enough information to identify and integrate our

568     RA and COVID-19 Benchmarking Datasets. Based on this observation, we propose that these

569     pre-defined ENCODE cCREs could be used as a gene-like reference set for future snATAC-seq

570     studies. This would enable easy and reproducible integration with a common feature set, rather

571     than defining bespoke feature sets for each unique application. Furthermore, the deep

572     characterization ENCODE has done of these cCREs provides a rich resource for understanding

573     the functional implications of open chromatin at that locus within the individual datasets. Similar

574     to gene annotation versions, the cCRE annotations could be periodically updated with additional

575     functional elements as datasets continue to increase. We note the possibility that there may be

576     specific biological instances where dataset-specific peaks capture biology missed by ENCODE

577     cCREs.

578         While we saw general trends across datasets, there were some dataset-specific caveats

579     when selecting appropriate pipelines. For example, datasets are routinely becoming even larger

580     than the 200K cells in Dataset 4; therefore, resource-intensive pipelines like rLSI may become

581     intractable for such datasets. Datasets with uneven fragment counts may lead to biased PeakVI

582     embeddings, particularly for the finer-grain cell substates whose cell phenotype variation is

583     smaller. Those cell state and substate datasets may prefer itLSI while broad cell type datasets

584     may want to avoid it.

585         We implemented our benchmarking strategy in an easy-to-use command line tool that

586     can be adapted to any dataset of interest, enabling easy benchmarking or inference. After

587     generating the input matrices with the provided feature files and scripts, a user only needs to

588     specify the dataset, feature, and method combinations to get a file of commands that can be run

589     locally or via a computing cluster. All the embeddings, NN graphs, metrics, and UMAPs

590     discussed here were generated in this way.

591         Our study has some limitations. First, we assume that the underlying cell phenotype will

592     be the same assayed via snRNA-seq and snATAC-seq. In our experience[21], we have generally

593     observed that cell states that are similar in transcriptional profiles also capture similar biology in

594     chromatin profiles. However, there may be important biological functions captured in one

595     modality, but missed by the other. For example, cell cycle usually affects gene expression more

596     dramatically than chromatin accessibility[30]. Second, our conclusions are based on a limited

597     number of datasets. We chose these datasets to span cell types, states, and substates across

598     both blood and tissue, and we generally observed consistent results across them. However,

599     there were situations where specific pipelines performed slightly better or worse, as discussed

600    above. We also restricted our analysis to the 10x multiome[TM] platform as it is among the most

601    popular and commercially available, but it is possible that alternative platforms like sci-CAR[31] or

602    SHARE-seq[30] may have different results. For the interested reader, we provide the code to

603    apply our framework to other contexts and datasets. Third, we note that our study focused on

604    integration performance. We were unable to quantify factors like software installation, data

605    structure preparation, and ease-of-use per method in an objective fashion. These requirements

606    varied dramatically across pipelines and may be an important factor in user choice.

607         We hope this benchmarking study will assist researchers decide which combination of

608    feature, method, and correction they will apply to their future snATAC-seq datasets.

609

## Conclusions

611    In conclusion, we benchmarked 58 snATAC-seq integration pipelines across 5 datasets utilizing

612    2 novel multimodal-guided metrics. Using our command-line tool to process these benchmarks,

613    we determined that SnapATAC2 + Harmony using cCRE features outperformed other pipelines

614    while also being generalizable and resource-efficient.

615

616

## Methods

**Benchmarking Metrics.**

619    *Bio-conservation K Nearest Neighbors (bioKNN).* To define a gold standard nearest-neighbor

620    graph, we used the snRNA-seq sample-harmonized[20] PCs (RNA-embedding1) as described in

621    the **Benchmarking Datasets** section. We also used Seurat[29] version 5 (RNA-embedding2) to

622    test our strategy's sensitivity to snRNA-seq embedding choice. We used Seurat functions:

623    CreateSeuratObject, split by sample, NormalizeData, FindVariableFeatures, ScaleData,

624    RunPCA, and IntegrateLayers with method=CCAIntegration.

625    In both cases, we calculated the NN graph from the resulting embedding matrix using

626    RANN::nn2 with K=50, 100, 200 NN. To evaluate each snATAC-seq integration strategy (see

627    **Integration Methods** section), we calculated the K=50, 100, 200 NN as done with the snRNA-

628    seq modality and evaluated what fraction ($bioKNN$) of the total NN tested ($K$) were not

629    represented in both modalities for each cell ($NN_{shared}$). We wanted lower values to be desirable

630    as in the batchKLD metric.

631
$$bioKNN = \frac{K - NN_{shared}}{K}$$

632

633    *Batch Kullback-Leibler Divergence (batchKLD).* To determine how well the samples integrated

634    in the batch-corrected snATAC-seq NN graph, we compared the sample ($x$) distribution in the

635    ATAC top K NN per cell ($Q$) to that of the gold standard batch-corrected snRNA-seq top K NN

636    for the same cell ($P$) using KL divergence (**Fig. 1d-e**).

637
$$D_{KL}(P \,||\, Q) = \sum_{x \in X} P(x) \log_2 \frac{P(x)}{Q(x)}$$

638
$$Q(x) = \frac{count(x)}{\sum_{x \in X} count(x)}$$

639    To avoid dividing by 0 and to ensure $KL = 0$ when $P = Q$, we added a pseudocount to $P$ and $Q$.

640    Instead of a uniform pseudocount, we added a dataset-driven pseudocount that accounted for

641    the overall sample distribution of all the cells ($A$) per dataset, scaled to the equivalent of one

642    additional cell.

643
$$Q_{ps}(x) = \frac{count(x) + A(x)}{\sum_{x \in X}(count(x) + A(x))}$$

644
$$= \frac{count(x) + A(x)}{\sum_{x \in X} count(x) + \sum_{x \in X} A(x)}$$

645
$$= \frac{count(x)}{\sum_{x \in X} count(x) + 1} + \frac{A(x)}{\sum_{x \in X} count(x) + 1}$$

33

646     We used the top K NN graphs as described in the *bioKNN* section. We used philentropy::KL to

647     calculate KL divergence; with the addition of the dataset-driven pseudocount, we did not require

648     the epsilon condition of the method, where the small epsilon value was added to 0s in the $Q$

649     distribution.

650

651     *LISI.* Local Inverse Simpson Index (LISI) scores measured the effective number of different

652     categories of a covariate ($C$) represented in the local neighborhood of each cell ($p_i$, proportional

653     abundances).

654
$$LISI_C = \frac{1}{\sum_{i=1}^{C} p_i^2}$$

655     We calculated the LISI scores via lisi::compute_lisi for each embedding matrix for both cell

656     phenotype (cLISI) and sample (iLISI). For all non-substate datasets, the cell type was defined

657     using the multimodal snRNA-seq as described in the **Benchmarking Datasets** section. For the

658     substate dataset (*Dataset 3*), we calculated the cLISI for only the run that had FACS-defined

659     cell types by subsetting both the embedding and metadata matrices to those cells before

660     calculations.

661

662     **ATAC Feature Sets and Matrices.** We tested four different genomic feature types: (1) peaks,

663     (2) tiles, (3) cCRE, and (4) two different GAS (**Fig. 1b**; **Table 1**; **Additional file 1: Table S1**).

664

665     *Dataset-specific Peaks.* We called peaks from each dataset's cells as we described

666     previously[21]. For each RA sample, we subsetted the BAM files to chromosomes 1-22XY and

667     QCed cells within each dataset (*Benchmarking Datasets 1-3*) and then converted the

668     subsequent files to MACS2 BEDPE files. For each COVID-19 sample, we started from fragment

669     files requested from the authors, subsetted to the BED-3 columns, chromosomes 1-22XY, and

670     QCed cells from each dataset (*Benchmarking Datasets 4-5*) to get BEDPE files. For each

671     dataset, we concatenated all BEDPE files across samples and called consensus peaks using

672     macs2 callpeak --call-summits with a control file where ATAC-seq was done on free DNA[32] to

673     account for Tn5's inherent cutting bias. We trimmed peaks to 500bp (summit +/- 250 bp) and

674     removed overlapping peaks iteratively, keeping the peak with the best q-value. These peaks

675     were overlapped with cell fragments via GenomicRanges::findOverlaps to get a peaks by cells

676     matrix.

677

678     *Tiles.* Tiles of 500bp were computed within ArchR[24] version 1.0.2 via createArrowFiles with

679     addTileMat = TRUE and exported via getMatrixFromProject. The mitochondrial and Y

680     chromosomes were excluded. To calculate fragment overlap, we used

681     GenomicRanges::findOverlaps between ArchR's tiles and the whole fragments, as done in the

682     other feature types.

683

684     *cCRE.* ENCODE SCREEN v3 candidate cis-regulatory elements (cCREs) were downloaded

685     from https://screen.wenglab.org/index/cversions. We also downloaded a set of cell-type-

686     agnostic cCRE with their max DNase-seq z-score from UCSC Data Browser

687     (https://hgdownload.soe.ucsc.edu/gbdb/hg38/encode3/ccre/encodeCcreCombined.bb). For

688     each cCRE, we extended it to 500 bp (midpoint +/- 250 bp) and removed overlapping cCREs

689     iteratively, keeping the cCRE with the highest cell-type-agnostic max DNase-seq z-score. If the

690     cCREs were not cell-type-agnostic or had the same max z-score, then we prioritized by cCRE

691     annotation in the following order: PLS, PLS/CTCF-bound, pELS, pELS/CTCF-bound, dELS,

692     dELS/CTCF-bound, DNase-H3K4me3, DNase-H3K4me3/CTCF-bound, CTCF-only/CTCF-

693     bound. We overlapped the final set of non-overlapping 500bp cCREs with each Benchmarking

694     Dataset's cell fragments to get a matrix.

695

696    *Gene Activity Scores.* We used two gene activity score (GAS) methods. The first was a basic

697    model ('basicGAS') that overlapped cell fragments with gene bodies with a 2kb promoter region

698    upstream of the gene. The second was an exponential decay model within ArchR[24]

699    ('archrGAS'), calculated via createArrowFiles with addGeneScoreMat = TRUE and exported via

700    getMatrixFromProject. It was then converted from a SummarizedExperiment to a named Sparse

701    Matrix for streamlined processing. In both cases, we used ArchR's default hg38 gene annotation

702    obtained via getGenes.

703

704    **Integration Methods.** We used the above feature matrices as input to a variety of algorithms

705    spanning linear and nonlinear methodologies (**Table 2**). All methods inputted these feature

706    matrices, except ArchR, which reconstructed the feature matrix internally from a user-defined

707    feature set. All non-GAS matrices used 500 bp features; they were also binarized, except for

708    Signac and SnapATAC2 as they specified counts. GAS matrices were processed as if they

709    were gene expression data. Note that not all snATAC-seq integration methods were suited for

710    GAS feature types; if excluded, a justification was given per method. We used defaults for all

711    methods unless otherwise stated. In total, we tested 58 strategies of feature types, integration

712    methods, and Harmony correction (**Additional file 1: Table S1**). For visualization purposes, we

713    calculated a UMAP via umap::umap for each embedding.

714

715    *Principal Component Analysis (PCA).* PCA defined orthogonal principal components (PCs) that

716    explained progressively smaller percentages of variation. It is widely used in genomics to

717    reduce data dimensionality, including for single-cell data[33]. We used the snATAC-seq

718    normalization strategy of term frequency-inverse document frequency (TF-IDF) as part of the

719    preprocessing.

720    We generated PC embeddings for each non-GAS feature type binary matrix via: subsetting

721    features by those with accessibility in at least 0.5% of cells, log(TFxIDF) normalization, variable

722     feature selection, center/scale features, and PCA via irlba::prcomp_irlba[34]. For the basic GAS

723     features, we log-normalized features; ArchR GAS features were already normalized to 10,000,

724     so we logged its matrix with log1p. In both GAS types, we then did variable feature selection,

725     center/scale features, and PCA via irlba::prcomp_irlba. We used 30 dimensions here since the

726     default 3 will likely be too small for this application and all other methods, except PeakVI,

727     defaulted to 30 dimensions. We used irlba R version 2.3.5.1, a package not specifically tailored

728     to snATAC-seq data analysis or pipelines.

729

730     *Signac.* Signac[27] is the commonly used[4,10,22] snATAC-seq extension of the Seurat package. It

731     used latent semantic indexing (LSI), a combination of TF-IDF normalization followed by Singular

732     Value Decomposition (SVD). For integration purposes, they recommended reciprocal LSI (rLSI),

733     iteratively projecting each dataset into a shared LSI space.

734     Using Signac version 1.14.0, we processed the snATAC-seq non-GAS non-binary feature

735     matrices with the LSI pipeline: CreateSeuratObject along with the cell metadata,

736     FindTopFeatures, RunTFIDF, and RunSVD. For rLSI in the non-GAS features, we created the

737     Seurat object as before, but split it by sample into a list of objects using SplitObject, and ran the

738     LSI pipeline per object. After merging objects and joining layers, we ran the LSI pipeline again

739     on the combined object. We then used FindIntegrationAnchors and IntegrateEmbeddings to get

740     the rLSI embedding. We calculated the default 30 dimensions, but only used the 2:30

741     dimensions as suggested by their tutorial as the first dimension was usually correlated to

742     sequencing depth. For GAS features, we used the related pipeline Seurat version 5.1.0,

743     developed primarily for scRNA-seq. However, since Seurat's default method is PCA, already

744     tested here, we only tested their anchor integration approach using rCCA as described in the

745     *bioKNN* section. For pre-normalized archrGAS features, we built a Seurat object using

746     CreateAssay5Object with the data layer set to the logged matrix before continuing the pipeline

747     outlined above.

748     We encountered two errors while running rLSI/rCCA. During the rLSI IntegrateEmbeddings

749     step, we would occasionally get FindWeights errors saying "Number of anchor cells is less than

750     k.weight. Consider lowering k.weight to less than XX or increase k.anchor", where XX was an

751     integer. We reran the same command with an additional argument k.weight=XX-1 until the error

752     no longer occurred. During rCCA with the pre-normalized archrGAS features, we got

753     FindVariableFeatures errors in match.arg. We reran the same command with

754     selection.method="mvp" instead of the default "vst".

755

756     *ArchR.* ArchR[24] is a popular[6,7,35] end-to-end pipeline for analyzing scATAC-seq data. It used

757     iterative LSI for dimensionality reduction and implicit batch correction. For the first iteration,

758     ArchR used the top accessible features (by default, 25,000 500 bp tiles). Subsequently, it

759     calculated LSI, defined clusters, sum-aggregated accessibility per cluster, log-normalized, and

760     identified most variable features across clusters to use in the next, and by default final, round of

761     LSI. ArchR did not explicitly correct for batch within this process, claiming that the first round of

762     LSI identified low resolution clusters that were not batch confounded. We did not use ArchR's

763     estimated LSI functionality as it was not set by default.

764     We inputted the per-sample post cell QC fragment files to ArchR version 1.0.2 createArrowFiles

765     with all additional QC flags nullified to avoid subsetting on cells; the tile and archrGAS matrices

766     were created in this step as mentioned previously. We then created an ArchRProject. Since

767     ArchR does not allow non-gene-expression matrices to be inputted directly, we used the feature

768     sets calculated in **ATAC Feature Sets and Matrices** to build the feature matrices within the

769     ArchRProject. The peak matrix was added to the project via addPeakSet and addPeakMatrix

770     while the cCRE matrix was added via addFeatureMatrix. Since the basicGAS matrix was count

771     data, it was converted to a SummarizedExperiment object and inputted with

772     addGeneExpression to allow for the calculation of gene-specific parameters. For each non-GAS

773     feature matrix, we applied the IterativeLSI process via addIterativeLSI with default parameters.

774    For the GAS feature matrices, we used addIterativeLSI with parameters to mimic scRNA-seq:

775    binarize = FALSE, firstSelection = "var", varFeatures = 2000. In both cases, we calculated the

776    default 30 dimensions, though of note, ArchR excluded dimensions correlated to sequencing

777    depth by default. ArchR did not require a batch covariate. The embedding matrices were

778    exported from the ArchR project via getReducedDims and converted from a

779    SummarizedExperiment to a named Sparse Matrix for streamlined postprocessing.

780

781    *SnapATAC2.* SnapATAC2[23] used a matrix-free spectral embedding for dimensionality reduction.

782    This nonlinear method was built on the Lanczos algorithm to implicitly use the Laplacian matrix

783    without storing it to decrease time and space complexity. SnapATAC2 utilized independent

784    batch correction algorithms in their paper, with Harmony[20] among its suggested set.

785    Since SnapATAC2 was written in python, we first converted the nonbinary R Sparse Matrix to

786    MatrixMarket format using writeMM. We then created an AnnData object using

787    snapatac2.read_10x_mtx with that matrix file, an observation file including cell metadata, and a

788    variable file of peak names. For non-GAS features, we selected 50,000 features with

789    snapatac2.pp.select_features, as suggested in their tutorial; this was the only post-IO step in

790    their pipeline before non-GAS embedding creation. Since basicGAS features are counts, we

791    used scanpy functions as suggested by their paper: scanpy.pp.highly_variable_genes with

792    flavor='seurat_v3' before subsetting to those variable genes, scanpy.pp.normalize_total,

793    scanpy.pp.log1p. As the archrGAS features were already normalized, we used scanpy.pp.log1p

794    and scanpy.pp.highly_variable_genes with flavor='seurat' before subsetting. We then ran their

795    dimensionality reduction method with snapatac2.tl.spectral with a seed and the default 30

796    dimensions set and if a GAS feature, features=None. Note that with default

797    weighted_by_sd=TRUE, the resulting matrix could be less than 30 dimensions. We used

798    SnapATAC2 version 2.5.3. We then extracted the embedding matrix to a text file that we loaded

799    back into R to save as a named Sparse Matrix. SnapATAC2 did not require a batch covariate.

800

801   *CellSpace.* CellSpace[28] co-embedded k-mers and cells to infer features and TF motifs from their

802   constituent k-mers using neural embedding modeling software StarSpace[36]. It randomly

803   sampled overlapping k-mers from inputted accessible features, peaks or tiles preferred, to

804   generate training examples of k-mers accessible per cell ("positive" cell). Negative cells without

805   accessibility for the k-mer set were randomly sampled since there were many more inaccessible

806   cells for any given feature; this also helped with false negatives and sparsity. During training, k-

807   mer and cell embeddings were updated to move the induced feature sequence embedding

808   closer to positive cells and further from negative cells. N-grams of flanking sequence around k-

809   mers provided additional context. Batch variables were not used here as the authors claim the

810   final embedding generated from small k-mers was less influenced by batch effects. We did not

811   run the GAS matrices here since randomly sampling 8-mers across an entire gene seemed

812   counterintuitive to the TF-centric theme of this method.

813   We downloaded and compiled CellSpace version 1.0.0 as their GitHub tutorial recommended.

814   CellSpace required a file of cell names, the fasta file of feature sequences, and a MatrixMarket

815   formatted matrix. CellSpace did not internally do feature selection, though its paper

816   recommended using top variable features to speed up runtime and possibly improve quality.

817   Therefore, we subsetted to variable features calculated as in the PCA pipeline discussed above

818   before converting from R Sparse Matrix to MatrixMarket using writeMM. We generated fasta

819   files of the feature regions via GenomicRanges::getSeq. We ran CellSpace with all default

820   parameters, except requesting 5 threads instead of 10. CellSpace did not require a batch

821   covariate. The resulting tsv file was then inputted into its corresponding R package CellSpace

822   function to access and save its embedding slot cell.emb as a 30-dimension named Sparse

823   Matrix.

824

40

825　*PeakVI.* As part of the scVI suite of tools, PeakVI[25] used a variational autoencoder (VAE) to

826　model a latent space conditioned on a user-specified batch variable to correct for those batch

827　effects and capture batch-independent biological variation. The batch variable was also used in

828　the decoding step to calculate the probability of accessibility, which was further modified to the

829　probability of observation by multiplying by the region-specific factors and cell-specific factors

830　calculated from the input. Because PeakVI modeled a Bernoulli distribution, we did not use the

831　non-binary GAS matrices as input.

832　Since PeakVI was written in python, we first converted the binary R Sparse Matrices to

833　MatrixMarket format using writeMM before converting into a python AnnData format using

834　scvi.data.read_10x_atac. Cell metadata was added to the AnnData object separately, verifying

835　that all observations and variables were unique. We then filtered the features to keep those in at

836　least 5% of cells as recommended by their tutorial. We setup the AnnData object using

837　scvi.data.setup_anndata with batch_key='sample' before training the PeakVI model with

838　scvi.model.PEAKVI with all default parameters. We used scvi-tools version 1.1.6. We then

839　extracted the embedding matrix to a text file that we loaded back into R to save as a named

840　Sparse Matrix.

841

842　**Harmony Correction.** To standardize the Harmony implementation and defaults, we applied

843　the stand-alone Harmony package version 1.2.1 to each integration method's output

844　embedding, using all dimensions given. We batch corrected by sample using default

845　parameters.

846

847　**Benchmarking Datasets.** An overview of the benchmarking datasets used in this study is in

848　**Table 3**. RNA UMAPs with gold standard biological cell phenotypes defined from RNA or

849　protein can be found in **Fig. 1g**; these annotations were used in the cLISI bio-conservation

850　metric (see **Benchmarking Metrics** section). Additionally, as part of the Seurat NN graphs

41

851    generated in the **Benchmarking Metrics** section, we removed samples with less than 100 cells

852    from the Benchmarking Datasets before any further processing since Seurat's anchor

853    integration had issues with very small cell count batches. **Additional file 2: Figs. S1-S5**

854    describe each dataset: panels **a**-**b** denote quality control statistics for each modality, panels **c-e**

855    denote snRNA-seq batch correction (iLISI) and bio-conservation (cLISI) metrics, and for the

856    RNA-based cell phenotypes, panels **f-i** denote comparisons between our reprocessing and the

857    original paper's processing. Full methodological details of the initial processing can be found in

858    the original studies[21,22], with summaries here.

859

860    *Benchmark Dataset 1. RA inflammatory tissue broad cell types*. This dataset[21] consisted of 12

861    synovial tissue samples, 11 from RA patients and 1 from an osteoarthritis (OA) patient, and

862    included 6 broad cell types (B/plasma, T, NK, myeloid, stromal [fibroblasts/mural], and

863    endothelial) across 31,547 cells (median 3,093 cells/sample). Each sample was processed

864    independently at the cell capture step and run through a 10x multiome experiment. To be

865    included, cells had to pass quality control metrics in both snRNA-seq and snATAC-seq

866    modalities (post-QC in **Additional file 2: Fig. S1a-b**) as well as have the same annotated broad

867    cell type. We reprocessed the snRNA-seq data using the updated software versions uniformly

868    used in this study. Briefly, we normalized, selected variable genes, centered/scaled genes,

869    computed 20 PCs, batch corrected by sample via Harmony[20], created a shared nearest

870    neighbor graph, clustered with Louvain clustering, and generated a UMAP. This reprocessing

871    generated the RNA-embedding1 used in the NN metrics. It showed good batch correction

872    (**Additional file 2: Fig. S1c-d**) and cell type bio-conservation (**Additional file 2: Fig. S1e**). We

873    chose a clustering resolution that closely mirrored the original processing (**Additional file 2:**

874    **Fig. S1f-i**).

875

42

876     *Benchmark Dataset 2. RA inflammatory tissue T cell states.* This dataset[21] used 8,069 T cells

877     from 11 of the RA/OA synovial tissue samples discussed above; one sample was dropped for

878     having less than 100 cells. As the original study did not define snRNA-seq cell states from the

879     multiome data alone, we clustered as above using 10 sample-harmonized[20] PCs. Samples were

880     well-mixed after Harmony (**Additional file 2: Fig. S2c-d**) with consistent clusters (**Additional**

881     **file 2: Fig. S2e**). We chose a clustering resolution that corresponded to the chromatin classes

882     defined by the original study (**Additional file 2: Fig. S2f-g**). We annotated cell states using

883     marker genes similar to those in the original study and via differential gene analysis using

884     presto::wilcoxauc on the normalized genes x cells matrix (**Additional file 2: Fig. S2h-i**).

885

886     *Benchmark Dataset 3. RA PBMC T cell substates.* This dataset[21] totaled 10,669 cells across

887     two runs of four pooled RA PBMC samples each. There were no donors shared between runs,

888     and the donors were not genotyped, disallowing genotype-based demultiplexing. After enriching

889     for CD4+ T cells, four populations were sorted using Fluorescence-Activated Cell Sorting

890     (FACS): $CD4^+CD127^-CD25^{hi}$ Tregs, $CD4^+CD127^-CD25^{int}$ Tregs, $CD4^+CD25^-PD1^+CXCR5^+$ TFH,

891     and $CD4^+CD25^-PD1^+CXCR5^-$ TPH. Each population was tagged with a hashing antibody. We

892     generated 10x multiome data from these cells with an HTO library generated to label the

893     specific populations. Cells had to pass quality control metrics in all modalities processed.

894     Overall, the quality was slightly worse than the RA tissue multiome datasets in *Benchmark*

895     *Datasets 1-2* with fewer fragments and genes detected and higher mitochondrial reads

896     (**Additional file 2: Fig. S3a-b**). We used 10 run-harmonized PCs to generate a snRNA-seq

897     UMAP (**Fig. 1g**) and a gene KNN graph for use in the bioKNN and batchKLD metrics as defined

898     above. Samples mixed reasonably well (**Additional file 2: Fig. S3c-d**), though there was a cell

899     imbalance between runs (2,998 in Run 1 vs 7,671 in Run 2). While both runs were used for all

900     snRNA-seq and snATAC-seq processing, only one run had processed FACS populations, so

901     the cLISI calculation was restricted to that run in both modalities (**Additional file 2: Fig. S3e**).

902

903 *Benchmark Dataset 4. COVID-19 PBMC broad cell types.* This study[22] profiled mature and

904 progenitor cell populations from peripheral blood samples using PBMC-PIE (peripheral blood

905 mononuclear cell analysis with progenitor input enrichment). This 10x multiome dataset was

906 comprised of 30 PBMC samples collected between March 2020 and March 2021 from 26

907 individuals, including subsets of healthy adults, adults recovering from a non-COVID-19 critical

908 illness, and adults at early (2-4 months) or late (4-12 months) post-infection stages of COVID-19

909 convalescence. In the original study, sample merging (using Seurat[29] for snRNA-seq and

910 Signac[27] for snATAC-seq), quality control, doublet removal (using Scrublet for snRNA-seq and

911 Amulet for snATAC-seq), batch correction by sample with Harmony[20], and iterative

912 clustering/annotation were preformed to get a final dataset of 197,360 cells (median 6,027

913 cells/sample). Overall, there were fewer genes and fragments found here with a much higher

914 mitochondrial percentage than in *Benchmark Datasets 1-2* (**Additional file 2: Fig. S4a-b**).

915 Using marker gene expression, the authors defined 10 major cell type clusters: NK, CD8+ T,

916 CD4+ T, B, plasma, CD16+ monocyte, CD14+ monocyte, DC, pDC, and hematopoietic stem

917 and progenitor cells (HSPC). For uniform processing and to generate a nearest neighbor graph

918 for use in the NN metrics, we reprocessed the multiome snRNA-seq data using the pipeline

919 denoted in *Benchmark Dataset 1* using 30 sample-harmonized PCs. We saw good sample

920 mixing (**Additional file 2: Fig. S4c-d**) and internal cell type consistency (**Additional file 2: Fig.**

921 **S4e**). Furthermore, we compared the reprocessed cell types to the originally-defined cell types

922 (**Additional file 2: Fig. S4f-i**) and saw general consistency.

923

924 *Benchmark Dataset 5. COVID-19 PBMC HSPC states.* In the original paper, the 28,069 HSPCs

925 in the previous dataset [22] were further subclustered into 6 states: hematopoietic stem cells and

926 multipotent progenitors (HSC_MPP), lymphoid-primed MPP (LMPP), granulocyte-monocyte

927 progenitors (GMP), megakaryocyte-erythroid progenitors (MEP), erythroid progenitors (Ery),

44

928    and basophil-eosinophil-mast (BEM) cell progenitors. As above, we dropped 2 samples with

929    fewer than 100 cells and re-annotated cells using 20 multimodal snRNA-seq sample-

930    harmonized PCs (**Additional file 2: Fig. S5c-e**). We concatenated the multiple HSC_MPP

931    clusters as the original authors did. The cluster borders were highlighted in a higher (lighter)

932    cLISI score in **Additional file 2: Fig. S5e**. We compared them to the original marker gene

933    annotations and saw common cell states (**Additional file 2: Fig. S5f-i**).

934

935

936    **Pipeline.** After creating all the feature matrices defined in **ATAC Feature Sets and Matrices** for

937    each dataset defined in **Benchmarking Datasets**, we used a command line tool we developed

938    to create all embeddings, Harmony-corrected embeddings, NN graphs, metrics, and UMAPs.

939    The datasets and features were inputted into this pipeline via a "dictionary file," which connected

940    dataset and feature keywords to file paths for the required types of files: feature matrix, features

941    for ArchR input, ArchR project, metadata, and gene NN graphs. These dataset and feature

942    keywords were then referenced in the pipeline script to generate dataset/feature/method-

943    specific command files that we piped into a SLURM scheduler. Therefore, any dataset or

944    feature combination future users desire can be used to generate the types of output for each

945    method we detailed here.

946

947    **Linear modeling.** We used two different linear models, combined and separated, to assess the

948    overall metric ranking for each of the three metric types: (1) NN metrics using RNA-

949    embedding1, (2) NN metrics using RNA-embedding2, (3) LISI metrics. The bio-conservation

950    metrics for NN and LISI were bioKNN and cLISI, respectively while the batch metrics were

951    batchKLD and iLISI, respectively. Each metric was mean-averaged across cells per dataset,

952    feature, method, and correction combination and subsequently ranked within dataset where 1 is

953    the best rank and 58 is the worst rank (**Additional file 1: Table S2-S4**). In both models, this

954    rank was used to create a score that was related to the dataset, feature, method, and correction

955    combination using stats::glm in R with family="gaussian".

956    $$score \sim dataset + feature + method + correction$$

957    In the combined model, using percentages based off Luecken et al.,[15], we ran 1 model using

958    score:

959    $$score = \ rank_{mKNN} * 60\% + rank_{lsKLD} * 40\%$$

960    In the separated model, we ran 2 models using scores:

961    $$score = \ rank_{mKNN}$$

962    $$score = \ rank_{lsKLD}$$

963    To avoid overdetermination in the model, we chose a reference value for each covariate as

964    follows: dataset – dataset1, feature – tile, method – LSI, and correction – No Harmony. Each

965    covariate's matrix was 1-hot encoded.

966

967    **Job Requirements.** We assessed time and memory requirements per dataset/feature/method

968    pipeline job and per step within a pipeline job using /usr/bin/time. We summed both system and

969    user CPU-minutes for time while using max resident set size (RSS) to assess memory. Steps

970    were separated into 5 categories. Firstly, pre-processing included variable peak/gene selection,

971    data type conversions from R to python, and feature matrix re-creation within ArchR projects.

972    Secondly, embeddings were generated and thirdly, they were corrected with Harmony. Fourthly,

973    post-processing was converting both uncorrected and corrected embeddings back into R

974    Matrices. Fifthly, NN graphs, metrics, and UMAPS were generated for metrics/visualization

975    purposes for both uncorrected and corrected embeddings.

976

977 **Declarations**

978 **Ethics approval and consent to participate**

979 Not applicable

980

981 **Consent for publication**

982 Not applicable

983

984 **Availability of data and materials**

985 RA datasets analyzed during this study are included in Weinand et al., Nat Commun, 2024[21]

986 and available from the Synapse repository using identifier syn53650034[37]

987 (https://doi.org/10.7303/syn53650034). COVID-19 datasets analyzed during this study are

988 included in Cheong et al., Cell, 2023[22] and available from the GEO repository using identifiers

989 GSE196987 and GSE196988

990 (https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE196987;

991 https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE196988).

992 Processed data generated in this study are available from the Zenodo repository using identifier

993 15073137 (https://doi.org/10.5281/zenodo.15073137). A website of UMAPs generated in this

994 study is available at https://immunogenomics.io/snATAC_benchmark/.

995 The command-line tool code used to generate the results and figures presented in this study

996 can be found on GitHub (https://github.com/immunogenomics/snATAC_benchmark/). We have

997 also listed there the conda and R packages we used.

998

999 **Competing Interests**

1000 The authors declare that they have no competing interests.

1001

1002 **Funding**

1005

**Author contributions**

1007    KW and SR conceptualized the study. KW defined the methodology, software, and visualization

1008    with input from EL and SR. KW processed and analyzed the Benchmarking Datasets. EL

1009    curated and initially processed and analyzed the COVID-19 datasets. MC designed the website.

1010    KW and SR drafted the original manuscript with edits by EL and MC. SR supervised the study

1011    and provided funding. All authors read and approved the final manuscript.

1012

**Acknowledgements**

## References

1016    **References**

1017  1.   Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory
1018       variation. *Nature* (2015) doi:10.1038/nature14590.
1019  2.   De Rop, F. V *et al.* Systematic benchmarking of single-cell ATAC-sequencing protocols.
1020       *Nat Biotechnol* **42**, 916–926 (2024).
1021  3.   Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human
1022       immune cell development and intratumoral T cell exhaustion. *Nat Biotechnol* **37**, 925–936
1023       (2019).
1024  4.   Domcke, S. *et al.* A human cell atlas of fetal chromatin accessibility. *Science* **370**, (2020).
1025  5.   Zhang, K. *et al.* A single-cell atlas of chromatin accessibility in the human genome. *Cell*
1026       **184**, 5985-6001.e19 (2021).
1027  6.   Sun, K., Liu, X. & Lan, X. A single-cell atlas of chromatin accessibility in mouse
1028       organogenesis. *Nat Cell Biol* **26**, 1200–1211 (2024).
1029  7.   Kim, H. *et al.* Single-Cell Transcriptional and Epigenetic Profiles of Male Breast Cancer
1030       Nominate Salient Cancer-Specific Enhancers. *Int J Mol Sci* **24**, (2023).
1031  8.   Li, Y. *et al.* Transcriptomics based multi-dimensional characterization and drug screen in
1032       esophageal squamous cell carcinoma. *EBioMedicine* **70**, 103510 (2021).
1033  9.   Carraro, C. *et al.* Decoding mechanism of action and sensitivity to drug candidates from
1034       integrated transcriptome and chromatin state. *Elife* **11**, (2022).
1035  10.  Zhang, B. *et al.* Multimodal single-cell datasets characterize antigen-specific CD8+ T
1036       cells across SARS-CoV-2 vaccination and infection. *Nat Immunol* **24**, 1725–1734 (2023).
1037  11.  Zhang, F. *et al.* Deconstruction of rheumatoid arthritis synovium defines inflammatory
1038       subtypes. *Nature* **623**, 616–624 (2023).
1039  12.  Korsunsky, I. *et al.* Cross-tissue, single-cell stromal atlas identifies shared pathological
1040       fibroblast phenotypes in four chronic inflammatory diseases. *Med (N Y)* **3**, 481-518.e14
1041       (2022).
1042  13.  Eraslan, G. *et al.* Single-nucleus cross-tissue molecular reference maps toward
1043       understanding disease gene function. *Science* **376**, eabl4290 (2022).
1044  14.  Sun, N. *et al.* Human microglial state dynamics in Alzheimer's disease progression. *Cell*
1045       **186**, 4386-4403.e29 (2023).
1046  15.  Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics.
1047       *Nat Methods* **19**, 41–50 (2022).
1048  16.  Li, Z. *et al.* Chromatin-accessibility estimation from single-cell ATAC-seq data with
1049       scOpen. *Nat Commun* **12**, 6386 (2021).
1050  17.  Chen, H. *et al.* Assessment of computational methods for the analysis of single-cell
1051       ATAC-seq data. *Genome Biol* (2019) doi:10.1186/s13059-019-1854-5.
1052  18.  Luo, S., Germain, P.-L., Robinson, M. D. & von Meyenn, F. Benchmarking computational
1053       methods for single-cell chromatin data analysis. *Genome Biol* **25**, 225 (2024).
1054  19.  ENCODE Project Consortium *et al.* Expanded encyclopaedias of DNA elements in the
1055       human and mouse genomes. *Nature* **583**, 699–710 (2020).
1056  20.  Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with
1057       Harmony. *Nature Methods 2019 16:12* **16**, 1289–1296 (2019).
1058  21.  Weinand, K. *et al.* The chromatin landscape of pathogenic transcriptional cell states in
1059       rheumatoid arthritis. *Nat Commun* **15**, 4650 (2024).
1060  22.  Cheong, J.-G. *et al.* Epigenetic memory of coronavirus infection in innate immune cells
1061       and their progenitors. *Cell* **186**, 3882-3902.e24 (2023).
1062  23.  Zhang, K., Zemke, N. R., Armand, E. J. & Ren, B. A fast, scalable and versatile tool for
1063       analysis of single-cell omics data. *Nat Methods* **21**, 217–227 (2024).
1064  24.  Granja, J. M. *et al.* ArchR is a scalable software package for integrative single-cell
1065       chromatin accessibility analysis. *Nature Genetics 2021 53:3* **53**, 403–411 (2021).

1066   25.   Ashuach, T., Reidenbach, D. A., Gayoso, A. & Yosef, N. PeakVI: A deep generative
1067         model for single-cell chromatin accessibility analysis. *Cell Reports Methods* **2**, 100182
1068         (2022).
1069   26.   Gupta, A. *et al.* Dynamic regulatory elements in single-cell multimodal data implicate key
1070         immune cell states enriched for autoimmune disease heritability. *Nat Genet* **55**, 2200–
1071         2210 (2023).
1072   27.   Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin
1073         state analysis with Signac. *Nat Methods* **18**, 1333–1341 (2021).
1074   28.   Tayyebi, Z., Pine, A. R. & Leslie, C. S. Scalable and unbiased sequence-informed
1075         embedding of single-cell ATAC-seq data with CellSpace. *Nat Methods* **21**, 1014–1022
1076         (2024).
1077   29.   Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell
1078         analysis. *Nat Biotechnol* **42**, 293–304 (2024).
1079   30.   Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and
1080         Chromatin. *Cell* **183**, 1103-1116.e20 (2020).
1081   31.   Cao, J. *et al.* Joint profiling of chromatin accessibility and gene expression in thousands
1082         of single cells. *Science (1979)* **361**, 1380–1385 (2018).
1083   32.   Martins, A. L., Walavalkar, N. M., Anderson, W. D., Zang, C. & Guertin, M. J. Universal
1084         correction of enzymatic sequence bias reveals molecular signatures of protein/DNA
1085         interactions. *Nucleic Acids Res* (2018) doi:10.1093/nar/gkx1053.
1086   33.   Tsuyuzaki, K., Sato, H., Sato, K. & Nikaido, I. Benchmarking principal component
1087         analysis for large-scale single-cell RNA-sequencing. *Genome Biol* **21**, 9 (2020).
1088   34.   Baglama, J. & Reichel, L. Augmented Implicitly Restarted Lanczos Bidiagonalization
1089         Methods. *SIAM Journal on Scientific Computing* **27**, 19–42 (2005).
1090   35.   Thakore, P. I. *et al.* BACH2 regulates diversification of regulatory and proinflammatory
1091         chromatin states in TH17 cells. *Nat Immunol* (2024) doi:10.1038/s41590-024-01901-1.
1092   36.   Wu, L. *et al.* StarSpace: Embed All The Things! (2017).
1093   37.   Weinand, K. *et al.* Dataset: The chromatin landscape of pathogenic transcriptional cell
1094         states in rheumatoid arthritis. *Synapse* (2024) doi:10.7303/syn53650034.
1095
1096

1097    **Additional files**

1098

1099    Additional file 1. Supplementary Tables S1-S4

1100    Excel Spreadsheet, .xlsx

1101

1102    Additional file 2. Supplementary Figures S1-S17

1103    Word Document, .docx

1104