

1 **Is single nucleus ATAC-seq accessibility a qualitative or quantitative measurement?**

2

3 Zhen Miao<sup>1,2</sup>, and Junhyong Kim<sup>1,2,†</sup>

4

5 <sup>1</sup>Graduate Group in Genomics and Computational Biology, Perelman School of Medicine,  
6 University of Pennsylvania, Philadelphia, PA, USA

7 <sup>2</sup>Department of Biology, University of Pennsylvania, Philadelphia, PA, USA

8 <sup>†</sup>Corresponding author

9

10 **Correspondence:**

11 **Junhyong Kim, Ph.D.**

12 Patricia M. Williams Term Professor, Department of Biology

13 Adjunct Professor, Department of Computer and Information Science

14 University of Pennsylvania

15 433 S. University Avenue

16 Philadelphia, PA 19104

17 [junhyong@sas.upenn.edu](mailto:junhyong@sas.upenn.edu)

18

19 **Abstract:**

20

21 Single nucleus ATAC-seq is a key assay for gene regulation analysis. Existing approaches to  
22 scoring feature matrices from sequencing reads are inconsistent with each other, creating  
23 differences in downstream analysis, and displaying artifacts. We show that even with sparse single  
24 cell data, quantitative counts are informative for estimating a cell's regulatory state, which calls  
25 for consistent treatment. We propose Paired-Insertion-Counting (PIC) as a uniform method for  
26 snATAC-seq feature characterization.

27

28 **Main:**

29

30 Single nucleus ATAC-seq (snATAC-seq) assays open chromatin profiles of individual cells.  
31 However, unlike RNA-seq where the counts estimate numbers of molecules, there is not a common  
32 agreement on what biological state is being estimated from snATAC-seq data. Existing snATAC-  
33 seq analysis methods create chromosomal domain features either by arbitrarily dividing the entire  
34 genome into fixed-width segments (features usually referred to as bins), or estimating discrete  
35 domains by peak-calling from aggregated pseudo-bulk data (features usually referred to as peaks).  
36 Using bins as features has problems associated with arbitrarily fixing length scales and phase (i.e.,  
37 starting positions of the bins) and the problem that many bins will contain no relevant information.  
38 Peaks subset functionally relevant genomic intervals, but there are technical challenges to resolve  
39 boundaries for heterotypic datasets and to identify functional elements for rare cells, and  
40 differences exist in numerical criterion for peak identification. After choosing bins or peaks, some  
41 methods assign the feature counts based on the number of fragments that overlap with a region  
42 (fragment-based counting; e.g., Signac<sup>1</sup> and snapATAC<sup>2</sup>), while others assign counts based on the  
43 number of insertions within the region (insertion-based counting; e.g., 10X cellranger ATAC<sup>3</sup> and  
44 ArchR<sup>4</sup>). After feature counting, most methods convert the counts into a binary state of “open” or  
45 “closed” (e.g., snapATAC<sup>2</sup>, SCALE<sup>5</sup>, scOPEN<sup>6</sup>, MASETRO<sup>7</sup>, and cisTopic<sup>8</sup>), while other retain  
46 quantitative count information, implying that single nucleus assays may contain quantitative  
47 information on nucleosome density or turnover (e.g., scABC<sup>9</sup>, chromVAR<sup>10</sup>, and ArchR<sup>4</sup>).

48

49 When considering counts, the configuration of fragment/insertion positions around the peak/bin  
50 interval can create different quantifications dependent on whether one uses fragments or insertions  
51 (Figure 1a-b). Histograms of counts for fragment-based or insertion-based counting applied to the  
52 same dataset (10X Genomics peripheral blood mononuclear cell dataset, PBMC-5k) show evident  
53 differences (**Figure 1c-f** and **Supplementary Table 1**). In particular, with insertion-based  
54 counting, there is an artifact of depleted odd numbers. In a standard ATAC-seq experiment, two  
55 Tn5 insertions in the appropriate directions are required to form one amplicon fragment, thus the  
56 unit of observation is pairs of insertions. Odd number of insertions only arise when rare fragments  
57 cross feature boundaries, artificially breaking up paired insertions of a fragment. Fragment-based  
58 counting also has problems because the entire interval of an amplicon from a pair of insertion is  
59 considered evidence of “openness”. However, longer the fragment, less likely the region away  
60 from the insertion sites is open. This is especially acute when there are long fragments with  
61 insertions completely outside the peak/bin of interest<sup>11,12</sup> (cell 1 in **Figure 1a**). The two counting  
62 strategies can result in discrepancies in downstream analysis. As an example, we analyzed a P0  
63 mouse kidney snATAC-seq dataset<sup>13</sup> for Differentially Accessible Region (DAR) identification  
64 between two most abundant cell types with ArchR<sup>4</sup> and Signac<sup>1</sup> (**Methods**). We found up to 4.7%  
65 peaks are only significant with one counting strategy, but not the other (**Supplementary Figure**  
66 **1a**).

67  
68 If the counts are binarized, both insertion and fragment counting are consistent with each other,  
69 except for rare cases (e.g., cell 1 in **Figure 1a**). Thus, the vagaries of counting only matter if  
70 snATAC-seq contains quantitative information about the chromosome state. While variable  
71 nucleosome density and turnover dynamics imply that “openness” is a quantitative state<sup>14</sup>, it is not  
72 clear whether sparse data in single cells contain quantitative information. We asked whether more  
73 fragments in a peak for a single cell indicates higher probability that a randomly selected cell of  
74 the same type would be in open state. That is, we asked whether within-cell insertion density is  
75 predictive of between-cell sampling of open states. We first analyzed a human cell line snATAC-  
76 seq dataset<sup>4</sup>. The cell-by-peak matrix was constructed with insertion-based counting. We retained  
77 166,142 peaks and 10,832 cells in ten cell types after stringent quality control (QC; see **Methods**).  
78 For each peak, we estimated the proportion of cells with the peak being accessible (hereafter we  
79 denote as open probability) in each of the ten cell types (**Methods**). With insertion-based counting

80 approach, a count greater or equal to three indicates at least two fragments (four insertion events)—  
81 we call such cases “high density peaks”. We calculated the relative proportion of cells with high  
82 density peaks for each of the ten cell types (i.e.,  $P(y \geq 3 | y > 0)$ ) and then compared their rank  
83 order with the rank order of cell type open probability by Spearman rank correlation. Among the  
84 peaks we tested, the great majority (>94.6%) showed positive correlation and 9.4% showed  
85 significant correlations at significance level of 0.05 after FDR p-value correction (34.5% without  
86 FDR correction, **Figure 2a**). We also investigated the relationship between open probability and  
87 the relative proportion of cells with counts equal to two given counts being either one or two (i.e.,  
88  $P(y = 2 | y = 1 \text{ or } 2)$ ) for the ten cell types. Consistent with our reasoning that the occurrence of  
89 one insertion mostly represents the boundary phasing artifact, we observed a symmetric  
90 distribution of Spearman correlation coefficients centered around 0 (**Figure 2b**), with only ~0.08%  
91 peaks showing significant correlations at significance level of 0.05 after FDR p-value correction.  
92 Example peaks are shown in **Figure 2c-d** and **Supplementary Figure 1b-c**. We next examined  
93 the P0 mouse kidney snATAC-seq dataset<sup>13</sup> we examined above. After QC, we retained 256,574  
94 peaks and 9,286 cells in seven most abundant cell types in the dataset. With both insertion-based  
95 and fragment-based counting matrices, we conducted the same analysis as above, and the results  
96 were consistent with the human cell line data (**Supplementary Figure 2a-c**) where we found high-  
97 density peaks provided significant information on greater probability of open peaks in the  
98 corresponding cell type.

99  
100 To investigate the potential relationship between snATAC-seq count and gene expression, we  
101 analyzed a 10X genomics PBMC multiome dataset with RNA and ATAC measured on the same  
102 cells. We quantified the cell-by-peak matrix with insertion-based counting approach. Because the  
103 regulatory structure of chromatin domains around a given gene may be complex and largely  
104 unknown, we considered only peaks that are close ( $\pm 100 \text{ bp}$ ) to Transcript Start Site (TSS) to  
105 focus on the most proximal relationship. We also focused on peaks that had a broad range of one  
106 to four counts across cells, filtering out those with too small number of cells within appropriate  
107 range (**Methods**). This resulted in 3,387 peak-gene pairs across 11,234 cells. We compared the  
108 gene expression levels with associated TSS peak insertion count = 1 or 2 (single fragment) against  
109 those with count  $\geq 3$  (more than two fragments) using Wilcoxon rank sum test. We found 199  
110 significant peak-gene pairs after FDR correction, 189 of which have positive log fold change

111 **(Figure 2e)**; 67.2% of peak-gene pairs showed higher non-zero expression proportion in the group  
112 of count  $\geq 3$ . When we compared gene expression levels associated with TSS peak insertion count  
113 = 1 against those with count = 2, we found only 18 significant peak-gene pairs after FDR correction,  
114 nine of which have positive log fold change **(Figure 2f)**. In addition, 52% peak-gene pairs showed  
115 higher non-zero expression proportion in the group of count = 2, suggesting no difference between  
116 the two groups. **Figure 2g-h** shows two examples of peak-gene pair where the distribution of RNA  
117 expression monotonically changes as a function of ATAC counts. We next analyzed a Bone  
118 Marrow Mononuclear Cells (BMMC) multiome dataset<sup>15</sup> which again indicated that peak density  
119 was informative for expression levels **(Supplementary Figure 3a-d)**.

120  
121 In sum, greater counts of snATAC-seq insertions are correlated with greater probability of peak  
122 open state and higher expression of proximal genes, suggesting that even with single nuclear data,  
123 quantitative counting provides important functional information about the epigenomic state of the  
124 cell. We noted above that insertion-based counting creates occasional artifacts and ignores the fact  
125 that, while insertions themselves may be random, the sequence evidence is always in terms of pairs  
126 of insertions. Fragment-based counting has the problem that direct evidence of open state is only  
127 at the insertion site and the evidence for open state decays as a function of distance from the  
128 insertion site. Ideally, it might be appropriate to estimate the quantitative open state of an interval  
129 as a function of fragment lengths and local chromosome features. However, such a model will  
130 need to be data-driven given the irregularities of locus-specific chromosome dynamics. Here, we  
131 propose a simple consistent counting strategy we call Paired-Insertion-Counting (PIC,  
132 <https://github.com/Zhen-Miao/PIC-snATAC>). With PIC, for a given chromosome interval, if an  
133 ATAC-seq fragment's pair of insertions are both within the interval, counted as one (pair); if only  
134 one insertion is within the interval also count one (pair).

135  
136 PIC is consistent with the fact that all fragments have two insertions. It also prevents counting a  
137 fragment when its ends are both outside the peak/bin interval. It has the drawback that when one  
138 insertion is in the peak/bin and the other insertion is far from this insertion, evidence is weak that  
139 both insertions provide information on the current peak/bin. However, in most datasets, long  
140 fragments are rare and unlikely to greatly distort the data **(Supplement Figure 4)**. We recommend  
141 treating snATAC-seq PIC count as a quantitative trait, wherever sensitivity is a critical factor.

142

143 In sum, snATAC-seq is increasingly an important tool for genomic analysis and despite sparse  
144 data at single cell resolution, we find evidence that it can be informative to consider “openness”  
145 as a quantitative trait. Existing approaches are inconsistent in how they quantify peak/bin openness  
146 and here we propose a new counting method that is consistent with the molecular basis of the  
147 assays.

148

149

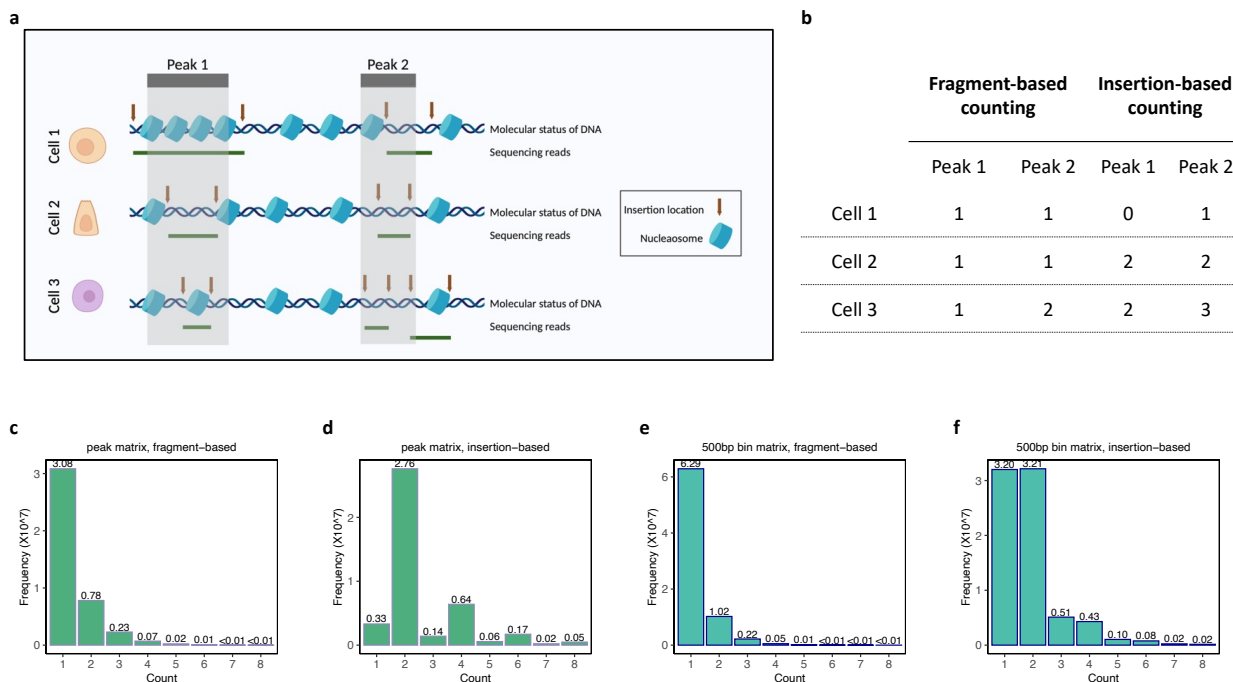
## 150 **References:**

- 151 1. Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin state  
152 analysis with Signac. *Nat. Methods* **18**, 1333–1341 (2021).
- 153 2. Fang, R. *et al.* Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat.*  
154 *Commun.* **12**, 1337 (2021).
- 155 3. Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human immune  
156 cell development and intratumoral T cell exhaustion. *Nat. Biotechnol.* **37**, 925–936 (2019).
- 157 4. Granja, J. M. *et al.* ArchR is a scalable software package for integrative single-cell chromatin  
158 accessibility analysis. *Nat. Genet.* **53**, 403–411 (2021).
- 159 5. Xiong, L. *et al.* SCALE method for single-cell ATAC-seq analysis via latent feature  
160 extraction. *Nat. Commun.* **10**, 4576 (2019).
- 161 6. Li, Z. *et al.* Chromatin-accessibility estimation from single-cell ATAC-seq data with scOpen.  
162 *Nat. Commun.* **12**, 6386 (2021).
- 163 7. Wang, C. *et al.* Integrative analyses of single-cell transcriptome and regulome using  
164 MAESTRO. *Genome Biol.* **21**, 198 (2020).
- 165 8. Bravo González-Blas, C. *et al.* cisTopic: cis-regulatory topic modeling on single-cell ATAC-  
166 seq data. *Nat. Methods* **16**, 397–400 (2019).

- 167 9. Zamanighomi, M. *et al.* Unsupervised clustering and epigenetic classification of single cells.  
168 *Nat. Commun.* **9**, 2410 (2018).
- 169 10. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring  
170 transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods*  
171 **14**, 975–978 (2017).
- 172 11. Schep, A. N. *et al.* Structured nucleosome fingerprints enable high-resolution mapping of  
173 chromatin architecture within regulatory regions. *Genome Res.* **25**, 1757–1770 (2015).
- 174 12. Lu, R. J.-H. *et al.* ATACgraph: Profiling Genome-Wide Chromatin Accessibility From  
175 ATAC-seq. *Front. Genet.* **11**, 618478 (2021).
- 176 13. Miao, Zhen *et al.* Single cell regulatory landscape of the mouse kidney highlights cellular  
177 differentiation programs and disease targets. *Nat. Commun.*  
178 doi:<https://doi.org/10.1038/s41467-021-222266-1>.
- 179 14. Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory  
180 epigenome. *Nat. Rev. Genet.* **20**, 207–220 (2019).
- 181 15. Luecken, M. *et al.* A sandbox for prediction and integration of DNA, RNA, and proteins in  
182 single cells. in *Proceedings of the Neural Information Processing Systems Track on Datasets*  
183 *and Benchmarks* (eds. Vanschoren, J. & Yeung, S.) vol. 1 (2021).
- 184 16. Thibodeau, A. *et al.* AMULET: a novel read count-based method for effective multiplet  
185 detection from single nucleus ATAC-seq data. *Genome Biol.* **22**, 252 (2021).

186

187 **Figures and Figure Legends**



188

189 **Figure 1. Two existing counting strategies for snATAC-seq data processing.**

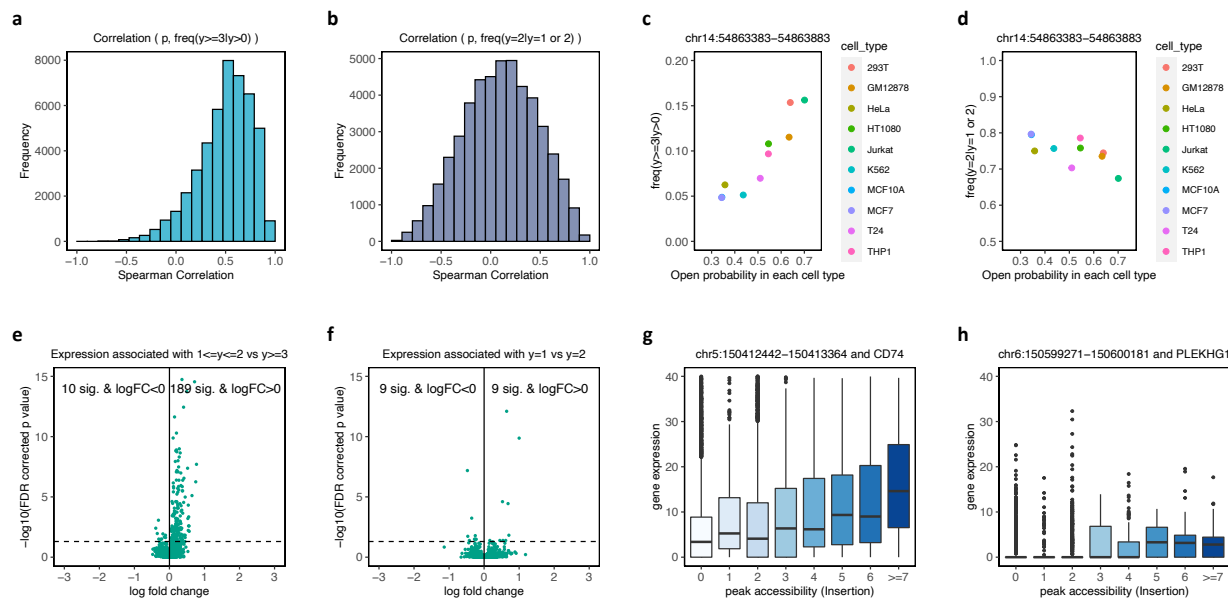
190 (a-b) Schematic example of how the same open chromatin profiles can result in different counts  
 191 with insertion-based or fragment-based counting strategies

192 (c-f) Histogram of count frequencies with two counting strategies and with peaks or bins as  
 193 features

194

195





196

197 **Figure 2. snATAC-seq data contain quantitative information of cellular states.**

198 (a) Histogram of Spearman correlation coefficients between open probability in each group and  
 199 the relative frequency of counts greater than or equal to 3 in human cell line data

200 (b) Histogram of Spearman correlation coefficients between open probability in each group and  
 201 the relative frequency of counts equal to 2 given counts being either 1 or 2 in human cell line data

202 (c-d) An example peak with different open probabilities across various cell types and the relative  
 203 frequency of peaks with counts greater than or equal to 3 or the relative frequency of counts equal  
 204 to 2 given counts were either 1 or 2 in human cell line data. Another example was displayed in

205 **Supplementary Figure 1b-c**

206 (e) Volcano plot showing the normalized gene expression levels between cells with TSS peak  
 207 insertion counts equal to 1 or 2 and cells with TSS peak insertion counts greater than or equal to 3  
 208 in PBMC data

209 (f) Volcano plot showing the normalized gene expression levels between cells with TSS peak  
 210 insertion counts equal to 1 and cells with TSS peak insertion counts equal to 2 in PBMC data

211 (g-h) Examples of peak-gene pairs where gene expression levels are related to the TSS peak  
 212 insertion counts in PBMC data

213

214

215

216 **Methods:**

217

218 Public Datasets

219 We downloaded the following snATAC-seq datasets from public repositories: mouse kidney data<sup>13</sup>  
220 (GEO accession number GSE157079), human cell line data<sup>4</sup> (GEO accession number GSE162690),  
221 and human BMMC data<sup>15</sup> (GEO accession number GSE194122). We downloaded the 10X  
222 Genomics human PBMC data (including a snATAC-seq dataset and a sn-multiome dataset) from  
223 10X Genomics website (<https://www.10xgenomics.com/resources/datasets>).

224

225 Data QC and pre-processing

226 To remove artifacts due to data processing, we conducted QC filtering for the datasets. First, we  
227 removed peaks with very high counts ( $\geq 7$  with fragment-based counting or  $\geq 14$  with insertion-  
228 based counting) across the entire dataset, which could be associated with repetitive or potentially  
229 uncharacterized blacklist regions<sup>2</sup>. We removed potential doublet cells by the number of regions  
230 with per-base coverage greater than 3 (Ref. <sup>16</sup>). We also removed fragments with interval length  
231 smaller than 10 that are likely to be misalignment.

232

233 Processing 10X Genomics PBMC snATAC-seq data (5k)

234 The 10X Genomics PBMC snATAC-seq data (ID: atac\_pbmc\_5k\_nextgem) were used to compare  
235 the count distribution obtained from different counting methods. The peak ranges and insertion-  
236 based peak-by-cell count matrices were obtained from cellranger pipeline. The insertion-based  
237 bin-by-cell matrix was constructed by ArchR<sup>4</sup>. Bins that are accessible in fewer than ten cells were  
238 filtered. To obtain the fragment-based peak or bin count matrix, we used Signac<sup>1</sup> pipeline.

239

240 Adjusting Open Probability

241 We define “open probability” as the probability that a given genomic region is accessible for a  
242 randomly sampled cell of a given cell type. Note that this open probability does not measure the  
243 degree of openness but the probability of capturing a cell in an open state accessible to ATAC-seq  
244 assay. This probability will be governed by the temporal dynamics of nucleosome-dependent  
245 accessibility of that region for that cell type. Typical snATAC-seq data have missing data issue  
246 and are very sparse. In order to unbiasedly estimate the chromatin open probability in each cell

247 type, we considered two sources of excessive zeros in the snATAC-seq data: biological  
248 inaccessibility and technical failure to capture open state in sequencing data. We developed the  
249 following model to estimate true open proportion.

250

251 Let  $\mathbf{Z}_{g,j}^c = (Z_{g,1}^c, \dots, Z_{g,J}^c)$  be a  $J \times 1$  binary vector denoting the open chromatin status of cell  $c$   
252 that depends on group label  $g$  (e.g., cell type label). Each element in the vector,  $Z_{g,j}^c \in \{1,0\}$   
253 represents the accessibility of  $j^{\text{th}}$  genomic region (e.g., bin or peak), where the value 1 indicates  
254 open and 0 indicates close. We consider  $Z_{g,j}^c$  to be sampled from a Bernoulli distribution  
255 parameterized by  $p_{g,j}$ , the probability that a random cell of  $g$  type will be open for  $j^{\text{th}}$  region:

$$256 \quad Z_{g,j}^c \sim \text{Bernoulli}(p_{g,j})$$

257

258 In practice, the true open chromatin status  $Z$  of cell  $c$  is unobserved. Instead, due to disparity of  
259 enzyme activity and sequencing depth across cells, an open state may not be observed in the data.  
260 We introduce  $\mathbf{T}_d^c$  as a  $J \times 1$  binary vector representing the capture state of different genomic  
261 regions in cell  $c$ . This status depends on sequencing depth  $d$  for cell  $c$ . Additional experimental  
262 factors and the particular chromosomal region may also affect the status, which we ignore here.  
263 We also drop index  $d$ , since every cell is associated with particular sequencing depth. We assume:

$$264 \quad \mathbf{T}^c \sim \text{Bernoulli}(q^c)$$

265

266 for some parameter vector  $q^c$  that is a function of the cell.

267

268 Let  $Y_g^{(1)}, Y_g^{(2)}, \dots, Y_g^{(C)}$  be a random vector representing observed data with  $g \in \{1, 2, \dots, G\}$  a  
269 priori assigned cell type labels.  $Y_g^{(c)} \in \{0,1\}$  where 1 indicates open and 0 indicates close. Then  
270  $\mathbf{Y}_g^c = \mathbf{Z}_g^c \otimes \mathbf{T}_d^c$  where  $\otimes$  denote element-wise direct product (Hadamard Product).

271

272 For a given dataset  $\mathbf{y}$ , we set the loss function  $\log L(\mathbf{p}, \mathbf{q}|\mathbf{y})$  as

273

$$274 \quad \log L(\mathbf{p}, \mathbf{q}|\mathbf{y}) = \sum_{j=1}^J \sum_{c=1}^C [y_{jc} \log(p_j q_c) + (1 - y_{jc}) \log(1 - p_j q_c)]$$

275

276 In order to compute both estimators for  $\mathbf{p}$  and  $\mathbf{q}$ , we implemented a coordinate descent  
277 algorithm. This iteration stops until convergence:

278 1. Start with an initial estimate of  $\mathbf{p}^{(0)}$

279 2. For  $t = 1, 2, \dots$

280 a. Compute  $q_c^{(t)}$  by:

281 
$$q_c^{(t)} = \frac{\sum_{j=1}^J y_{jc}}{\sum_{j=1}^J p_j^{(t-1)}}$$

282 b. Update  $p_j^{(t)}$  by moment estimator:

283 
$$p_c^{(t)} = \frac{\sum_{c=1}^C y_{jc}}{\sum_{c=1}^C q_c^{(t)}}$$

284

285 *Analysis of count frequency and open probability in human cell line data*

286 The cell line data matrix was constructed by insertion-based counting method, and the maximum  
287 count was 4 in this matrix. The open probability for each cell type,  $p_g$ , was estimated with the  
288 method described above. Since the count 2 and 1 mainly represent the boundary phasing issue, we  
289 estimated the probability of observing count greater or equal to 3 given observing a non-zero count,  
290  $P_g[y \geq 3 | y > 0]$

291 
$$P_g[y \geq 3 | y > 0] = \frac{f_3 + f_4}{f_1 + f_2 + f_3 + f_4}$$

292

293 Since some peaks do not have counts that are greater than three, we only retained peaks with at  
294 least five count greater than 3, and 46,499 peaks were left. The Spearman correlation was  
295 computed between the open probability and frequency of counts greater than three. In addition, we  
296 also computed the probability of observing a count equal to 2 given the count being 1 or 2,  
297  $P_g[y = 2 | y > 0]$

298 
$$P_g[y = 2 | y = 1 \text{ or } 2] = \frac{f_2}{f_1 + f_2}$$

299 and its correlation with open probability.

300

### 301 Analysis of differentially accessible regions (DAR) in P0 mouse kidney data

302 The peak information as well as cell type annotations were obtained from the original publication<sup>13</sup>.

303 The peak-by-cell matrix was then constructed by both insertion-based and fragment-based  
304 approaches. The count correspondence is summarized in the **Supplementary Table 2**. We then  
305 picked the two most abundant cell types, nephron progenitor cells and stroma cells for the DAR  
306 analysis. Two DAR approaches, Signac<sup>1</sup> and ArchR<sup>4</sup>, were used to identify DARs. Peaks with  
307 FDR-adjusted p value  $\leq 0.05$  were regarded as DARs.

308

### 309 Analysis of count frequency and open probability in P0 mouse kidney data

310 We retained cell types with more than 600 cells to get accurate estimations of the parameters,  
311 which resulted in seven cell types. The open probability for each cell type,  $p_g$ , was estimated with  
312 the method described above. Within a cell type, assuming there are  $f_1$  cells with count 1,  $f_2$  cells  
313 with count 2 and so on, the probability of observing counts greater than or equal to 3 given  
314 observing a non-zero count is estimated by

$$315 \quad P_g[y \geq 3 | y > 0] = \frac{f_3 + \dots + f_n}{f_1 + f_2 + f_3 + \dots + f_n}$$

316

317 Spearman correlation was computed between the two quantities, and results were shown in  
318 **Supplementary Figure 2a-b**. We observed the same pattern with fragment-based counting when  
319 we compare the rank correlation between open probability and  $P_g[y \geq 2 | y > 0]$ .  
320 **(Supplementary Figure 2c)**.

321

### 322 Analysis of gene expression and different counts for PBMC data

323 The 10X Genomics PBMC sn-multiome data (ID: pbmc\_granulocyte\_sorted\_10k) were used to  
324 study the relationship between the number of insertions around TSS and its associated gene  
325 expression. We first retained peaks that overlap with  $\pm 100$  bp region around TSS and with at  
326 least five instances of counts greater than or equal to two. Then, we linked these peaks with their  
327 associated genes to form peak-gene pairs. The peak-gene pairs were then filtered by requiring the

328 non-zero expression proportion with chromatin insertion counts greater than zero to be at least  
329 10%. 3,387 such peak-gene pairs were kept for the downstream analysis.

330

331 For each peak-gene pair, we grouped the normalized gene expression levels by the insertion count  
332 in the TSS peak. Mean expression level and non-zero expression proportion were calculated for  
333 each group. Two-sided Wilcoxon Rank Sum test was then conducted between the two groups and  
334 log fold change was computed by comparing the mean expression differences.

335

336 *Analysis of gene expression and different counts for BMMC data*

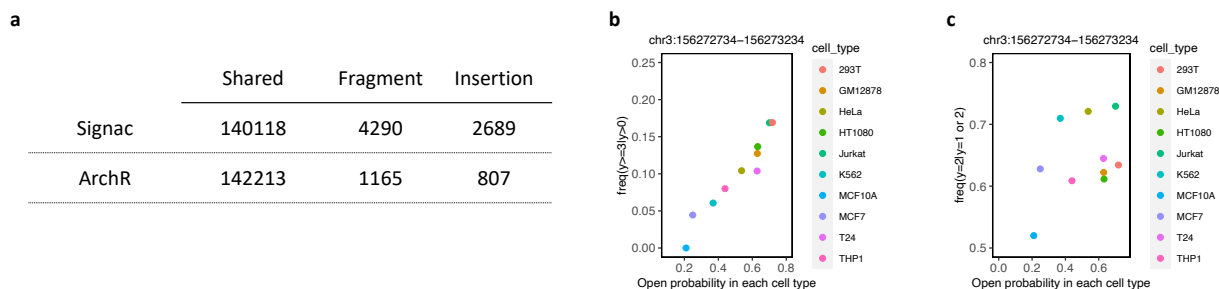
337 The BMMC dataset<sup>15</sup> was collected across multiple institutes and multiple donors with batch effect.  
338 To prevent batch effect, we focused on one donor sample that was collected at one institute (donor  
339 #2 collected from institute #1). There are 6,740 cells across multiple cell types. With the same  
340 filtration criteria as above, we retained 2,488 peak-gene pairs for our analysis. The same analyses  
341 were conducted as above and were shown in **Supplementary Figure 3a-c**.

342

343

344

345 **Supplementary Information:**



346

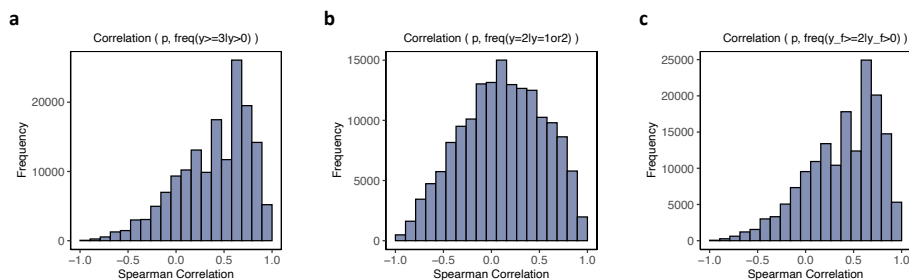
347 **Supplementary Figure 1**

348 (a) Number of significant Differentially Accessible Regions between the two most abundant cell  
 349 types, nephron progenitor cells and stroma cells with two different counting approaches and two  
 350 different pipelines

351 (b-c) An example of a peak with different open probabilities across various cell types and the  
 352 relative frequency of peaks with counts greater than or equal to 3 or the relative frequency of  
 353 counts equal to 2 given counts were either 1 or 2. Another example was displayed in **Figure 2c-d**

354

355



356

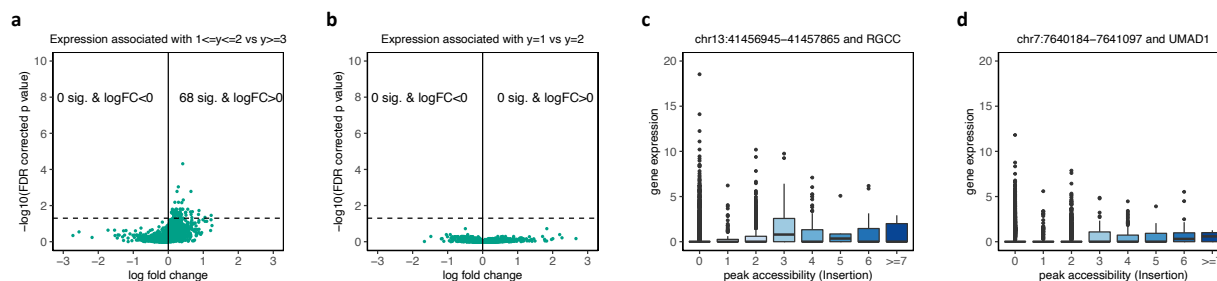
357 **Supplementary Figure 2**

358 (a) Histogram of Spearman correlation coefficients between open probability in each group and  
 359 the relative frequency of counts greater than or equal to 3 in P0 mouse kidney data

360 (b) Histogram of Spearman correlation coefficients between open probability in each group and  
 361 the relative frequency of counts equal to 2 given counts being either 1 or 2 in P0 mouse kidney  
 362 data

363 (c) Histogram of Spearman correlation coefficients between open probability in each group and  
 364 the relative frequency of counts greater than or equal to 2 with fragment-based counting in P0  
 365 mouse kidney data

366  
367  
368



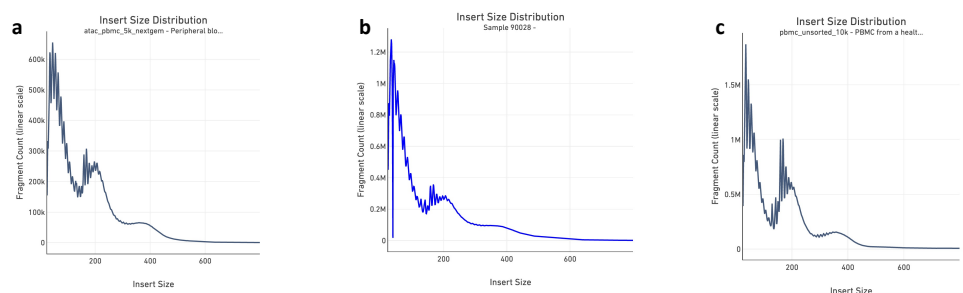
369  
370  
371  
372  
373  
374  
375  
376  
377  
378

### Supplementary Figure 3

(a) Volcano plot showing the normalized gene expression levels between cells with TSS peak insertion counts equal to 1 or 2 and cells with TSS peak insertion counts greater than or equal to 3 in BMMC data

(b) Volcano plot showing the normalized gene expression levels between cells with TSS peak insertion counts equal to 1 and cells with TSS peak insertion counts equal to 2 in BMMC data

(c-d) Examples of peak-gene pairs where gene expression levels are related to the TSS peak insertion counts in BMMC data



379  
380  
381  
382  
383  
384  
385  
386

### Supplementary Figure 4

(a) Tn5 Insert size distribution in 10X Genmoics PBMC-5k snATAC-seq dataset

(b) Tn5 Insert size distribution in P0 mouse kidney snATAC-seq dataset

(c) Tn5 Insert size distribution in 10X Genmoics PBMC-10k snMultiome dataset

### Supplementary Table 1: Frequency of counts with different counting strategies (PBMC-5k data)



387 **Supplementary Table 2: Correspondence between different counting strategies (kidney P0**  
388 **data)**

389

390

391 **Acknowledgment:**

392 This work was supported by NIDDK grant 5UC2DK126024-02 as part of the ReBuilding a Kidney  
393 (RBK) consortium.

394

395 **Code Availability:**

396 All codes used in this project including PIC algorithm are in the GitHub repository:  
397 <https://github.com/Zhen-Miao/PIC-snATAC>

398

399 **Competing interests:**

400 The authors declare no competing interests.

401

402

403

404

405

406

407