1  **Systematic evaluation of single-cell multimodal data integration for comprehensive**
2  **human reference atlas.**

3  Mario Acera-Mateos[1,2,*], Xian Adiconis[3,4,*], Jessica-Kanglin Li[1,*], Domenica Marchese[5],
4  Ginevra Caratù[5], Chung-Chau Hon[6], Prabha Tiwari[7], Miki Kojima[7], Beate Vieth[8], Michael A.
5  Murphy[9,10,11], Sean K. Simmons[3,4], Thomas Lefevre[12,13], Irene Claes[12,13], Christopher L.
6  O'Connor[14], Rajasree Menon[14,15], Edgar A. Otto[14], Yoshinari Ando[7], Katy Vandereyken[12,13],
7  Matthias Kretzler[14,15], Markus Bitzer[14], Ernest Fraenkel[9], Thierry Voet[12,13,$], Wolfgang
8  Enard[8,$], Piero Carninci[7,16,$], Holger Heyn[2,5,17,$], Joshua Z. Levin[3,4,$,#], Elisabetta Mereu[1,$,#].


9  [1]Josep Carreras Leukemia Research Institute, Barcelona, Spain.
10  [2]University of Barcelona (UB), Barcelona, Spain.
11  [3]Klarman Cell Observatory, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA.
12  [4]Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA.
13  [5]Centro Nacional de Análisis Genómico (CNAG), Barcelona, Spain.
14  [6]Laboratory for Regulatory Genomics, RIKEN Center for Integrative Medical Sciences, Yokohama, Kanagawa,
15  Japan.
16  [7]Laboratory for Transcriptome Technology, RIKEN Center for Integrative Medical Sciences, Yokohama,
17  Kanagawa, Japan.
18  [8]Anthropology and Human Genomics, Faculty of Biology, Ludwig-Maximilians Universität München, 82152
19  Planegg, Germany.
20  [9]Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA 02139, USA.
21  [10]Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge,
22  MA 02139, USA.
23  [11]Current affiliation: Osmo; New York, NY 10016, USA.
24  [12]Department of Human Genetics, University of Leuven, KU Leuven, Leuven, Belgium.
25  [13]KU Leuven Institute for Single Cell Omics (LISCO), University of Leuven, KU Leuven, Leuven, Belgium.
26  [14]Division of Nephrology, Department of Internal Medicine, University of Michigan, Ann Arbor, MI 48109,
27  USA.
28  [15]Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI 48109,
29  USA.
30  [16]Human Technopole, Milano, Italy.
31  [17]ICREA, Barcelona, Spain.
32
33  [*]These authors contributed equally.
34  [$]Senior authors.
35  [#]Correspondence: Joshua Levin (jlevin@broadinstitute.org) and Elisabetta Mereu
36  (emereu@carrerasresearch.org).
37

38  **Abstract.**
39  The integration of multimodal single-cell data enables comprehensive organ reference atlases,
40  yet its impact remains largely unexplored, particularly in complex tissues. We generated a
41  benchmarking dataset for the renal cortex by integrating 3' and 5' scRNA-seq with joint
42  snRNA-seq and snATAC-seq, profiling 119,744 high-quality nuclei/cells from 19 donors. To
43  align cell identities and enable consistent comparisons, we developed the interpretable machine
44  learning tool scOMM (single-cell Omics Multimodal Mapping) and systematically assessed
45  integration strategies. "Horizontal" integration of scRNA and snRNA-seq improved cell-type
46  identification, while "vertical" integration of snRNA-seq and snATAC-seq had an additive

47      effect, enhancing resolution in homogeneous populations and difficult-to-identify states.
48      Global integration was especially effective in identifying adaptive states and rare cell types,
49      including WFDC2-expressing Thick Ascending Limb and Norn cells, previously undetected in
50      kidney atlases. Our work establishes a robust framework for multimodal reference atlas
51      generation, advancing single-cell analysis and extending its applicability to diverse tissues.
52

53      **Introduction.**
54      Single-cell genomics is a fast-evolving field and provides tools to understand complex organs
55      in incredible detail. Each single-cell method, whether for studying RNA or open chromatin
56      (ATAC), offers a unique perspective into cell identity and function. While large-scale scRNA-
57      seq datasets are the most prevalent, transcriptional profiles can be distorted during single-cell
58      isolation, and cellular representation can be biased by incomplete tissue preparation[1,2].
59      Complementing scRNA-seq with protocols that provide full-length and nuclear RNA profiles
60      is often adopted to minimize biases and improve data accuracy. Additionally, transcriptome
61      profiling alone may not capture the full spectrum of cellular diversity while incorporating other
62      modalities, like epigenetic measurements (e.g., open chromatin and DNA methylation) can
63      enhance the completeness and resolution of an organ atlas. As such, researchers have
64      increasingly employed multimodal, both paired and unpaired, single-cell measurements to
65      generate reference atlases of human samples, as exemplified by large-scale efforts like the
66      Human Cell Atlas (HCA)[3]. These approaches have significantly advanced the characterization
67      of cellular heterogeneity and molecular diversity by enabling cross-validation of findings and
68      elucidating functional relationships between different molecular layers[4,5,6,7]. However, the need
69      for systematic benchmarking to evaluate the unique contributions of distinct modalities in
70      resolving cell types and states has also emerged. While recent efforts in benchmarking single-
71      cell approaches have primarily focused on comparing experimental methods within specific
72      classes of single-cell assays, such as single-cell RNA sequencing (scRNA-seq)[8,9,10] or single-
73      cell ATAC sequencing (scATAC-seq)[11], they have not addressed their integration. These
74      studies have provided valuable insights into the efficiency and technical biases of each method
75      in detecting molecules within cell types, offering a relevant but partial view of their unique
76      ability to characterize complex tissues. Additionally, computational tools have been compared
77      in their performance to correct batch effects and preserve biological variability in the
78      integration of scRNA-seq datasets alone[12] or in combination with scATAC-seq data[13].
79      However, these comparisons have not examined the advantages of joint integrations for
80      defining cellular types based on multimodal data.
81      In the present study, we extend prior investigations by exploring the power of integrating
82      multimodal single-cell omics data, offering a broader perspective on their combination in
83      representing tissue complexity accurately and comprehensively. Using the kidney as an
84      emblematic example of a complex organ, we conducted a multicenter comparative study,
85      analyzing 33 samples, including 25 from donors with multiple samples (paired) and 8 from
86      donors with only one sample (unpaired). This included single-cell 3' and 5' transcriptomic data,
87      as well as multiomic single-nucleus RNA (snRNA-seq) and chromatin accessibility (snATAC-
88      seq) data from the same cells. Additionally, we processed a subset of samples with Smart-seq2
89      and single-cell nucleosomal occupancy, DNA methylation and transcriptome sequencing
90      (scNMT-seq) to further support the multimodal characterization. We systematically evaluated:

91  i) how each modality contributes to the identification and characterization of specific cell types
92  and ii) how the combination of multimodal data deepens our understanding of distinct cell types
93  compared to any single modality. Our analysis includes three main components: 1) data
94  harmonization, which standardizes and preprocesses single-cell data from different sources to
95  ensure compatibility and minimize biases; 2) multimodal integration, which integrates and
96  analyzes diverse data types; and 3) evaluation and benchmarking, which assesses the
97  performance and efficiency of the integrated data in defining a wide range of cell types. Here,
98  we introduce scOMM, a novel analytical framework designed to ensure a consistent projection
99  of diverse data types onto a reference dataset while providing interpretable feature importance
100  scores for cell type classification. Built on a supervised neural network strategy, scOMM learns
101  cell identities and facilitates benchmarking across modalities, offering insights into their
102  relative contributions to cell type identification and marker feature detection. We combine
103  scOMM with unsupervised approaches, such as graph embedding followed by clustering, to
104  explore the intrinsic structure of the data, offering supplementary insights into cell type
105  diversity and relationships. To fully dissect distinct scenarios in cell atlas projects, we
106  leveraged the established concept of anchors in single-cell multiomics integration[14], which use
107  shared elements to align and integrate different datasets. Using this framework, we compare
108  different types of integrations: matching the same type of data from different sources (i.e.,
109  horizontal), combining different kinds of data from the same cells (i.e. vertical), and mixing
110  different types of data from different cells, as well as incorporating paired cells when available
111  (i.e. diagonal or mosaic).
112  Altogether, our work elucidates the synergistic value of integrating distinct single-cell data
113  types, offering an integration framework that results in a more robust definition of kidney cell
114  types and states. This framework can serve as a set of standards and guidelines for integrating
115  multimodal data, highlighting the specific features of integration depending on the
116  characteristics of the cell types being analyzed, and providing a roadmap for similar analyses
117  in other complex tissues.
118
119  **Results.**
120
121  **Generation of a benchmarking dataset for single-cell multimodal characterization of the**
122  **human renal cortex.**
123  Our study started with kidney tissue collection in the operating room from donors undergoing
124  nephrectomy. Multiple aliquots from each normal tissue sample were frozen under different
125  conditions selected to be appropriate for a variety of assays to profile them (see Methods).
126  These aliquots were thawed and processed for 3' scRNA-seq, 5' scRNA-seq, and multiomics
127  (snRNA-seq and snATAC-seq), which were central to our integrative framework (Fig. 1, Supp.
128  Fig. 1). In addition, we also processed a small number of samples with Smart-seq2[15] and
129  scNMT-seq[16], which are reported here separately because these assays were performed on a
130  limited scale, generating insufficient data for our main analyses.
131  This generated a multimodal benchmarking dataset for renal cortex (mBDRC) characterization,
132  encompassing 119,744 high-quality cells/nuclei (see Methods) from 19 healthy donors (Fig.
133  2A-B, Supp. Fig. 2), representing diverse sex, age, BMI, and other clinical renal-associated
134  characteristics (Table 1, Table 2). The mBDRC was anchored on previously established human

135   kidney references[5,17] through which we have obtained two main layers of annotations, 12 broad
136   cell types (referred to as L1 annotation) and 39 cell types/states (referred to as L2 annotation),
137   encompassing epithelial, endothelial, immune and other stromal subtypes (Fig. 2C, Supp. Fig.
138   3A-B, Table 3). This comprehensive dataset enabled us to perform a comparison and
139   integration of multiple single-cell genomics data types, assessing technical biases, estimating
140   statistical power, and exploring complementarity and redundancy in cell type/state
141   identification, described in the following sections. We first harmonized and annotated cells
142   from each technology individually (see Methods). Subsequently, we integrated the distinct
143   protocols and modalities using the different integration strategies: horizontal, vertical, and
144   diagonal/mosaic. For the mBDRC mosaic integration (Fig. 2B), we utilized MultiVI[18] to
145   generate a lower-dimensional latent space and correct for batch effects. We conducted a series
146   of iterative benchmarking investigations throughout the analysis pipeline (see Methods),
147   assessing the integration of individual technology samples and across sn/scRNA technologies,
148   and ultimately the integration of snRNA with snATAC data in both paired and unpaired
149   settings. These steps are discussed in greater detail in subsequent sections. This stepwise
150   evaluation provided immediate feedback on the performance and consistency of the integrated
151   datasets at each stage, ensuring that MultiVI was a well-validated choice. Additionally, to
152   harmonize cell identities between technologies, we developed scOMM (see Methods), a
153   machine learning tool specifically designed for supervised cell-type annotation and
154   benchmarking of multimodal single-cell data. ScOMM mapped each dataset onto the external
155   references by leveraging their curated cell annotations. Unlike existing tools such as Harmony[19]
156   and Seurat[20], which focus on general reference mapping or integration, scOMM also addresses
157   the unmet need for systematic benchmarking by facilitating direct comparison of distinct data
158   types, offering flexible parameter tuning, and evaluating feature importance for cell-type
159   predictability.
160   Analyzing differentially expressed genes across the dataset revealed a consensus set of markers
161   for each cell type (Fig. 2D), which demonstrated consistent co-expression patterns across
162   platforms and modalities. Specifically, homogeneous populations, such as podocytes (POD),
163   showed consistent detection and co-expression of canonical markers like *NPHS1* and *NPHS2*
164   across all platforms, highlighting the robustness of these markers in defining well-
165   distinguishable cell types. By contrast, populations such as distal convoluted tubules (DCT),
166   connecting tubules (CNT), and principal cells (PC) displayed overlapping or gradient-like
167   expression patterns, reflecting a spectrum of cellular states. This heterogeneity adds complexity
168   to capturing these markers consistently across different technologies and emphasizes the value
169   of integrating them. Indeed, comparing cell type-specific markers after downsampling cells
170   and unique molecular identifiers (see Methods), showed that different technologies captured
171   different markers more effectively (Fig. 2E, Supp. Fig. 3C, Table 4). For instance, *PARD3* and
172   *COL4A5*, which are crucial for maintaining POD structure and function[21,22], were identified as
173   significant markers exclusively in the snRNA-seq data. Likewise, *SLC25A4*, a marker vital for
174   the metabolic activity of CNT cells, was uniquely detected in the 5' scRNA-seq dataset.
175   Additionally, we performed Smart-seq2 and scNMT-seq analyses on four and two samples,
176   respectively, yielding 597 and 1245 high-quality single-cell transcriptomes (Supp. Fig. 4A).
177   Despite the limited sample size, these datasets align well with and reinforce the computational
178   and manually curated markers identified using other protocols (Supp. Fig. 4B-C). For scNMT-

179 seq, a small set of PT and DCT-CNT cells identified from transcriptomic data (Supp. Fig. 4A)
180 was selected to match the DNA methylomes of these same cells (Supp. Fig. 5A). Since GpC
181 methylation in scNMT-seq marks open chromatin, we further investigated the list of genes with
182 open regions based on snATAC-seq data (see Methods; Table 3), which confirmed increased
183 DNA accessibility at transcription start sites of those genes (Supp. Fig. 5B). Next, we
184 investigated GpC methylation in promoter regions in genes detected as transcribed or non-
185 transcribed on a pseudo-bulk level per cell type and provided evidence (see Methods) that non-
186 transcribed genes were less accessible than transcribed genes ($Pr(>\chi^2) < 0.01$, Supp. Fig. 5C).
187 Endogenous CpG methylation showed an overall reduction in promoter methylation which was
188 more pronounced in transcribed than non-transcribed genes ($Pr(>\chi^2) < 0.01$, Supp. Fig. 5D-E).
189

**Empirical power analysis quantifies cell-type specific advantages of scRNA-seq and snRNA-seq.**

192 Before proceeding to the integration of different data types, we performed an empirical power
193 analysis to quantify the advantages of scRNA-seq versus snRNA-seq, two widely used
194 techniques in generating human organ atlases. Our mBDRC dataset was particularly well-
195 suited for this analysis due to the unique advantage of profiling several donors using both
196 protocols, enabling quantitative comparisons and specific recommendations for these assays.
197 Using the scOMM broad cell type (L1) annotations and additional filtering steps to achieve
198 balanced datasets for the different comparisons (see Methods), we first analyzed how much
199 variance in expression levels can be explained by cell type (n=14), donor (n=5) and assay (3'
200 scRNA-seq and snRNA-seq). While assay (median 24%) and cell type (13%) explained most
201 of the variance, also the interaction of assay and cell type explained a considerable proportion
202 (6%; Supp. Fig. 3D). Hence, the different assays measured considerably different expression
203 levels in different cells, exceeding differences among donors (3%). Nevertheless, gene
204 expression differences among cell types were large and consistent enough to lead to similar
205 distances among them, as reflected by the congruent cell type trees of scRNA-seq and snRNA-
206 seq (Supp. Fig. 3E).
207 To compare the two assays on the gene level, we estimated mean, variance and detection rate
208 of expression levels weighted by the number of cells. On the cell level, we estimated the
209 proportion of genes per cell with non-zero expression levels weighted by number of reads and
210 cells as well as the homogeneity of the cell type cluster based on its purity and silhouette (see
211 Methods). We calculated values for all 14 cell types (Supp. Fig. 6) but restricted the
212 visualization to the four main cell types for simplicity (Supp. Fig. 3F-G). In addition, we
213 assessed reproducibility of metrics across our sampling cohort by calculating the Kolmogorov-
214 Smirnov distance between assays within the same donor and between donors per assay (Supp.
215 Fig. 6). We opted to rank the two assays based on eight summary statistics for the four main
216 cell types (Fig. 2G). Intriguingly, in most cases snRNA-seq outcompeted scRNA-seq in terms
217 of summary statistics as well as reproducibility. This is quantitative evidence for kidney tissue
218 and probably also for tissues with similar cellular fragility that a) more information is gained
219 when using snRNA-seq compared to scRNA-seq at both gene and cell levels; and b) that
220 scRNA-seq and snRNA-seq modalities yields complementary information. Hence, their
221 integration can improve cell-type characterization and provide a more comprehensive
222 understanding of kidney complexity.

**Horizontal integration of sn/scRNA-seq data highlights differences in cell type identification accuracy across assays.**

To achieve a comprehensive transcriptional characterization of the cell types in the human kidney, we performed horizontal integration of sn/scRNA-seq datasets (Fig. 3A) using scVI[23], which demonstrated the highest benchmarking scores (see Methods) for integration both within and across protocols (Supp. Fig. 7A, Supp. Fig. 8A). This integration approach effectively reduced technical noise (as indicated by the highest batch correction scores) and enhanced the detection of true biological signals (reflected in the highest biological conservation scores). Then, clustering produced cell-type annotations that were consistent with scOMM across datasets at both broad (L1) and fine (L2) annotation levels (Supp. Fig. 7B-C, Supp. Fig. 8B-C). To evaluate the relative contribution of each protocol to the transcriptional definition of each cell type, we first harmonized cell type annotations (see Methods) through the inspection of cluster-specific markers (Supp. Fig. 7B). The integrated data then served as a biological reference or "ground truth" for each cell type/subtype in subsequent comparison. We, thus, compared the annotations to the reference-based (scOMM) cell mapping obtained independently and from each sn/scRNA-seq dataset (Supp. Fig. 7D). By evaluating area under the curve (AUC) scores for each cell type/subtype (see Methods), stratified by protocol (Fig. 3B), we quantified how each dataset independently contributed to cell type definition, leveraging the strengths of integrated data while minimizing biases. High AUC scores for a specific cell type suggest strong agreement between the unsupervised clustering-based annotations and the supervised reference-based annotations (i.e., obtained by scOMM), indicating effective capture of the transcriptional cell type signature. Conversely, low AUC scores indicate discrepancies, reflecting a lower effectiveness for cell type identification. We focused on high specificity values (>0.9; see Methods) to prioritize the accurate identification with minimal false positive rates. We found that PT subtypes (e.g., PT-S1, aPT, PT-S3) were more accurately characterized by snRNA-seq, as indicated by higher AUC scores. Conversely, TAL groups (C-TAL, aTAL) were better identified with 3' scRNA-seq data, while transcriptionally continuous and low-abundant populations (e.g., MD, DCT1, CNT) were better captured by 5' scRNA-seq data (Fig. 3B).

To further assess the resolution of the cell-type representations generated by the graph embedding and latent space approximation, we used the Mahalanobis distance (see Methods), which accounts for cell-type distribution variance and correlation structure, providing a robust metric of cluster separation in high-dimensional data. Calculating Mahalanobis distances independently for each dataset (before and after sample integration) revealed that individual protocols showed similar cell-type relationship, with consistently higher distance scores observed across all populations following the integration of samples, indicating improved performance in separating distinct cell types (Fig. 3C). While integration across protocols did not necessarily outperform individual datasets in all cases, it mitigated protocol-specific deficits, resulting in a more balanced and robust representation of cell populations, such as in IC, PT, and DCT, where distances were consistently greater than the minimum observed in individual protocols or near their average (Fig. 3D). Complementing this, the binary cell-type Local Inverse Simpson's Index (b-cLISI; see Methods), which measures biological conservation and local diversity after integration, provided higher scores post-integration (Fig.

267   3E, values above the diagonal). This improvement was particularly notable in continuous and
268   transitional populations, such as CNT, DCT and PC. Here, the DCT population, exhibited a
269   clear improvement in both Mahalanobis (Fig. 3D, summary score 0.98) and b-cLISI (Fig. 3E,
270   with values >0.8 after integration observed in snRNA and 5' scRNA), reflecting its accurate
271   separation from CNT after integration despite being a challenging case due to their biological
272   proximity. Together, these results show how integration enhances both global separation (as
273   assessed by Mahalanobis distances) and local diversity (as captured by b-cLISI), resulting in a
274   comprehensive and nuanced view of cell population relationships.
275
276   **Integration of sn/scRNA-seq modalities boosts marker discovery and cell type**
277   **predictability.**
278   Identifying marker genes that exhibit variable expression across populations is essential for
279   classifying cell types in sn/scRNA-seq. However, this process can be biased by technical noise,
280   dropouts, and batch effects, which may mask true biological signals and hinder accurate marker
281   identification. By combining data from different sn/scRNA-seq protocols, we investigated
282   whether this could improve the detection of cell type-specific markers. Also, we examined if
283   this improvement varied depending on the characteristics of each cell type. We trained a
284   scOMM model for each sn/scRNA-seq dataset, downsampling to 6,000 nuclei/cells per dataset
285   to reduce biases associated with varying cell numbers. The models were trained using all
286   potential marker genes (see Methods) identified independently within each dataset among their
287   differentially expressed genes. To determine which genes were most important for
288   distinguishing cell types, we applied a feature importance (FI) approach (see Methods) that
289   perturbs the data and evaluates the impact on the overall cell type prediction. This method
290   simulated the elimination of a gene's contribution to the model by replacing the expression
291   values of each tested gene with zeros. Genes with higher FI values were those that increased
292   the model's accuracy. We set a threshold of 10 for FI to indicate significance, meaning that
293   when perturbing a gene resulted in a reduction of at least a relative 10 points in the model's
294   accuracy, the gene was considered critical for distinguishing cell types. Next, we also trained
295   a model (see Methods) to assess the impact of combining data from different protocols on
296   marker detection. The detected markers were compared to the cell-type markers listed in the
297   HCA kidney reference, which integrates analyses from both single-cell and single-nucleus
298   RNA-seq data, thereby minimizing assay-specific bias. We observed substantial differences in
299   marker detection rates across protocols and cell types (Fig. 3F), with snRNA-seq consistently
300   demonstrating the highest detection rates for several cell types, particularly POD, PC, and
301   DCT, at around 0.5. This may be due to snRNA-seq avoiding issues associated with cell
302   dissociation. The integrated dataset demonstrated strong performance across most cell types,
303   again balancing the performance of individual dataset (Supp. Fig. 9A). Notably, the 5' scRNA-
304   seq dataset performed comparably well in several cell types, particularly in POD and TAL,
305   where it showed the highest detection rates. On the other hand, the 3' scRNA-seq dataset
306   showed lower rates overall, with a very low score surprisingly observed in POD. The marker
307   detection performance also highlights the unique, common, and partially shared markers across
308   the different protocols (Supp. Fig. 9B-C), where snRNA-seq and the integrated dataset resulted
309   in a higher proportion of unique markers, particularly in POD (e.g., *PCOLCE2* significantly
310   detected only in snRNA-seq), epithelial cells (e.g., *SLC7A13* and *SLC5A8* specific to PT-S3,

311   *PTP4A1* as a marker of CNT-PC significantly detected exclusively in the integrated data), and
312   endothelial cells (e.g., *TBX1* in the integrated dataset and *PECAM1* in snRNA-seq) (Table 5).
313   Importantly, assay-specific markers that were also detected as significant in the integrated
314   model (Fig. 3G, Supp. Fig. 9D-E, referred to as "both" and highlighted by the yellow line in
315   the figures) exhibited higher FI scores compared to those from their original protocol,
316   indicating their relevance in distinguishing cell types. Furthermore, we compared the AUC
317   scores of the integrated model across all cell types when predicting cells from each standalone
318   assay-specific dataset (Fig. 3H, Supp. Fig. 9F-G), confirming that integrating data from
319   multiple protocols not only improved marker detection but also provided a more robust
320   framework for cell-type classification.
321
322   **Vertical integration of snRNA- and snATAC-seq data enhances cell subtypes**
323   **identification beyond single-modality approaches.**
324   To perform vertical integration, we analyzed 11 multiome samples, generating simultaneous
325   snRNA-seq and snATAC-seq data from 37,717 high-quality nuclei (Supp. Fig. 2; see
326   Methods). To harness the full potential of both modalities, we employed two approaches
327   specifically designed for single-cell multiomics analysis: multimodal spectral (multi-spectral)
328   (Fig. 4A) as implemented in SnapATAC2[24] and the Weighted Nearest Neighbor (WNN) from
329   Seurat[7]. While both methods enable the generation of a joint embedding of RNA and ATAC
330   data, WNN assigns cell-specific modality weights, prioritizing the most informative data type
331   for each cell. This resulted in predominantly RNA-weighted data across all cell types (Supp.
332   Fig. 10A), diminishing the contribution of snATAC-seq data to overall cell profiling. In light
333   of this limitation, where RNA data was predominantly prioritized over ATAC data, we
334   conducted additional analyses to evaluate the specific contribution of snATAC-seq data in
335   renal cortex samples. We harmonized snATAC-seq samples (see Methods) similarly to RNA
336   horizontal integration and assessed the performance of four widely used computational
337   pipelines: SnapATAC2[24] (referred to as Spectral MNN), Signac[25] (referred to as LSI
338   Harmony), PeakVI[26], and PoissonVI[27] (Supp. Fig. 11A). In agreement with recent
339   benchmarking studies[28], Spectral MNN best captured the complex structure of the kidney
340   cortex, even at finer L2 resolution (Supp. Fig. 11B). Unsupervised clustering further supported
341   this, showing a high level of agreement between snATAC cluster annotations and their RNA
342   counterparts, with minimal mismatches in rare populations such as macula densa (MD) cells
343   and transitional states like CNT-PC and PT-S1/2 (Supp. Fig. 10B-C). Based on these findings,
344   we first optimized WNN by combining the best snRNA-seq and snATAC-seq embeddings
345   derived from their respective horizontal integrations (Supp. Fig. 10D; See Methods). We then
346   compared this WNN result with multi-spectral for vertical integration. Despite this
347   optimization, a comparison of silhouette scores revealed that multi-spectral consistently
348   achieved higher scores across several subtypes, resulting in greater average deviance, which
349   indicates better-separated cell groups (Fig. 4B, Supp. Fig. 10E). This trend was particularly
350   evident in populations such as POD and PEC, which are typically well-resolved through
351   clustering. We hypothesize that WNN underrepresented complementary ATAC data, limiting
352   its ability to capture subtle cellular differences, whereas multi-spectral enhanced sensitivity to
353   chromatin accessibility, resulting in stronger cluster separability. Indeed, a comparison of
354   average silhouette width scores from individual assays (RNA and ATAC) with those from

vertically integrated data (i.e., joint RNA and ATAC) revealed a complementary effect, approximating an additive trend with the multi-spectral method (Fig. 4C, top table). This pattern was especially pronounced in homogeneous populations and extended to more complex subtypes, such as C-TAL and IC-B, underscoring the value of integrating both modalities. Rather than simply balancing the contributions of RNA and ATAC, the multi-spectral approach appears to optimally combine them, enhancing cell type identification by preserving the complementary strengths of each technology. At the local level, b-cLISI scores further demonstrated multi-spectral superiority in preserving neighborhood structure (Fig. 4C, bottom table). Higher b-cLISI values in challenging-to-identify cell states, such as aPT, CNT-PC, and DCT1, confirmed the benefits of multimodal integration. Altogether, these findings highlight the importance of leveraging both transcriptomic and chromatin accessibility data in parallel to achieve a more nuanced understanding of cellular heterogeneity, particularly in complex tissues.

**Evaluation of the contribution of chromatin accessibility to cell identity prediction and profiling consistency with snRNA data.**

To further examine the similarity between RNA and ATAC data and assess how well snATAC-seq captures regulatory features reflective of transcriptional identities, we trained two scOMM models (see Methods). These models transferred cell identities from snRNA to snATAC using either gene activity (i.e., peaks-associated genes are used as features for the model, also referred to as Cross-Modality, CM) as a proxy for gene expression or through bridge integration (referred to as Bridge), which directly incorporates the accessible peaks as model features. Additionally, we trained a third model (referred to as Integrated) that uses both types of features as input data to determine if combining the distinct features derived from each modality improves cell-type detection. In all models, cell type annotations from snRNA data served as the ground truth for performance evaluation. Model accuracy, which was measured by AUC, was highest for the Bridge model across most cell types, while the Integrated model displayed comparable performance for specific subtypes, such as PT, C-TAL, IC-B, DCT1, and PEC (Fig. 4D). Notably, both the Bridge and Integrated models improved sensitivity and classification rates compared to the CM model, whereas specificity remained similar across all three approaches (Supp. Fig. 11C-D). The drop in performance for the CM model may be attributed to discrepancies between gene activity (reflecting chromatin accessibility) and actual gene expression (Supp. Fig. 11E), as well as the lack of non-coding region information in snRNA-seq data. We also examined the contribution of the two different modalities to cell identity prediction by comparing the significance of peaks and/or genes across models. In the unimodal models (CM and Bridge), roughly half of the features showed high importance for cell-type prediction (Fig. 4E). In the Integrated model, which used both genes and peaks, peaks appeared to play a more dominant role than genes, likely due to the greater number of peaks relative to genes in the feature space. Interestingly, some features that were not significant in the unimodal models became relevant, exhibiting high feature importance in the integrated model (represented in red in the Genes + Peaks bar, Fig. 4E). However, despite the prominence of peaks and newly significant features in this model, overall cell type prediction accuracy did not outperform that of the Bridge model (Fig. 4D), suggesting that the added complexity, including the higher dimensionality introduced by peaks, may introduce redundancy or noise

399  without providing substantial improvements in predictive power. Therefore, while integrating
400  both data types uncovers new relationships, it does not necessarily enhance cell type
401  predictability.
402
403  **Global data integration improves detection and characterization of rare cell types by**
404  **refining signals across modalities.**
405  The combination of distinct single-cell data types can be approached through diagonal and
406  mosaic integrations. In diagonal integration, distinct features (e.g., RNA and ATAC) are
407  combined from unpaired cells, while in mosaic integration, these features are measured in
408  paired or unpaired cells, providing a more comprehensive view. It is important to investigate:
409  i) whether paired data significantly improves cell-type identification and characterization
410  within global mosaic integration, and ii) whether integrating all data types in an unpaired
411  fashion outperforms simpler combinations. To integrate all the data, we utilized MultiVI[18] and
412  GLUE[29] (Supp. Fig. 12A-B), two models designed for single-cell multiomics integration.
413  MultiVI is suited for mosaic integration, where paired data are leveraged for optimized
414  integration, while GLUE handles diagonal integration using unpaired data. Interestingly, no
415  significant differences were observed between these two approaches at both the local level,
416  referring to the b-cLISI score (Fig. 5A), and the global level, corresponding to the silhouette
417  score (Supp. Fig. 12C), with cell annotations obtained using scOMM for each individual
418  dataset. Furthermore, global multimodal integration (mosaic and diagonal) did not outperform
419  other multimodal combinations (Fig. 5A) in cell type identification, as reflected by the
420  comparable cell-type b-cLISI scores in the renal epithelial compartment. This suggests that,
421  depending on the heterogeneity of cell types or states, vertical or horizontal integrations may
422  already offer sufficient resolution for accurate identification. For example, in POD, PEC, and
423  IC subtypes, silhouette and b-cLISI values were comparable or slightly higher in the horizontal
424  RNA and vertical integrations, indicating subtle but meaningful improvements in resolving
425  these cell subtypes. However, in adaptive TAL cells, particularly in aTAL1, global integrations
426  demonstrated better local resolution, as evidenced by increased b-cLISI scores, which
427  improved significantly from 0.45 in vertical integration to 0.61 in the mosaic integration. We
428  hypothesized that this adaptive state is better resolved in 3' scRNA-seq datasets, displaying
429  higher b-cLISI score in the unimodal integration (b-cLISI=0.55, Fig. 5A), suggesting that 3'
430  scRNA-seq captures biological variations that are less effectively represented in other
431  modalities. Therefore, we re-clustered the TAL population and identified a subtype of aTAL1,
432  referred to as aTAL1_0, which was predominantly enriched in cells from the 3' scRNA-seq
433  data (Fig. 5B). By analyzing genes enriched in this cluster, we found *WFDC2* to be the top
434  marker for this population (Fig. 5C, Supp. Fig. 13A). *WFDC2* was also expressed in aTAL2
435  and MD cells, though at lower levels in both 3' and 5' scRNA-seq, but it was absent in snRNA-
436  seq data in these populations (Supp. Fig. 13B). The identification of this *WFDC2*-marked
437  subpopulations is intriguing, as *WFDC2*, which encodes Human Epididymis Protein 4 (HE4),
438  has been proposed as a serum biomarker for lupus nephritis and chronic kidney disease[30].
439  Recent studies also suggest that *WFDC2* is a suitable pan-distal nephron marker in the human
440  kidney[31]. We further compared additional markers of the aTAL1_0 population to better
441  characterize this subpopulation with respect to the other aTAL1 groups, aTAL1_1 and
442  aTAL1_2 (Fig. 5B). The differential detection of markers across these subtypes and protocols

(Supp. Fig. 13C) suggests underlying biological, and potentially functional, differences rather than solely technical biases. *WFDC2* and *B2M,* top markers of aTAL1_0*,* are often linked to active immune or stress responses[30,32].

Applying a similar approach to the stromal compartment, b-cLISI scores highlighted higher values in the 3' scRNA-seq data for the adaptive state of fibroblasts (aFIB; Supp. Fig. 12D), where reclustering of the stromal cells identified a rare subpopulation of erythropoietin (EPO)-producing cells, known as Norn cells[33], characterized by the expression of DCN, TIMP1, and CFD (Fig. 5D). Typically, EPO is produced by peritubular fibroblast-like cells in the kidney, which respond to low oxygen levels by increasing EPO production to stimulate red blood cell generation, a process crucial for maintaining oxygen delivery throughout the body[34]. Norn cells were predominantly detected in 3' and 5' scRNA-seq datasets (Fig. 5E) but were not observed in the same samples analyzed by snRNA-seq, likely due to differences in method sensitivity. Identifying EPO-producing cells is especially relevant in the context of kidney disorders, where chronic conditions such as renal hypoxia can impair EPO production, leading to anemia. Understanding these cells' behavior and regulation could provide valuable insights into treating anemia associated with chronic kidney disease and other conditions characterized by reduced oxygen levels. The differences in marker detection across these protocols highlight the biological complexity of these subpopulations, with each technology offering complementary insights into their cellular states.

**Integrative multimodal analysis reveals enhanced trait heritability in kidney cell-type-specific cis-regulatory elements.**

To further understand the impact of multimodal analysis in uncovering the molecular mechanisms driving cellular function and disease states, we focused on gene regulatory networks (Supplementary Material) and cis-regulatory elements (CREs) and their enrichment in Genome-wide association studies (GWAS) traits. GWAS have identified genetic variants associated with traits and diseases[35], many of which are enriched within CREs, highlighting their critical role in gene regulation and trait heritability. Chromatin accessibility assays are commonly used to identify accessible CREs[36]. However, many distal CREs lack epigenomic features characteristic of active enhancers[37]. While some of these elements may function as insulators[38] or silencers[39], their broader roles in gene regulation remain poorly understood, making it challenging to annotate trait-associated variants based solely on chromatin accessibility. Emerging evidence indicates that transcriptional activity at distal CREs, as identified through 5′-end RNA-seq methods[40], can serve as a marker of enhancer activity[41]. This suggests that transcriptional activity may offer a more informative and interpretable metric than chromatin accessibility for functional annotation of trait-associated variants. To demonstrate this, we systematically compared chromatin accessibility and transcription as metrics for assessing trait heritability in kidney cell-type-specific CREs. Specifically, we assessed heritability enrichment for hypertension[42] and estimated glomerular filtration rate (eGFR)[43] using accessibility-defined, cell-type-specific CREs (Fig. 5F; Table 6**;** see Methods). MD cells exhibited strong heritability enrichment for both traits, consistent with their role in regulating glomerular filtration rate and systemic blood pressure through the renin-angiotensin-aldosterone system[44]. Proximal tubular cells (PT-S1/2 and PT-S3) also showed enrichment for

487 both traits, while adaptive proximal tubular cells (aPT) were enriched only for eGFR. This may
488 reflect that genes activated in adaptive states (a potential failed-repair population)[5] are enriched
489 in eGFR-associated loci. Endothelial cells (EC-PTC, EC-GC, EC-AEA) and vascular smooth
490 muscle/pericytes (VSMC, VSMC/P, MC) were enriched for hypertension heritability but not
491 eGFR, emphasizing the role of blood vessels in hypertension[45]. These results imply that cell-
492 type-specific CREs defined by accessibility are critical for interpreting cell-type-specific
493 heritability enrichment. Next, we compared the utility of accessibility and transcription signals
494 in defining cell-type-specific CREs for heritability enrichment (see Methods; Table 5). We
495 found that either incorporating transcription signals in CRE selection (Set 2) or ranking CREs
496 by transcription (Set 3) yielded significantly higher heritability enrichment (Fig. 5G, Wilcoxon
497 rank-sum test, $p < 0.05$) compared to using accessibility alone (Set 1). This aligns with evidence
498 suggesting that transcribed enhancers are more likely to validate in functional assays compared
499 to epigenetically defined enhancers[41]. Altogether, these results highlight the advantage of
500 integrating accessibility and transcription signals to enhance the sensitivity and interpretability
501 of trait heritability analyses.
502
503 **Discussion.**
504 To fully harness the potential of multimodal data integration for understanding cellular and
505 tissue function, benchmarking experimental methods and computational tools is crucial.
506 Setting reliable standards helps ensure accuracy and shows the strengths and weaknesses of
507 each modality. In this study, we investigated the integration of distinct single-cell protocols
508 and multimodal data to create a detailed human organ reference focused on the kidney's cortex.
509 We aimed to address technical biases, improve tissue characterization, and assess how each
510 single-cell data type contributes to identify cell types and states in both matched and unmatched
511 samples. We also evaluated their complementarity and robustness, along with the
512 reproducibility of data in samples from the same donors. The resulting mBDRC serves as a
513 detailed reference for this complex organ, capturing a range of heterogeneous populations,
514 from well-defined, homogeneous cell types to continuous, challenging-to-distinguish states.
515 This diversity enabled us to explore key features relevant to cell type identification and marker
516 detection across multiple integrative scenarios by using two complementary classes of
517 approaches: supervised (i.e., label transfer-based, via scOMM) and unsupervised (i.e., graph
518 embedding-based) methods (Fig. 6). Not all integration strategies were evaluated using both
519 label transfer-based and clustering-based approaches. For instance, diagonal and mosaic
520 integrations were only assessed with embedding-based methods, reflecting methodological
521 constraints and the differing applicability of evaluation metrics across integration scenarios. In
522 embedding-based approaches, both protocol-specific and multimodal integrations consistently
523 demonstrated a clear positive impact on cell type identification, with RNA modalities
524 outperforming ATAC in most cell compartments. Among the integration strategies, mosaic
525 integration (i.e., including paired or unpaired data types) achieved the highest overall scores,
526 while diagonal (i.e., all data are used but multiomics data are used as unpaired data) and
527 horizontal integrations (same data type across samples) outperformed vertical integration
528 (different data types within the same samples). Supervised classification using scOMM further
529 emphasized the critical role of external references for accurate cell type prediction, as
530 evidenced by the higher performance scores achieved for cell type identification compared to

531  embedding-based approaches. This underscores the importance of reference-guided annotation
532  in achieving robust and consistent cell type classifications. In this regard, the integration of
533  diverse sn/scRNA-seq data demonstrated significant advantages in precisely identifying cell
534  states, as different protocols contribute in varying ways to this process. For instance, snRNA-
535  seq data provided higher resolution for proximal tubules (i.e., PT-S1/2, aPT, PT-S3), 3'
536  scRNA-seq excelled in resolving TAL subtypes, while 5' scRNA-seq data were more effective
537  in identifying smaller subpopulations, such as MD cells and subtypes within distal convoluted
538  tubules and collecting ducts (Fig. 3C). Employing both horizontal and vertical integration
539  strategies substantially improved data resolution, allowing for more accurate cellular mapping
540  (Fig. 4C). Methods like multi-spectral integration preserved unique information from both
541  RNA and chromatin accessibility data (Fig. 4B-C), offering refined insights into specific cell
542  subtypes, such as podocytes (Fig. 4C, Fig. 5A), which exhibit distinct transcriptional and
543  epigenetic profiles that make them particularly well-suited to benefit from multimodal
544  integration.

545  A key outcome of multimodal global integration was its ability to detect rare and clinically
546  significant cell populations. For instance, we identified a unique TAL cell subpopulation
547  expressing *WFDC2*, previously suggested as a lupus biomarker, detectable only in scRNA-seq
548  data (Fig. 5B-C, Supp. Fig. 13A). Notably, we uncovered the very rare erythropoietin-
549  producing Norn cells within the adaptive fibroblast population (Fig. 5D-E), a finding with
550  potential therapeutic relevance for addressing anemia in kidney disease. These results clearly
551  show the potential of combining distinct protocols to identify functionally and clinically
552  relevant cell types that might be missed using a single modality.

553  Beyond cell type identification, integrating transcriptional and chromatin accessibility data
554  revealed key regulatory mechanisms underlying trait heritability. We showed that
555  transcriptional activity at CREs enriched in GWAS-associated traits more accurately reflects
556  the functional relevance of trait-associated variants than chromatin accessibility alone (Fig.
557  5G). This was particularly evident for hypertension and eGFR, where heritability was enriched
558  in CREs specific to proximal tubular, endothelial, and vascular smooth muscle cells. Overall,
559  this work provides a systematic framework for creating comprehensive tissue atlases,
560  establishing a valuable precedent for multimodal approaches that enable biologically
561  meaningful insights into complex organs. The integration methodologies applied here
562  demonstrate that multiomics data integration can significantly enhance the resolution and
563  accuracy of cell type identification. Additionally, we are confident that our mBDRC will serve
564  as an important resource for the scientific community, advancing both kidney research and
565  next-generation protocols and tools for large-scale single-cell studies, such as those in the
566  Human Cell Atlas.

567

**Study limitations.**

569  This study establishes a robust framework for integrating multimodal single-cell data and
570  evaluating their individual and combined contribution to kidney cell biology. However, there
571  are limitations that need to be considered. Some experimental protocols, including Smart-seq2
572  and scNMT-seq were not fully optimized for frozen kidney tissue, potentially affecting data
573  quality and restricting their use to exploratory analyses. Similarly, we were unable to generate
574  any useful data with MERFISH even with multiple attempts to optimize this method for our

575   samples (unpublished data). The limited number of matched samples across protocols also
576   introduces variability from donor heterogeneity, which could confound computational analyses
577   such as gene regulatory network inference and transcription factor identification, detailed in
578   the Supplementary material. Moreover, the exclusive focus on kidney samples could limit the
579   extent of the generalizability of these findings to other tissues with distinct cellular
580   compositions and regulatory mechanisms. On the other hand, computational challenges remain,
581   particularly due to the lack of a ground truth for cell-type identification and marker detection,
582   which limits the precision of evaluations. We minimized this bias as much as possible by
583   consulting kidney experts for annotations and simulating realistic analysis scenarios to create
584   a robust framework for evaluating cell type prediction performances.
585

586 **Figure legends**

587 **Figure 1. Experimental design and computational workflow.**
588 Overview of the study, integrating multimodal single-cell data from 33 kidney cortex samples
589 from 19 matched and unmatched donors. Samples were analyzed using multiome snRNA-seq,
590 multiome snATAC-seq, 3' scRNA-seq, and 5' scRNA-seq to create a comprehensive, multi-
591 layered resource. Single-cell multiomics integration was conducted at different levels (with
592 shared and unshared features) using unsupervised and supervised methods, such as graph
593 embedding, clustering, and cross-modality bridging, to harmonize cell annotations and evaluate
594 integration performance across cell types. Quantitative metrics, including latent space- and
595 neighborhood graph-based, AUC, and feature importance, were employed to assess each
596 modality's contribution to cell type identification and the improved resolution achieved through
597 multimodal integration.

598 **Figure 2. Overview of the multimodal Benchmarking Dataset for Renal Cortex**
599 **(mBDRC).**
600 A) Color legend of the renal cortex cell types at the two levels of annotations. Level 1 (L1)
601 represents 14 broad cell types, while Level 2 (L2) provides detailed annotations of 39 specific
602 cell subtypes or states within these categories. Broad cell types and their subtypes are as
603 follows: PT (Proximal Tubule) includes aPT, PT-S1/S2, and PT-S3; TAL (Thick Ascending
604 Limb) includes cTAL, aTAL1, aTAL2, and MD (Macula Densa); PC (Principal Cells) includes
605 PC; CNT (Connecting Tubule Cells) includes CNT and CNT-PC; DCT (Distal Convoluted
606 Tubule) includes DCT1 and DCT2; IC (Intercalated Cells) includes CCD-IC-A, CNT-IC-A,
607 and IC-B; PEC (Parietal Epithelial Cells) includes PEC; POD (Podocytes); EC (Endothelial
608 Cells) includes EC-AEA, EC-GC, EC-LYM, and EC-PTC; FIB (Fibroblasts) includes MYOF,
609 and aFIB; VSM/P (Vascular Smooth Muscle/Pericytes) includes MC, VSMC, VSM/P, and
610 REN (Renin-producing cells); and IMM (Immune Cells) includes B, MAC-M2, MAST, MDC,
611 NK/T, PL, T, cDC, ncMON, and pDC. B) UMAP plot of 119,744 nuclei/cells across different
612 single-cell data modalities, including multiome snRNA, multiome snATAC, 3' and 5' scRNA.
613 Colors indicate the refined cell-type annotations (L2), with broad cell-type annotations (L1)
614 shown in the top-left squares.C) Alluvial plots displaying the distribution of cells analyzed
615 across different protocols, with colors representing L1/L2 cell-type annotations. D) Dot plot
616 showing consensus cell-type markers for the main renal cortex populations as detected by the
617 different single-cell data types. E) Upset plot illustrating the overlap of detected markers
618 identified by each technology for the primary epithelial populations, including POD, PT, DCT,
619 and TAL. F) Summary Statistics for CNT, PT, POD and TAL, illustrating the best performing
620 assay (snRNA-seq represented as a dark teal nucleus and scRNA-seq represented as a light teal
621 cell) in terms of gene metrics: mean (the minimal median standard error of mean expression),
622 variance (the minimal median variance of gene expression) and dropout (the minimal median
623 dropout of gene expression); in terms of cell metrics: detection (the maximum median
624 proportion of genes with nonzero expression per cell), purity ( maximum median purity of k
625 neighborhoods in a cell type), purity > 0.5 (maximal proportion of cells with 50% or more pure
626 neighborhoods), silhouette (maximum median silhouette per cell type) and silhouette > 0

627 (maximal proportion of cells with silhouette index above 0). Both symbols are shown when
628 they are essentially equivalent. See Supp. Fig. 6 for details underlying this summary.
629

630 **Figure 3. Horizontal integration of sn/scRNA-seq data.**
631 A) Horizontal integration of 97,125 nuclei/cells from 33 renal cortex samples. Colors represent
632 the inferred annotations from HCA reference mapping, using broad L1 (left) and fine L2 cell
633 types/subtypes (right). B) AUC scores from scOMM L2 annotations compared to clusters from
634 the horizontal integration in A. Colored dots indicate modalities, while bars represent L2 cell
635 groups, with color bars stratified by L1 cell groups. C) Pairwise Mahalanobis distances (see
636 Methods) between epithelial populations within each technology before and after integration,
637 where higher values indicate greater separation between populations. D) Mahalanobis distances
638 between epithelial-derived cell-types in the embedding space (left). Summary of the cell type-
639 specific Mahalanobis distances as overlap ratios derived from a Chi-squared distribution
640 (right), with values closer to 1.0 indicating greater separability (lower overlap) between cell
641 types, and values closer to 0 reflecting higher overlap (less separability). E) Binary cLISI (b-
642 cLISI; see Methods) scores for L1 epithelial groups calculated before and after integration. F)
643 Marker detection rate comparison based on the feature importance of L1 cell types across data
644 models. G) Feature importance score comparison between snRNA-seq and the integrated
645 dataset for markers that showed significant FI in the snRNA model but did not exhibit
646 significant FI in the 3' and 5' scRNA models. H) AUC scores for the predictability of each data
647 model in predicting the snRNA-seq data type for L1 cell-type annotations. The background is
648 divided into shaded regions that correspond to different sn/scRNA technologies. The dots,
649 connected by lines, represent the AUC scores for individual cell types, with distinct colors
650 assigned to each cell type.
651

652 **Figure 4. Vertical integration of snRNA/snATAC-seq data.**
653 A) Multiomics integration of 37,717 paired nuclei obtained with the multimodal spectral
654 (multi-spectral) approach. B) Scatter plots compare silhouette scores from optimized WNN (X-
655 axis) and multi-spectral (Y-axis) integration methods for various kidney cell types and
656 subtypes. Each blue dot represents a single cell, while the red line indicates the trend, with the
657 pink shading showing the confidence interval. Marginal histograms illustrate the distribution
658 of silhouette scores along each axis. The "Avg Dev (y)" value represents the average deviation
659 of the Y-axis (multi-spectral scores) for each cell type or subtype. C) The heatmaps compare
660 four data integration methods across kidney L2 cell populations. Average Width Silhouette
661 (AWS) are shown at the top, with darker blue indicating worse performance. B-cLISI scores
662 are shown at the bottom, where darker green represents worse clustering consistency. Red
663 boxes highlight key cell populations (C-TAL, aPT, CNT-PC, IC-B, POD) that exhibit
664 improved values with vertical integration using the multi-spectral approach. D) Bar charts
665 compare AUC scores across L2 cell types using three classification methods for predicting cell
666 types in the paired snATAC-seq data: Cross Modality (yellow), Bridge (cyan), and Integrated
667 (red). E) Stacked bar charts show the percentage of significant features (Genes, Peaks, and
668 Genes + Peaks) across the three classification models with significance defined by a score
669 threshold of 10. The colors represent the percentage of significant features for Genes (yellow),

670 Peaks (cyan), and Genes + Peaks (red), while gray indicates the percentage of features that are
671 not significant (Score ≤ 10).

672

673 **Figure 5. Global multimodal integration of paired and unpaired sn/scRNA/ATAC data.**
674 A) Heatmap showing the comparison of binary cell type LISI scores (L2 level) across unimodal
675 (pink) and multimodal (blue) integration strategies for different kidney cell types. Each cell's
676 value represents the accuracy of cell type annotation, with red indicating higher accuracy and
677 blue indicating lower accuracy. Bar plots at the top show the relative number of L1 cells per
678 cell type. Cell types are listed at the bottom, and the right side provides a color key for the
679 corresponding modalities and cell types. B) UMAP plots showing the re-clustering of the TAL1
680 subtype within the TAL compartment, split by snRNA, 3' scRNA, and 5' scRNA modalities.
681 Notably, aTAL1_0 is enriched in the scRNA-seq protocols (3' scRNA and 5' scRNA), while
682 aTAL1_2 is more specific to the snRNA-seq dataset, highlighting modality-specific differences
683 in cell subpopulation detection. C) Dot plot illustrating the expression levels of key marker
684 genes (*WFDC2*, *B2M*, *TMSB10*, *ITM2B*, *CLU*, *TACSTD2*) across aTAL1 subpopulations. Dot
685 size represents the fraction of cells expressing the gene, and color intensity indicates log-
686 normalized mean expression on each group. D) UMAP plot showing the re-clustering of
687 stromal compartment cell types, including FIB (Fibroblasts), MC (Mesangial Cells), MYOF
688 (Myofibroblasts), REN (Renin-producing cells), VSMC (Vascular Smooth Muscle Cells), and
689 others. The red circle highlights the rare population of aFIB cells identified as the Norn cells.
690 E) Violin plots showing the Norn signature scores across stromal subtypes splitter by
691 sn/scRNA-seq protocols. F) Enrichment of hypertension and eGFR heritability in accessibility-
692 defined cell type-specific CREs. Set 1 cell type-specific CREs were used for heritability
693 enrichment analyses, and only cell types with an enrichment p-value < 0.1 for either trait are
694 displayed. G) Comparison of cell type-specific CREs defined by accessibility and transcription.
695 Only cell types present in all three sets and showing an enrichment p-value < 0.05 in at least
696 one set are included. Asterisks indicate significance based on the Wilcoxon rank-sum test,
697 whith * corresponding to p-values < 0.05 and ** corresponding to p-values < 0.01.

698 **Figure 6. Benchmarking summary of the integrative scenarios for distinct single-cell data**
699 **types.** Different unimodal and multimodal integration strategies were evaluated based on
700 customized scores (see Methods) for cell type identification and marker detection based on
701 supervised (transfer learning-based via scOMM) and unsupervised (embedding-based)
702 approaches. Performance scores are displayed for 6 kidney populations/compartments: PT
703 (Proximal Tubule), TAL (Thick Ascending Limb), POD (Podocytes), ST (Stromal), EC
704 (Endothelial Cells), and IMM (Immune Cells). AVG (Average) represents the arithmetic mean
705 across shown cell groups in each evaluated approach and integration scenario. Shapes (i.e.,
706 circles for individual cell types and squares for the average across cell groups) and colors
707 display columns with normalized scores, making the results comparable across different types
708 of integration, while values are the original performance scores before normalization.

709 **Authors contribution.**
710 E.M. and J.Z.L. coordinated the overall study. E.M., J.Z.L., H.H., P.C., W.E., and T.V.
711 designed the original project, secured funding, and supervised the project, with contributions

712 from E.F., M.B., and M.Kretzler. C.L.O. procured and processed human kidney tissue and
713 performed pathological assessments. E.A.O. prepared human kidney tissue for scRNA-seq and
714 oversaw quality control of its data. Single-cell experiments were conducted by X.A., D.M.,
715 G.C., R.M., T.L., P.T., and M.Kojima. Y.A. and K.V. coordinated parts of the study. E.M.
716 conceived and guided the analysis framework for the study. M.A.-M. and J.-K.L. performed
717 most of the computational analyses, with contributions from S.K.S., C.-C.H.H., B.V., M.A.M.,
718 R.M., T.L., and I.C. R.M., M.B., and M.Kretzler contributed to data interpretation. All authors
719 contributed to writing the manuscript, read, and approved the final version.

**Data and code availability.**

721 All data analysis code used in this study is available at
722 https://github.com/mereulab/Multimodal_single-cell_Benchmarking. The integrated
723 multimodal single-cell data has been deposited and will be accessible via the CELLxGENE
724 platform. ScOMM is publicly available as an R package at
725 https://github.com/mereulab/scOMM. The empirical power analysis is available at
726 https://github.com/bvieth. Raw sequencing data, including FASTQ files for gene expression
727 datasets (scRNA-seq and snRNA-seq), will be released on GEO. Additional datasets, including
728 chromatin accessibility (snATAC-seq) and multiome data, will be made available through their
729 respective repositories upon publication.

**Competing interests.**
747 H.H. is co-founder and Chief Scientific Officer of Omniscope, a Scientific Advisory Board
748 member at Nanostring and Mirxes, a consultant for Moderna and Singularity, and has received
749 honoraria from Genentech. T.V. is co-inventor on licensed patents WO/2011/157846 (Methods
750 for haplotyping single cells), WO/2014/053664 (High-throughput genotyping by sequencing
751 low amounts of genetic material) and WO/2015/028576 (Haplotyping and copy number typing
752 using polymorphic variant allelic frequencies). M.Kretzler reports grants and contracts through

**Online Methods.**

**Human Tissue Procurement.**
Fresh normal tissue from the unaffected part of surgically removed kidney of patients undergoing total nephrectomy at the University of Michigan was obtained directly from the operating theater at Michigan Medicine. Patients were enrolled in the PRECISE cohort at the University of Michigan. The PRECISE study was approved by the institutional review board (IRB) of the University of Michigan (HUM00165536, HUM00052918) as previously described[46]. Patient data were obtained from each participant's electronic medical record. Procured tissue specimens were a minimum of 1.5 cm x 1.0 cm x 1.0 cm and included all the kidney compartments (capsule, cortex, and medulla). In order to minimize cold ischemic time, tissue was obtained directly from the operating room and immediately processed on site using a 3D printed device ("PRECISE Pyramid"). Utilizing this device, nephrectomy specimens were simultaneously cut into approximately 25 cores of equal dimension that were closely related in space and resemble 16-gauge clinical kidney biopsy cores (Supp. Fig. 1). Biopsy cores were subsequently preserved across a variety of media, including CryoStor, RNALater, OCT, and LN$_2$[47,48]. Biopsy cores were stored for later use at -80°C and shipped to interrogation sites on dry ice.

**3' scRNA-seq.**
We utilized 3' scRNA-seq data sets from tumor-free kidney cortical tissue of nephrectomy specimens within the PRECISE cohort at the University of Michigan (IRB: HUM165536). Detailed tissue processing, single-cell isolation, and scRNA-seq generation protocols are available in the KPMP scRNA-Seq protocol (https://www.protocols.io/view/single-cell-rna-

795  sequencing-scrna-seq-7dthi6n). Briefly, kidney biopsies preserved in CryoStor and frozen in
796  liquid nitrogen were rapidly thawed and enzymatically dissociated into a single-cell solution
797  using Liberase TL digestion for 12 minutes at 37°C. The solution was then filtered through a
798  30 µm strainer (Miltenyi Biotec) and washed in DMEM/F12 medium supplemented with 10%
799  fetal calf serum and HEPES buffer. The cells (average 40,000 cells) were processed using the
800  droplet-based platform with Chromium Single Cell 3'chemistry (v3.1, 10x Genomics). The
801  cDNA libraries were prepared and sequenced on an Illumina NovaSeq 6000 platform,
802  generating over 200 million reads (average 25,000 reads per cell) per sample using one of the
803  following two conditions: 28 bases (Read 1), 10 bases (Index read), and 151 bases (Read 2) or
804  151 bases (Read 1), 8 bases (Index read), and 151 bases (Read 2) – see Table 2 for details. The
805  Cell Ranger (v3, 10x Genomics) pipeline was employed to extract the cell x gene matrix from
806  FASTQ files aligned to the GRCh38 genome reference (version 2020-A).

807

808  **Tissue dissociation and single-cell plate sorting.**
809  In this section, we describe the preparation of samples for 5' scRNA-seq, Smart-seq2, and
810  scNMT-seq. We placed the cryopreserved tissue containing tube in a 37°C water bath for 1
811  minute for quick thawing. We transferred the tissue to a plastic petri dish (diameter 5 cm) filled
812  with 1 ml DMEM/F12/10% fetal bovine serum (FBS) prepared with 90% DMEM/F12/HEPES
813  (Thermo Fisher Scientific, 11330057) and 10% heat inactivated FBS (Thermo Fisher
814  Scientific, 16140-071) for 10 seconds to wash off the remaining DMSO at room temperature,
815  followed by transferring the tissue to a new petri dish containing 1 ml DMEM/F12/10% FBS
816  and incubating for 10 minutes at room temperature. We prepared 1 ml dissociation media by
817  mixing 900 µl DMEM/F12 and 100 µl Liberase TL (Millipore Sigma, 5401020001, 2.5 mg/ml
818  in $H_2O$). We then minced the tissue in another new petri dish filled with 500 µl dissociation
819  media for 1 minute (or 2 minutes for 5' scRNA-seq) with a razor blade. We transferred the
820  minced tissue to a 1.5 ml LoBind Eppendorf tube and rinsed the petri dish with the remaining
821  500 µl media and collected to the same tube. We incubated the tube at 37°C in a thermomixer
822  for 12 minutes at 500 rpm, during which we triturated after 6 minutes 15 times with a wide
823  bore 1 ml pipette tip. We stopped the reaction by adding 1 ml (or 500 µl for 5' scRNA-seq)
824  DMEM/F12/10% FBS (room temperature) and gentle mixing and incubating at room
825  temperature for 1 minute. We passed the dissociated tissue through a 30 µm filter (Sysmex, 04-
826  004-2326 or Miltenyi Biotec, 130-098-458 for 5' scRNA-seq) and into a 15 ml tube on ice. We
827  rinsed the filter with 10 ml cold DMEM/F12/10% FBS. We then filtered the flowthrough again
828  through a new 30 µm filter and collected into a new 15 ml tube. We rinsed the previous 15 ml
829  tube with 1 ml cold DMEM/F12/10% FBS and passed through the second filter and collected
830  into the same tube. We spun down the collection for 10 minutes at 200 g and 4°C. After
831  supernatant removal, we resuspended the pellet with 200 µl (or 40 µl for 5' scRNA-seq)
832  DMEM/F12/10% FBS.

833

834  For 5' scRNA-seq, we added 1 ml PBS/1% UltraPure BSA (Thermo Fisher Scientific,
835  AM2616) and transferred the cells to a new 1.5 ml tube. We centrifuged for 10 minutes at 100
836  g at 4°C. After supernatant removal, we resuspended the pellet with 400 µl PBS/1% BSA. We
837  then filtered the cells through a 20 µm filter (pluriSelect, 43-10020-40) and collected into a
838  new 1.5 ml tube.

839

840     For Smart-seq2 and scNMT-seq, after cell counting with a Cellometer K2 (Nexcelom
841     Bioscience) and AOPI stain (Nexcelom Bioscience, CS2-0106-5ML), we diluted the cell
842     suspension to between 1 to 10 cells/µl with Phosphate Buffered Saline (PBS) containing 0.04%
843     BSA (Miltenyi Biotec, 130-091-376) and 0.2 U/µl Recombinant RNase Inhibitor (Takara,
844     2313A). We stained the cells with 2 µl Propidium Iodide (PI, BioLegend, 421301) and 1 µl
845     (1:10) CalceinAM (LifeTechnologies, C3099) per ml of cell suspension. To isolate single cells
846     for the Smart-seq2 assay, we used a Hana cell sorter (Namocell) and gated on PE$^-$/FITC$^+$ to
847     exclude dead cells and debris and sorted live single cells into 96-well plates with each well
848     preloaded with 1 µl lysis buffer (Takara, 635013) containing 0.4 U Recombinant RNase
849     Inhibitor. For the scNMT-seq assay, we sorted the above single cells into another set of special
850     96-well plates (Thomas Scientific, 1149V59 (4ti-0970/C)), in which each well was preloaded
851     with 1.5 µl reaction buffer containing 0.25 µl 10x M.CviPI reaction buffer (New England
852     Biolabs, M0227L), 2 U of GpC methyltransferase M.CviPI (New England Biolabs, M0227L),
853     0.5 µl of 800 µM SAM (New England Biolabs, M0227L), 0.025 µl 10% IGEPAL (Sigma,
854     I8896-50ml ), 2.5 U of Recombinant RNase Inhibitor. We spun down the Smart-seq2 plates at
855     2000 g for 2 minutes at 4°C and stored them at -80°C for later processing. We also spun down
856     the scNMT-seq plates with the same conditions before a 37°C incubation for 15 minutes
857     followed by adding 5 µl RLT buffer (Qiagen, 1053393) and freezing at -80°C.

858

859 **5' scRNA-seq.**
860     We did a final counting with Trypan Blue and loaded aiming to recover 10,000 cells per sample
861     preparation on the Chromium Next GEM Chip G (10x Genomics, 1000127). We performed
862     the subsequent library preparation with a Chromium Next GEM Single Cell 5' v1.1 kit (10x
863     Genomics, 1000167) following the vendor protocol with single indexing for the 5' gene
864     expression libraries. The 5' gene expression libraries were sequenced on a NovaSeq6000 S2
865     200 cycle flow cell (Illumina) with 111 bases for read 1, 91 bases for read 2 and 8 bases for
866     index 1.

867

868 **Smart-seq2.**
869     For four samples (lib_34, lib_36, lib_29, lib_676), we prepared cDNA using the SMART-Seq
870     Single Cell Kit (Takara, 634471) following the vendor manual with 1/5x reduced reaction
871     volumes and purified the cDNA with AMPureXP beads (Beckmann Coulter, A63881) with a
872     0.8x volume ratio. We quantified cDNA with a Wallac EnVision plate reader and normalized
873     concentrations across all wells. We used 0.075 ng cDNA for each and prepared Illumina
874     sequencing libraries with the NexteraXT kit (Illumina, FC-131-1096) with a 1/4x reduced
875     reaction volume. We pooled individual libraries with equal volumes and sequenced on a
876     NextSeq 500 flow cell (Illumina) with 50 bases for read 1, 25 bases for read 2, and 8 bases
877     each for index 1 and index 2, and aimed for about 1 million reads per cell.
878     For one (lib_34) of those four samples, full-length single-cell RNA-seq libraries were also
879     prepared separately using the SMART-Seq v5 Ultra Low Input RNA Kit for Sequencing
880     (Takara Bio). All reactions were downscaled to one quarter of the original protocol and
881     performed following the manufacturer's thermal cycling conditions. Briefly, reverse
882     transcription was performed using 2.5 µl of the RT MasterMix (SMART-Seq v5 Ultra Low

883    Input RNA Kit for Sequencing, Takara Bio). cDNA was amplified using 8 µl of the PCR
884    MasterMix (SMART-Seq v5 Ultra Low Input RNA Kit for Sequencing, Takara Bio) with 23
885    cycles of amplification. Following purification with Agencourt Ampure XP beads (Beckmann
886    Coulter), product size distribution and quantity were assessed on a Bioanalyzer using a High
887    Sensitivity DNA Kit (Agilent Technologies). A total of 140 pg of the amplified cDNA was
888    fragmented using Nextera XT (Illumina) and amplified with double indexed Nextera PCR
889    primers (IDT). Products of each well of the 96-well plate were pooled and purified twice with
890    Agencourt Ampure XP beads (Beckmann Coulter). Final libraries were quantified and checked
891    for fragment size distribution using a Bioanalyzer High Sensitivity DNA Kit (Agilent
892    Technologies). Pooled sequencing of Nextera libraries was carried out using a NovaSeq 6000
893    (Illumina) to an average sequencing depth of 0.5 million reads per cell. Sequencing was carried
894    out as paired-end (PE150) reads with library indexes corresponding to cell barcode.

895

896    **Nuclei isolation and single-cell multiome sequencing.**
897    A) For eight of the samples (lib_09, lib_10, lib_23, lib_29, lib_51, lib_54, lib_56, lib_57), we
898    isolated nuclei from snap frozen human kidney issue based on a vendor (10x Genomics,
899    CG000366_DemonstratedProtocol_SingleCellMultiome_Nuclei_EmbMouseBrain_Rev)
900    provided protocol with these modifications. Briefly, we chopped the tissue in 0.75 ml cold lysis
901    buffer containing 10 mM Tris-HCl pH 7.4, 10 mM NaCl, 3 mM $MgCl_2$, 0.01% Tween-20
902    (Sigma, 11332465001), 0.01% NP40 (Sigma, 11332473001), 0.001% Digitonin, 1% BSA
903    (Miltenyi Biotec, 130-091-376), 1mM DTT and 1 U/µl Recombinant RNase Inhibitor in a 1.5
904    ml Eppendorf tube with surgical scissors for about 3 minutes. We then transferred the chopped
905    tissue together with 1.25 ml rinse with lysis buffer into a 2 ml dounce tissue grinder tube
906    (Sigma, T2690). We ground the tissue with pestle A for about 15 strokes until resistance went
907    away followed by 15 strokes with pestle B until smooth. We passed the tissue lysate through a
908    50 µm filter (Sysmex, 04-004-2327) together with another 2 ml tube rinse with the lysis buffer.
909    We then passed the flowthrough through another 35 µm blue cap filter (Thermo Fisher
910    Scientific, 08-771-23). We spun down the flowthrough at 500 g for 5 minutes at 4°C and
911    resuspended the nuclei pellet in 1 ml cold wash buffer containing 10 mM Tris-HCl pH 7.4, 10
912    mM NaCl, 3 mM $MgCl_2$, 0.1% Tween-20, 1% BSA, 1 mM DTT and 0.04 U/µl Recombinant
913    RNase Inhibitor and repeated the centrifugation one more time. We resuspended the pellet in
914    1 ml cold wash buffer and counted the nuclei using the Cellometer K2 with the AO stain. We
915    spun down the nuclei suspension at 500 g for 5 minutes at 4°C and resuspended the pellet with
916    the diluted nuclei buffer containing 1x Nuclei Buffer (10x Genomics, 2000207), 1 mM DTT
917    and 1 U/µl Recombinant RNase Inhibitor in a volume based on the previous counting and
918    aiming for around 4,840 nuclei/µl. We did a final counting and loaded aiming to recover 6,000
919    nuclei per sample preparation.

920    We performed the tagmentation reaction followed by loading on the Chromium Next GEM
921    Chip J (10x Genomics, 1000230) and subsequent library preparation with Chromium Next
922    GEM Single Cell Multiome ATAC + Gene Expression kit (10x Genomics, 1000285) following
923    the vendor protocol with single indexing for the gene expression libraries. We sequenced the
924    ATAC-seq libraries on a NovaSeq 6000 SP 100 cycle flow cell (Illumina) with 34 bases each
925    for read 1 and read 2, 8 bases for index 1 and 24 bases for index 2 loading together with 1%

926 PhiX. We sequenced the gene expression libraries on a separate NovaSeq 6000 SP 100 cycle
927 flow cell with 28 bases for read 1, 55 bases for read 2 and 8 bases for index 1.

929 B) For five of the samples (lib_15, lib_34, lib_36, lib_38, lib_55), nuclei were isolated from
930 snap-frozen human kidney samples by following the Demonstrated Protocol for Single Cell
931 Multiome Nuclei Isolation from Embryonic Mouse Brain (10x Genomics, CG000366) with the
932 following modifications. Briefly, frozen tissues were transferred from –80°C to a petri dish on
933 dry ice and cut into 2–3 smaller pieces. The samples were placed in a dounce homogenizer
934 with 2 ml of 0.1x Lysis Buffer containing 10 mM Tris-HCl (pH 7.4), 10 mM NaCl (Invitrogen,
935 AM9759), 3 mM MgCl$_2$ (Invitrogen, AM9530G), 0.01% Tween-20 (Sigma, 11332465001),
936 0.01% NP40 (Sigma, 11332473001), 0.001% Digitonin (Invitrogen, 10636033), 1% BSA
937 (Miltenyi Biotec, 130-091-376), 1 mM DTT (Sigma, 646563-10X), and 1 U/µl RNase Inhibitor
938 (Roche, 03335402001). The tissue was disaggregated with 30 strokes (approximately 15
939 strokes per pestle), then transferred to a 2 ml Eppendorf tube and incubated on ice for 5 minutes.
940 A 70 µm strainer was pre-wetted with 0.5 ml of wash buffer containing 10 mM Tris-HCl, 10
941 mM NaCl, 3 mM MgCl$_2$, 0.1% Tween-20, 1% BSA, 1 mM DTT, and 1 U/µl RNase Inhibitor.
942 The tissue lysate was then passed through the filter, and 1.5 ml of wash buffer was added to
943 rinse the filter. The sample was centrifuged at 550 g for 5 minutes at 4°C using a swinging
944 bucket rotor to pellet the nuclei. The pellet was resuspended in 1 ml of cold wash buffer, and
945 the nuclei were counted using a TC20™ Automated Cell Counter (Bio-Rad) after staining with
946 Trypan Blue (Gibco, 15250-061). Nuclei were washed a total of three times and finally
947 resuspended in the appropriate volume of chilled Diluted Nuclei Buffer (10x Genomics)
948 supplemented with 1 mM DTT and 1 U/µl of RNase Inhibitor to achieve a nuclei concentration
949 suitable for a target recovery of 5000–7000 nuclei. If large clumps or debris were observed in
950 the final suspension, the nuclei were additionally filtered with a 40 µm Flowmi Cell Strainer
951 (Bel-ART, H13680-0040), before staining with Trypan Blue and performing final manual
952 counting using a Neubauer chamber.

954 Nuclei transposition and library preparation were performed following the Chromium Next
955 GEM Single Cell Multiome ATAC + Gene Expression User Guide (10x Genomics,
956 CG000338). Transposed nuclei were partitioned into GEMs using the Chromium Controller
957 with Chip J, aiming for a target recovery of 7000 nuclei per sample. After GEM incubation for
958 mRNA reverse transcription and transposed DNA barcoding, the resulting cDNA and barcoded
959 gDNA were purified and pre-amplified with 7 cycles, following the 10x Genomics protocol.
960 After clean-up, 35 µl of the pre-amplified cDNA was amplified with 7 additional PCR cycles.
961 The resulting cDNA was quantified on an Agilent Bioanalyzer High Sensitivity chip (Agilent
962 Technologies), and 100 ng was used for library preparation. GEX libraries were indexed with
963 13 cycles of amplification using the Dual Index Plate TT Set A (10x Genomics; PN-3000431).
964 In parallel, 40 µl of the pre-amplified DNA was indexed with 7 cycles of amplification using
965 the Sample Index N Set A (10x Genomics; PN 3000427). The size distribution and
966 concentration of full-length GEX and ATAC-seq libraries were verified on an Agilent
967 Bioanalyzer High Sensitivity chip. Finally, sequencing of GEX libraries was carried out on a
968 NovaSeq 6000 sequencer (Illumina) using the following sequencing conditions: 28 bases (Read
969 1), 8 bases (i7 index), 8 bases (i5 index), and 90 bases (Read 2), to obtain approximately 40,000

970 paired-end reads per nucleus. The ATAC-seq libraries were also sequenced on a NovaSeq 6000
971 sequencer using the following conditions: 50 bases (Read 1), 8 bases (i7 Index), 16 bases (i5
972 Index), and 49 bases (Read 2), aiming for a sequencing depth of >20,000 reads/nucleus.
973

974 C) For one sample (lib_34), we also followed a nuclei isolation method indicated for the
975 snDropSeq method, as described in the Kidney Precision Medicine Project Tissue Interrogation
976 Site Manual of Procedures (V7.0 – 2019). Briefly, a piece of tissue was placed in 1 ml of Nuclei
977 Extraction Buffer (NEB) containing 20 mM Tris (pH 8), 320 mM sucrose, 5 mM $CaCl_2$, 3 mM
978 $MgAcetate_2$, 0.1 mM EDTA, 0.1% Triton X-100, and 1 U/µl of RNase Inhibitor. The tissue
979 was roughly disaggregated by pipetting 10–15 times using a wide-bore 1000 µl tip. The sample
980 was then transferred to a dounce homogenizer (Sigma, D8938-1SET) and stroked 5 times with
981 the loose pestle and 20 times with the tight pestle before being transferred to a 5 ml tube. An
982 additional 2 ml of NEB buffer was used to rinse the dounce and added to the rest of the sample.
983 After incubation on ice for 10 minutes, the sample was filtered through a 40 µm strainer
984 (Pluristrainer, 43-10040-70). Then, 8 ml of PBS supplemented with 1 mM EDTA and 1 U/µl
985 of RNase Inhibitor was added, and the sample was centrifuged for 10 minutes at 900 g at 4°C.
986 The nuclei pellet was then resuspended in the appropriate volume of 1x NB buffer (10x
987 Genomics), filtered with a Flowmi cell strainer, and stained with Trypan Blue for manual
988 counting using a Neubauer chamber. Once the concentration was determined, nuclei were
989 transposed and loaded onto the Chromium system following the manufacturer's instructions as
990 described above.
991

992 **scNMT-seq: gDNA and mRNA separation.**
993 The genomic DNA (gDNA) and mRNA were separated using oligo(dT)$_{30}$VN beads on an
994 automated liquid-handling robotics platform (Hamilton STAR) as described for single cell
995 genome-plus-transcriptome sequencing[49] with minor modifications. More specifically, three
996 washes instead of two were executed using 15 µl instead of 10 µl G&T wash buffer. After
997 separation, the gDNA plate was centrifuged for 1 min at 1,000 rpm at room temperature and
998 stored at -80 °C until further processing, while the mRNA was immediately subjected to reverse
999 transcription and PCR amplification.
1000

1001 **Modified Smart-seq2 and cDNA library preparation.**
1002 Following separation, the mRNA was reverse transcribed and PCR amplified with an adapted
1003 Smart-seq2 protocol[49]. Single-cell samples were amplified for 24 cycles and subsequently
1004 purified with SPRI beads (Beckman Coulter) according to the manufacturer's instructions at
1005 0.8:1 bead:sample ratio and eluted in 25 µl nuclease-free water. Sample concentrations were
1006 measured using Quantifluor (Promega) according to the manufacturer's instructions and
1007 fragment size was assessed by Bioanalyzer (Agilent). The cDNA length ranged from 500 to
1008 2,000 bp, peaking around 1-1.5 kb, with a concentration of approximately 1-2 ng/µl. Libraries
1009 were prepared following the Nextera XT (Illumina) library preparation kit according to the
1010 manufacturer's instructions, using quarter volumes. 96 single-cell libraries were pooled
1011 together and SPRI-purified according to the manufacturer's instructions at a 0.65:1
1012 bead:sample ratio and eluted in 50 µl elution buffer (Qiagen). The concentration of the library
1013 pool was measured using Qubit (Thermo Fisher Scientific) and the size was measured using a

1014 Bioanalyzer (Agilent). Expected library pool concentrations were approximately 30 ng/µl with
1015 an average size between 400 and 600 bp and smooth profiles. Libraries with expected profiles
1016 were 384-plex equimolarly pooled to 2 nM and 50 base paired-end sequenced on an NextSeq
1017 2000 (Illumina) aiming for 1-2 million reads per cell.
1018
1019 **scBS-seq.**
1020 The gDNA was first SPRI-purified at a 0.65:1 bead:sample ratio. After adding the beads, the
1021 samples were incubated 20 min at room temperature. Next, the plate was spun down and placed
1022 on a magnet for 20 min. The supernatant was discarded and the beads were washed twice with
1023 80 % ethanol. Finally, the beads were resuspended in a 10 µl elution buffer (Qiagen). To the
1024 resuspended beads, 65 µl CT conversion reaction buffer (Zymo, Methylation-Direct MagPrep
1025 kit, D5044) was added and samples were incubated as follows: 8 min at 98°C, 3 hours at 65 °C
1026 and at most 20 hours at 4 °C. Bisulfite-converted samples were purified using the Methylation-
1027 Direct MagPrep kit according to the manufacturer's instructions using half volumes. Samples
1028 were then subjected to single-cell bisulfite sequencing library preparation as described by Clark
1029 et al.[7] Amplified libraries were pooled together and twice SPRI-purified at 0.8:1 bead:sample
1030 ratio and eluted in 50 µl elution buffer (Qiagen). The quality of the library pool was assessed
1031 using Qubit and Bioanalyzer (Agilent). The concentrations ranged from 10 to 20 ng/µl with an
1032 average fragment length of 400-600 bp. Libraries were 44-plex sequenced on an NextSeq 2000
1033 in 150 base paired-end mode aiming for at least 5 million reads per cell.
1034
1035 **3' and 5' scRNA-seq and snRNA/ATAC-seq multiome raw data processing.**
1036 FASTQ files per sample from 10x Chromium 3' scRNA-seq, 5' scRNA-seq, and
1037 snRNA/ATAC-seq multiome protocols were processed using the 10x Genomics pipelines: Cell
1038 Ranger (v.5)[50] for scRNA-seq data and Cell Ranger Arc (v.2) for multiome data. Both analyses
1039 utilized the GRCh38 (v. 2020-A) reference genome (hg38). Data processing included barcode
1040 processing, alignment to the reference genome, and quantification of gene expression for
1041 scRNA-seq data using standard parameters. For multiome data, processing also included
1042 quantification of accessible chromatin and sample aggregation. Specifically, for snATAC-seq
1043 data, sample aggregation was performed to identify a shared set of peaks across samples, with
1044 depth normalization set to *None*.
1045
1046 **Smart-Seq2 data and scNMT-seq transcriptomics analysis.**
1047 FASTQ files for Smart-seq2 and scNMT-seq transcriptome data were processed with the
1048 Nextflow DSL1[51] bulk pipeline, available at
1049 https://github.com/seanken/BulkIsoform/tree/main/Pipeline/BulkPipeline.strand.nf. Reads
1050 were mapped to the reference genome (used the GTF and FASTA file from the Cell Ranger
1051 refdata-cellranger-arc-GRCh38-2020-A reference to construct STAR and Salmon references)
1052 with STAR v 2.7.9a[52], with arguments "--outSAMattributes NH HI AS nM  --outSAMtype
1053 BAM SortedByCoordinate --readFilesCommand zcat --outStd BAM_SortedByCoordinate --
1054 outSAMunmapped Within". The resulting BAM file was used by PICARD v2.27.5[53] to extract
1055 QC information with CollectRnaSeqMetrics and MarkDuplicates (default settings). Salmon
1056 v1.6.0[54] was run with the arguments "-l A --posBias --seqBias --gcBias –validateMapping"

1057 which was then used for downstream analysis. Other tools were run by the pipeline though not
1058 used in this manuscript, so are not described.

1060 Transcripts per million from Salmon were loaded into R with tximport v1.18[55] and processed
1061 with Seurat v4.0.0[56]. Cells with <500 genes were filtered out. NormalizeData was used to
1062 normalize to TP10K, and variable genes were selected with FindVariableFeatures, both with
1063 default arguments. Data were scaled with the ScaleData command with default parameters and
1064 RunPCA was run with npcs = 60. Harmony v1.0[19] was then used with the RunHarmony
1065 command with default parameters to correct for lab of origin differences. The UMAP was
1066 generated with the RunUMAP command, while the clustering was calculated with the
1067 FindNeighbors command followed by FindClusters. Defaults were used for UMAP and
1068 clustering, except with reduction="harmony" and dims=1:25. Cell types were labeled based on
1069 marker genes and manual inspection based on the HCA reference atlas[5].

**scNMT-seq epigenome data analysis.**
1072 After sequencing and demultiplexing, adapter sequences were trimmed as described by Clark
1073 et al.[57] using cutadapt v2.103. Trimmed sequences were aligned using Bismark v0.23.18 [58] in
1074 non-directional mode and the methylation state at individual CpG and GpC sites was extracted
1075 using Bismark Methylation extractor using the coverage2cytosine script with --NOMe and --
1076 gc options. Endogenous DNA methylation was examined using CpG methylation, while GpC
1077 methylation was used to analyse chromatin accessibility. To eliminate potential biases, CpG
1078 methylation was assessed only in the ACG and TCG contexts, while GpC methylation was
1079 analysed in the GCH context, where H represents A, C, or T. Methylation levels were linked
1080 to genomic features, including gene bodies and promoters. Gene bodies were defined as the
1081 region from the transcription start site (TSS) to the transcription end site (TES), extended by
1082 15 kb upstream and downstream. Promoter regions were defined as 1,200 bp upstream and 300
1083 bp downstream of the TSS. Gene body and promoter regions were transformed into a
1084 percentage-based scale, where detected positions within a feature were normalised according
1085 to the length of the feature, taking strand orientation into account. This enabled conversion of
1086 specific genomic locations into relative percentages, representing their position within the
1087 feature. For each relative position, the average methylation percentage was calculated across
1088 overlapping features within a single cell. Subsequently, per-cell methylation values were
1089 pooled across all cells, and a best-fit was generated using a loess regression model to visualise
1090 the trends within promotor and gene body endogenous CpG methylation (%CpG) and
1091 chromatin accessibility (%GpC).

1093 To determine whether accessibility and methylation levels at the transcription start site differed
1094 significantly, two nested regression models were compared. Both models employ a restricted
1095 cubic spline approach, with 15 knots for %GpC and 14 knots for endogenous %CpG. The
1096 number of knots was selected based on the Akaike Information Criterion (AIC)[59]. Model 1,
1097 without interaction effect, assumes that the distance (around TSS and gene body) ($X_1$) and
1098 transcription status ($X_2$, 1 = transcribed, 0 = non-transcribed) have independent effects on the
1099 average methylation across different cells ($Y$, measured as %CpG or %GpC methylation). No
1100 interaction effect between these variables was considered. The model was the following:

1101
$$Y = \beta_0 + f(X_1) + \beta_2 X_2 + \epsilon$$

1102 where $f(X_1)$ was a restricted cubic spline function describing the nonlinear relationship
1103 between distance and methylation, and where $\epsilon$ were i.i.d. from a normal distribution with
1104 mean 0 and residual variance $\sigma^2$.
1105 In model 2, with interaction effect, the effect of distance $(X_1)$ on the average methylation across
1106 different cells $(Y)$ was assumed to depend on transcription status $(X_2)$, and vice versa. This was
1107 modelled by including an interaction term between distance and transcription.

1108
$$Y = \beta_0 + f(X_1) + \beta_2 X_2 + f(X_1) X_2 + \epsilon$$

1109 where the term $f(X_1)X_2$ modeled the interaction between distance and transcription, allowing
1110 for each transcription type to have its own specific curve. The number of knots for the restricted
1111 cubic spline models were determined based on the AIC[60]. The two models were compared
1112 using the likelihood ratio test (LRT), to determine whether adding the interaction term in Model
1113 2 provided a significant improvement over Model 1, under the null hypothesis $(H_0)$ that Model
1114 1 was sufficient and the interaction term did not contribute to the model, while the alternative
1115 hypothesis $(H_A)$ stated that there was a significant interaction term between distance and
1116 transcription on methylation levels.
1117

1118 **Quality control procedure for sn/scRNA-seq data.**
1119 A consistent quality control (QC) procedure was applied across scRNA-seq and snRNA-seq
1120 protocols to ensure reproducibility. We began with cell barcodes that passed the 10x Genomics
1121 Cell Ranger (or Cell Ranger Arc) filters, followed by additional filtering steps for low-quality
1122 cell barcodes and doublet removal, as detailed below.

1123 **Low-quality cell barcode removal.**
1124 Low-quality cell barcodes were filtered in two stages:
1125     1) Stage 1 - General QC Thresholds:
1126         ● Cell barcodes were excluded if they had fewer than 200 or more than 7,500
1127           unique features.
1128         ● Mitochondrial content exceeding 70% and ribosomal content exceeding 40%
1129           were also excluded.
1130     2) Stage 2 - Cell-Type Specific Filtering:
1131         ● Cell-type classifications were informed by scOMM label transfer from the
1132           reference data at L1 resolution.
1133         ● For immune cells, a mitochondrial content threshold of 10% was applied.
1134         ● For non-immune cells: A 20% mitochondrial threshold was applied for cell
1135           barcodes with fewer than 500 unique genes.
1136         ● Proximal tubule cells were an exception, with a relaxed mitochondrial threshold
1137           of 30%.
1138

1139 **Doublet detection.**
1140 Doublets were identified using the DoubletDetection[61] software (v.4.2), employing a
1141 BoostClassifier model with 30 components and Louvain clustering. Predictions were made on
1142 raw counts with a p-value threshold of $1 \times 10^{-16}$ and a voter threshold of 0.5. Barcodes classified
1143 as doublets were excluded.

**Quality control procedure for snATAC-seq data.**

Similarly, for snATAC-seq data, cell barcodes passing the 10x Genomics Cell Ranger Arc filters were retained as an initial filter before additional low-quality cell barcode and doublet removal, as detailed below.

**Low-quality cell barcode removal.**

To assess cell barcode quality, metrics specific to snATAC-seq data were applied. Fragment size distribution was used to inspect nucleosome binding patterns and a nucleosome ratio value was computed. Transcriptional start site (TSS) enrichment score was used to confirm the accessibility enrichment around TSS compared to the accessibility on flanking regions. Atac module from Muon[62] (v0.1.2) was employed to compute was employed to compute QC metrics (muon.atac.tl.nucleosome_signal with n=1e6, muon.atac.tl.tss_enrichment with n_tss=1000). QC threshold:

- Cell barcodes were excluded if they had fewer than 2500 or more than 30000 total features.
- Cell barcodes were excluded if the nucleosome ratio exceeded 4.
- Cell barcodes we filtered if the TSS enrichment score was lower than 2.

**Doublet detection.**

Doublets were identified using the AMULET[63] (v1.1) with algorithm standard parameters and specifying the blacklist file for hg38 genome. Barcodes classified as multiplets were filtered out.

**Variance partitioning and cell type trees.**

For the quantitative comparison of sn/scRNA-seq, eight donors profiled with both technologies were considered (lib_09, lib_10, lib_15, lib_29, lib_36, lib_51, lib_55, lib_56) (Table 2). Broad cell type annotation of scOMM was used. After initial filtering for genes with at least 10% non-zero counts in any cell type and cells with reasonable UMI counts and number of detected genes (i.e. no more than 5 MADs computed with function *isOutlier()* in package scuttle version 1.14.0) per assay, the expression of 22,707 genes in 63,711 cells was kept.

Out of the initial list of eight donors, five donors (lib_09, lib_10, lib_15, lib_51, lib_55) had enough cells across the 14 cell types to achieve a balanced design for the gene-wise partitioning of variance. The 14 cell types were annotated using Muto et al.'s granularity, with alternative names from Lake e.t al.[5] shown in parentheses: CNT, DCT, ENDO, FIB, ICA (CCD-IC-A), ICB (IC-B), LEUK (IMM), MES (VSM/P), PC, PEC, PODO, PT, TAL, and aPT. For that, pseudobulk samples per cell type, assay and donor were constructed by summing up all counts per cell type with function *aggregateAcrossCells()* in package scuttle (version 1.14.0). The contributions of assay, cell type, interaction of cell type and assay while controlling for donor identity in explaining the variance of normalized gene expression (function *voom()* in limma package version 3.60.4) were estimated with function *fitExtractVarPartModel()* in package variancePartition (version 1.34.0).

The same pseudobulk samples were used to draw dendrograms of cell types per assay based on hierarchical clustering using Euclidean distances and complete agglomeration method with

1187 function *buildClusterTreeFromPB()* in package dreamlet (version 1.2.1). The assay-specific
1188 dendrograms were compared using stepwise greedy forward selection algorithm by rotating the
1189 dendrograms until a local optimal solution of entanglement score (0 = no entanglement, i.e.
1190 similar; 1 = full entanglement, i.e. dissimilar) was found (function *untangle(method =*
1191 *"step2side")* in package dendextend[64] version 1.19.0).

**Metrics and statistics per cell type and assay.**
1194 The following metrics were estimated per library: mean, variance and detection rate of gene
1195 expression as well as cellular detection rate (CDR, i.e. proportion of genes with nonzero
1196 expression per cell), purity (i.e. proportion of k neighboring cells with the same cell type
1197 annotation in a Nearest Neighborhood graph based on 500 highly variable genes and 50
1198 principal components using function in package bluster[65] version 1.14.0) and silhouette of cell
1199 type clusters (silhouette based on 500 highly variable genes and 50 principal components using
1200 function in package bluster version 1.14.0). For comparison of metrics, the metrics were
1201 weighted by number of cells per library for gene metrics, purity and silhouette and weighted
1202 by number of cells and UMI reads for CDR, respectively. To ensure comparability of the
1203 metrics, we restricted the representation of gene metrics distributions (Fig. 3G, Supp. Fig. 3F-
1204 G) to genes with nonzero values in both assays per cell type as well as to the 22,707 genes
1205 deemed to be expressed for the cellular detection rate, respectively.

**Summary statistics per cell type.**
1208 To provide recommendations for researchers to choose a particular assay, we derived the
1209 following criteria of performance:
   1. Gene Metrics
       a. minimal standard error of mean gene expression ("Mean" in Figure 2G).
       b. minimal variance of gene expression ("Variance" in Figure 2G).
       c. minimal dropout of gene expression ("Dropout" in Figure 2G).
   2. Cell Metrics
       a. maximal cellular detection rate ("Detection" in Figure 2G).
       b. maximal purity ("Purity" in Figure 2G).
       c. maximal proportion of cells with a purity value > 0.5 ("Purity > 0.5" in Figure 2G).
       d. maximal silhouette ("Silhouette" in Figure 2G)
       e. maximal proportion of cells with a silhouette value > 0 ("Silhouette > 0" in Figure 2G).
1222 In addition, we used the Kolmornov-Smirnov distance as a measure of reproducibility. In
1223 general, this metric quantified how (dis-)similar a pair of univariate distributions are, e.g., the
1224 mean gene expression between one snRNA-seq library compared to one scRNA-seq library,
1225 where a value close to zero indicated a high similarity. Given our experimental design where
1226 the expression was profiled in several patients using both scRNA-seq and snRNA-seq, we
1227 could compare the reproducibility of above-mentioned metrics per cell type within donors
1228 between scRNA-seq and snRNA-seq and across donors within one assay.

**Horizontal integration of sc/snRNA-seq and snATAC-seq modalities.**

1231    Individual sn/scRNA protocols were preprocessed independently prior to integration,
1232    following consistent procedures as explained below. The Scanpy[66] (v1.9.5) standard pipeline
1233    was employed to generate PCA-based latent spaces, using the top 50 Principal Components
1234    (PCs) as the latent representation. Sample batch correction within protocols was applied
1235    through the following methods:

1236        ● Harmony[19] integration: Performed on the top 50 PCs using *run_harmony* function from
1237          *harmonypy* (v.0.4.7).
1238        ● Scanorama[67] integration: Conducted using the *correct_scanpy* function (v1.7.3)
1239          according to the standard parameters as suggested in authors' guidelines.
1240        ● scVI integration[23]: The scVI model (v.1.0) was configured with 256 nodes in the first
1241          hidden layer (*n_hidden*), two hidden layers in the encoder-decoder architecture
1242          (*n_layers*), and 30 nodes in the bottleneck layer (*n_latent*). For all latent spaces, the
1243          same set of 5,000 genes was used, identified using Scanpy's '*cell_ranger*' algorithm.
1244

1245    For sn/scRNA integration between protocols, 5,000 highly variable genes were selected using
1246    '*seurat_v3*', setting the protocol of origin as a variable to regress out. For tools supporting
1247    multiple batch-effect variables, samples of origin were also included beyond suspension type
1248    (nucleus/cell). This approach ensured consistency and effective integration across different
1249    sn/scRNA sequencing protocols.

1250    For the snATAC-seq modality, the open chromatin regions data were represented in two ways
1251    for latent space computations:

1252      1. Peaks (Variable-Length Windows): A peaks x nuclei matrix with a homogenous set of
1253        peaks across human samples was directly obtained from the *Cell Ranger ARC*
1254        *aggregate* pipeline as *filtered_feature_bc_matrix.h5*.

1255      2. Bins (Fixed-Length Tiles): A bin feature matrix was created using the *add_tile_matrix*
1256        function from SnapATAC2[24] (v0), based on the *atac_fragments.tsv* file provided by
1257        Cell Ranger ARC. The matrix consisted of fixed-size tiles for downstream
1258        dimensionality reduction tasks.

1259    Depending on the requirements of the dimensionality reduction method, either the peak or bin
1260    matrix was selected.

1261        ● Latent Semantic Indexing (LSI):
1262          LSI was computed using Signac[25] (v,1.14.0) with standard parameters. The first
1263          component of the LSI space was discarded, as typically associated with sequencing
1264          depth, and batch correction was performed using Harmony (via Seurat[20] v5).
1265        ● scVI-Based Models:
1266          PeakVI[26] and PoissonVI[27] models from scVI (v.1.0) were configured with 512 nodes
1267          in the first hidden layer (n_hidden) and 30 nodes in the bottleneck layer (n_latent).
1268          Features were filtered to retain those present in at least 1% of cells. Training continued
1269          until the elbow loss converged (early_stopping=True). For additional analysis, the

1270           PoissonVI model was trained on a fragment matrix approximation derived from the
1271           peak matrix.
1272     ●   Spectral Embedding:
1273           Based on Laplacian Eigenmaps, unimodal spectral dimensionality reduction was
1274           computed following SnapATAC2 guidelines, using a tile matrix with a 500-bin size
1275           and the top 200,000 features. Batch effect correction was evaluated using Mutual
1276           Nearest Neighbor[68] (MNN) and Harmony with standard parameters.

1277

**1278 Benchmarking of computational tools for horizontal integrations in sn/scRNA and**
**1279 snATAC-seq data.**

1280 For each protocol and modality integration described in the horizontal integration section, a
1281 benchmarking pipeline from scIB[12] (v1.1.3) was executed for the computational tools tested.
1282 The benchmarking was performed using the two levels of annotations defined via label transfer.
1283 For individual sn/scRNA protocols and data harmonization across RNA protocols, all metrics
1284 included in the standard *scIB* Benchmarker function were applied. In the snATAC-seq
1285 scenario, the benchmarking pipeline was modified by removing the *Isolated_labels* metric
1286 from the *BioConservation* set and the *pcr_comparison* metric from the *BatchCorrection* set.
1287 These metrics were excluded as they were not informative and produced similar values across
1288 the different tools. Additionally, for the snATAC-seq horizontal integration, nuclei failing
1289 quality control in the snRNA modality were filtered out to include only high-quality nuclei
1290 from both modalities in the multiome experiment. Labels transferred from snRNA were
1291 propagated and used in the snATAC-seq data. For each protocol, the sample of origin was set
1292 as the batch variable. For integration across sn/scRNA-seq protocols, the batch variable was
1293 defined as the protocol of origin.

1294

**1295 Cell-type annotations across single-cell data types with scOMM.**

1296 To achieve consistent and robust cell annotations across datasets and modalities, we developed
1297 scOMM (https://github.com/mereulab/scOMM/tree/master), a versatile, reference-based
1298 classifier designed for the automatic classification of single-cell multimodal data. Built within
1299 a unified architectural framework, scOMM features implementations in Keras, R, and Python,
1300 ensuring accessibility and reproducibility across diverse computational workflows. ScOMM
1301 employs a sequential neural network architecture with adjustable parameters, including the
1302 number of layer nodes, activation functions, dropout rates, and training settings, providing
1303 flexibility for various applications. The basic workflow of the scOMM annotation is the
1304 following:

1305

**1306 Data Preparation**. The *ds_prepare_data* function aligns the reference ($R$) and query ($Q$)
1307 datasets and extracts relevant features. Marker genes ($M$) are either user-specified or
1308 determined using the *FindAllMarkers* function (Seurat v.5.1.0). A gene $g$ is retained if it
1309 satisfies:

1310 $$g \in M \cap G_R \cap G_Q$$

1311 where $G_R$ and $G_Q$ represent the gene sets of the reference and query datasets, respectively.

1312    The resulting expression matrices $X_R \in \mathbb{R}^{|M| \times n_R}$ and $X_Q \in \mathbb{R}^{|M| \times n_R}$, where $n_R$ and $n_Q$ are the

1313    number of cells in $R$ and $Q$, are processed alongside their associated cell-type annotations

1314    $(C_R, C_Q)$ for downstream tasks.

1315    *Data Splitting.* The *ds_dplit_data_dnn* function divides the reference dataset $X_R$ into training

1316    $(X_{train})$ and test $(X_{test})$ subsets while maintaining class balance. For each celltype $c$ in $C_R$:

1317 $$|X_{train}^c| = p \times |X_R^c|, |X_{test}^c| = (1 - p) \times |X_R^c|$$

1318    Where $p$ is the proportion of data allocated to training, and $|X_R^c|$ is the number of cells of type

1319    $c$.

1320    Labels $C_{train}$ and $C_{test}$ are one-hot encoded into matrices $Y_{train} \in \{0,1\}^{n_{train} \times k}$ and

1321    $Y_{test} \in \{0,1\}^{n_{test} \times k}$, where $k$ denotes the number of cell types.

1322

1323    **Model Architecture and Training.** The deep neural network (DNN) model is constructed and

1324    trained using the *ds_dnn_model* function. Its architecture includes:

1325      1. An input layer with $|M|$ nodes.

1326      2. $L$ hidden layers, each with user-defined nodes ($h_l$ for layer $l$), *ReLU* activation, and

1327        optional dropout. The transformation at layer $l$ is expressed as:

1328 $$h_l = \max(0, W_l h_{l-1} + b_l), l \in \{1, \dots, L\}$$

1329        where $W_l$ and $b_l$ are the weight matrix and bias vector, respectively.

1330      3. An output layer with $k$ nodes and softmax activation:

1331 $$\hat{y_i} = softmax(W_o h_L + b_o)$$

1332        where $\hat{y_i}$ represents the predicted probability vector for cell $i$.

1333    The model optimizes the categorical cross-entropy loss:

1334 $$\mathcal{L} = -\frac{1}{n_{train}} \sum_{i=1}^{n_{train}} \sum_{j=1}^{k} y_{i,j} \log(\widehat{y_{i,j}})$$

1335    Where $y_{i,j}$ and $\widehat{y_{i,j}}$ are the true and predicted probabilities for class $j$ in cell $i$.

1336    The Adam optimizer is used, with early stopping employed to prevent overfitting.

1337    Hyperparameters such as the learning rate ($\eta$), batch size ($b$), and dropout rate ($r$) are

1338    configurable.

1339

1340    **Classification.** The *ds_dnn_classify* function applies the trained model to the query dataset $X_Q$,

1341    producing a probability matrix $P_Q \in [0,1]^{n_Q \times k}$:

1342 $$P_Q = f(X_Q; \Theta)$$

1343    where $f$ denotes the DNN and $\Theta$ its learned parameters. Classification is determined by:

1344 $$class(i) = \begin{cases} \arg max_j P_{Q,i,j} & if\ max_j P_{Q,i,j} \geq \tau \\ unclassified & otherwise \end{cases}$$

1345    where $\tau$ is a user-defined threshold.

1346

1347    While originally developed for scRNA-seq data, scOMM is highly adaptable and extends its

1348    functionality to snATAC-seq data and transformed gene activity values, making it suitable for

1349    a wide range of omics datasets. This extension enables its application to cross-modality models,

1350    where the snRNA-seq data serve as reference for training the model to predict the query gene

activity values inferred from snATAC-seq, and vice versa. A notable strength of scOMM is its bridge-like classification approach, which leverages the relationships between snRNA and snATAC profiles within the same individual nuclei. This approach enables the transfer of cell labels from chromatin accessibility profiles in multiomics datasets to unpaired ATAC data using the same basic workflow, replacing gene expression with chromatin accessibility features. By preserving consistency in cell annotations across both modalities and datasets, scOMM facilitates comprehensive, scalable, and integrated single-cell multiomics analysis.

**Label transfer from external references.**

For the cell-type annotation, we relied on scOMM. Feature selection for model training was based on cell-type markers, identified using the Wilcoxon test on log-normalized data. The model architecture shared several common elements, including *ReLU* as the activation function, batch normalization, weight regularization, and dropout. The initial learning rate for training was set to 0.001.

Two reference single-cell datasets from human kidney studies were used in order to guarantee precise cell type labeling and enable cross-referencing. One dataset contained snRNA-seq data, whereas the other contained both snRNA-seq and scRNA-seq data. One dataset captured finer cell states (Lake et al.), while the other provided broader annotations (Muto et al.).

### 1. snRNA/snATAC-seq reference (Muto et al., 2021)

The reference dataset was obtained through the public data in Muto et al.[17]. Features were selected using the "*author_cell_type*" metadata via Scanpy[66] software. For training, the top 400 non-overlapping gene markers per cell type were used. The model architecture consisted of three hidden layers with 256, 128, and 68 nodes, respectively, and a dropout probability of 0.5.

### 2. Human Kidney Atlas (Lake et al., 2023)

The reference dataset was obtained through the public data in Lake et al.[5], and was subsetted to include only cortex-associated states, either healthy or adaptive, from the reference samples.

- Level 1 Label Transfer (Broad Annotation):
  The top 500 non-overlapping features per cell type were selected as input for the model. The model architecture included two hidden layers with 1024 and 256 nodes, respectively, and training was performed with a dropout probability of 0.1. Predictions were made without applying an unclassified threshold, across all RNA-based protocols.
- Level 2 Label Transfer (Fine Annotation):
  Level 2 label transfer was conducted iteratively from Level 1, with one model trained per Level 1 cell type to predict the Level 2 sub-states. Each model used two hidden layers with 512 and 128 nodes, respectively. The top 300 non-overlapping features per cell type were selected as input. Training was performed with a dropout probability of 0.3.

**External references annotation matching.**

1389    The Human Kidney Atlas (Lake et al., 2023) provides hierarchical annotations for kidney
1390    biology. From the publicly available dataset, we selected *subclass.l1* metadata as Level 1 (L1)
1391    annotation and *subclass.l3* metadata as Level 2 (L2) annotation. The kidney cortex was
1392    resolved into the following stratified categories: For PT (L1), this included PTS1/2, aPT, and
1393    PT-S3 (L2). TAL (L1) was annotated as C-TAL, aTAL1, MD, and aTAL2 (L2). PC (L1)
1394    comprised PC (L2). CNT (L1) included CNT and CNT-PC (L2). DCT (L1) was further
1395    resolved into DCT1 and DCT2 (L2). IC (L1) was divided into CCD-IC-A, IC-B, and CNT-IC-
1396    A (L2). PEC (L1) was annotated as PEC (L2). POD (L1) included POD (L2). EC (L1) was
1397    categorized into EC-PTC, EC-GC, EC-AEA, and EC-LYM (L2). FIB (L1) comprised FIB,
1398    MYOF, and aFIB (L2). VSM/P (L1) included VSMC/P, VSMC, MC, and REN (L2). Finally,
1399    IMM (L1) encompassed T, MDC, NKC/T, ncMON, MAC-M2, B, N, PL, cDC, pDC, and
1400    MAST (L2).
1401    The snRNA/snATAC-seq reference (Muto et al., 2021) provides a single annotation layer,
1402    defined in the publicly available dataset as *author_cell_type* metadata. This annotation broadly
1403    matches the Level 1 annotation from the Human Kidney Atlas, with increased resolution for
1404    certain populations. Here, we provide the annotation match between these two external
1405    datasets: CNT (L1), DCT (L1), ENDO (L1, corresponding to EC), FIB (L1), ICA (L2,
1406    corresponding to CCD-IC-A under IC), ICB (L2, corresponding to IC-B under IC), LEUK (L1,
1407    corresponding to IMM), MES (L1, corresponding to VSM/P), PC (L1), PEC (L1), PODO (L1,
1408    corresponding to POD), PT (L1), TAL (L1), and PT_VCAM1 (L2, corresponding to aPT under
1409    PT).
1410

1411    **Cell identity harmonization following integration based on clustering.**
1412    To harmonize cell identities after data integration, we performed clustering followed by manual
1413    annotation using known cell type-specific markers. Clustering was performed to define two
1414    main annotation levels, L1 and L2, as outlined in the supervised label transfer approach with
1415    scOMM, which was applied independently for each assay to minimize biases. Leiden[69]
1416    clustering was conducted on a k-nearest neighbors (k=50) graph, using a range of resolution
1417    values (0.3, 0.6, 1, 1.4, 2, 6, 10). Final cluster labeling was determined by majority voting based
1418    on labels transferred at each annotation level, with resolution level 6 used as the primary
1419    reference. Results were then visually inspected for cluster consistency, and marker-based
1420    supervision was applied for accurate cell type annotation. If needed, cluster assignments were
1421    refined by adjusting the resolution to balance cluster structure and cell identity.
1422

1423    **Population reclustering.**
1424    Reclustering is performed with Leiden[69] algorithm on a k-nearest neighbors (k=10) graph
1425    derived from the embedding computed on the whole dataset. No recomputation of highly
1426    variable genes or embedding was performed.
1427

1428    **Contribution of sn/scRNA-seq protocols to cell type identification.**
1429    To evaluate the contribution of each sn/scRNA-seq protocol to cell type annotations, partial
1430    area under the receiver operating characteristic curve (pAUC) by specificity was employed
1431    using the R package pROC[70] (v.1.18.5), focusing on the high-specificity range (values > 0.9)
1432    to minimize false-positive assignments. Partial AUCs were calculated separately for Level 1

1433 (broad annotation) and Level 2 (fine annotation) label transfers. The curated manual clustering
1434 annotation, derived after sn/scRNA-seq horizontal integration, was used as the ground truth for
1435 cell-type annotation.
1436

**Vertical integration of snRNA- and snATAC-seq data.**
1438 Integration of paired snRNA and snATAC modalities from the multiome dataset uses nuclei
1439 (observations) as anchors to generate a joint embedding across modalities. A joint embedding
1440 can be generated either directly from the raw feature spaces of both data types (snRNA-seq and
1441 snATAC-seq) or by using pre-computed embeddings specific to each modality. Preprocessing
1442 settings and highly variable feature selections are imported and fixed from horizontal unimodal
1443 integration. Final vertical embedding representation was performed with the SnapATAC2
1444 module through Laplacian Eigenmap computation on both modalities simultaneously (i.e.,
1445 *snapatac2.tl.multi_spectral*), with linear space and time complexity, to obtain the joint latent
1446 space. MNN was applied to mitigate the batch effect. UMAP projection was computed on a 50
1447 kNN graph.
1448

**Mosaic integration.**
1450 Integration of paired and unpaired sn/scRNA-seq and snATAC-seq datasets was achieved
1451 using scVI tool (v.1.0) through the MultiVI[18] model. This model was trained on the previously
1452 selected highly variable features, with protocol and sample origin included as categorical
1453 covariates for batch effect correction. Clustering of the joint embedding was computed on a k-
1454 nearest neighbors (k=50) graph using the Leiden[69] algorithm. Differential expression analysis
1455 (Wilcoxon test) on the sn/scRNA-seq data was used to detect and exclude low-quality cells and
1456 outlier populations inconsistent with the kidney cortex biology. The final UMAP projection
1457 was performed on a 50 kNN graph.
1458

**Diagonal integration.**
1460 To explore global integration potentials, we simulated fully unpaired datasets by decoupling
1461 snMultiome modalities, thereby testing diagonal integration. For this purpose, we utilized
1462 GLUE[29] as a representative method for state-of-the-art unpaired multimodal integration.
1463 Preprocessing and feature selection from both sn/scRNA- and snATAC-seq horizontal
1464 integrations. The GLUE model was trained according to the authors' guidelines. During dataset
1465 configuration (via *scglue.models.configure_dataset*), cells from each modality were assigned
1466 a batch effect variable defined as a combinatorial factor of sample origin and protocol. The
1467 best-performing embedding from horizontal integration task was used as the guidance
1468 embedding for each modality.
1469

**Systematic evaluation of vertical integration embedding.**
1471 To evaluate vertical integration embedding, we employed the Weighted Nearest Neighbor
1472 (WNN) algorithm from Seurat as a baseline. WNN leverages unimodal latent spaces to assign
1473 weights to each cell and modality, generating a final multimodal latent space. We computed
1474 WNN embeddings twice: (1) using standard unimodal embeddings as per the authors' vignette
1475 (Standard WNN), and (2) using the best-performing embedding from each modality obtained
1476 during horizontal integration (referred as to optimized WNN).

The performance of the optimized WNN embedding was systematically compared to that of the multi-spectral MNN-corrected embedding (described in Vertical Integration of snRNA-seq and snATAC-seq Data) across every L2 population. Silhouette width and b-cLISI scores were computed for each cell in unimodal and multimodal embeddings using scOMM labels as the cell-type reference. The distributions of cell-wise Silhouette values for multimodal embeddings were compared via scatter plots (generated with Seaborn[71] v0.13.0). Average deviations from the diagonal (referred to as Avg. Dev.) toward a specific axis was used to quantify the performance gap between embeddings in detecting and isolating cell types or states. It was computed by averaging pairwise distances between the two distributions. To provide a broader comparison perspective, a heatmap was generated showing the average metric values per L2 population for multiome cells with high-quality data in both RNA and ATAC modalities.

**Systematic evaluation of diagonal and mosaic integration embeddings.**
To systematically evaluate the ability of unimodal and multimodal approaches to detect and isolate cell types and states, silhouette width and b-cLISI scores were computed using scOMM labels. For each population category, cell-wise metric values were averaged and displayed in a heatmap, with integration methods as rows and populations as columns.

**Systematic evaluation of single-cell data integrations using scOMM for supervised label transfer.**
ScOMM was employed as a benchmark tool to systematically evaluate the integration of the distinct data types, focusing on how different strategies preserved biological information, maintained data quality, and achieved consistent cell-type annotations. Horizontal and vertical integration were assessed by quantifying modality-specific contributions and performance using metrics such as cell-type identification accuracy and marker detection efficiency.

**Cell type predictability of sn/scRNA-seq assays before and following horizontal integration.**
To quantify cell type predictability (L1 granularity) for each individual sn/scRNA-seq protocol, we trained three distinct scOMM models, one for each protocol. For model training, ~6,000 nuclei/cells per dataset were used, representing the maximum number that could be equally included across all three protocols. Additionally, a combined model was constructed by randomly sampling ~2,000 nuclei or cells from the training data of each protocol, resulting in a combined dataset of ~6,000 nuclei/cells. The remaining cells were allocated as testing data. Each model was subsequently used to predict cell types in the other datasets. Model performance was evaluated using the total area under the curve (AUC), with AUC scores computed using the *bench_calcAUC* function in scOMM. Predictions from the *ds_dnn_classify* scOMM function were compared to the original protocol-specific cell annotations. The visualization of this comparison was generated via the *bench_plotTileAUC* function in scOMM.

**Comparative analysis of snRNA- and snATAC-seq data modalities in cell type identification (label transfer-based).**
To compare snRNA- and snATAC-seq data in cell type identification under a supervised, label

1520  transfer-based scenario, three scOMM models were evaluated using cell type annotations from
1521  snRNA-seq data as the ground truth:

1522  1.  Cross-Modality model: trained on snRNA-seq data as the reference and tested on gene
1523      activity scores derived from snATAC-seq.
1524  2.  Bridge model: both training and testing were conducted within the snATAC-seq
1525      modality, using peak data.
1526  3.  Integrated model: combined features from the training data of both previous models to
1527      assess whether integrating RNA and chromatin features improves cell type resolution.

1529  To ensure comparable cell composition across models, 70% of nuclei (~26,400 cells) were
1530  sampled as the training set for the gene-based, peak-based, and integrated models, with the
1531  remaining 30% used for testing. Sampling was performed using the *ds_split_data_dnn*
1532  function, and training was conducted using the *ds_dnn_model* function, both implemented in
1533  scOMM. Performance evaluation focused on two key metrics:
1534  i) Cell type prediction accuracy: to assess the relative contributions of RNA and ATAC
1535  modalities in resolving cell types, AUC scores were computed using the *bench_calcAUC*
1536  function. Predictions from the *ds_dnn_classify* function were compared against the ground
1537  truth.
1538  ii) Identification of cell type-specific features: significant markers were identified through
1539  feature importance scores and marker detection rates, evaluated using the
1540  *ds_feature_importance* and *bench_MarkerDetect* functions, respectively, in scOMM.

**Feature importance in scOMM**.
1543  Feature importance (FI) scores were calculated from the scOMM models using a permutation
1544  importance algorithm as implemented in the *ds_feature_importance* function. This approach
1545  evaluates the relative contribution of individual features, such as genes, peaks, or their
1546  combination, to the model's overall cell type prediction performance. The algorithm works by
1547  systematically setting the values of a single feature to zero across all cells and measuring the
1548  resulting impact on model accuracy. Features that cause a drop of more than 10% in accuracy
1549  are assigned the highest importance scores.

**Marker detection rates.**
1552  The marker detection rate comparison is conducted using the function *bench_MarkerDetect* (in
1553  scOMM), where each marker detection rate was calculated based on the feature importance
1554  scores generated by each model. For each cell type, marker detection rate was defined as the
1555  proportion of significant features (importance scores ≥ 10) relative to the predefined reference
1556  list of markers derived from the Human Kidney Atlas . This allowed the quantification of each
1557  model's capacity to detect biologically meaningful markers across sn/scRNA-seq, snATAC-
1558  seq, and their integrated datasets. For models using genomic sites (peak-based and peak-
1559  integrated models), the *ClosestFeature* function from Signac[25] (v1.14) was employed to map
1560  each genomic site to its nearest gene.

**Embedding-based metrics for systematic evaluation of data integration.**

1563 The primary metrics employed in this evaluation include the Mahalanobis distance and width
1564 silhouette, which serve as global measures of cell type separability, and the binary version cell
1565 type Local Inverse Simpson's Index (b-cLISI), which provides a local assessment of cell-type
1566 preservation and accuracy. For all these metrics, calculations were performed at the level of
1567 cell types, where cell types were inferred using scOMM.

1568 **Multidimensional Mahalanobis Distance**

1569 This metric calculates the distance between two distributions, one as a reference ($V_i$) and the
1570 other as a query ($V_j$), in an N-dimensional embedding space. The distance ($D_{ij}$) is computed
1571 as the separation from the reference distribution $V_i$ to the centroid of $V_j$, measured in units of
1572 dispersion of $V_i$. The formula is as follows:

1573
$$D_{ij} = \sqrt{\Delta^T \, \Sigma_{V_i}^{-1} \, \Delta}$$

1574 where:

1575 • $\Delta = \mu_{V_i} - \mu_{V_j}$

1576 • $\Sigma_{V_i}^{-1}$ is the inverse of the covariance matrix of $V_i$, computed using the Ledoit-Wolf
1577 shrinkage adjustment (Scikit-learn v1.4.0)

1578 • $\mu_{V_w} = \frac{1}{n_{\mathrm{obs_w}}} \sum_{k=1}^{n_{\mathrm{obs_w}}} V_{wk}$ is the centroid of $V_w$

1579 Each cell type defines a distribution in the N-dimensional embedding space, allowing for
1580 Mahalanobis distances to be calculated between cell types. Distances greater than 3 indicate
1581 good separation in low-dimensional embeddings. For larger embeddings or precise evaluations,
1582 the chi-squared distribution can serve as a guide ($D^2 = \chi_{p,\alpha}^2$), where $p$ is the dimensionality of
1583 the embedding and α is the overlap ratio). As this measure is not symmetric ($D_{ij} \neq D_{ji}$), the
1584 harmonic mean is used to summarize distances between two populations.

1585 **Accuracy of cell-type identification in the embedding-based (unsupervised) scenario.**
1586 Cell type identification performance was assessed using the following embedding-based
1587 metrics. The evaluation combined the median values of b-cLISI and Silhouette Width (SW)
1588 scores for each cell type, averaging these values across all observations in the category to
1589 produce a final score. The SW metric evaluates label identity assignments within N-
1590 dimensional embeddings. It calculates the ratio between intra-category distances and inter-
1591 category distances (to the nearest distinct category) for individual observations. Implemented
1592 in the scIB[12] package, the *silhouette_samples* function was used to compute sample-wise
1593 scores, providing insights into how well clusters are separated in the embedding space.

1594 **Binary Cell-Type Local Inverse Simpson's Index (b-cLISI).**
1595 The Binary Cell-Type Local Inverse Simpson's Index (b-cLISI) score is a variant of the cLISI
1596 to assess local neighbor diversity in a neighborhood graph, considering only the label of the

central node's category. This modification enhances interpretability by focusing on intra-category consistency. For a node $i$ with neighborhood $N(i)$ and category $j$, the formula is:

$$\text{b-cLISI}_i = \frac{n_j}{\sum_{k \in N(i)} n_k}$$

where $n_j$ is the number of nodes in $N(i)$ belonging to category $j$. The b-cLISI metric ranges between 0 and 1, where higher scores indicate better local clustering and consistency for the same cell type.

**Differential expression and accessibility analyses.**

Differential analyses were conducted for each individual protocol and modality using annotation labels derived from the label transfer tasks at the two levels of granularity. For sn/scRNA-seq protocols, differential expression analysis was performed using Scanpy's[66] *rank_genes_groups* function with the Wilcoxon method on log-normalized data. For the snATAC-seq modality, cell type-specific narrow peak calling was conducted using the *macs3* function in SnapATAC2 with standard settings. Differentially accessible regions (DARs) were identified on a cell type basis from peak calling results using the *diff_test* function in SnapATAC2, with thresholds set to *min_lfc=0* and *min_pct=0.01*. Accessible regions were mapped to "accessible genes" by extending 2000 bp upstream and 500 bp downstream from the transcriptional start site, based on the GRCh38 genome reference.

Differential accessibility results were propagated to genes in the following way:

1. P-values for accessible regions were aggregated to corresponding genes using Stouffer's Z-score method. This was implemented via custom Python scripts utilizing the *ppf* and *cdf* functions from the *scipy.stats.norm* module (v1.11.4).
2. Log Fold change is approximated by averaging the linear fold changes of individual features assuming all features contribute equally to the aggregate:

$$LFC_g = log_2\left(\frac{\sum_{i=1}^{n} 2^{LFC_{r_i}}}{n}\right); \forall r \in g$$

A significance score ($w_r(g)$) was calculated as:

$$w_r(g) = \text{LFC } (r) \times -\log_{10}\big(Pval \ (r)\big)$$

These scores were averaged across regions associated with the same gene. Statistically significant expressed genes were defined using thresholds of a minimum log fold change of 0.25 and an FDR-corrected p-value of 0.05. For accessible genes, significance scores were thresholded at the 0.8 quantile, enabling comparisons with statistically significant differentially expressed genes across protocols        .

**Identification of cell type-specific markers.**

We implemented a postprocessing step on differential expression results to identify markers that are specific to individual populations within a given batch (i.e., protocol or modality).

These markers were further validated to ensure they were not statistically significant in any other population within the same batch. Filtering was performed using thresholds for adjusted p-value, log-fold change, and mean expression level, which could be customized to meet the user's desired level of statistical rigor.

For each batch:

 1.Candidate markers for each population were compared against genes identified in other populations within the same batch, excluding overlapping or non-significant genes.

 2.Markers unique to each population were identified, ranked by significance scores, and compiled into a sorted list for each population.

This approach enabled the precise identification of population-specific markers, facilitating robust and reliable downstream analyses. This was used for the discovery of consensus sets of markers across protocols.

**Norn signature scoring.**

A Norn signature score was calculated in each nucleus/cell using a reference set of marker genes (DCN, PDGFRA, GSN, TIMP1, CFD) by using the *score_genes* function from Scanpy[66] with default parameters. This method assigns a score to each cell based on the average expression of the provided marker genes subtracted with the average of a reference set of genes drawn randomly from all expressed genes. This score enables the identification and quantification of cell-type-specific activity across the single-cell RNA-seq dataset.

**Downsampling for the comparison of assay-specific cell type markers.**

Custom downsampling workflow was developed for sn/scRNA-seq data to ensure proportional representation of cell types and preservation data distributions while standardizing number cells and counts depth across dataset, seeking to avoid technical biases on the comparison of marker features. Data from each protocol were first downsized to 7,500 cells. Its implementation on Python was achieved through *StratifiedShuffleSplit* function from sklearn[72] (v1.4.1) model_selection module, providing the cell type L1 labels to maintain the population proportion on the downsized object. Checks were done for the maintenance of total counts and number of features per cell distribution quantiles within their original values. Once objects were downsized, *downsample_counts* function from Scanpy[66] preprocessing module was employed, setting an upper threshold of 10,000 total counts per cell.

**Benchmarking summary scores.**

We focused on two primary analytical tasks: cell type identification and marker detection, each implemented through label transfer–based (supervised) and embedding-based approaches (unsupervised). To facilitate a concise comparison of results, a single metric was computed derived from the analyses within each analytical arm, providing an integrated perspective on the overall performance of the different methodologies and cell types.

- Cell-type identification in label transfer-based scenario: population wise AUC predictions for each scenario scOMM model were averaged across testing sets to obtain the final population score.

- Cell-type identification in embedding-based scenario: cell-wise width silhouette and b-cLISI scores were averaged per cell group, then their values were averaged across populations to obtain the average score.
- Marker detection in label transfer-based scenario: final score was drawn proportion of reference marker sets provided in Lake et al.[5] detected as significant features for each scenario scOMM model based on feature importance score.
- Marker detection in embedding-based scenario: marker detection scores were based on differential expression analysis statistical results from each individual dataset/modality.

**Marker detection metric.**

A composite marker detection score was used to evaluate the uniqueness and statistical significance of marker features across cell types and protocols. For a given technology ($t$) and cell type ($c$), a set of marker genes were computed ($G_{t,c} = \{g_1, g_2, \dots, g_n\}$), with their associated statistics ($LFC_{tc}(g), Pval_{t,c}(g)$), the following scores were computed:

1) Uniqueness score ($U_{tc}$): Is computed as the summatory over the $LFC_{tc}$ of markers ($G_{tc}'$) for $c$, found only in $t$ and not in any other technology ($t'$). In the case that a gene is considered marker ($G_{tc}''$) for any other cell-type ($c'$) in any other technology, a penalty ($P_{t,c}(g)$) is applied by subtracting $LFC_{t'c'}$.

$$U_{tc} = \sum_{g \in G'} S_{tc}(g)$$

Where:

- $S_{t,c}(g) = \begin{cases} LFC_{t,c}(g) & g \notin G_{t',c'} \; ; \; \forall (t', c') \neq (t, c) \\ LFC_{t,c}(g) - P_{t,c}(g) & otherwise \end{cases}$
- $P_{t,c}(g) = \sum_{g \in G''} LFC_{t'c'}(g)$
- $G_{tc}' = g \in G_{t,c} \cap g \notin G_{t',c} \; ; \forall t' \neq t$
- $G_{tc}'' = G_{tc}' \cap g \in G_{t',c'} \; ; \; \forall (t', c') \neq (t, c)$

2) Power score ($W_{tc}$): For each gene, a significance score ($w_{tc}(g)$) was computed based on ($LFC_{tc}(g), Pval_{t,c}(g)$). Power score is computed as the summatory of the significance score over markers ($\overline{G_{tc}'}$) for c, found in at least more than one technology.

$$W_{tc} = \sum_{g \in \overline{G_{tc}'}} w_{tc}(g)$$

Where:

- $\overline{G_{tc}'} = G_{tc} - G_{tc}'$

1707 $\quad \bullet \quad w_{tc}(g) = \text{LFC}_{t,c}(g) \times - \log_{10} \left( Pval_{t,c}(g) \right)$

1708 The final marker detection score was a weighted sum of the uniqueness and power scores, with
1709 adjustments for snATAC-associated features on the power score computation (weighted by 0.3
1710 to account for differences in feature statistics between snRNA and snATAC data).

1711 **Defining Cell-Type-Specific CREs Using Accessibility and Transcription.**
1712 CREs were defined using scATAC-seq as described. Accessibility at CREs in single cells was
1713 quantified as described. Transcription within CREs in individual cells was quantified using
1714 SCAFE[73] v1.0. Briefly, the 5' scRNA-seq read alignments were converted to capped TSS
1715 (CTSS) BED files and the number of UMIs within each CRE was counted in a strand-agnostic
1716 manner. CTSS signals from all 5' scRNA-seq samples were clustered into TSS clusters using
1717 default SCAFE parameters. These clusters were subsequently used to identify transcribed
1718 CREs. Cell types in scATAC-seq nuclei and 5' scRNA-seq cells were annotated as described.
1719 The specificity of CREs in each level 2 cell type were calculated using Wilcoxon-rank sum
1720 tests implemented in the *scanpy.tl.rank_genes_groups* function in Scanpy[66], based on either
1721 accessibility or transcription. Three sets of cell-type-specific CREs were defined, including Set
1722 1: CREs ranked by accessibility, Set 2: Transcribed CREs ranked by accessibility, and Set 3:
1723 Transcribed CREs ranked by transcription. CREs were ranked by Wilcoxon p-values ($\log_2$ fold
1724 change $> 0$) and the top 3000 CREs were defined as cell-type-specific. Cell types with fewer
1725 than 3,000 CREs in each set were discarded. Set 1, 2 and 3 contains 34, 27 and 13 cell-types,
1726 respectively.
1727

1728 **Enrichment of Trait Heritability in Cell-Type-Specific CREs.**
1729 Enrichments of heritability of two traits, i.e., hypertension[42] and estimated glomerus filtration
1730 rate (eGFR)[43], were assessed in the three sets of cell type specific CREs using partitioning
1731 linkage disequilibrium score regression implemented in LDSC[74] v1.0.1. GWAS summary
1732 statistics for hypertension were obtained from UK Biobank[75] (phenotype code 6150), and those
1733 for eGFR were obtained from the GWAS catalog (study GCST90428446). The summary
1734 statistics of eGFR were munged using the "munge_sumstats.py" script and that of hypertension
1735 were used as it is. Annotation files and LD score files for each set of CREs were generated
1736 using the "make_annot.py" and "ldsc.py" scripts with default parameters. Each set of CREs
1737 was added to the 97 annotations of the baseline-LD model v2.2, and heritability enrichment
1738 (i.e., the ratio of the proportion of heritability to the proportion of SNPs) for each trait in each
1739 set of CREs was estimated using the "ldsc.py" script with the "--h2" flag under default
1740 parameters. All cell types from Set 1 were analyzed in Fig. 5F and overlapping cell types across
1741 all sets were analyzed in Fig. 5G.

1742

1743 **Supplementary material.**

1744 **Comparison of gene regulatory network inference across single-cell platforms.**

1745   Understanding cell-type-specific gene regulatory networks (GRNs) is essential for unraveling
1746   the mechanisms that drive cellular identity and function[76]. Single-cell sequencing technologies,
1747   such as 3' and 5' scRNA-seq as well as joint snRNA-seq and snATAC-seq data, provide
1748   distinct insights into transcriptional regulation[77,78]. However, the robustness, reproducibility,
1749   and biological relevance of GRNs inferred from these platforms are not well understood. To
1750   address this, we developed a systematic framework to compare GRNs across platforms based
1751   on predictive accuracy, reproducibility of transcription factor (TF) identification, cross-
1752   platform complementarity, and biological validation with independent datasets.
1753   We analyzed data from three donors ("lib_09," "lib_10," and "lib_36") across all three
1754   platforms, encompassing 25,669 nuclei/cells (4,393 from 3' scRNA, 6,416 from 5' scRNA,
1755   and 14,860 from multiome). After filtering for features expressed in at least 30 cells, gene
1756   expression was log(1+x)-scaled, and peak accessibility was binarized. The top 1,000 most
1757   highly variable genes were selected per platform using Scanpy's[66] *sc.pp.highly_variable_genes*
1758   with the "batch_key" set to the biological sample identifier.
1759

1760   **Regulatory network inference.**
1761   We inferred gene regulatory networks within each sequencing platform and biological sample
1762   independently. This entailed fitting a gradient-boosting-machine (GBM) tree regression model
1763   to each gene to predict its expression, similarly to the approach used by SCENIC[77]. We used
1764   the LightGBM[79] package, with an ensemble of 20 estimators per model, learning rate of 0.5
1765   and early stopping. For each gene, we used as input features: (1) gene expression data for all
1766   transcription factors (selected using JASPAR 2024[80]); and for only the multiome data, (2)
1767   chromatin accessibility of peaks proximal to the gene (within +/- 5kb), using the PyRanges[81]
1768   package and TSS loci from RefTSS[82] v4.1). This approach allowed us to treat "snRNA" as an
1769   additional modality, which is simply the multiome data excluding the chromatin accessibility
1770   features. To identify regulators for each gene, we assigned TFs based on their importance
1771   scores (total information gains from tree splits using that TF). A TF was considered a regulator
1772   if its importance score for a given gene exceeded the 95th percentile of scores globally across
1773   all genes within the sample and platform. Since the models were trained on cells pooled from
1774   all cell types within each sample, the resulting GRNs did not inherently include "cell type
1775   specificity" for network edges. To address this, we assigned cell type specificity to TF-target
1776   links that passed the importance filtering by evaluating the Pearson correlation between the TF
1777   and target gene's ground-truth expressions in the training set. A link was deemed active in a
1778   particular cell type if the Pearson correlation exceeded 0.2.
1779

1780   **Cross-platform prediction accuracy.**
1781   To assess the cross-sample generalization ability of models trained on each platform, we
1782   evaluated each model trained for that platform on the remaining held-out biological replicates.
1783   Specifically, for each pairing of training and held-out samples, we computed the Pearson
1784   correlation between the predicted and true log-transformed expression values for each gene
1785   (across cells), stratified by cell type (Supp. Fig. 14A). The 5' scRNA-seq model demonstrated
1786   slightly higher but statistically significant predictive accuracy compared to 3' scRNA-seq and
1787   multiome-derived models (Supp. Fig. 14B, left). However, the overall distributions of
1788   correlations revealed considerable variability across cell types, suggesting that biological and

1789 technical heterogeneity and cell-type-specific differences within patients may influence model
1790 performance for the gene expression inference through GRN (Supp. Fig. 14B, right).

1791

1792 **Reproducibility of TF identification.**
1793 To estimate robustness of TF identifications within platforms, we calculated, per cell type, two
1794 metrics: (1) mean number of TFs identified per cell type across replicates; (2) consensus TF
1795 count: defined as TFs identified consistently in all replicates within each cell type. We defined
1796 a "reproducibility ratio" for each (platform, cell type) combination, as the latter quantity
1797 divided by the former. Overall, snMultiome exhibited the highest reproducibility (Supp. Fig.
1798 14C, right), underscoring its robustness in detecting cell-type-specific regulators. While 5'
1799 scRNA-seq also performed well in certain cell types, it exhibited greater variability compared
1800 to snMultiome. By contrast, 3' scRNA-seq demonstrated the lowest reproducibility ratios,
1801 suggesting that RNA-only approaches may be less reliable for consistent TF inference.
1802

1803 **Cross-Platform complementarity in TF-target associations.**
1804 To evaluate potential complementarity of cell-type-specific transcriptional regulators detected
1805 across sequencing platforms, we represented each consensus TF-target regulatory network as
1806 a binary indicator vector, where each entry denoted the presence or absence of a TF-target link.
1807 Pairwise Jaccard similarity scores were then computed between these vectors for each
1808 (platform, cell-type) pair, representing the intersection-over-union of shared regulatory
1809 elements across platforms. Hierarchical clustering revealed that platform-specific effects were
1810 more prominent than biological differences, with GRNs clustering primarily by platform rather
1811 than cell type (Supp. Fig. 14D). To characterize cross-platform overlap at the aggregated (cell
1812 type, TF) level, we also calculated the number of (cell type, TF) identifications shared by
1813 consensus networks between each pair of platforms. As expected, multiome and snRNA
1814 exhibited significant overlap in TFs due to their common RNA-sequencing foundation and
1815 potential representation of the same cells (Supp. Fig. 14E). However, their overlap with 3' and
1816 5' scRNA, was considerably lower (Supp. Fig. 14E). These differences appear to be influenced
1817 by the RNA sequencing protocols rather than the inclusion of chromatin accessibility data, as
1818 multiome and snRNA exhibit much higher log-odds ratios (LORs) compared to the others
1819 scRNA-seq platforms (Supp. Fig. 14E, right). This indicates that complementary regulatory
1820 elements can be identified by scRNA-seq, with 5' scRNA-seq detecting the highest number of
1821 unique TF across cell types (Supp. Fig. 14E, left). Altogether, this underscores the role of
1822 platform-specific technical factors in shaping GRN architecture.
1823

1824 **Validation using independent datasets.**
1825 A TF needs to localize to the nucleus to act as a regulator. To test whether regulatory inferences
1826 from each platform satisfy this property, we cross-referenced consensus TFs identified per
1827 platform against localization data from v23 of the Human Protein Atlas[83] (HPA). Specifically,
1828 we queried the HPA for genes with immunohistochemical staining patterns annotated as
1829 "detected" in the kidney, with a reliability score of "Enhanced", and intersected these against
1830 TFs listed in JASPAR 2024. Intersecting the inferred TFs from each consensus network yielded
1831 a set of TFs per platform with independent evidence of nuclear localization in kidney tissue.

1832    With this approach, we found that a subset of inferred TFs was supported by experimental
1833    evidence (Supp. Fig. 14F). While the proportion of nuclear-localized TFs was consistent across
1834    platforms, 3' scRNA-seq exhibited the highest validation fraction.
1835    We also validated whether predicted regulatory TF-target associations were corroborated by
1836    independent evidence of TF binding or chromatin accessibility. To do so, we cross-referenced
1837    our consensus networks with bulk epigenomic data in kidney tissue. Using ChIPAtlas 2021[84],
1838    we retrieved bulk ChIP-Seq peaks for transcription factors in all kidney cell lines, and
1839    intersected those with bulk ATAC-Seq peaks from kidney cortex (All peaks were filtered at
1840    reported significance level >= 50). We then linked each peak to the closest promoter in RefTSS
1841    v4.1 within a ±500 bp window. This yielded a set of potential TF-gene associations based on
1842    promoter binding in bulk samples. Altogether, this validation of TF-target associations using
1843    bulk ChIP-Seq data revealed low validation rates across all platforms (Supp. Fig. 14G).
1844    However, while the 3' scRNA-seq platform seems to recover more consensus edges in ChIP-
1845    Seq promoters, the 5' scRNA-seq platform, along with multiome, demonstrated the highest
1846    total average TF detection, suggesting that 5' scRNA-seq and multiome may capture a broader
1847    repertoire of transcription factors, potentially enhancing the scope of inferred regulatory
1848    networks despite lower promoter-level validation in bulk ChIP-Seq data.
1849
1850    **Supplementary Figures.**
1851    **Supplementary Figure 1.** A) Live normal kidney tissue is obtained directly from the operating
1852    room at Michigan Medicine. Ideal procurements are a minimum of 1.5 cm x 1.0 cm x 1.0 cm
1853    in dimension and represent both renal cortex and medulla. B) Cortex (blue), medulla (black),
1854    capsule (green), corticomedullary junction (blue-black), and large arteries (white) are
1855    annotated. Tissue sections are chosen which contain complete transitions from renal capsule to
1856    corticomedullary junction. C) Using a 3D printed device (PRECISE Pyramid), nephrectomy
1857    specimens are cut into mock 16 gauge biopsy cores. D) The Precise Pyramid creates 25
1858    identically proportioned mock 16 gauge biopsy cores per procurement. E) Biopsy cores are
1859    preserved across a variety of media, including but not limited to FFPE, CryoStor, RNALater,
1860    OCT, and liquid nitrogen (LN$_2$). F) Biopsy core preserved in hypothermosol/CryoStor10. Note
1861    capsule, cortex, corticomedullary junction, and medulla. G) FFPE preserved core, stained with
1862    Masson Trichrome and digitally scanned at 40x from the same case.
1863
1864    **Supplementary Figure 2.** A) Violin plots showing quality control metrics for RNA-based
1865    single-cell technologies: number of unique features (genes) detected per cell (top row),
1866    percentage of mitochondrial reads (MT %) (middle row), and percentage of ribosomal protein
1867    genes (RB %) (bottom row). Metrics are displayed for snRNA from the Multiome assay (left
1868    column), 5' scRNA-seq (middle column), and 3' scRNA-seq (right column). Upper and lower
1869    thresholds for quality filtering are indicated by solid and dashed lines, respectively. B) Quality
1870    control metrics for snATAC-seq related to fragment count and length. Boxplot (box center
1871    represents the median, box size represents the interquartile range (IQR), and whiskers extend
1872    the smallest and largest value within 1.5 times the IQR) of the number of fragments per cell for
1873    each sample (left), fragment size distribution per sample (middle), and boxplot of nucleosome
1874    signal ratios per cell for each sample (right). C) Quality control metrics for scATAC-seq related
1875    to transcription start site (TSS) enrichment. Boxplot of TSS enrichment score distributions per

1876 sample (left), average TSS enrichment score distribution centered around the TSS site for cells
1877 with high TSS enrichment (middle), and average TSS enrichment score distribution across TSS
1878 sites for all high-enrichment cells (right). D) Bivariate histogram of TSS enrichment scores and
1879 the number of unique fragments per cell, illustrating the relationship between these two quality
1880 control parameters.

1881 **Supplementary Figure 3.** A) Stacked barplot of L1 cell-type composition across single cell
1882 data modalities. B) Alluvial plots displaying the distribution of label transfer annotation from
1883 external references Lake et al.[5] - L1 and Muto et al.[17] C) Upset plot of statistically significant
1884 genes on downsampled data across mBDRC protocols for Parietal Epithelial cells (PEC),
1885 Intercalated cells (IC), Distal Connecting Tubule cells (DCT), and Principal cells (PC),
1886 Vascular Smooth Muscle/Pericytes (VSM-P), Endothelial cells (EC) and Immune cells (IMM).
1887 D) The partial variance explained (PVE) by assay, cell type and donor using pseudobulk
1888 samples of 14 cell types with matched 3' single-cell and single-nucleus libraries of 5 donors.
1889 The distributions of PVEs for the 22,707 genes are plotted as violin plots with median (bar),
1890 IQR (box), and outliers (dots). E) Tanglegram of Cell Type Dendrograms. Dendrograms are
1891 derived from hierarchical clustering using Euclidean distances and complete agglomeration of
1892 cell types per assay whereas thicker lines represent more confident clustering. The connecting
1893 lines between the dendrograms represent the local optimal solution of similarity based on a
1894 stepwise greedy forward selection algorithm (entanglement = 0.01), resulting in a nearly
1895 congruent tree structure. Some lines are not connected due to differences in clustering
1896 structures and the algorithm not forcing matches when similarity is insufficient. Dashed vs. full
1897 thick lines indicate slight differences in internal branch node splitting from the greedy
1898 algorithm, but the branch tips maintain the same similarity order. F) Mean gene expression,
1899 variance of mean expression and detection rate weighted by number of cells per library for
1900 single cell and single nuclei assays. The distributions of all genes with nonzero values in both
1901 assays (average and variance) as well as all 22,707 genes detected in either assay (detection)
1902 are plotted as violin plots with median (bar), interquartile range (box) and outliers (dots). G)
1903 Cellular detection rate per assay weighted by number of reads and cells per library for single
1904 cell and single nuclei assay, purity and silhouette index per assay. The dot represents the
1905 weighted median per assay.
1906

1907 **Supplementary Figure 4.** A) UMAP embedding of 1,842 cells from the Smart-seq2 and
1908 scNMT-seq transcriptomics dataset, showing L1 cell type populations, high mitochondrial
1909 (hiMT) cells, and injured state (INJ) cells. The top-right inset displays the same UMAP with
1910 colors indicating the scRNA-seq technology used for the experiment. B) Average expression
1911 values across Smart-seq2 and scNMT-seq cell populations in the mBDRC consensus marker
1912 analysis, consistent with Fig. 2D. C) Upset plots displaying statistically significant genes
1913 identified in downsampled data across mBDRC protocols, including Smart-seq2/scNMT-seq
1914 data, for podocytes (POD), proximal tubule (PT), thick ascending limb (TAL), and connecting
1915 tubule (CNT) populations markers.
1916

1917 **Supplementary Figure 5.** A) Percentage CpG and GpC sites covered by reads vs mapping
1918 efficiency (%). B) Promoter accessibility confirmed by GpC methylation of genes with open

1919    chromatin based on snATAC data common between PT and DCT. C) Promotor accessibility
1920    determined by GpC methylation of genes transcribed vs non-transcribed. D) CpG promoter and
1921    gene body methylation of genes with open chromatin based on ATAC data common between
1922    PT and DCT. E) CpG promoter and gene body methylation of genes transcribed vs non-
1923    transcribed.

1924

1925    **Supplementary Figure 6.** For each cell-type panel: A) Mean gene expression. A.1
1926    Logarithmic weighted mean gene expression of snRNA-seq against scRNA-seq using log-
1927    normalized counts with dashed bisecting line and Pearson correlation coefficient in red. A.2
1928    Logarithmic weighted mean gene expression values of snRNA-seq and scRNA-seq. A.3 KS
1929    Statistic comparing the distribution of mean gene expression values within snRNA-seq,
1930    between sc- and sn-RNA-seq and within scRNA-seq. Blue indicates within the same donor,
1931    orange across donors. B) Variance of gene expression. B.1 Logarithmic weighted gene
1932    expression variance of snRNA-seq against scRNA-seq using log-normalized counts with
1933    dashed bisecting line and Pearson correlation coefficient in red. B.2 Logarithmic weighted gene
1934    expression variance of snRNA-seq and scRNA-seq. B.3 KS Statistic comparing the distribution
1935    of gene expression variance within snRNA-seq, between sc- and sn-RNA-seq and within
1936    scRNA-seq. Blue indicates within the same donor, orange across donors. C) Detection of gene
1937    expression. C.1 Gene expression detection of snRNA-seq against scRNA-seq using log-
1938    normalized counts with dashed bisecting line and Pearson correlation coefficient in red. C.2
1939    Gene expression detection of snRNA-seq and scRNA-seq. C.3 KS Statistic comparing the
1940    distribution of gene expression detection within snRNA-seq, between sc- and sn-RNA-seq and
1941    within scRNA-seq. Blue indicates within the same donor, orange across donors. D) Cellular
1942    Detection Rate. D.1 Proportion of nonzero gene counts in snRNA-seq and scRNA-seq. D.2 KS
1943    Statistic comparing the distribution of cellular detection rate within snRNA-seq, between sc-
1944    and snRNA-seq and within scRNA-seq. Blue indicates within the same donor, orange across
1945    donors. E) Purity. E.1 Purity of cell type cluster neighborhoods in snRNA-seq and scRNA-seq,
1946    dashed line indicates 50% pure neighborhood so half of the cells belong to the same cell type
1947    cluster. E.2 KS Statistic comparing the distribution of purity values within snRNA-seq,
1948    between sc- and snRNA-seq and within scRNA-seq. Blue indicates within the same donor,
1949    orange across donors. F) Silhouette. F.1 Silhouette values of cells within a cell type cluster in
1950    snRNA-seq and scRNA-seq, dashed line indicates 0 so below this value cells are closer to
1951    another cell type cluster or suited for forming their own cluster. F.2 KS Statistic comparing the
1952    distribution of silhouette values within snRNA-seq, between sc- and sn-RNA-seq and within
1953    scRNA-seq. Blue indicates within the same donor, orange across donors.

1954

1955    **Supplementary Figure 7.** A) Benchmarking tables for Horizontal integration across mBDRC
1956    RNA protocols on the different cell type annotation granularity levels (L1 left, L2 right)
1957    showing methods ranked from best to poorest performance for bio conservation (preserving
1958    biological differences) and batch correction. B) scVI associated UMAP embedding (97,125
1959    cells) of RNA protocols integration showing clustering-based annotation for the two cell type
1960    granularity levels L1 and L2. C) Heatmap showing the confusion matrix between Label transfer
1961    from HCA external resource (X axis) and Clustering-based annotation on L2 cell type
1962    annotation (Y axis). D) scVI associated UMAP embedding for individual RNA protocols

(53,799 / 35,513 / 7813 cells respectively) showing transferred labels from HCA external resources for the two cell type granularity levels L1 and L2. E) scVI associated UMAP embedding (97,125 cells) of sn/scRNA protocol integration showing cells' respective protocol. F) AUC score between label transfer annotation from HCA external resource and clustering-based annotation on L1 annotation level.

**Supplementary Figure 8.** A) Benchmarking tables for Horizontal Integration of each mBDRC RNA protocols (Top to bottom: snRNA, 3' scRNA, 5' scRNA) on the different cell type annotation levels (L1 left, L2 right) showing methods ranked from best to poorest performance. B) Violin plots of silhouette width cell wise values by L1 transferred labels from HCA on scVI RNA Integration embedding. C) Violin plots of silhouette width cell wise values by L2 transferred labels from HCA on scVI RNA Integration embedding.

**Supplementary Figure 9.** A) Marker detection rate comparison based on the feature importance of L1 groups across data models. B) Bar plots for Integrated model, snRNA, 5' scRNA, and 3' scRNA (left to right) showing fraction of markers composing detection rate metric in L1 groups that are common for every data model, shared by some, or uniquely detected by one data model. C) Bar plots for Integrated model, snRNA, 5' scRNA, and 3' scRNA (right to left) showing fraction of markers composing detection rate metric over every L1 group that are common for every data model, shared by some, or uniquely detected by one data model. D) Feature importance score comparison between 3' scRNA-seq (left) / 5' (right) and the integrated dataset for markers in Fig. 3G, which were found significant in snRNA-seq and integration models. E) Dot plot showing the log-normalized expression of highlighted genes in panel D across L1 cell-types per RNA protocol F) Heatmap showing AUC score between label transfer annotation from HCA external resource and clustering-based annotation on L1 annotation level across protocols. G) AUC scores for the predictability of each data model in predicting the 3' scRNA-seq (top) / 5' (bottom) data type for L1 cell-type annotations. Colored bars indicate distinct models for each technology, and dots represent L1 cell types.

**Supplementary Figure 10.** A) Violin plots of WNN modality weights per L1 cell type, obtained using the baseline embedding method implemented in Seurat. B) Spectral MNN-associated UMAP embedding (37,717 cells) showing snATAC-seq data clustering-based annotation at L2 resolution. C) Heatmap displaying the confusion matrix comparing clustering-based annotations from the best-performing horizontal integrations of both snRNA-seq (X-axis) and snATAC-seq (Y-axis) datasets using the same nuclei. D) WNN-associated UMAP embedding (37,717 cells) showing scOMM annotations obtained in the snRNA-seq data by projecting onto the HCA external reference data at L2 cell type resolution. E) Scatter plots comparing silhouette scores from optimized WNN (X-axis) and multi-spectral (Y-axis) integration methods across various kidney cell types and subtypes. Each blue dot represents a single cell, while the red line indicates the trend, with pink shading showing the confidence interval. Marginal histograms illustrate the distribution of silhouette scores along each axis. The "Avg Dev (y)" value represents the average deviation of Y-axis (multi-spectral scores) for each cell type or subtype. F) Marker detection rates from the scOMM model across different classification scenarios.

**Supplementary Figure 11.** A) Benchmarking of methods associated with UMAPs embeddings (60,336 cells) of snATAC-seq showing label transferred annotations from HCA external resource for the L2 cell type granularity. B) Benchmarking tables for Vertical Integration of Multiome ATAC-RNA modalities on the different cell type annotation granularity levels (L1 & L2) showing methods ranked from best to poorest performance (top to bottom) C) Radar chart for model performance metrics (Sensitivity, Error Rate, Classification Rate, Accuracy and Sensitivity) across classification scenarios. D) Heatmaps showing the confusion matrix between clustering-based annotation on L2 cell type annotation and scOMM predictions for cross-modality classification (left) and bridge classification (right). E) Pearson correlation values across L2 groups between paired snRNA-seq gene expression profiles and gene activity profiles derived from snATAC-seq data from the same cells.

**Supplementary Figure 12.** A) MultiVI associated UMAP embedding (119,744 cells) for mosaic integration showing scOMM annotations from HCA external resource with L2 cell type resolution. B) GLUE associated UMAP embedding (134,842 cells) for diagonal integration showing scOMM annotations from HCA external resource with L2 cell type resolution. C) Summary heatmap showing L2 cell type silhouette score over L1 cell type populations across every embedding scenario, unimodal and multimodal integrations. D) Summary heatmap showing L2 cell type silhouette score over L1 cell type populations for the stromal compartment across every embedding scenario, unimodal and multimodal integrations; aFIB population is highlighted with a red box.

**Supplementary Figure 13.** A) TAL compartment recomputed UMAP embedding (16,301 cells) as in Fig. 5B showing the log-normalized expression of top markers for aTAL1_0 cluster B) Violin plot of *WFDC2* log-normalized expression on TAL substates split by sn/scRNA-seq data. C) Average expression values of top marker genes across aTAL1 subclusters, split by RNA technology. D) Stromal compartment on UMAP embedding (3386 cells) of Horizontal Integration across mBDRC RNA protocols displaying scOMM L2 labels, together with Norn cells associated markers and Norn signature score.

**Supplementary Figure 14.** A) Cell-type annotation legend. B) Left panel: Pearson correlations between observed and out-of-sample predicted log-expression per gene, stratified by cell type. Right panel: Distribution of Pearson correlations, aggregated across cell types. Asterisks indicate significant differences in predictive accuracy across platforms (independent two-tailed T-test; *** = $p < 0.001$; * = $p < 0.05$). C) Left panel: Median number of transcription factors (TFs) identified per cell type across biological replicates, represented by the top of each bar, with whiskers indicating the minimum and maximum TF counts. The bottom of each bar represents the number of consensus TFs, found consistently in all replicates within each cell type. Right panel: Distribution of cell-type-specific "reproducibility ratios" per platform, calculated as the ratio of the number of TFs in the cross-sample consensus network to the mean number of TFs per replicate (independent two-tailed T-test; *** = $p < 0.001$). D) Ward's hierarchical clustering of consensus regulatory networks per platform and cell type, based on Jaccard similarity of TF-target indicator vectors. Row and column annotations indicate platform (batch) and cell type. E) Left panel: number of cell type-TF assignments identified in

each consensus network (diagonal) and shared between consensus networks across platforms (off-diagonal). Right panel: log-odds ratios of shared cell type-TF assignments across platforms (Fisher's exact test; *** = p < 0.001). F) Left panel: Fraction of consensus TFs per platform that are also nuclear-localized in kidney tissue, per HPA immunohistochemistry annotations. Right panel: Distribution of these fractions across platforms, aggregated over cell types; no significant differences were observed between platforms (independent two-tailed T-test). G) Left panel: Fraction of consensus TF-target links per cell type with supporting evidence from intersected bulk ChIP-Seq and ATAC-Seq peaks at kidney promoters. Right panel: Distribution of these fractions across platforms, aggregated over cell types, with statistical significance between platforms indicated (independent two-tailed T-test; *** = p < 0.001).

**References**

1. Denisenko, E. *et al.* Systematic assessment of tissue dissociation and storage biases in single-cell and single-nucleus RNA-seq workflows. *Genome Biol* **21**, (2020).

2. Jiménez-Gracia, L. *et al.* FixNCut: single-cell genomics through reversible tissue fixation and dissociation. *Genome Biol* **25**, (2024).

3. Regev, A. *et al.* Science Forum: The Human Cell Atlas. *Elife* (2017).

4. Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin. *Cell* **183**, (2020).

5. Lake, B. B. *et al.* An atlas of healthy and injured cell states and niches in the human kidney. *Nature* **619**, (2023).

6. Wang, G. *et al.* Integrating genetics with single-cell multiomic measurements across disease states identifies mechanisms of beta cell dysfunction in type 2 diabetes. *Nat Genet* **55**, (2023).

7. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, (2021).

8. De Simone, M. *et al.* Comparative Analysis of Commercial Single-Cell RNA Sequencing Technologies. *bioRxiv* 2024–2026 (2024).

9. Mereu, E. *et al.* Benchmarking single-cell RNA-sequencing protocols for cell atlas projects. *Nat Biotechnol* **38**, (2020).

10. Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing methods. *Nat Biotechnol* **38**, (2020).

11. De Rop, F. V. *et al.* Systematic benchmarking of single-cell ATAC-sequencing protocols. *Nat Biotechnol* **42**, (2024).

12. Luecken, M. D. *et al.* Benchmarking atlas-level data integration in single-cell genomics. *Nat Methods* **19**, (2022).

13. Lee, M. Y. Y., Kaestner, K. H. & Li, M. Benchmarking algorithms for joint integration of unpaired and paired single-cell RNA-seq and ATAC-seq data. *Genome Biol* **24**, (2023).

14. Argelaguet, R., Cuomo, A. S. E., Stegle, O. & Marioni, J. C. Computational principles and challenges in single-cell data integration. *Nature Biotechnology* vol. 39 Preprint at https://doi.org/10.1038/s41587-021-00895-7 (2021).

15. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods* **10**, (2013).

16. Clark, S. J. *et al.* ScNMT-seq enables joint profiling of chromatin accessibility DNA methylation and transcription in single cells e. *Nat Commun* **9**, (2018).

17. Muto, Y. *et al.* Single cell transcriptional and chromatin accessibility profiling redefine cellular heterogeneity in the adult human kidney. *Nat Commun* **12**, (2021).

18. Ashuach, T. *et al.* MultiVI: deep generative model for the integration of multimodal data. *Nat Methods* **20**, (2023).

19. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat Methods* **16**, (2019).

20. Hao, Y. *et al.* Dictionary learning for integrative, multimodal and scalable single-cell analysis. *Nat Biotechnol* **42**, (2024).

21. Koehler, S. *et al.* Scaffold polarity proteins Par3A and Par3B share redundant functions while Par3B acts independent of atypical protein kinase C/Par6 in podocytes to maintain the kidney filtration barrier. *Kidney Int* **101**, (2022).

22. Daga, S. *et al.* New frontiers to cure Alport syndrome: COL4A3 and COL4A5 gene editing in podocyte-lineage cells. *European Journal of Human Genetics* **28**, (2020).

23. Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-cell transcriptomics. *Nat Methods* **15**, (2018).

24. Zhang, K., Zemke, N. R., Armand, E. J. & Ren, B. A fast, scalable and versatile tool for analysis of single-cell omics data. *Nat Methods* **21**, (2024).

25. Stuart, T., Srivastava, A., Madad, S., Lareau, C. A. & Satija, R. Single-cell chromatin state analysis with Signac. *Nat Methods* **18**, (2021).

26. Ashuach, T., Reidenbach, D. A., Gayoso, A. & Yosef, N. PeakVI: A deep generative model for single-cell chromatin accessibility analysis. *Cell Reports Methods* **2**, (2022).

27. Martens, L. D., Fischer, D. S., Yépez, V. A., Theis, F. J. & Gagneur, J. Modeling fragment counts improves single-cell ATAC-seq analysis. *Nat Methods* **21**, (2024).

28. Luo, S., Germain, P. L., Robinson, M. D. & von Meyenn, F. Benchmarking computational methods for single-cell chromatin data analysis. *Genome Biol* **25**, (2024).

29. Cao, Z. J. & Gao, G. Multi-omics single-cell data integration and regulatory inference with graph-linked embedding. *Nat Biotechnol* **40**, (2022).

30. Yang, Z. *et al.* Human Epididymis Protein 4: A Novel Biomarker for Lupus Nephritis and Chronic Kidney Disease in Systemic Lupus Erythematosus. *J Clin Lab Anal* **30**, (2016).

31. Sivakamasundari, V. *et al.* Comprehensive Cell Type Specific Transcriptomics of the Human Kidney. *bioRxiv* (2017).

32. Bingle, L. *et al.* WFDC2 (HE4): A potential role in the innate immunity of the oral cavity and respiratory tract and the development of adenocarcinomas of the lung. *Respir Res* **7**, (2006).

33. Kragesteen, B. K. *et al.* The transcriptional and regulatory identity of erythropoietin producing cells. *Nat Med* **29**, (2023).

34. Haase, V. H. Regulation of erythropoiesis by hypoxia-inducible factors. *Blood Rev* **27**, (2013).

35. Cerezo, M. *et al.* The NHGRI-EBI GWAS Catalog: standards for reusability, sustainability and diversity. *Nucleic Acids Res* **53**, D998–D1005 (2024).

36. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, (2015).

37. Thibodeau, A., Uyar, A., Khetan, S., Stitzel, M. L. & Ucar, D. A neural network based model effectively predicts enhancers from clinical ATAC-seq samples. *Sci Rep* **8**, (2018).

38. Kim, T. H. *et al.* Analysis of the Vertebrate Insulator Protein CTCF-Binding Sites in the Human Genome. *Cell* **128**, (2007).

39. Pang, B. & Snyder, M. P. Systematic identification of silencers in human cells. *Nat Genet* **52**, (2020).

40. Murata, M. *et al.* Detecting expressed genes using CAGE. *Methods in Molecular Biology* **1164**, (2014).

41. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues. *Nature* **507**, (2014).

42. Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, (2018).

43. Loeb, G. B. *et al.* Variants in tubule epithelial regulatory elements mediate most heritable differences in human kidney function. *Nat Genet* **56**, 2078–2092 (2024).

44. Peti-Peterdi, J. & Harris, R. C. Macula densa sensing and signaling mechanisms of renin release. *Journal of the American Society of Nephrology* vol. 21 Preprint at https://doi.org/10.1681/ASN.2009070759 (2010).

45. Mayet, J. & Hughes, A. Cardiac and vascular pathophysiology in hypertension. *Heart* vol. 89 Preprint at https://doi.org/10.1136/heart.89.9.1104 (2003).

46. Schaub, J. A. *et al.* Quantitative morphometrics reveals glomerular changes in patients with infrequent segmentally sclerosed glomeruli. *J Clin Pathol* **75**, (2022).

47. Pippin, J. W. *et al.* Upregulated PD-1 signaling antagonizes glomerular health in aged kidneys and disease. *Journal of Clinical Investigation* **132**, (2022).

48. Atchison, D. K. *et al.* Hypertension induces glomerulosclerosis in phospholipase C-ε1 deficiency. *Am J Physiol Renal Physiol* **318**, (2020).

49. Macaulay, I. C. *et al.* Separation and parallel sequencing of the genomes and transcriptomes of single cells using G&T-seq. *Nat Protoc* **11**, (2016).

50. Zheng, G. X. Y. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat Commun* **8**, (2017).

51. DI Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nature Biotechnology* vol. 35 Preprint at https://doi.org/10.1038/nbt.3820 (2017).

52. Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, (2013).

53. Picard. *http://broadinstitute.github.io/picard/*.

54. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. *Nat Methods* **14**, (2017).

55. Soneson, C., Love, M. I. & Robinson, M. D. Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Res* **4**, (2015).

56. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, (2019).

57. Clark, S. J. *et al.* Genome-wide base-resolution mapping of DNA methylation in single cells using single-cell bisulfite sequencing (scBS-seq). *Nat Protoc* **12**, (2017).

58. Krueger, F. & Andrews, S. R. Bismark: A flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* **27**, (2011).

59. Harrel Jr., F. E. Regression Modeling Strategies - With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis. *R Software* **70**, (2015).

60. Akaike, H. A New Look at the Statistical Model Identification. in (1974). doi:10.1007/978-1-4612-1694-0_16.

2188    61.    Gayoso, A., Shor, J., Carr, A. J., Sharma, R. & Pe'er, D. DoubletDetection. Preprint at
2189            https://doi.org/10.5281/zenodo.2678041 (2020).
2190    62.    Bredikhin, D., Kats, I. & Stegle, O. MUON: multimodal omics analysis framework.
2191            *Genome Biol* **23**, (2022).
2192    63.    Thibodeau, A. *et al.* AMULET: a novel read count-based method for effective multiplet
2193            detection from single nucleus ATAC-seq data. *Genome Biol* **22**, (2021).
2194    64.    Galili, T. dendextend: An R package for visualizing, adjusting and comparing trees of
2195            hierarchical clustering. *Bioinformatics* **31**, (2015).
2196    65.    Lun, A. bluster: Clustering Algorithms for Bioconductor. Preprint at
2197            https://bioconductor.org/packages/bluster (2024).
2198    66.    Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: Large-scale single-cell gene
2199            expression data analysis. *Genome Biol* **19**, (2018).
2200    67.    Hie, B., Bryson, B. & Berger, B. Efficient integration of heterogeneous single-cell
2201            transcriptomes using Scanorama. *Nat Biotechnol* **37**, (2019).
2202    68.    Haghverdi, L., Lun, A. T. L., Morgan, M. D. & Marioni, J. C. Batch effects in single-cell
2203            RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat
2204            Biotechnol* **36**, (2018).
2205    69.    Traag, V. A., Waltman, L. & van Eck, N. J. From Louvain to Leiden: guaranteeing well-
2206            connected communities. *Sci Rep* **9**, (2019).
2207    70.    Robin, X. *et al.* pROC: An open-source package for R and S+ to analyze and
2208            compare ROC curves. *BMC Bioinformatics* **12**, (2011).
2209    71.    Waskom, M. seaborn: statistical data visualization. *J Open Source Softw* **6**, (2021).
2210    72.    Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *Journal of Machine
2211            Learning Research* **12**, (2011).
2212    73.    Moody, J. *et al.* SCAFE: a software suite for analysis of transcribed cis-regulatory
2213            elements in single cells. *Bioinformatics* **38**, (2022).
2214    74.    Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-
2215            wide association summary statistics. *Nat Genet* **47**, (2015).
2216    75.    UK Biobank — Neale lab. *http://www.nealelab.is/uk-biobank/*.
2217    76.    Badia-i-Mompel, P. *et al.* Gene regulatory network inference in the era of single-cell
2218            multi-omics. *Nature Reviews Genetics* vol. 24 Preprint at
2219            https://doi.org/10.1038/s41576-023-00618-5 (2023).
2220    77.    Aibar, S. *et al.* SCENIC: Single-cell regulatory network inference and clustering. *Nat
2221            Methods* **14**, (2017).
2222    78.    Bravo González-Blas, C. *et al.* SCENIC+: single-cell multiomic inference of enhancers
2223            and gene regulatory networks. *Nat Methods* **20**, (2023).
2224    79.    Ke, G. *et al.* LightGBM: A highly efficient gradient boosting decision tree. in *Advances
2225            in Neural Information Processing Systems* vols 2017-December (2017).
2226    80.    Rauluseviciute, I. *et al.* JASPAR 2024: 20th anniversary of the open-access database
2227            of transcription factor binding profiles. *Nucleic Acids Res* **52**, (2024).
2228    81.    Stovner, E. B. & Sætrom, P. PyRanges: Efficient comparison of genomic intervals in
2229            Python. *Bioinformatics* **36**, (2020).
2230    82.    Abugessaisa, I. *et al.* refTSS: A Reference Data Set for Human and Mouse
2231            Transcription Start Sites. *J Mol Biol* **431**, (2019).
2232    83.    Uhlén, M. *et al.* Tissue-based map of the human proteome. *Science (1979)* **347**,
2233            (2015).

2234    84.    Zou, Z., Ohta, T., Miura, F. & Oki, S. ChIP-Atlas 2021 update: a data-mining suite for
2235           exploring epigenomic landscapes by fully integrating ChIP-seq, ATAC-seq and
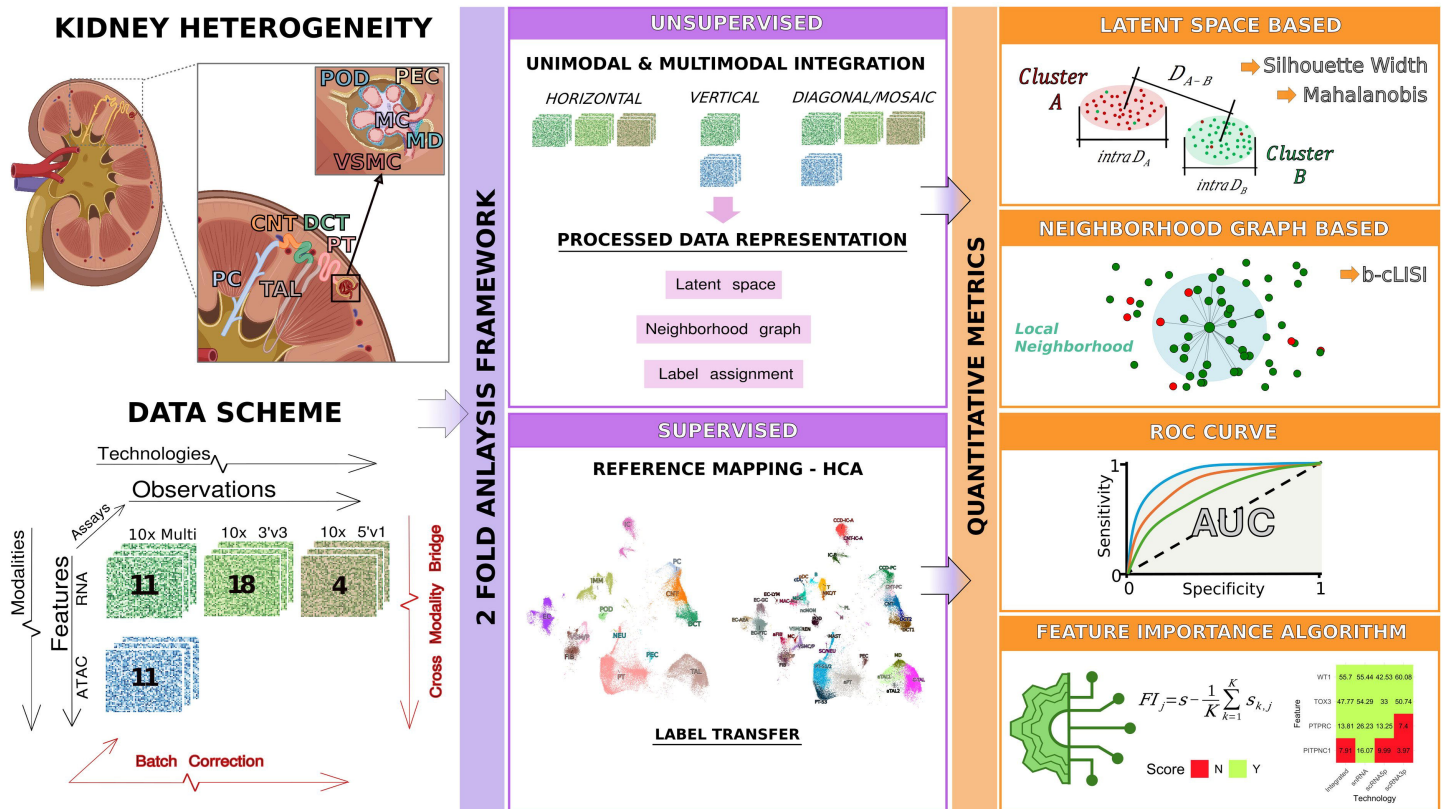2236           Bisulfite-seq data. *Nucleic Acids Res* **50**, (2022).
2237
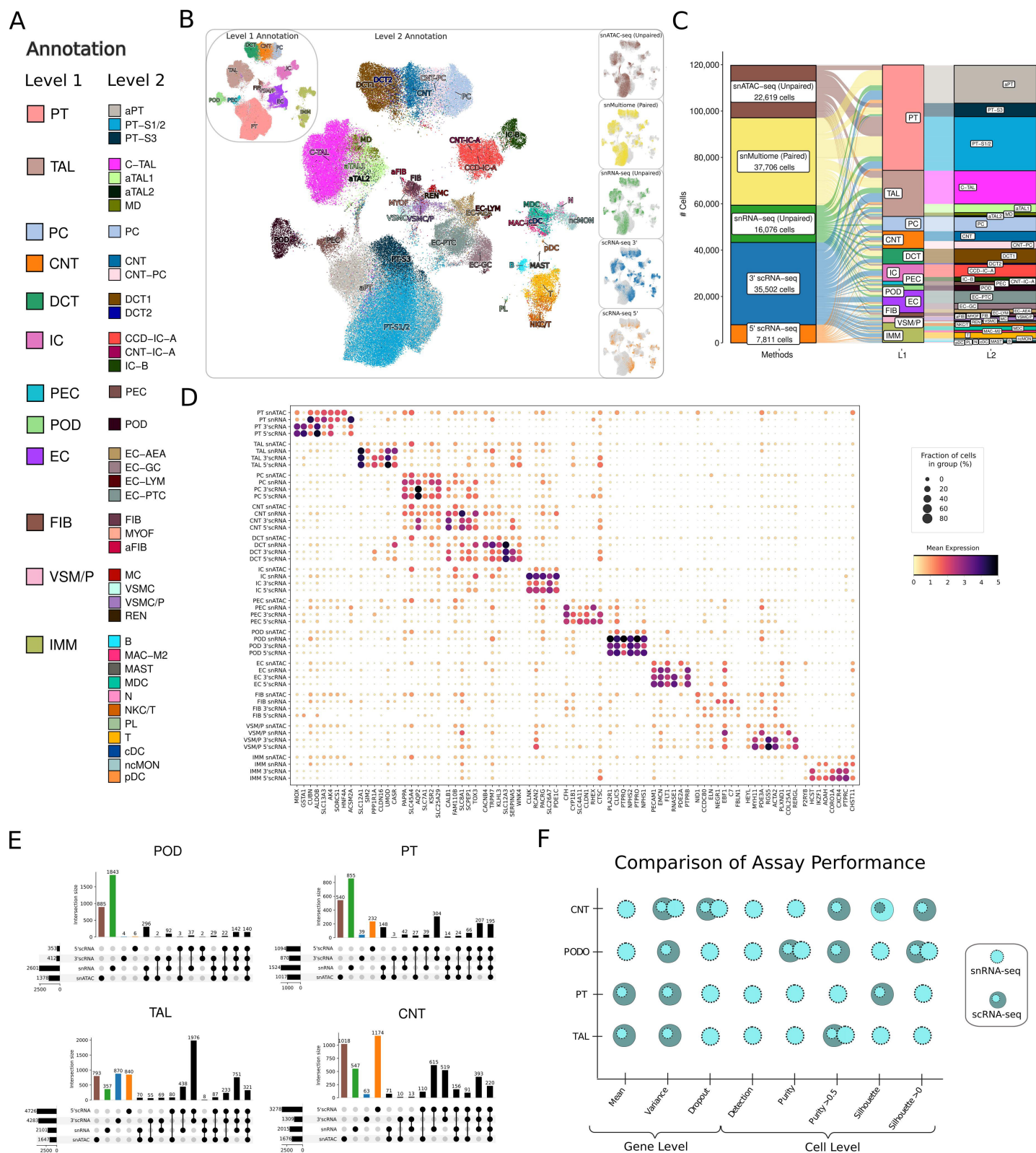2238
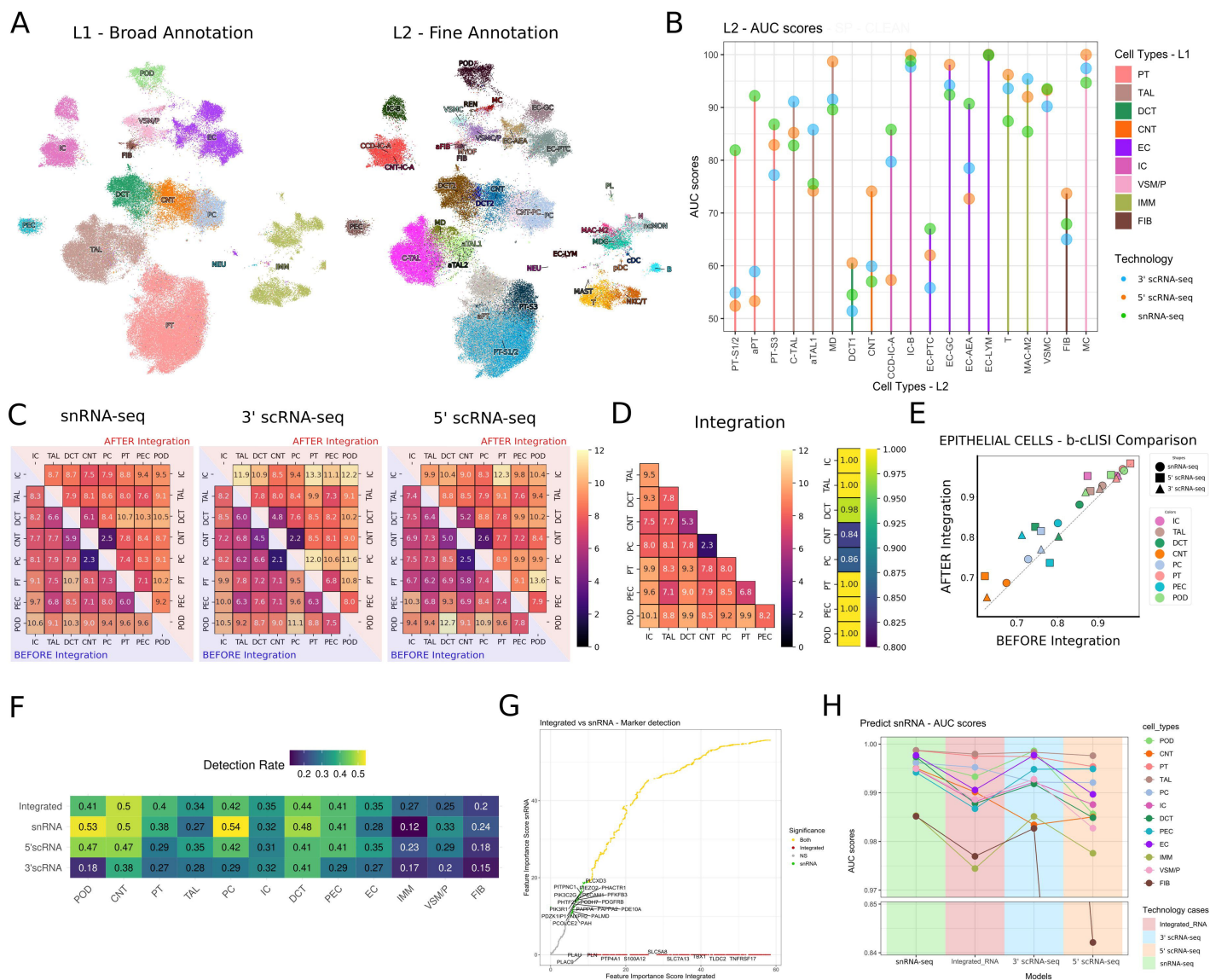2239
2240
2241

**Figure 1**

**Figure 2**

**Figure 3**

**Figure 4**

**Figure 5**

**Figure 6**

**Supplementary Figure 1**

**Supplementary Figure 2**

**Supplementary Figure 4**

**Supplementary Figure 6**

**Supplementary Figure 8**

A

**snRNA-seq (L1)**

| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| Harmony | 0.50 | 0.62 | 0.41 | 0.57 | 1.00 | 0.90 | 0.23 | 0.27 | 0.66 | 0.75 | 0.56 | 0.62 | 0.60 |
| scVI | 0.48 | 0.68 | 0.42 | 0.52 | 1.00 | 0.89 | 0.26 | 0.15 | 0.72 | 0.71 | 0.55 | 0.62 | 0.59 |
| Scanorama | 0.50 | 0.63 | 0.37 | 0.56 | 1.00 | 0.91 | 0.17 | 0.12 | 0.68 | 0.22 | 0.42 | 0.61 | 0.56 |
| Unintegrated | 0.52 | 0.59 | 0.35 | 0.56 | 1.00 | 0.89 | 0.07 | 0.05 | 0.70 | 0.00 | 0.34 | 0.60 | 0.50 |

**snRNA-seq (L2)**

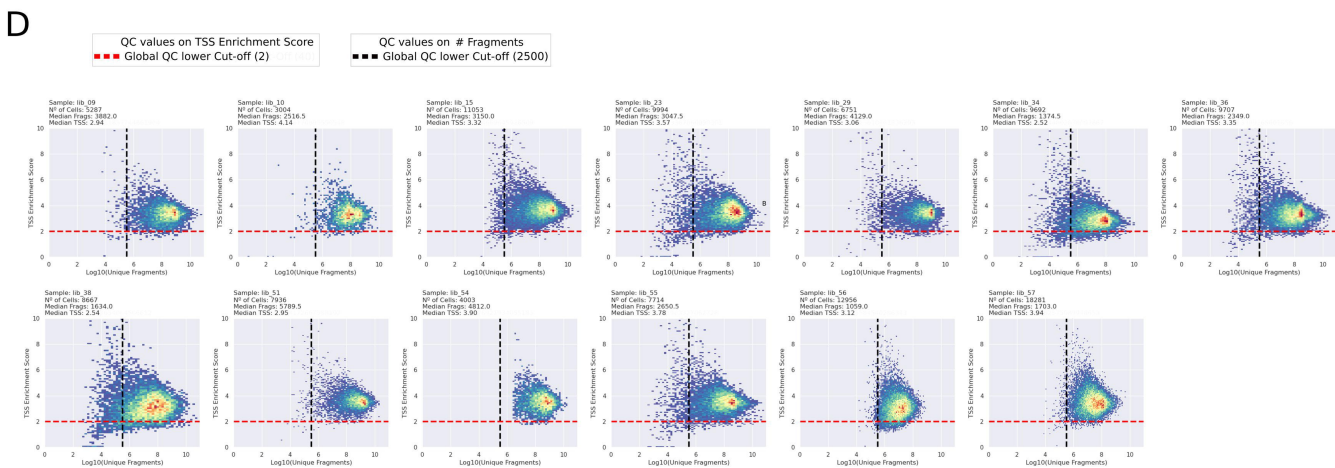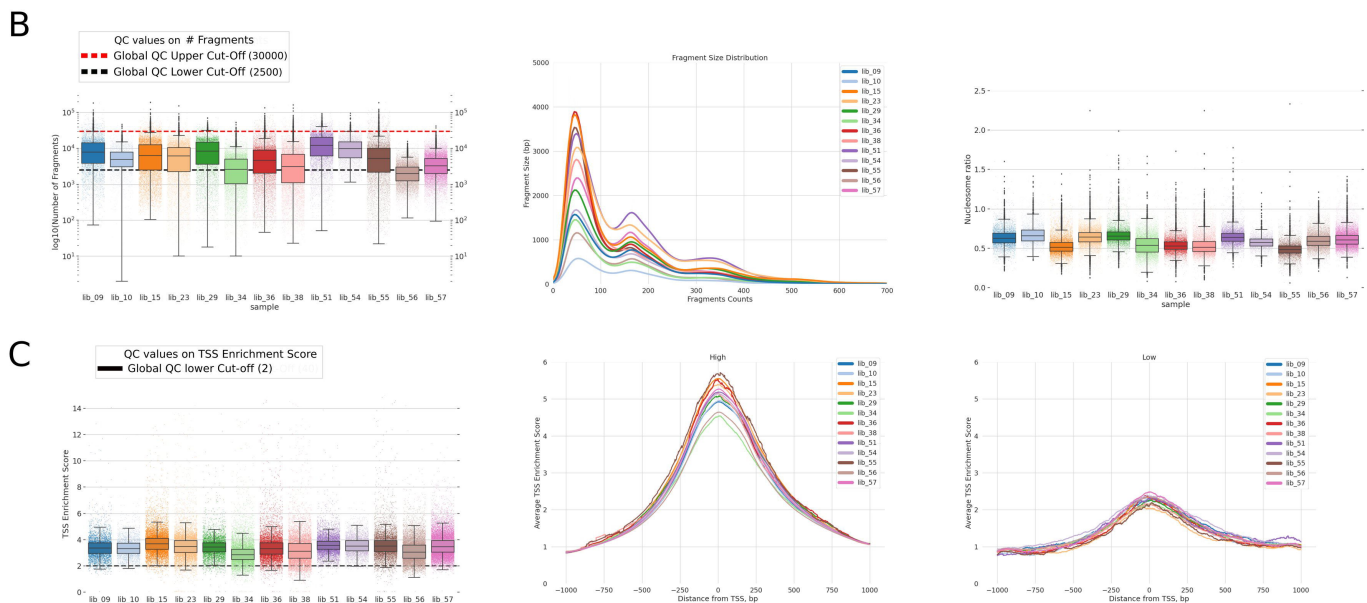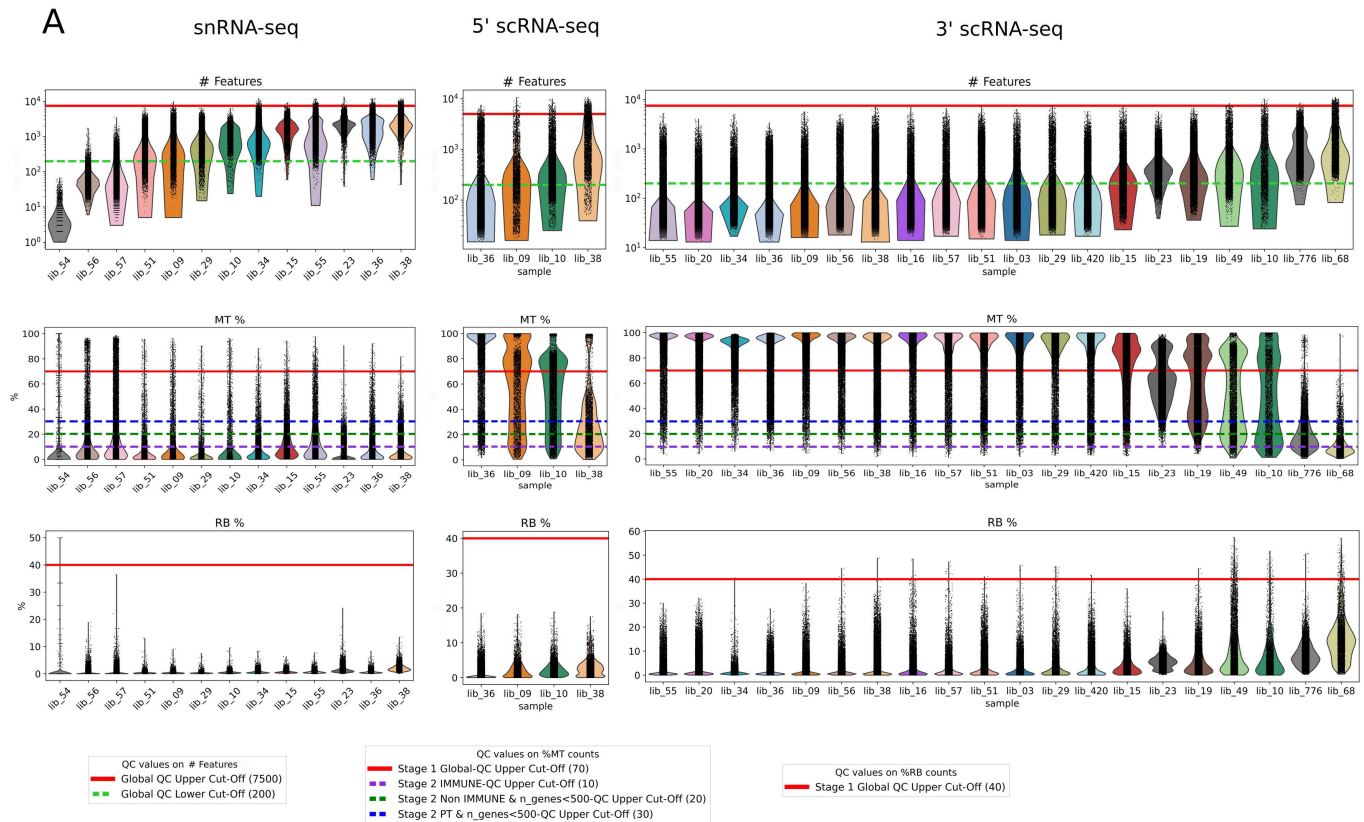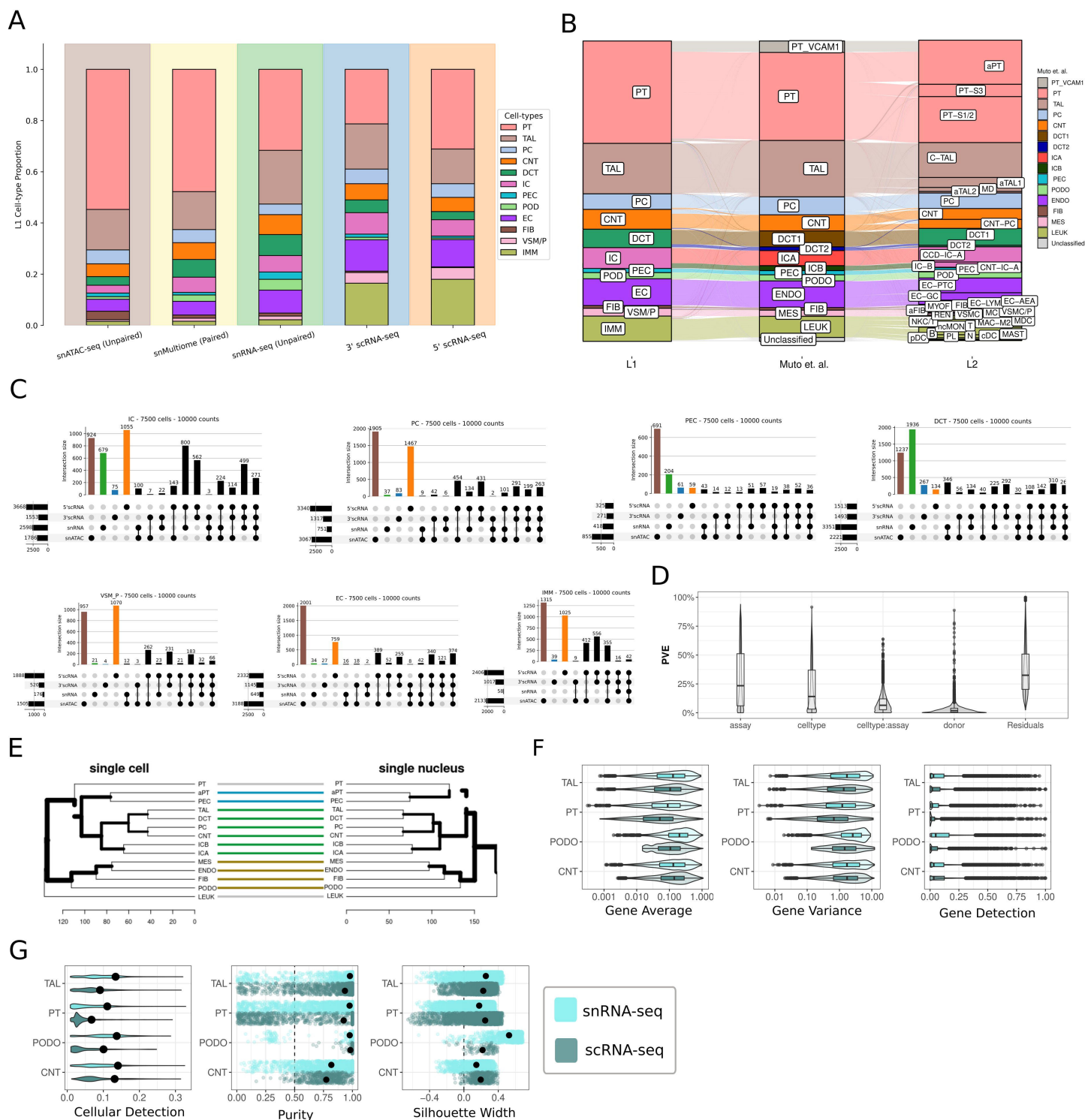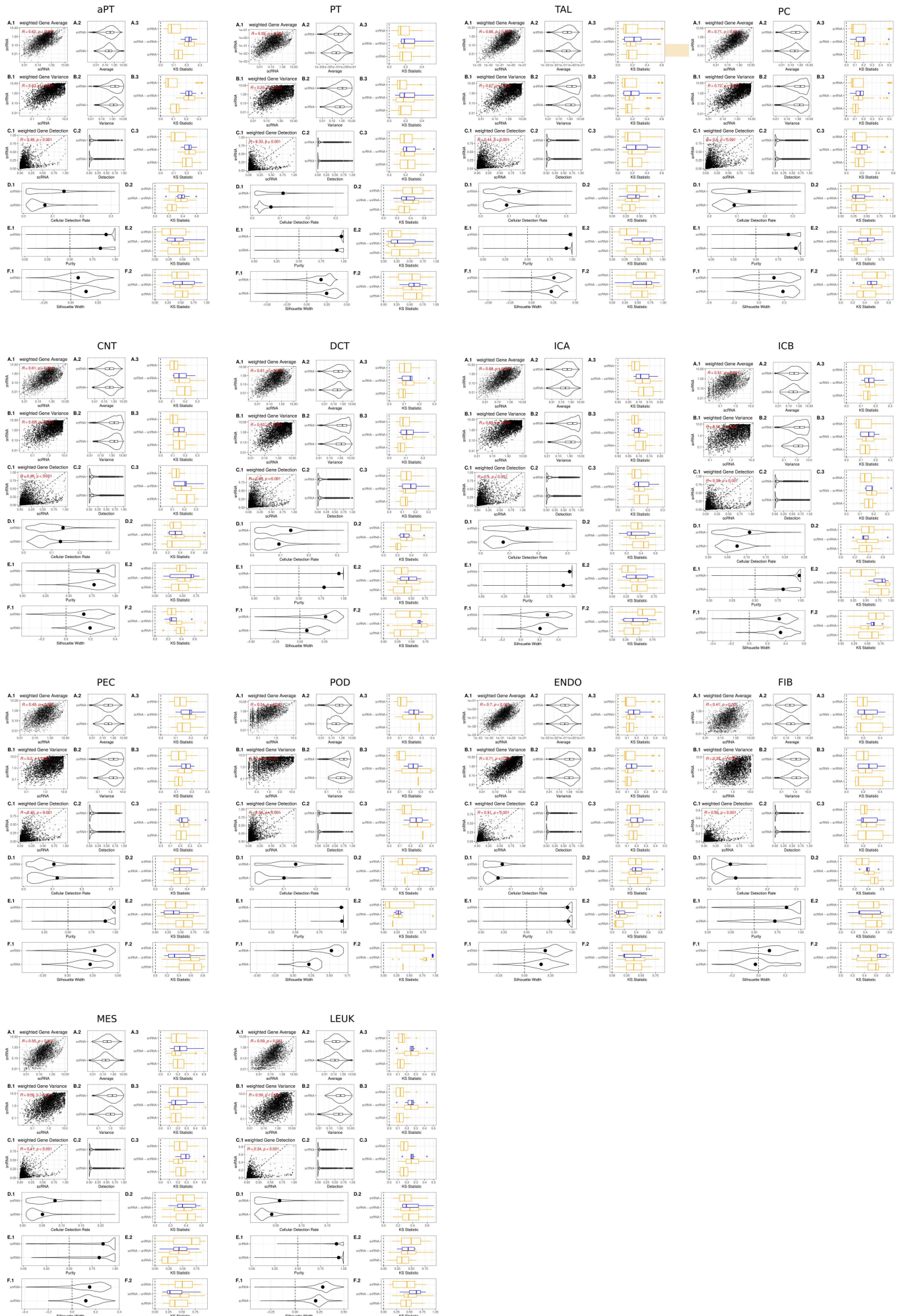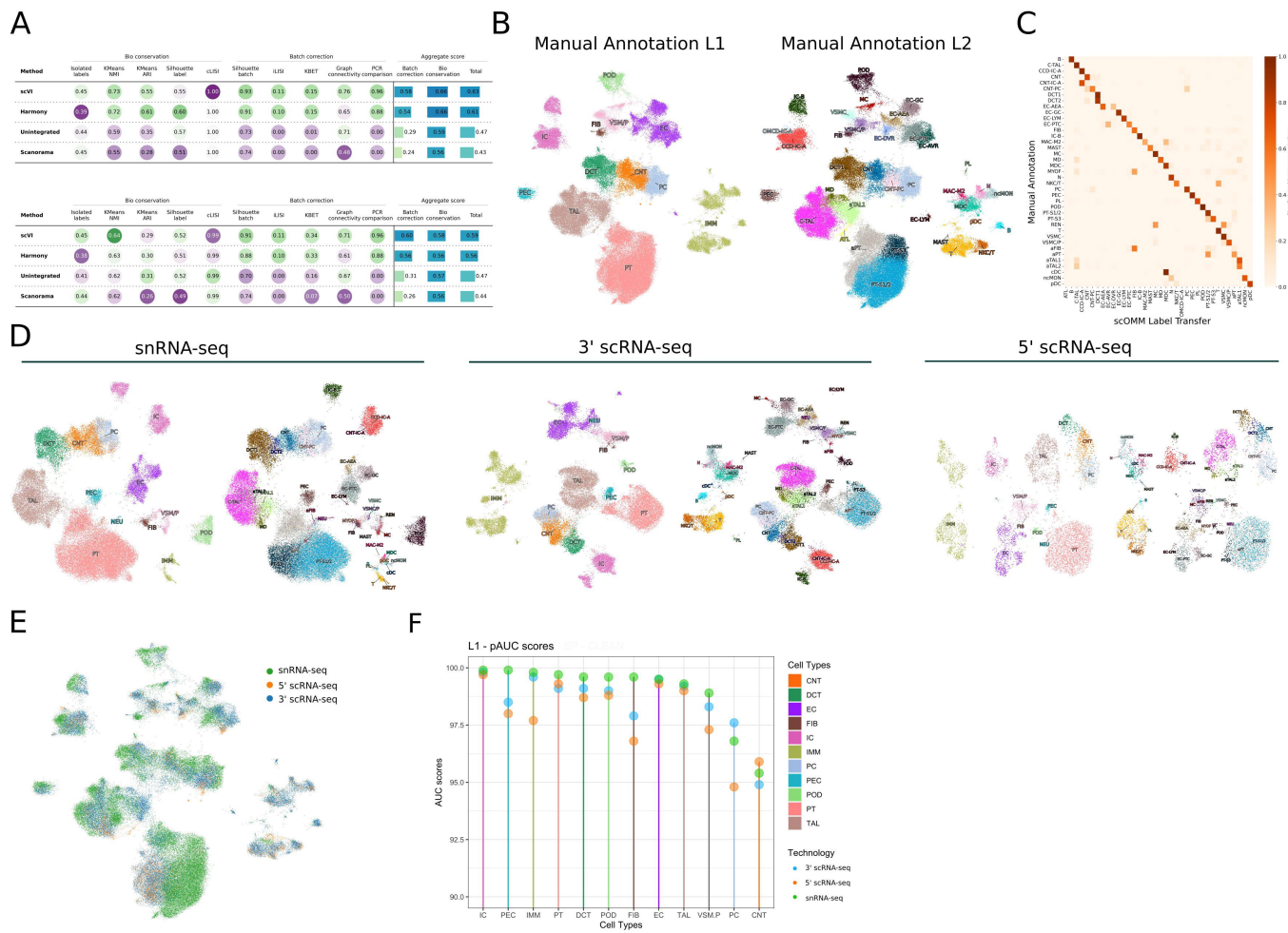| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| scVI | 0.45 | 0.60 | 0.24 | 0.51 | 0.99 | 0.87 | 0.26 | 0.38 | 0.59 | 0.71 | 0.56 | 0.56 | 0.56 |
| Harmony | 0.43 | 0.54 | 0.21 | 0.52 | 0.99 | 0.88 | 0.23 | 0.34 | 0.50 | 0.75 | 0.54 | 0.54 | 0.54 |
| Scanorama | 0.45 | 0.57 | 0.21 | 0.53 | 0.99 | 0.90 | 0.17 | 0.33 | 0.52 | 0.22 | 0.43 | 0.55 | 0.50 |
| Unintegrated | 0.44 | 0.56 | 0.21 | 0.53 | 0.99 | 0.87 | 0.07 | 0.23 | 0.56 | 0.00 | 0.35 | 0.55 | 0.47 |

**3' scRNA-seq (L1)**

| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| scVI | 0.44 | 0.76 | 0.61 | 0.55 | 1.00 | 0.92 | 0.18 | 0.25 | 0.75 | 0.66 | 0.55 | 0.67 | 0.62 |
| Harmony | 0.35 | 0.71 | 0.56 | 0.61 | 1.00 | 0.87 | 0.20 | 0.42 | 0.71 | 0.26 | 0.49 | 0.65 | 0.58 |
| Scanorama | 0.34 | 0.75 | 0.58 | 0.60 | 1.00 | 0.90 | 0.14 | 0.24 | 0.59 | 0.00 | 0.37 | 0.65 | 0.54 |
| Unintegrated | 0.37 | 0.71 | 0.55 | 0.59 | 1.00 | 0.88 | 0.10 | 0.12 | 0.76 | 0.00 | 0.37 | 0.65 | 0.54 |

**3' scRNA-seq (L2)**

| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| scVI | 0.57 | 0.70 | 0.41 | 0.52 | 0.99 | 0.90 | 0.18 | 0.57 | 0.70 | 0.66 | 0.60 | 0.64 | 0.62 |
| Harmony | 0.60 | 0.64 | 0.32 | 0.51 | 0.99 | 0.87 | 0.20 | 0.64 | 0.61 | 0.26 | 0.52 | 0.61 | 0.57 |
| Scanorama | 0.68 | 0.70 | 0.38 | 0.52 | 0.99 | 0.90 | 0.14 | 0.54 | 0.55 | 0.00 | 0.43 | 0.63 | 0.56 |
| Unintegrated | 0.66 | 0.68 | 0.36 | 0.52 | 0.99 | 0.89 | 0.10 | 0.41 | 0.70 | 0.00 | 0.42 | 0.64 | 0.55 |

**5' scRNA-seq (L1)**

| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| scVI | 0.56 | 0.77 | 0.55 | 0.55 | 1.00 | 0.91 | 0.34 | 0.43 | 0.84 | 0.80 | 0.64 | 0.69 | 0.68 |
| Harmony | 0.55 | 0.74 | 0.51 | 0.64 | 1.00 | 0.89 | 0.33 | 0.53 | 0.82 | 0.00 | 0.51 | 0.69 | 0.62 |
| Scanorama | 0.51 | 0.78 | 0.58 | 0.64 | 1.00 | 0.89 | 0.32 | 0.42 | 0.77 | 0.00 | 0.48 | 0.70 | 0.61 |
| Unintegrated | 0.55 | 0.78 | 0.61 | 0.64 | 1.00 | 0.86 | 0.14 | 0.25 | 0.81 | 0.00 | 0.41 | 0.72 | 0.59 |

**5' scRNA-seq (L2)**

| Method | Bio conservation | | | | | Batch correction | | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Isolated labels | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | PCR comparison | Batch correction | Bio conservation | Total |
| scVI | 0.51 | 0.68 | 0.29 | 0.52 | 0.99 | 0.91 | 0.34 | 0.65 | 0.74 | 0.80 | 0.69 | 0.60 | 0.63 |
| Scanorama | 0.51 | 0.69 | 0.33 | 0.53 | 0.99 | 0.90 | 0.32 | 0.63 | 0.68 | 0.00 | 0.51 | 0.61 | 0.57 |
| Harmony | 0.47 | 0.66 | 0.30 | 0.53 | 0.99 | 0.90 | 0.33 | 0.68 | 0.71 | 0.00 | 0.52 | 0.59 | 0.56 |
| Unintegrated | 0.48 | 0.67 | 0.30 | 0.53 | 0.99 | 0.87 | 0.14 | 0.47 | 0.74 | 0.00 | 0.44 | 0.60 | 0.54 |

B

RNA Integration - Annotation L1

C

RNA Integration - Annotation L2

A. SnapAtac2 (Spectral MNN), PoissonVI, PeakVI, Signac (LSI Harmony) — UMAP embeddings.

B.

L1

| Method | Bio conservation | | | | Batch correction | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | Batch correction | Bio conservation | Total |
| Spectral_MNN | 0.66 | 0.54 | 0.62 | 1.00 | 0.86 | 0.18 | 0.18 | 0.78 | 0.50 | 0.71 | 0.62 |
| Spectral_Harmony | 0.63 | 0.42 | 0.62 | 1.00 | 0.87 | 0.17 | 0.17 | 0.77 | 0.50 | 0.67 | 0.60 |
| PeakVI | 0.61 | 0.38 | 0.64 | 1.00 | 0.85 | 0.22 | 0.19 | 0.76 | 0.50 | 0.66 | 0.60 |
| PoissonVI_fragme | 0.58 | 0.35 | 0.60 | 1.00 | 0.88 | 0.24 | 0.24 | 0.75 | 0.53 | 0.66 | 0.59 |
| Spectral | 0.62 | 0.40 | 0.61 | 1.00 | 0.86 | 0.09 | 0.13 | 0.78 | 0.46 | 0.66 | 0.58 |
| PoissonVI | 0.57 | 0.29 | 0.59 | 1.00 | 0.90 | 0.23 | 0.22 | 0.75 | 0.52 | 0.61 | 0.58 |
| LSI_Harmony | 0.45 | 0.19 | 0.51 | 0.99 | 0.87 | 0.29 | 0.41 | 0.59 | 0.54 | 0.54 | 0.54 |

L2

| Method | Bio conservation | | | | Batch correction | | | | Aggregate score | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | KMeans NMI | KMeans ARI | Silhouette label | cLISI | Silhouette batch | iLISI | KBET | Graph connectivity | Batch correction | Bio conservation | Total |
| Spectral_MNN | 0.55 | 0.25 | 0.54 | 0.99 | 0.80 | 0.18 | 0.23 | 0.60 | 0.45 | 0.58 | 0.53 |
| Spectral_Harmony | 0.54 | 0.23 | 0.54 | 0.99 | 0.82 | 0.17 | 0.21 | 0.60 | 0.45 | 0.57 | 0.53 |
| PoissonVI | 0.54 | 0.21 | 0.53 | 0.99 | 0.84 | 0.23 | 0.21 | 0.57 | 0.45 | 0.57 | 0.52 |
| PoissonVI_fragme | 0.52 | 0.19 | 0.53 | 0.99 | 0.83 | 0.24 | 0.28 | 0.50 | 0.46 | 0.56 | 0.52 |
| PeakVI | 0.54 | 0.23 | 0.54 | 0.99 | 0.79 | 0.22 | 0.21 | 0.54 | 0.43 | 0.57 | 0.52 |
| Spectral | 0.54 | 0.22 | 0.54 | 0.99 | 0.81 | 0.09 | 0.21 | 0.59 | 0.42 | 0.57 | 0.51 |
| LSI_Harmony | 0.48 | 0.17 | 0.51 | 0.98 | 0.82 | 0.29 | 0.37 | 0.43 | 0.48 | 0.53 | 0.51 |

C. Radar plot with axes: Sensitivity, Error Rate, Classification Rate, Accuracy, Specificity. Method: Cross Modality, Bridge, Integrated.

D. Cross-Mod. cell-type classification; Bridge cell-type classification — confusion matrices (Clusters vs Predicted).

E. RNA/ATAC correlation dot plot.