1  **Single-cell transcriptomic characterization of 20 organs and tissues from individual mice**
2  **creates a *Tabula Muris***
3
4
5  The *Tabula Muris* Consortium
6
7  **We have created a compendium of single cell transcriptome data from the model**
8  **organism *Mus musculus* comprising more than 100,000 cells from 20 organs and**
9  **tissues.   These data represent a new resource for cell biology, revealing gene**
10 **expression in poorly characterized cell populations and allowing for direct and**
11 **controlled comparison of gene expression in cell types shared between tissues, such**
12 **as T-lymphocytes and endothelial cells from distinct anatomical locations. Two**
13 **distinct technical approaches were used for most tissues: one approach, microfluidic**
14 **droplet-based 3'-end counting, enabled the survey of thousands of cells at relatively**
15 **low coverage, while the other, FACS-based full length transcript analysis, enabled**
16 **characterization of cell types with high sensitivity and coverage. The cumulative**
17 **data provide the foundation for an atlas of transcriptomic cell biology.**
18
19  The cell is a fundamental unit of structure and function in biology, and multicellular
20  organisms have evolved a wide variety of different cell types with specialized roles.
21  Although cell types have historically been characterized on the basis of morphology and
22  phenotype, the development of molecular methods has enabled ever more precise
23  defining of their properties, typically by measuring protein or mRNA expression
24  patterns[1].   Technological advances have enabled increasingly greater degrees of
25  multiplexing of these measurements[2-7], and it is now possible to use highly parallel
26  sequencing to enumerate nearly every mRNA molecule in a given single cell[7,8]. This
27  approach has provided many novel insights into cell biology and the composition of
28  organs from a variety of organisms[9-18].  However, while these reports provide valuable
29  characterization of individual organs, it is challenging to compare data taken with varying
30  experimental techniques in independent labs from different animals. It therefore remains
31  an open question whether data from individual organs can be synthesized and used as a
32  more general resource for biology.
33
34  Here we report a compendium of cell types from the mouse *Mus musculus*. We analyzed
35  multiple organs and tissues from the same animal, thereby generating a data set
36  controlled for age, environment and epigenetic effects.   This enables the direct
37  comparison of cell type composition between organs as well as comparison of shared cell
38  types across the entire organism. The compendium is comprised of single cell
39  transcriptome sequence data from 100,605 cells isolated from 20 organs and tissues (Fig.
40  1). Those were collected from 3 female and 4 male, C57BL/6 NIA, 3 month old mice
41  (10-15 weeks), whose developmental age is roughly analogous to humans at 20 years of
42  age. All data, protocols, and analysis scripts from the *Tabula Muris* are shared as a public
43  resource (http://tabula-muris.ds.czbiohub.org/), gene counts and metadata from all single
44  cells are accessible on Figshare (https://figshare.com/account/home#/projects/27733),
45  raw data are available on GEO (GSE109774), and all code used for analysis is available
46  on GitHub (https://github.com/czbiohub/tabula-muris). While these data are by no means

47    a complete representation of all mouse organs and cell types, they provide a first draft
48    attempt to create an organism-wide representation of cellular diversity and a comparative
49    framework for future studies using the large variety of murine disease models.
50

51    We developed a procedure to collect 20 organs and tissues from the same mouse in which
52    aorta, bladder, bone marrow, brain (cerebellum, cortex, hippocampus, striatum),
53    diaphragm, fat (brown, gonadal, mesenteric, subcutaneous), heart, kidney, large intestine,
54    limb muscle, liver, lung, mammary gland, pancreas, skin, spleen, thymus, tongue, and
55    trachea were immediately dissected and processed into single cell suspensions, which in
56    turn were either single cell sorted into plates with FACS or loaded into microfluidic
57    droplets (see Extended Data and Methods). Single cell transcriptomes were sequenced to
58    an average depth of 814,488 reads per cell for the plate data and 7,709 unique molecular
59    identifiers (UMI) per cell for the microfluidic droplet data. After quality control filtering,
60    44,949 FACS sorted cells and 55,656 microfluidic droplet processed cells were retained
61    for further analysis. A comparison of the two methods shows differences for each organ
62    in the number of cells analyzed (Fig. 1b,c), reads per cell (Supp. Fig. 1a,c) and genes per
63    cell (Supp. Fig. 1b,d).
64

65    We performed unbiased graph-based clustering of the pooled set of transcriptomes across
66    all organs, and visualized them using tSNE (Fig. 2 and Supp. Fig. 2). The majority of
67    clusters contain cells from only one organ (n=29/54), but a number of clusters (n=25/54)
68    (Supp. Fig. 2) contained cells from multiple organs. To further dissect these clusters we
69    analyzed each organ independently, first by performing principal component analysis
70    (PCA) on the most variable genes in the organ, followed by nearest-neighbor graph-based
71    clustering. We then used cluster-specific gene expression of known markers as well as
72    genes differentially expressed between clusters to assign cell type annotations (Fig. 3,
73    Supp.Fig.3, TableS1). A detailed description of the cell types and defining genes for each
74    organ and tissue is available in the Supplementary Information. We used a standardized
75    analysis approach for all organs and tissues and an example using liver can be found in
76    the Organ Annotation Vignette. For each cell, we provide annotations in the controlled
77    vocabulary of a cell ontology[19] to facilitate comparisons with other experiments. Many of
78    these cell clusters have not previously been obtained in pure populations and our data
79    provide a wealth of new information on their characteristic gene expression profiles.
80    Initial annotation of the cellular diversity of each organ and tissue can be found in the
81    extended data, and a detailed discussion of each cell type on an organ by organ basis can
82    be found in the supplement. Some unexpected discoveries include a potential new role
83    for genes *Neurog3*, *Hex3*, and *Prss53* in the adult pancreas, a cell population expressing
84    *Chodl* in limb muscle, transcriptional heterogeneity of brain endothelial cells, the
85    expression of MHCII genes by adult mouse T cells, and sets of transcription factors that
86    can specifically distinguish between similar cell types across multiple organs and tissues.
87

88    Any individual single-cell sequencing experiment offers a partial view of the diversity of
89    cell types within an organism and the gene expression within each cell type. We illustrate
90    the variability to be expected between methods and experiments by comparing our two
91    measurement approaches to one another, and to data from Han *et al.*[20] generated using a
92    third method, microwell-seq. One striking feature is the variability in the number of

93   genes detected per cell between organs and tissues and between methods. For example,
94   the median number of genes detected per cell in bladder is about 4900 in the FACS data,
95   2900 in the droplet data, and 900 in the microwell-seq data, while the number detected in
96   kidney is about 1400 in the FACS data, 1900 in the droplet data, and 500 in the
97   microwell-seq data. The bladder, liver, lung, mammary gland, trachea, tongue, and spleen
98   all show nearly twice as many genes detected per cell in the FACS data as compared to
99   the microfluidic data, whereas heart and marrow show comparable numbers detected in
100  both methods (Supp. Fig. 4a). This difference does not appear to be due to sequencing
101  depth, as the microfluidic droplet libraries are nearly saturated (Supp. Fig. 4b) and deeper
102  sequencing of the FACS libraries could only increase the number of genes detected. In
103  every organ, there are fewer genes detected per cell in microwell-seq data than either
104  droplet or FACS data. In these comparisons, a gene is considered detected if a single read
105  maps to it, as that is the only standard for expression at which reads and UMIs can be
106  treated equally. We also looked at how the number of detected genes across each organ
107  changes with different thresholds on the number of reads or UMIs (Supp. Fig. 5). We
108  found that the number of detected genes decreases monotonically with increasing
109  thresholds at similar rates across different organs and tissues within each method. We
110  observed that in the droplet data more than half of the detected genes are represented by
111  only a single UMI; this is to be expected given that only a few thousand UMIs are
112  captured per cell. The FACS data are sampled much more deeply and one needs to set a
113  relatively high threshold of 40 reads to see a comparable reduction in gene detection
114  sensitivity.
115
116  Next, we investigated whether the three methods demonstrate concordance on the genes
117  which define each of the cell clusters. To do so, we computed lists of genes (see Methods
118  "Differential expression overlap analysis") that differentiate between each cell cluster and
119  the rest of the cell clusters in each organ across all three methods, focusing on common
120  organs and cell clusters for the three methods. As expected, data from FACS and
121  microfluidic droplet are in better agreement due to the fact that cells originated from the
122  exact same organ or tissue and were prepared in parallel. For each cell cluster there
123  appears to be a core of a few hundred defining genes on which all three methods agree
124  (Supp. Fig. 6 and Table S2). This comparison suggests that independent datasets
125  generated from the various tissue atlases that are beginning to arise can be combined and
126  collectively analyzed to generate more robust characterizations of gene expression.
127
128  To understand the relationships between cell types, we mapped the annotations of organ-
129  specific cell types onto the unbiased clustering of all cells. It is evident that the clusters in
130  Figure 2 (also Supp. Fig. 2) containing cells from multiple organs generally represent
131  shared cell types common to those organs (Fig. 4).  For example, B cells from fat, limb
132  muscle, diaphragm, lung, spleen and marrow cluster together, as do T cells from spleen,
133  marrow, lung, limb muscle, fat and thymus. Interestingly, while endothelial cells from
134  fat, heart, and lung cluster together, they are segregated from endothelial cells from the
135  mammary gland, kidney, trachea, limb muscle, aorta, diaphragm, and pancreas. Such
136  differences could be caused by true differential gene expression signatures across
137  different organs, but could also potentially be influenced by organ-specific batch effects.
138  The fact that many cells cluster together across organs and biological replicates is

139  evidence that batch effects are not the main source of variance in the dataset. Our
140  findings show that manual annotation of cell types is consistent with unbiased
141  transcriptomic clustering, and that most cell types are unique enough to enable their
142  unbiased identification across organs and tissues. We expect that further refinements of
143  comparison algorithms will facilitate the discovery of finer, organ-specific distinctions
144  between these shared cell types.
145
146  To investigate common cell types across all organs, we pooled all cells annotated as T
147  cells and analyzed them collectively (Fig. 5). Our analysis revealed 5 clusters. Cluster 0
148  comprises cells from the thymus that are undergoing VDJ recombination characterized by
149  the expression of RAG (*Rag1*, *Rag2*) and TdT (*Dntt*), and includes uncommitted double
150  positive T-cells (*Cd4*$^+$, *Cd8a*$^+$). Cluster 4 contains proliferating T cells, predominantly
151  from the thymus. We hypothesize that these are pre-T cells expanding after the
152  completion of VDJ recombination. Clusters 1-3 contain predominantly single positive T
153  cells (*Cd4*$^+$ or *Cd8a*$^+$). Cluster 3 contains *Cd5*$^{high}$ thymic T cells possibly undergoing
154  positive selection while Cluster 2 contains mostly non-thymic T cells expressing the high
155  affinity IL2 receptor (*Il2ra*, *Il2rb*), suggesting they are activated. Interestingly, they also
156  express MHC type II genes (*H2-Aa*, *H2-Ab1*). While this is known to occur in human T
157  cells, MHCII was previously thought restricted to professional antigen presenting cells in
158  mice[11]. Finally, Cluster 1 also represents mature T cells, but primarily from the spleen.
159
160  A key challenge for many single cell studies is understanding the potential changes to the
161  transcriptome caused by handling, dissociation and other experimental manipulation.  A
162  previous study in limb muscle showed that quiescent satellite cells tend to become
163  activated by dissociation and consequently express immediate early genes among other
164  genes[21]. We found that expression of these dissociation-related markers was also clearly
165  observed in our limb muscle data, as well as in mammary gland and bladder (Supp. Fig.
166  7), but that many organs and tissues showed little evidence of similar cellular activation.
167  Therefore the dissociation-related activation markers found in limb muscle are not
168  universal across all organs and tissues.  This is not to say that other organs lack
169  dissociation-related gene expression changes, but that some of the genes involved are
170  specific to a given organ.  Importantly, the presence of such gene expression changes
171  does not prevent the identification of cell type or the comparison of cell types across
172  organs and tissues.
173
174  One major goal of defining cell identities is to understand the transcription factor (TF)
175  regulatory networks that underlie them. We first investigated the combinatorial
176  specificity of TF expression across all cell types (defined as unique combinations of cell
177  ontology annotation and tissue)  (**Fig. 6**). We searched for the combination of four (n=4)
178  enriched TFs that best specified each target cell type over all others. For each
179  combination of TFs, we counted every cell expressing all four TFs as a positive, and
180  anything else as a negative. We then calculated cell type-specificity by the precision
181  (ratio of number of positive target cells to total number of positive cells) and recall (ratio
182  of number of positive target cells to total number of target cells) of each combination of
183  TFs for the target cell type over the rest of the cells (**Table S3**). We found 41 cell types
184  with TF combinations with precision > 0.3 and recall > 0.3. We noted that the

185  combinatorial nature of TF expression was critical to specificity; for example, *Ctnnb1,*
186  combined with one of two TF sets, specified either skin keratinocyte stem cells or lung
187  type II pneumocytes (**Fig. 6a**). We found many TF combinations for cell types with
188  challenging *in vitro* differentiation protocols[22] (e.g., hepatocytes; *Creb3l3*, *Nr1h3*, *Hnf4a,*
189  and *Klf15*) and cell types with no established direct differentiation protocol (e.g.,
190  microglia; *Mafb, Sall1, Irf5,* and *Maf*) (**Fig. 6a**).
191
192  We then analyzed organ-specific TFs by isolating a set of closely-related, cross-organ
193  cell groups (epithelial cells and endothelial cells). We performed TF correlation analysis,
194  similar to [15] within the cell groups (**Fig. 6b-g**). We found many TFs within epithelial cells
195  that clustered strongly by organ and were enriched in organ-specific epithelial clusters
196  (**Fig. 6b**). For example, *Sox4* (mammary basal cells), *Foxq1* (bladder basal cells of the
197  urothelium), *Pax9* (tongue basal cells of the epidermis), and *Lhx2* (skin keratinocyte stem
198  cells) were highly organ-specific (**Fig. 6c,d**). Within endothelial cells, liver, brain,
199  mammary gland/limb muscle, and lung-specific clusters of TFs were evident (**Fig. 6e-g**).
200  *Gata4,* known to specify liver endothelium, appeared in a cluster of liver-enriched TFs
201  (**Fig. 6g**). Another cluster of TFs, including *Pbx1*, were enriched in kidney endothelial
202  cells (**Fig. 6g**). The roles of *Pbx1* in kidney endothelial development are not explored,
203  and could aid in tissue engineering for kidney regeneration. A highly distinct cluster of
204  cells specified the heart endocardium, including *Plagl1*, a TF whose role in endocardial
205  specification is unknown (**Fig. 6g**). These results illustrate how single cell data taken
206  across many organs and organs can identify the transcriptional regulatory programs
207  which are specific to cell types of interest.
208
209  In conclusion, we have created a compendium of single-cell transcriptional
210  measurements across 20 organs and tissues of the mouse. This *Tabula Muris*, or "Mouse
211  Atlas", has many uses, including the discovery of new putative cell types, the discovery
212  of novel gene expression in known cell types, and the ability to compare cell types across
213  organs and tissues. It will also serve as a reference of healthy young adult organs and
214  tissues which can be used as a baseline for current and future mouse models of disease.
215  While it is not an exhaustive characterization of all organs of the mouse, it does provide a
216  rich data set of the most highly studied organs and tissues in biology. The *Tabula Muris*
217  provides a framework and description of many of the most populous and important cell
218  populations within the mouse, and represents a foundation for future studies across a
219  multitude of diverse physiological disciplines.
220
221  **Supplementary Information** is available in the online version of the paper.
222

230

**The *Tabula Muris* Consortium:**

**Overall Coordination:** Nicholas Schaum[1], Jim Karkanias[2], Norma F Neff[2], Andrew P. May[2], Stephen R. Quake[2,3]*, Tony Wyss-Coray[4-6]*, and Spyros Darmanis[2]*

* Correspondence to: quake@stanford.edu, twc@stanford.edu, spyros.darmanis@czbiohub.org

**Logistic Coordination:** Joshua Batson[2], Olga Botvinnik[2], Michelle B. Chen[3], Steven Chen[2], Foad Green[2], Robert Jones[3], Ashley Maynard[2], Lolita Penland[2], Rene V. Sit[2], Geoffrey M. Stanley[3] , James T. Webber[2], Fabio Zanini[3]

**Organ and Tissue collection and processing:** Ankit S. Baghel[1], Isaac Bakerman[1,7,8], Ishita Bansal[2], Daniela Berdnik[4], Biter Bilen[4], Douglas Brownfield[9], Corey Cain[10], Michelle B. Chen[3], Steven Chen[2], Min Cho[2], Giana Cirolia[2], Stephanie D. Conley[1], Spyros Darmanis[2], Aaron Demers[2], Kubilay Demir[1,11], Antoine de Morree[4], Tessa Divita[2], Haley du Bois[4], Laughing Bear Torrez Dulgeroff[1], Hamid Ebadi[2], F. Hernán Espinoza[9], Matt Fish[1,11,12], Qiang Gan[4], Benson M. George[1], Astrid Gillich[9], Foad Green[2], Geraldine Genetiano[2], Xueying Gu[12], Gunsagar S. Gulati[1], Yan Hang[12], Shayan Hosseinzadeh[2], Albin Huang[44], Tal Iram[4], Taichi Isobe[1], Feather Ives[2], Robert Jones[3], Kevin S. Kao[1], Guruswamy Karnam[13], Aaron M. Kershner[1], Bernhard Kiss[1,14], William Kong[1], Maya E. Kumar[15,16], Jonathan Lam[12], Davis P. Lee[6], Song E. Lee[4], Guang Li[17], Qingyun Li[18], Ling Liu[4], Annie Lo[2], Wan-Jin Lu[1,9], Anoop Manjunath[1], Andrew P. May[2], Kaia L. May[2], Oliver L. May[2], Ashley Maynard[2], Marina McKay[2], Ross J. Metzger[19,20], Marco Mignardi[3], Dullei Min[21], Ahmad N. Nabhan[9], Norma F Neff[2], Katharine M. Ng[3], Joseph Noh[1], Rasika Patkar[13], Weng Chuan Peng[12], Lolita Penland[2], Robert Puccinelli[2], Eric J. Rulifson[12], Nicholas Schaum[1], Shaheen S. Sikandar[1], Rahul Sinha[1,22-24], Rene V Sit[2], Krzysztof Szade[1,25], Weilun Tan[2], Cristina Tato[2], Krissie Tellez[12], Kyle J. Travaglini[9], Carolina Tropini[26], Lucas Waldburger[2], Linda J. van Weele[1], Michael N. Wosczyna[4], Jinyi Xiang[1], Soso Xue[3], Justin Youngyunpipatkul[2], Fabio Zanini[3], Macy E. Zardeneta[6], Fan Zhang[19,20], Lu Zhou[18]

**Library preparation and sequencing:** Ishita Bansal[2], Steven Chen[2], Min Cho[2], Giana Cirolia[2], Spyros Darmanis[2], Aaron Demers[2], Tessa Divita[2], Hamid Ebadi[2], Geraldine Genetiano[2], Foad Green[2], Shayan Hosseinzadeh[2], Feather Ives[2], Annie Lo[2], Andrew P. May[2], Ashley Maynard[2], Marina McKay[2], Norma F. Neff[2], Lolita Penland[2], Rene V. Sit[2], Weilun Tan[2], Lucas Waldburger[2], Justin Youngyunpipatkul[2]

**Computational Data Analysis:** Joshua Batson[2], Olga Botvinnik[2], Paola Castro[2], Derek Croote[3], Spyros Darmanis[2], Joseph L. DeRisi[2,27], Jim Karkanias[2], Angela Pisco[2], Geoffrey M. Stanley[3], James T. Webber[2], Fabio Zanini[3]

**Cell Type Annotation:** Ankit S. Baghel[1], Isaac Bakerman[1,7,8], Joshua Batson[2], Biter Bilen[4], Olga Botvinnik[2], Douglas Brownfield[9], Michelle B. Chen[3], Spyros Darmanis[2], Kubilay Demir[1,11], Antoine de Morree[4], Hamid Ebadi[2], F. Hernán Espinoza[9], Matt Fish[9,11,12], Qiang Gan[4], Benson M. George[1], Astrid Gillich[9], Xueying Gu[12], Gunsagar S.

277 Gulati[1], Yan Hang[12], Albin Huang[4], Tal Iram[4], Taichi Isobe[1], Guruswamy Karnam[13],
278 Aaron M. Kershner[1], Bernhard M. Kiss[1,14], William Kong[1], Christin S. Kuo[9,11,21], Jonathan
279 Lam[12], Benoit Lehallier[4], Guang Li[17], Qingyun Li[18], Ling Liu[4], Wan-Jin Lu[1,9], Dullei
280 Min[21], Ahmad N. Nabhan[9], Katharine M. Ng[3], Patricia K. Nguyen[1,7,8,17], Rasika Patkar[13],
281 Weng Chuan Peng[12], Lolita Penland[2], Eric J. Rulifson[12], Nicholas Schaum[1], Shaheen S.
282 Sikandar[1], Rahul Sinha[1,22-24], Krzysztof Szade[1,25], Serena Y. Tan[22], Krissie Tellez[12], Kyle
283 J. Travaglini[9], Carolina Tropini[26], Linda J. van Weele[1], Bruce M. Wang[13], Michael N.
284 Wosczyna[4], Jinyi Xiang[1], Hanadie Yousef[4], Lu Zhou[18]
285

286 **Writing Group:** Joshua Batson[2], Olga Botvinnik[2], Steven Chen[2], Spyros Darmanis[2],
287 Foad Green[2], Andrew P. May[2], Ashley Maynard[2], Angela Pisco[2], Stephen R. Quake[2,3],
288 Nicholas Schaum[1], Geoffrey M. Stanley[3], James T. Webber[2], Tony Wyss-Coray[4-6], Fabio
289 Zanini[3]
290

291 **Supplemental Text Writing Group:** Philip A. Beachy[1,9,11,12], Charles K. F. Chan[28],
292 Antoine de Morree[4], Benson M. George[1], Gunsagar S. Gulati[1], Yan Hang[12], Kerwyn
293 Casey Huang[2,3,26], Tal Iram[4], Taichi Isobe[1], Aaron M. Kershner[1], Bernhard M. Kiss[1,14],
294 William Kong[1], Guang Li[17], Qingyun Li[18], Ling Liu[4], Wan-Jin Lu[1,9], Ahmad N. Nabhan[9],
295 Katharine M. Ng[3], Patricia K. Nguyen[1,7,8,17], Nicholas Schaum[1], Shaheen S. Sikandar[1],
296 Rahul Sinha[1,22-24], Krzysztof Szade[1,25], Kyle J. Travaglini[9], Carolina Tropini[26], Bruce M.
297 Wang[13], Kenneth Weinberg[21], Michael N. Wosczyna[4], Sean Wu[17], Hanadie Yousef[4]
298

299 **Principal Investigators:** Ben A. Barres[18], Philip A. Beachy[1,9,11,12], Charles K. F. Chan[28],
300 Michael F. Clarke[1], Spyros Darmanis[2], Kerwyn Casey Huang[2,3,26], Jim Karkanias[2], Seung
301 K. Kim[12,29], Mark A. Krasnow[9,11], Christin S. Kuo[9,11,21], Andrew P. May[2], Norma Neff[2],
302 Roel Nusse[9,11,12], Patricia K. Nguyen[1,7,8,17], Thomas A. Rando[4-6], Justin Sonnenburg[26],
303 Bruce M. Wang[13], Kenneth Weinberg[21], Irving L. Weissman[1,22-24], Sean M. Wu[1,7,17],
304 Stephen R. Quake[2,3], Tony Wyss-Coray[4,5,6]
305

306 [1] Institute for Stem Cell Biology and Regenerative Medicine, Stanford University School
307 of Medicine, Stanford, California, USA
308 [2] Chan Zuckerburg Biohub, San Francisco, California, USA
309 [3] Department of Bioengineering, Stanford University, Stanford, California, USA
310 [4] Department of Neurology and Neurological Sciences, Stanford University School of
311 Medicine, Stanford, California, USA
312 [5] Paul F. Glenn Center for the Biology of Aging, Stanford University School of
313 Medicine, Stanford, California, USA
314 [6] Center for Tissue Regeneration, Repair, and Restoration, V.A. Palo Alto Healthcare
315 System, Palo Alto, California, USA
316 [7] Stanford Cardiovascular Institute, Stanford University School of Medicine, Stanford,
317 California, USA
318 [8] Department of Medicine, Division of Cardiology, Stanford University School of
319 Medicine, Stanford, California, USA
320 [9] Department of Biochemistry, Stanford University School of Medicine, Stanford,
321 California, USA
322 [10] Flow Cytometry Core, V.A. Palo Alto Healthcare System, Palo Alto, California, USA

323    [11] Howard Hughes Medical Institute, USA

324    [12] Department of Developmental Biology, Stanford University School of Medicine,
325    Stanford, California, USA

326    [13] Department of Medicine and Liver Center, University of California San Francisco, San
327    Francisco, California, USA

328    [14] Department of Urology, Stanford University School of Medicine, Stanford, California,
329    USA

330    [15] Sean N. Parker Center for Asthma and Allergy Research, Stanford University School
331    of Medicine, Stanford, California, USA

332    [16] Department of Medicine, Division of Pulmonary and Critical Care, Stanford University
333    School of Medicine, Stanford, California, USA

334    [17] Department of Medicine, Division of Cardiovascular Medicine, Stanford University,
335    Stanford, California, USA

336    [18] Department of Neurobiology, Stanford University School of Medicine, Stanford, CA
337    USA

338    [19] Vera Moulton Wall Center for Pulmonary and Vascular Disease, Stanford University
339    School of Medicine, Stanford, California, USA

340    [20] Department of Pediatrics, Division of Cardiology, Stanford University School of
341    Medicine, Stanford, California, USA

342    [21] Department of Pediatrics, Stanford University school of Medicine, Stanford,
343    California, USA

344    [22] Department of Pathology, Stanford University School of Medicine, Stanford,
345    California, USA

346    [23] Ludwig Center for Cancer Stem Cell Research and Medicine, Stanford University
347    School of Medicine, Stanford, California, USA

348    [24] Stanford Cancer Institute, Stanford University School of Medicine, Stanford,
349    California, USA

350    [25] Department of Medical Biotechnology, Faculty of Biophysics, Biochemistry and
351    Biotechnology, Jagiellonian University, Poland

352    [26] Department of Microbiology & Immunology, Stanford University School of Medicine,
353    Stanford, California, USA

354    [27] Department of Biochemistry and Biophysics, University of California San Francisco,
355    San Francisco, California USA

356    [28] Department of Surgery, Division of Plastic and Reconstructive Surgery, Stanford
357    University, Stanford, California USA

358    [29] Department of Medicine and Stanford Diabetes Research Center, Stanford University,
359    Stanford, California USA

360

**References**

1. Alberts, B. *et al.* *Essential Cell Biology*. (Garland Pub, 2014).

2. Guo, G. *et al.* Resolution of cell fate decisions revealed by single-cell gene expression analysis from zygote to blastocyst. *Dev. Cell* **18,** 675–685 (2010).

3. Dalerba, P. *et al.* Single-cell dissection of transcriptional heterogeneity in human colon tumors. *Nat. Biotechnol.* **29,** 1120–1127 (2011).

4. Thorsen, T., Roberts, R. W., Arnold, F. H. & Quake, S. R. Dynamic pattern formation in a vesicle-generating microfluidic device. *Phys. Rev. Lett.* **86,** 4163–4166 (2001).

5. Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161,** 1202–1214 (2015).

6. Klein, A. M. *et al.* Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell* **161,** 1187–1201 (2015).

7. Ramsköld, D. *et al.* Full-length mRNA-Seq from single-cell levels of RNA and individual circulating tumor cells. *Nat. Biotechnol.* **30,** 777–782 (2012).

8. Wu, A. R. *et al.* Quantitative assessment of single-cell RNA-sequencing methods. *Nat. Methods* **11,** 41–46 (2014).

9. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509,** 371–375 (2014).

10. Enge, M. *et al.* Single-Cell Analysis of Human Pancreas Reveals Transcriptional Signatures of Aging and Somatic Mutation Patterns. *Cell* **171,** 321–330.e14 (2017).

11. Halpern, K. B. *et al.* Single-cell spatial reconstruction reveals global division of labour in the mammalian liver. *Nature* **542,** 352–356 (2017).

12. Haber, A. L. *et al.* A single-cell survey of the small intestinal epithelium. *Nature* **551,** 333–339 (2017).

13. Villani, A.-C. *et al.* Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors. *Science* **356,** eaah4573 (2017).

14. Darmanis, S. *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proc. Natl. Acad. Sci. U.S.A.* **112,** 7285–7290 (2015).

15. Gokce, O. *et al.* Cellular Taxonomy of the Mouse Striatum as Revealed by Single-Cell RNA-Seq. *Cell Rep* **16,** 1126–1137 (2016).

16. Usoskin, D. *et al.* Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing. *Nat. Neurosci.* **18,** 145–153 (2015).

17. Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347,** 1138–1142 (2015).

18. Li, H. *et al.* Classifying Drosophila Olfactory Projection Neuron Subtypes by Single-Cell RNA Sequencing. *Cell* **171,** 1206–1220.e22 (2017).

19. Smith, B. *et al.* The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nat. Biotechnol.* **25,** 1251–1255 (2007).

20. Han, X. *et al.* Mapping the Mouse Cell Atlas by Microwell-Seq. *Cell* **172,** 1091–1107.e17 (2018).

21. Holling, T. M., Schooten, E. & van Den Elsen, P. J. Function and regulation of MHC class II molecules in T-lymphocytes: of mice and men. *Hum. Immunol.* **65,** 282–290 (2004).

415  22.    Soldatow, V. Y., Lecluyse, E. L., Griffith, L. G. & Rusyn, I. In vitro
416         models for liver toxicity testing. *Toxicol Res (Camb)* **2,** 23–39 (2013).
417  23.    Reichardt, J. & Bornholdt, S. Statistical mechanics of community detection.
418         *Phys Rev E Stat Nonlin Soft Matter Phys* **74,** 016110 (2006).
419  24.    van den Brink, S. C. *et al.* Single-cell sequencing reveals dissociation-
420         induced gene expression in tissue subpopulations. *Nat. Methods* **14,** 935–936
421         (2017).
422

423    **Methods**
424
425    **Mice and Tissue Collection**
426    Four 10-15 week old male and four virgin female C57BL/6 mice were shipped from the
427    National Institute on Aging colony at Charles River to the Veterinary Medical Unit
428    (VMU) at the VA Palo Alto (VA). At both locations, mice were housed on a 12-h
429    light/dark cycle, and provided food and water *ad libitum*. The diet at Charles River was
430    NIH-31, and Teklad 2918 at the VA VMU. Littermates were not recorded or tracked, and
431    mice were housed at the VA VMU for no longer than 2 weeks before euthanasia. Prior to
432    tissue collection, mice were placed in sterile collection chambers for 15 minutes to collect
433    fresh fecal pellets. Following anesthetization with 2.5% v/v Avertin, mice were weighed,
434    shaved, and blood drawn via cardiac puncture before transcardial perfusion with 20 ml
435    PBS. Mesenteric adipose tissue (MAT) was then immediately collected to avoid exposure
436    to the liver and pancreas perfusate, which negatively impacts cell sorting. Isolating viable
437    single cells from both pancreas and liver of the same mouse was not possible, therefore, 2
438    males and 2 females were used for each. Whole organs were then dissected in the
439    following order: large intestine, spleen, thymus, trachea, tongue, brain, heart, lung,
440    kidney, gonadal adipose tissue (GAT), bladder, diaphragm, limb muscle (*tibialis*
441    *anterior*), skin (dorsal), subcutaneous adipose tissue (SCAT, inguinal pad), mammary
442    glands (fat pads 2, 3, and 4), brown adipose tissue (BAT, interscapular pad), aorta, and
443    bone marrow (spine and limb bones). Following single cell dissociation as described
444    below, cell suspensions were either used for FACS sorting of individual cells into 384-
445    well plates, or for microfluidic droplet library preparation. All animal care and
446    procedures were carried out in accordance with institutional guidelines approved by the
447    VA Palo Alto Committee on Animal Research.
448
449    **Tissue dissociation and sample preparation**
450    Specific protocols for each tissue are described in the supplement.
451
452    **Single Cell Methods**
453
454    **Lysis plate preparation**
455    Lysis plates were created by dispensing 0.4 μl lysis buffer (0.5 U Recombinant RNase
456    Inhibitor (Takara Bio, 2313B), 0.0625% Triton$^{TM}$ X-100 (Sigma, 93443-100ML), 3.125
457    mM dNTP mix (Thermo Fisher, R0193), 3.125 μM Oligo-dT$_{30}$VN (IDT,
458    5'AAGCAGTGGTATCAACGCAGAGTACT$_{30}$VN-3') and 1:600,000 ERCC RNA
459    spike-in mix (Thermo Fisher, 4456740)) into 384-well hard-shell PCR plates (Biorad
460    HSP3901) using a Tempest liquid handler (Formulatrix). 96-well lysis plates were also
461    prepared with 4 μl lysis buffer. All plates were sealed with AlumaSeal CS Films (Sigma-
462    Aldrich Z722634) and spun down (3,220 x g, 1 minute) and snap frozen on dry ice. Plates
463    were stored at -80°C until sorting.
464
465    **FACS sorting**
466    After dissociation, single cells from each organ and tissue were isolated into 384- or 96-
467    well plates via Fluorescence Activated Cell Sorting (FACS). Most organs were sorted
468    into 384-well plates using SH800S (Sony) sorters. Heart and liver were sorted into 96-

469   well plates and cardiomyocytes were hand-picked into 96-well plates. Limb muscle and
470   diaphragm were sorted into 384-well plates on an Aria III (Becton Dickinson) sorter. The
471   last two columns of each 384 well plate were intentionally left as blanks.  For most
472   organs, single cells were selected with forward scatter, and dead cells and common cell
473   types were excluded with a single color channel. Combinations of fluorescent antibodies
474   were used for most organs to enrich for rare cell populations (see supplemental text), but
475   some were stained only for viable cells. Color compensation was used whenever
476   necessary. On the SH800, the highest purity setting ("Single cell") was used for all but
477   the rarest cell types, for which the "Ultrapure" setting was used. Sorters were calibrated
478   using FACS buffer every day before collecting any cells, and also after every 8 sorted
479   plates. For a typical sort, 1-3 ml of pre-stained cell suspension was filtered, vortexed
480   gently, and loaded onto the FACS machine. A small number of cells were flowed at low
481   pressure to check cell and debris concentrations. The pressure was then adjusted, flow
482   paused, the first destination plate unsealed, loaded and sorting started. If a cell suspension
483   was too concentrated, it was diluted using FACS buffer or 1X PBS. For some cell types
484   like hepatocytes, 96-well plates were used because it was not possible to sort individual
485   cells accurately into 384-well plates. Immediately after sorting, plates were sealed with a
486   pre-labeled aluminum seal, centrifuged, and flash frozen on dry ice. On average, each
487   384-well plate took 8 minutes to sort.
488
489   **cDNA synthesis and library preparation**
490   cDNA synthesis was performed using the Smart-seq2 protocol[2,3]. Briefly, 384-well plates
491   containing single-cell lysates were thawed on ice followed by first strand synthesis. 0.6 μl
492   of reaction mix (16.7 U/μl SMARTScribe Reverse Transcriptase (Takara Bio, 639538),
493   1.67 U/μl Recombinant RNase Inhibitor (Takara Bio, 2313B), 1.67X First-Strand Buffer
494   (Takara      Bio,      639538),      1.67      μM      TSO      (Exiqon,      5'-
495   AAGCAGTGGTATCAACGCAGAGTGAATrGrGrG-3'), 8.33 mM DTT (Bioworld,
496   40420001-1), 1.67 M Betaine (Sigma, B0300-5VL), and 10 mM $MgCl_2$ (Sigma, M1028-
497   10X1ML)) was added to each well using a Tempest liquid handler. Reverse transcription
498   was carried out by incubating wells on a ProFlex 2 x 384 thermal-cycler (Thermo Fisher)
499   at 42°C for 90 minutes, and stopped by heating at 70°C for 5 minutes.
500
501   Subsequently, 1.5 μl of PCR mix (1.67X KAPA HiFi HotStart ReadyMix (Kapa
502   Biosystems,      KK2602),      0.17      μM      IS      PCR      primer      (IDT,      5'-
503   AAGCAGTGGTATCAACGCAGAGT-3'), and 0.038 U/μl Lambda Exonuclease (NEB,
504   M0262L)) was added to each well with a Mantis liquid handler (Formulatrix), and second
505   strand synthesis was performed on a ProFlex 2x384 thermal-cycler by using the
506   following program: 1) 37°C for 30 minutes, 2) 95°C for 3 minutes, 3) 23 cycles of 98°C
507   for 20 seconds, 67°C for 15 seconds, and 72°C for 4 minutes, and 4) 72°C for 5 minutes.
508
509   The amplified product was diluted with a ratio of 1 part cDNA to 10 parts 10mM Tris-
510   HCl (Thermo Fisher, 15568025), and concentrations were measured with a dye-
511   fluorescence assay (Quant-iT dsDNA High Sensitivity kit; Thermo Fisher, Q33120) on a
512   SpectraMax i3x microplate reader (Molecular Devices). Sample plates were selected for
513   downstream processing if the mean concentration of blanks (ERCC-containing, non-cell
514   wells) was greater than 0 ng/μl, and, after linear regression of the values obtained from

515  the Quant-iT dsDNA standard curve, the $R^2$ value was greater than 0.98. Sample wells
516  were then selected if their cDNA concentrations were at least one standard deviation
517  greater than the mean concentration of the blanks. These wells were reformatted to a new
518  384-well plate at a concentration of 0.3 ng/μl and final volume of 0.4 μl using an Echo
519  550 acoustic liquid dispenser (Labcyte).
520
521  Illumina sequencing libraries were prepared as described in Darmanis et al. 2015.[4]
522  Briefly, tagmentation was carried out on double-stranded cDNA using the Nextera XT
523  Library Sample Preparation kit (Illumina, FC-131-1096). Each well was mixed with 0.8
524  μl Nextera tagmentation DNA buffer (Illumina) and 0.4 μl Tn5 enzyme (Illumina), then
525  incubated at 55°C for 10 minutes. The reaction was stopped by adding 0.4 μl "Neutralize
526  Tagment Buffer" (Illumina) and centrifuging at room temperature at 3,220 x g for 5
527  minutes. Indexing PCR reactions were performed by adding 0.4 μl of 5 μM i5 indexing
528  primer, 0.4 μl of 5 μM i7 indexing primer, and 1.2 μl of Nextera NPM mix (Illumina).
529  PCR amplification was carried out on a ProFlex 2x384 thermal cycler using the following
530  program: 1) 72°C for 3 minutes, 2) 95°C for 30 seconds, 3) 12 cycles of 95°C for 10
531  seconds, 55°C for 30 seconds, and 72°C for 1 minute, and 4) 72°C for 5 minutes.
532
533  **Library pooling, quality control, and sequencing**
534  Following library preparation, wells of each library plate were pooled using a
535  Mosquito liquid handler (TTP Labtech). Pooling was followed by two purifications using
536  0.7x AMPure beads (Fisher, A63881). Library quality was assessed using capillary
537  electrophoresis on a Fragment Analyzer (AATI), and libraries were quantified by qPCR
538  (Kapa Biosystems, KK4923) on a CFX96 Touch Real-Time PCR Detection System
539  (Biorad). Plate pools were normalized to 2 nM and equal volumes from 10 or 20 plates
540  were mixed together to make the sequencing sample pool. A PhiX control library was
541  spiked in at 0.2% before sequencing.
542
543  **Sequencing libraries from 384-well and 96-well plates**
544  Libraries were sequenced on the NovaSeq 6000 Sequencing System (Illumina) using 2 x
545  100bp paired-end reads and 2 x 8bp or 2 x 12bp index reads with either a 200- or 300-
546  cycle kit (Illumina, 20012861 or 20012860).
547
548  **Microfluidic droplet single cell analysis**
549  Single cells were captured in droplet emulsions using the GemCode Single-Cell
550  Instrument (10x Genomics, Pleasanton, CA, USA), and SC RNA-seq libraries were
551  constructed as per the 10X Genomics protocol using GemCode Single-Cell 3′ Gel Bead
552  and Library V2 Kit. Briefly, single cell suspensions were examined using an inverted
553  microscope, and if sample quality was deemed satisfactory, the sample was diluted in
554  PBS with 2% FBS to a concentration of 1000 cells/μl.  If cell suspensions contained cell
555  aggregates or debris, two additional washes in PBS with 2% FBS at 300 x g for 5 minutes
556  at 4°C were performed. Cell concentration was measured either with a Moxi GO II (Orflo
557  Technologies) or a hemocytometer. Cells were loaded in each channel with a target
558  output of 5,000 cells per sample. All reactions were performed in the Biorad C1000
559  Touch Thermal cycler with 96-Deep Well Reaction Module. 12 cycles were used for
560  cDNA amplification and sample index PCR. Amplified cDNA and final libraries were

561 evaluated on a Fragment Analyzer using a High Sensitivity NGS Analysis Kit (Advanced
562 Analytical). The average fragment length of 10x cDNA libraries was quantitated on a
563 Fragment Analyzer (AATI), and by qPCR with the Kapa Library Quantification kit for
564 Illumina. Each library was diluted to 2 nM, and equal volumes of 16 libraries were
565 pooled for each NovaSeq sequencing run. Pools were sequenced with 100 cycle run kits
566 with 26 bases for Read 1, 8 bases for Index 1, and 90 bases for Read 2 (Illumina
567 20012862). A PhiX control library was spiked in at 0.2 to 1%. Libraries were sequenced
568 on the NovaSeq 6000 Sequencing System (Illumina)
569
570 **Data Processing**
571 Sequences from the Novaseq were de-multiplexed using bcl2fastq version 2.19.0.316.
572 Reads were aligned using to the mm10plus genome using STAR version 2.5.2b with
573 parameters TK. Gene counts were produced using HTSEQ version 0.6.1p1 with default
574 parameters, except "stranded" was set to "false", and "mode" was set to "intersection-
575 nonempty".
576
577 Sequences from the microfluidic droplet platform were de-multiplexed and aligned using
578 CellRanger, available from 10x Genomics with default parameters.
579
580 **Clustering**
581 Standard procedures for filtering, variable gene selection, dimensionality reduction, and
582 clustering were performed using the Seurat package. A detailed worked example,
583 including the mathematical formulae for each operation, is in the Tissue Annotation
584 Vignette. The parameters that were tuned on a per-tissue basis (resolution and number of
585 PCs can be viewed in the tissue-specific Rmd files available on GitHub). For each tissue
586 and each sequencing method (FACS and microfluidic droplet), the following steps were
587 performed:
588
589 1. Cells were lexicographically sorted by cell ID to ensure reproducibility.
590 2. Cells with fewer than 500 detected genes were excluded. (A gene counts as
591 detected if it has at least one read mapping to it). Cells with fewer than 50,000
592 reads (FACS) or 1000 UMI (microfluidic droplet) were excluded.
593 3. Counts were log-normalized for each cell using the natural logarithm of 1 +
594 counts per million (for FACS) or 1 + counts per ten thousand (for microfluidic
595 droplet).
596 4. Variable genes were selected using a threshold (0.5) for the standardized log
597 dispersion, where the standardization was done in separately according to binned
598 values of log mean expression.
599 5. The variable genes were projected onto a low-dimensional subspace using
600 principal component analysis. The number of principal components was selected
601 based on inspection of the plot of variance explained.
602 6. A shared-nearest-neighbors graph was constructed based on the Euclidean
603 distance in the low-dimensional subspace spanned by the top principal
604 components. Cells were clustered using a variant of the Louvain method that
605 includes a resolution parameter in the modularity function[23].

606  7. Cells were visualized using a 2-dimensional t-distributed Stochastic Neighbor
607     Embedding of the PC-projected data.
608  8. Cell types were assigned to each cluster using the abundance of known marker
609     genes. Plots showing the expression of the markers for each tissue appear in the
610     extended data.
611  9. When clusters appeared to be mixtures of cell types, they were refined either by
612     increasing the resolution parameter for clustering or subsetting the data and
613     rerunning steps 3-7.
614
615 A similar analysis was done globally for all FACS processed cells and for all microfluidic
616 droplet processed cells to produce an unbiased clustering.
617
618 **Differential expression overlap analysis**
619
620 For FACS and microfluidic droplet data differential expression analysis for each organ
621 was performed using a Wilcox rank test as implemented in the "FindAllMarkers"
622 function of the Seurat package. Differential expression was performed between cell
623 ontology groups and resulted in a list of differentially expressed genes ($\log_e$FoldChange >
624 0.25) between each cell ontology group and all other ontology groups of the same organ.
625 For the microwellSeq we used the corresponding published lists for each cell type and for
626 every organ. We then assessed the overlap (Supp. Fig. 6) of those lists between the three
627 methods. As the nomenclature is not identical, the analysis was performed between cell
628 types that could be matched with a certain degree of confidence between the three
629 methods (TableS2).
630
631 **Calculation of dissociation scores**
632
633 For each organ, gene expression matrices were subset to 140 genes[24], and principal
634 component analysis was performed on this gene subset. The first principal component
635 was used as the "dissociation score" as it corresponds to the variance within these genes.
636
637 **Defining cell type-enriched transcription factors**
638
639 Transcription factors were defined as the 1140 genes annotated by the Gene Ontology
640 term "DNA binding transcription factor activity", downloading from the Mouse Genome
641 Informatics    database    (http://www.informatics.jax.org/mgihome/GO/project.shtml,
642 accessed on 2017-11-10).  Cell types were defined as unique combinations of cell
643 ontology and organ annotation (e.g. Lung__Endothelial_cell). All analysis was performed
644 on the full 3 month dataset, subsampled by randomly selecting 60 cells from each cell
645 type. Enriched TFs were defined by the Seurat FindMarkers function with the
646 "Wilcoxon" significance test for the target cell type against the all of rest of the cell types
647 combined. These were filtered by p_val < 10-3, avg_diff > 0.2, pct.1 – pct.2 > 0.1
648 (percent detected difference > 0.1), and pct.1 > 0.3 (detected in > 30% of target cells).
649
650 **Discovering cell type-specific TF combinations**
651

652  For each cell type that contained at least 6 cells, and had at least 4 enriched TFs, the top
653  30 TFs or all that passed filter, whichever was smaller, were selected by highest avg_diff.
654  The specificity of each four-TF combination (up to 27405 combinations for 30 TFs) was
655  assessed by a score defined from two standard metrics, precision and recall:

$$\text{Precision} = \frac{TP}{TP + FP}$$
$$\text{Recall} = \frac{TP}{TP + FN}$$
$$\text{Score} = 2 * \text{Precision} + \text{Recall}$$

656
657  Where TP (true positive) is the number of cells in the target cell type expressing all 4
658  TFs, FP (false positive) is the number of cells not in the target cell type expressing all 4
659  TFs, and TN (true negative) is the number of cells in the target cell type not expressing
660  all 4 TFs. The top TFs by this score for several cell types was plotted in Figure 6a.
661
662  **Defining TF networks by correlation analysis**
663
664  Organ-specific TF regulatory networks were measured by the correlations of TFs. TFs
665  were selected by enrichment in a cell type over all other cell type with the test described
666  in "Defining cell type-enriched transcription factors", filtered by p_val < $10^{-8}$, avg_diff >
667  0.3, and pct.1-pct.2 > 0.1. The top 8 markers per cell type (or however many passed the
668  filters) were selected by avg_diff. The Pearson correlations between genes were
669  calculated, and genes ordered by hierarchical clustering with optimal ordering (hclust and
670  cba::optimal). For analysis of TFs within single broad cross-organ cell types, endothelial
671  cells were defined as cell ontology annotations containing "endothelial" or "capillary"
672  (Fig. 6e-g). Epithelial cells were defined as cell ontology annotations containing
673  "epithelial", "basal", "keratinocyte", or "epidermis" (Fig. 6b-d). Exemplary organ-
674  specific TFs were visualized on t-SNE plots. t-SNE was computed for a single cell
675  annotation across all organs, by the top variable genes (Seurat FindVariableGenes,
676  RunPCA with 10 PCs, and RunTSNE with perplexity = 30).
677

678
679 **Figure captions**
680
681 **Figure 1.** Overview of *Tabula Muris*
682 a) 20 organs and tissues from 4 male and 3 female mice were analyzed. After
683 dissociation, cells were either sorted by FACS or captured in microfluidic oil droplets,
684 after which they were lysed and their transcriptomes amplified, sequenced, and reads
685 mapped, followed by data analysis. b) Barplot showing number of sequenced cells
686 prepared by FACS sorting from each organ (n=20). c) Barplot showing number of
687 sequenced cells prepared by microfluidic droplets from each organ (n=12).
688
689 **Figure 2**. tSNE visualization of all FACS sorted cells.
690 tSNE plot of all cells sorted by FACS, color coded by organ.
691
692 **Figure 3**. tSNE visualization of individual organs.
693 a) tSNE plots for each organ of cells sorted by FACS. Color coding indicates distinct
694 clusters. b) Barplots of annotated cell types based on differential gene expression across
695 all organs. Coloring of clusters within each organ is consistent between panels a and b.
696
697 **Figure 4**. Comparison of cell type determination.
698 Comparison of cell type determination as done by unbiased whole transcriptome
699 comparison versus manual annotation by organ-specific experts. The x-axis represents
700 clusters from Figure 2 and Figure S2 with multiple organs contributing, while the y-axis
701 represents manual expert annotation of cell types in an organ-specific fashion. The
702 unbiased method discovers relationships between similar cell types found in different
703 organs (highlighted regions); in particular it groups T cells from different organs into a
704 single cluster, B cells from different organs into a different single cluster, and endothelial
705 cells from different organs into a single cluster.
706
707 **Figure 5**. Analysis of all sorted T-cells.
708 a) tSNE plot of all T cells colored by cluster membership. Five clusters were identified.
709 b) Dotplot showing level of expression (color scale) and number of expressing cells
710 (point diameter) within each cluster of T cells. c) tSNE plot of all T cells colored by
711 organ of origin (Fat, Lung, Marrow, Limb Muscle, Spleen or Thymus). d) tSNE plot of
712 all T cells colored by classification of T cells to 4 categories based on expression of Cd4
713 and Cd8 ($Cd4^+$/ $Cd8^+$/ $Cd4^+Cd8^+$ / $Cd4^-Cd8^-$).
714
715 **Figure 6**. Transcription factor (TF) expression analysis.
716 a) Visualization of the precision (ppv) and recall of combinations of 4 TFs. Red bars
717 indicate the number of cells expressing all 4 TFs in the target cell type (true positive) in
718 both the ppv and recall columns. Other colored bars in the ppv column represent the
719 number of cells in the non-target cell types expressing all 4 TFs (false positives). The
720 height of the grey bar in the recall column is the number of cells in the target cell type not
721 expressing all 4 TFs (false negatives). The legend indicates the target cell type next to the
722 red square and all non-target cell types with coexpression. Data shown is the entire
723 dataset subsampled to at most 60 cells per cell type. b) Correlogram of top organ-specific

724    TFs for epithelial cells. Row colors correspond to organ of the most-enriched cell type. c)
725    tSNE visualization of epithelial cells, colored by organ. d) tSNE visualization of
726    endothelial cell expression of select TFs. (grey/low to red/high).  e) Correlogram of top
727    organ-specific TFs for epithelial cells. Row colors correspond to organ of the most-
728    enriched cell type. f) tSNE visualization of epithelial cells, colored by organ. g) tSNE
729    visualization of epithelial cell expression of select TFs.
730
731
732

733 **Supplementary Figure Captions**
734
735 **Supplementary Figure 1** a) Histogram of number of reads per cell for each organ from
736 FACS sorted cells. b) Histogram of number of genes detected per cell for each organ
737 from FACS sorted cells. c) Histogram of number of unique molecular identifiers (UMI)
738 sequenced per cell for each organ from cells prepared by microfluidic droplets. d)
739 Histogram of number of genes detected per cell for each organ for cells prepared by
740 microfluidic droplets.
741
742 **Supplementary Figure 2**. tSNE visualization of all FACS sorted cells annotated by
743 cluster. Clusters are discussed in the text and further analyzed in Figure 4.
744
745 **Supplementary Figure 3** a) tSNE plot of all cells captured by microfluidic droplets
746 color coded by organ. b) Dimensionally reduced tSNE plots for each organ of cells sorted
747 by microfluidic droplets. Color coding indicates distinct clusters. c) Barplots of
748 manually annotated cell types based on differential gene expression across all organs.
749 Coloring of clusters within each organ is consistent between panels b and c.
750
751 **Supplementary Figure 4** a) Number of genes detected by FACS (red), microfluidic
752 droplets (green) and microwell-Seq (blue) (Han *et al.*). b) library saturation fraction for
753 all 10x libraries included in the study. Dotted horizontal line demarcates the median
754 (=0.86).
755
756 **Supplementary Figure 5** Fraction of all detectable genes, for each cell across all organs,
757 (UMI/read threshold is >0) detected at increasing UMI/read thresholds for FACS (left),
758 microfluidic droplet (middle) and microwell-Seq (right).
759
760 **Supplementary Figure 6** Venn diagrams showing the overlap between differentially
761 expressed genes for each common cell type and organs across three methods (FACS,
762 droplet, microwell-Seq). Plotted data are provided in tabular form in Table S2.
763
764 **Supplementary Figure 7** Analysis of dissociation induced gene expression scores
765 across organs.
766
767 **Supplementary Tables**
768
769 **Supplementary Table 1** Number of cells belonging to each annotated cell type across all
770 organs for FACS and microfluidic droplets.
771
772 **Supplementary Table 2** Cell type comparisons and lists of differentially expressed
773 genes across three methods (FACS, droplet, microwell-Seq) and all common organs and
774 tissues.
775
776 **Supplementary Table 3** Combinatorial specificity of transcription factors (TFs) to single
777 cell types. Three combinations of 4 TFs with the highest combinatorial specificity score

778    are presented. The precision (ppv) and recall of each 4-TF combination and cell type is
779    calculated as described in the Methods and main text.

**a**

Circulatory System
Respiratory System
Digestive System
Urinary System
Muscular System
Integumentary System
Adipose Tissues
Immune System
Nervous System

♂ N=4
♀ N=3

Single-cell suspension

FACS

Microfluidic

Sequencing

Data processing

Data analysis

**b**

organ

| | Number of cells |
|---|---|
| Aorta | |
| Bladder | |
| Brain Myeloid | |
| Brain Non-Myeloid | |
| Diaphragm | |
| Fat | |
| Heart | |
| Kidney | |
| Large Intestine | |
| Limb Muscle | |
| Liver | |
| Lung | |
| Mammary Gland | |
| Marrow | |
| Pancreas | |
| Skin | |
| Spleen | |
| Thymus | |
| Tongue | |
| Trachea | |

Number of cells
0    1000    2000    3000    4000    5000

**c**

organ

| | Number of cells |
|---|---|
| Bladder | |
| Heart and Aorta | |
| Kidney | |
| Limb Muscle | |
| Liver | |
| Lung | |
| Mammary Gland | |
| Marrow | |
| Spleen | |
| Thymus | |
| Tongue | |
| Trachea | |

Number of cells
0    2000    4000    6000    8000    10000

**a**

Aorta
Bladder
Brain Myeloid
Brain Non–Myeloid
Diaphragm
Fat
Heart
Kidney
Large Intestine
Limb Muscle
Liver
Lung
Mammary Gland
Marrow
Pancreas
Skin
Spleen
Thymus
Tongue
Trachea

**Aorta**
endothelial cell
erythrocyte
fibroblast
professional antigen presenting cell

**Bladder**
bladder cell
bladder urothelial cell

**Brain Myeloid**
macrophage
microglial cell

**Brain Non-Myeloid**
astrocyte of the cerebral cortex
Bergmann glial cell
brain pericyte
endothelial cell
neuron
oligodendrocyte
oligodendrocyte precursor cell

**Diaphragm**
endothelial cell
lymphocyte
macrophage
mesenchymal stem cell
skeletal muscle satellite stem cell

**Fat**
B cell
endothelial cell
mesenchymal stem cell of adipose
myeloid cell
natural killer cell
T cell

**Heart**
cardiac muscle cell
cardiac neuron
endocardial cell
endothelial cell
fibroblast
leukocyte
myofibroblast cell
smooth muscle cell

**Kidney**
endothelial cell
epithelial cell of proximal tubule
kidney collecting duct epithelial cell
macrophage
natural killer cell

**Large Intestine**
Brush cell of epithelium proper of large intestine
enterocyte of epithelium of large intestine
enteroendocrine cell
epithelial cell of large intestine
large intestine goblet cell

**Limb Muscle**
B cell
endothelial cell
macrophage
mesenchymal stem cell
skeletal muscle satellite cell
T cell

**Liver**
B cell
endothelial cell of hepatic sinusoid
hepatocyte
Kupffer cell
natural killer cell

**Lung**
B cell
ciliated columnar cell of tracheobronchial tree
classical monocyte
epithelial cell of lung
leukocyte
lung endothelial cell
monocyte
myeloid cell
natural killer cell
stromal cell
T cell

**Mammary Gland**
basal cell
endothelial cell
luminal epithelial cell of mammary gland
stromal cell

**Marrow**
B cell
basophil
common lymphoid progenitor
granulocyte
granulocyte monocyte progenitor cell
hematopoietic precursor cell
immature B cell
immature natural killer cell
immature NK T cell
immature T cell
late pro-B cell
macrophage
mature natural killer cell
megakaryocyte-erythroid progenitor cell
monocyte
naive B cell
pre-natural killer cell
pro-B cell
regulatory T cell
Slamf1-negative multipotent progenitor cell
Slamf1-positive multipotent progenitor cell

**Pancreas**
endothelial cell
leukocyte
pancreatic A cell
pancreatic acinar cell
pancreatic D cell
pancreatic ductal cell
pancreatic PP cell
pancreatic stellate cell
type B pancreatic cell

**Skin**
basal cell of epidermis
epidermal cell
keratinocyte stem cell
leukocyte
stem cell of epidermis

**Spleen**
B cell
macrophage
T cell

**Thymus**
DN1 thymic pro-T cell
immature T cell
professional antigen presenting cell

**Tongue**
basal cell of epidermis
keratinocyte

**Trachea**
blood cell
endothelial cell
epithelial cell
mesenchymal cell

$10^3$  $10^2$  $10^1$
Number of cells

myofibroblast cell (Heart)
brain pericyte (Brain Non-Myeloid)
smooth muscle cell (Heart)
basal cell of epidermis (Tongue)
keratinocyte (Tongue)
cardiac muscle cell (Heart)
endothelial cell (Brain Non-Myeloid)
endothelial cell of hepatic sinusoid (Liver)
epithelial cell of proximal tubule (Kidney)
keratinocyte stem cell (Skin)
leukocyte (Skin)
neuron (Brain Non-Myeloid)
skeletal muscle satellite cell (Limb Muscle)
skeletal muscle satellite stem cell (Diaphragm)
type B pancreatic cell (Pancreas)
bladder urothelial cell (Bladder)
bladder cell (Bladder)
pancreatic A cell (Pancreas)
pancreatic PP cell (Pancreas)
pancreatic D cell (Pancreas)
macrophage (Spleen)
immature natural killer cell (Marrow)
mature natural killer cell (Marrow)
natural killer cell (Lung)
pre-natural killer cell (Marrow)
natural killer cell (Fat)
immature B cell (Marrow)
granulocyte (Marrow)
leukocyte (Lung)
pro-B cell (Marrow)
late pro-B cell (Marrow)
pancreatic stellate cell (Pancreas)
pancreatic ductal cell (Pancreas)
endothelial cell (Pancreas)
pancreatic acinar cell (Pancreas)
hepatocyte (Liver)
Kupffer cell (Liver)
Slamf1-negative multipotent progenitor cell (Marrow)
Slamf1-positive multipotent progenitor cell (Marrow)
common lymphoid progenitor (Marrow)
megakaryocyte-erythroid progenitor cell (Marrow)
hematopoietic precursor cell (Marrow)
classical monocyte (Lung)
monocyte (Lung)
granulocyte monocyte progenitor cell (Marrow)
astrocyte of the cerebral cortex (Brain Non-Myeloid)
Bergmann glial cell (Brain Non-Myeloid)
ciliated columnar cell of tracheobronchial tree (Lung)
T cell (Fat)
DN1 thymic pro-T cell (Thymus)
T cell (Limb Muscle)
T cell (Lung)
immature T cell (Marrow)
natural killer cell (Liver)
regulatory T cell (Marrow)
T cell (Spleen)
immature NK T cell (Marrow)
B cell (Marrow)
immature T cell (Thymus)
enteroendocrine cell (Large Intestine)
cardiac neuron (Heart)
kidney collecting duct epithelial cell (Kidney)
erythrocyte (Aorta)
Brush cell of epithelium proper of large intestine (Large Intestine)
epithelial cell (Trachea)
luminal epithelial cell of mammary gland (Mammary Gland)
epithelial cell of lung (Lung)
microglial cell (Brain Myeloid)
oligodendrocyte (Brain Non-Myeloid)
oligodendrocyte precursor cell (Brain Non-Myeloid)
epithelial cell of large intestine (Large Intestine)
enterocyte of epithelium of large intestine (Large Intestine)
large intestine goblet cell (Large Intestine)
endothelial cell (Fat)
endothelial cell (Trachea)
endothelial cell (Diaphragm)
endothelial cell (Limb Muscle)
endothelial cell (Kidney)
endothelial cell (Mammary Gland)
endothelial cell (Aorta)
endothelial cell (Heart)
lung endothelial cell (Lung)
endocardial cell (Heart)
basal cell of epidermis (Skin)
stem cell of epidermis (Skin)
epidermal cell (Skin)
basal cell (Mammary Gland)
mesenchymal stem cell of adipose (Fat)
mesenchymal stem cell (Diaphragm)
mesenchymal stem cell (Limb Muscle)
stromal cell (Mammary Gland)
mesenchymal cell (Trachea)
stromal cell (Lung)
fibroblast (Aorta)
fibroblast (Heart)
myeloid cell (Fat)
macrophage (Brain Myeloid)
professional antigen presenting cell (Aorta)
leukocyte (Heart)
blood cell (Trachea)
macrophage (Limb Muscle)
macrophage (Diaphragm)
leukocyte (Pancreas)
macrophage (Kidney)
monocyte (Marrow)
myeloid cell (Lung)
basophil (Marrow)
macrophage (Marrow)
naive B cell (Marrow)
B cell (Fat)
B cell (Lung)
B cell (Spleen)
B cell (Limb Muscle)
B cell (Liver)
lymphocyte (Diaphragm)
professional antigen presenting cell (Thymus)

cell ontology (organ)

clusters

Log10(nCells+1)

0    1    2    3

**a**

0
1
2
3
4

**b**

T cell cluster

4
3
2
1
0

Aurkb
Kif11
Pbk
Rrm2
Top2a
Eng
Pacsin1
Ckb
Acvrl1
Itm2a
Gem
Klrk1
S100a4
H2.Aa
S100a6
Il6ra
Ly6c1
Fam101b
Dapl1
S1pr1
Tdrd5
Tctex1d1
Wdr78
Rag2
Rag1

Fraction
of expressing cells
· 0.00
· 0.25
● 0.50
● 0.75

Scaled average
expression level
High
Low

**c**

Fat
Lung
Marrow
Limb Muscle
Spleen
Thymus

**d**

Cd4+
Cd8+
Double
negative
Double
positive

**a**

**Myt1l / Bcl11b / Purb / Tcf25**

**Creb3l3 / Nr1h3 / Hnf4a / Klf15**

**Nfib / Lhx2 / Ctnnb1 / Foxp1**

**Etv5 / Nkx2.1 / Tfcp2l1 / Ctnnb1**

**Isl1 / Neurod1 / Nkx6.1 / Xbp1**

**Runx2 / Tcf4 / Spib / Mef2c**

Legend:
- ■ Brain_NonMyeloid__neuron
- ■ Neg_Brain_NonMyeloid__neuron

- ■ Liver__hepatocyte
- ■ Liver__endothelial.cell.of.hepatic.sinusoid
- ■ Neg_Liver__hepatocyte

- ■ Skin__keratinocyte.stem.cell
- ■ Skin__epidermal.cell
- ■ Skin__stem.cell.of.epidermis
- ■ Neg_Skin__keratinocyte.stem.cell

- ■ Lung__type.II.pneumocyte
- ■ Neg_Lung__type.II.pneumocyte

- ■ Pancreas__type.B.pancreatic.cell
- ■ Pancreas__pancreatic.PP.cell
- ■ Pancreas__pancreatic.A.cell
- ■ Pancreas__pancreatic.D.cell
- ■ Neg_Pancreas__type.B.pancreatic.cell

- ■ Marrow__Fraction.A.pre.pro.B.cell
- ■ Fat__myeloid.cell
- ■ Thymus__mesenchymal.stem.cell
- ■ Liver__hepatocyte
- ■ Trachea__leukocyte
- ■ Liver__Kupffer.cell
- ■ Lung__B.cell
- ■ Neg_Marrow__Fraction.A.pre.pro.B.cell

**Organ**
- ● Aorta
- ● Bladder
- ● Brain Myeloid
- ● Brain Non-Myeloid
- ● Large Intestine
- ● Diaphragm
- ● Fat
- ● Heart
- ● Kidney
- ● Liver
- ● Lung
- ● Mammary Gland
- ● Marrow
- ● Limb Muscle
- ● Pancreas
- ● Skin
- ● Spleen
- ● Thymus
- ● Tongue
- ● Trachea

**b** Epithelial cells — Foxq1, Lhx2, Cited1, Pitx1

**c**

**d** Sox4, Pax9, Lhx2, Foxq1

**e** Endothelial cells — Meox2, Foxq1, Plagl1, Gata4, Pbx1, Sox11

**f**

**g** Foxq1, Plagl1, Pbx1, Gata4