**Supplementary Materials: Multi-modal profiling of peripheral blood cells across the human lifespan reveals distinct immune cell signatures of aging and longevity**

**Supplementary Methods**
**Experimental Procedure:**

Recruitment of human subjects. Centenarian study participants were recruited via recruitment mailings based upon statewide voter registration lists throughout the United States. A healthy volunteer effect made it more likely that the enrolled centenarian sample was healthier than centenarians generally. Participants with capacity provided informed written consent and otherwise next of kin acted as a legally authorized representative in providing informed written consent on behalf of the participant. The Boston University Medical Campus Institutional Review Board (BUMC IRB) reviewed and approved this minimum risk study of centenarians and their family members. For participants asked to provide a blood sample for the creation of induced pluripotent stem cell lines, they provided informed written consent for a separate study, again reviewed and approved by the BUMC IRB that is conducted by the Center for Regenerative Medicine also based on the Boston University Medical Campus.


Processing of blood samples. For each centenarian and younger individual involved in this study, 8 mLs of peripheral blood was drawn into each of two BD Vacutainer Cell Preparation Tubes with sodium citrate (BD Biosciences catalog #362760). The tubes were centrifuged at 1,800 x g for 30 minutes at room temperature (RT) and cellular fraction was processed for cell isolation. The cell layer containing peripheral blood mononuclear cells (PBMCs) isolated by Ficoll gradient centrifugation was transferred into a sterile 15 mL conical centrifuge tube. The PBMC sample was brought to 10 mLs with sterile Dulbecco's phosphate buffered saline (DPBS, Invitrogen catalog #14190-144) and centrifuged at 300 x g for 15 min at RT. The supernatant was aspirated and the pellet resuspended in 10 mL sterile DPBS and a cell count performed via hemocytometer. The sample was centrifuged at 300 x g for 10 min at RT and the supernatant aspirated. The pellet was resuspended in chilled (4℃) resuspension medium (40% fetal bovine serum (FBS) hyclone defined, Cytiva catalog #SH30070.03 in Iscove's Modified Dulbecco's Medium (IMDM), catalog #12440053) to achieve a cell concentration of 4 x 10^6 cells/ mL. An equal amount of chilled 2X freezing medium (30% Dimethyl Sulfoxide (DMSO), Sigma catalog #D2650 in IMDM/ 40% FBS) was added to achieve a cell concentration of 2 x 10^6 cells / mL. This mixture was then aliquoted into 1.2 mL cryovials (Corning catalog #430487) at 1 mL / vial. These vials were then brought to -80℃ before being transferred to a -150℃ deep freezer.

CITE-seq of PBMCs of centenarians. PBMC samples (2 x 10^6 cells/sample) from the centenarian cohort were thawed rapidly and mixed with 15 mL StemSpan SFEM II medium (CAT#09655) with L-glutamine (1:1000 conc). These samples were then brought to 50 mL with sort buffer (2% Bovine Serum Albumin (BSA), Millipore Sigma catalog #EM-2930 in DPBS) and centrifuged for 5 min at 400 x g at RT. The supernatant was aspirated, and the pellet of PBMCs resuspended in 30 mL sort buffer and centrifuged for 5 min at 400 x g at RT. The pellet was resuspended in 2 mL StemSpan medium (+L-glutamine) and filtered through a 40 uM filter. The samples were then incubated in this medium for 1 hour at 37℃/5% CO2. Following this

1   incubation, the samples were centrifuged for 5 min at 400 rcf at RT and the supernatant
2   aspirated. Each pellet was then resuspended in 50 uL labeling buffer (1% BSA in PBS) and 5 uL
3   Human TruStain FcX(Biolegend) was added. The samples were incubated at 4℃ for 10
4   minutes. During this incubation, the TotalSeq-C antibody pool (Biolegend) containing 1 ug of
5   each antibody was prepared and centrifuged for 10 min at 14,000 x g at 4℃. The supernatant
6   was then transferred and used as the antibody mix for each sample. The following TotalSeq-C
7   antibodies were used: CD274, CD3, CD8, CD19, CD33, CD4, CD14, CD16, CD56, and CD279
8   (BioLegend). 20 uL of the antibody mix was added to each sample and the samples were
9   brought to 100 uL with labeling buffer. The samples were incubated for 30 minutes at 4℃.
10  Following this incubation, the samples were washed with 1.3 mL labeling buffer and centrifuged
11  at 400 x g for 5 minutes at RT. This washing step was repeated for a total of 3 washes. The
12  pellets were then resuspended in 500 uL labeling buffer with calcein blue AM (1:1000) (Thermo
13  Fisher, catalog #C1429). 1 x 10^5 calcein blue positive cells (live cells) were then sorted on a
14  Beckman Coulter MoFlo Astrios cell sorter. The sorted samples were centrifuged for 5 min at
15  400 rcf at room temperature and resuspended in 120 uL resuspension buffer (0.04% BSA in
16  DPBS). The samples were counted via hemocytometer and diluted to 600 cells/uL with
17  resuspension buffer.
18
19  Flow cytometry analysis. Frozen PBMCs from centenarian and younger samples were banked
20  and thawed and following the same protocol as described above and below for the scRNA-seq.
21  Cytometry panel design and validation, sample staining and sample acquisition were performed
22  closely following OMIP-069 protocol[1]. Briefly, cells were stained with a live-dead dye (Live Dead
23  Blue, Thermo Fisher), blocked with FcBlock reagent (Biolegend) and Monocyte Blocker
24  (Biolegend), stained with the fluorescent antibody cocktail and Brilliant Buffer Plus (BD
25  Biosciences), washed and analyzed on the Cytek Aurora spectral cytometer (Cytek
26  Biosciences). At least 500,000 cells were recorded for each PBMC sample. Younger PBMCs
27  and Ultracomp eBeads plus (Thermo Fisher) were used as single stain unmixing controls. Data
28  were processed in SpectoFlo 2.2 (Cytek Biosciences) and OMIQ data analysis platform. We
29  reproduced the exact gating strategy as shown in Alpert et al [2] and extracted cell counts and
30  population proportions for downstream statistical analysis. The following antibodies were used
31  to identify cell populations: CD45RA BUV395, CD16 BUV496, CD56 BUV737, CD8 BUV805,
32  CCR7 BV421 (BD Biosciences), CD123 Super Bright 436 (Thermo Fisher), CD33 BV510, CD14
33  BV570, CD3 Spark Blue 550, CD19 Spark NIR 685 (Biolegend), CD4 cFluor YG584 (Cytek
34  Biosciences).
35
36
37  **Single cell analysis:**
38  New England Centenarian dataset:
39
40  *CITE-seq and CellRanger Preprocessing.* Cellular Indexing and epitopes sequencing (CITE-
41  seq) was performed on the 7 centenarians and 2 younger age individuals using a commercial
42  droplet-based platform (10x Chromium). We constructed 5' gene expression libraries (GEX), as
43  well as surface protein libraries (antibody derived tags, ADT) following the manufacturer's user
44  guide. These libraries were sequenced on two runs of an Illumina NextSeq 2000 instrument

1    generating 438 and 535 million reads respectively. Raw sequencing files were converted to
2    fastq and demultiplexed using bcl2fastq v.2.20 and Cellranger v.3.0.2. Counts for the
3    expression and antibody capture libraries were derived by simultaneously mapping the
4    respective fastq files to the human genome (GRCh38) and to the feature reference of the
5    TotalSeq-C antibodies used using the corresponding parameters in cellranger count (v.3.0.2).
6    This pipeline includes the alignment, barcode and UMI counting
7
8    *Filtering, PCA Analysis, Batch Correction, and Clustering.* After processing the samples through
9    CellRanger, we performed filtering, normalization, and principal component analysis using
10   Seurat v.3 [3].
11        First, we performed quality control steps based on the number of genes and UMIs
12   detected per cell, and percent of mitochondrial genes expressed per cell. To remove poor
13   quality cells with low RNA content, we removed cells with less than 200 genes detected. To filter
14   out outlier cells and doublets, we filtered out cells with greater than 3,000 detected genes, as
15   well as cells with greater than 15,000 UMIs. To account for cells that are damaged or dying, we
16   removed cells with greater than 15 percent mitochondrial counts expressed.
17        After filtering, we normalized the RNA-level expression data for each cell to compare
18   gene expression between sample cells; Gene counts for each cell were normalized by total
19   expression, multiplied by a scale factor of 10,000 and transformed to a log scale. We
20   normalized the protein-level expression data by applying a centered log ratio (CLR)
21   transformation for each cell to account for differences in total protein ADT counts that make up
22   each cell.
23        Further downstream analyses were performed on the RNA-level expression data. PCA
24   based on the top 2000 highly variable genes was performed for dimensionality reduction, and
25   the top 20 significant PCs were selected that explain the most variability in the data. The top
26   significant PCs were used as an input for clustering the cells and for nonlinear dimension
27   methods mentioned below to identify populations of cells with similar expression profiles. To
28   account for technical variations between samples from different experimental batches, we
29   corrected the PCA embeddings using the Harmony algorithm [4], a method that iteratively clusters
30   and corrects the PC coordinates to adjust for batch specific effects. We assessed the integration
31   of these datasets by employing PCA visualizations of batches of cells and calculating the
32   average silhouette width (ASW) score[5,6] for each cell type population based on the top 20
33   principal components before correction and the top 20 harmony components after batch
34   correction reported with Wilcoxon rank sum test and p-value significance threshold of 0.05. We
35   clustered cells based on graph-based methods (SNN and Louvain community detection
36   method) using the top 20 Harmony-adjusted components and used the Unifold Manifold
37   Approximation and Projection (UMAP) algorithm[7] to visualize clusters of cells and other known
38   annotations.
39
40    *Identification and classification of cell types.* We used a multi-modal approach to identify
41   immune subpopulations in the NECS dataset. First, we used the 10 cell-surface protein immune
42   cell marker panel of expression to identify main immune cell types. We then further partitioned
43   the main immune cell types into immune subtypes using graph-based clustering and based on
44   the expression of immune cell type signatures from literature [8,9]. The average expression score

1    of each signature was calculated for a single cell by calculating the average scaled expression
2    of all genes within a signature, with the scaling based on the expression of a control set of
3    genes (AddModuleScore function in Seurat[3]), and by taking the absolute value of the average
4    scaled expression to compare scores across signatures within a cell population.
5
6    Publicly available datasets:
7    *Data collection, filtering, PCA analysis, and clustering.* We downloaded the raw UMI matrix for
8    the scRNA-seq dataset of PBMCs from 45 younger age individuals of European descent [10],
9    which we will refer to as NATGEN. We also downloaded the raw UMI matrix for the scRNA-seq
10   dataset of PBMCS from 7 supercentenarians and 5 younger age individuals of Japanese
11   descent[11], which we will refer to as PNAS. For both PBMC datasets, we performed all
12   downstream processing including filtering, normalization, and scaling of data using the Seurat
13   v.3 [3]. For both datasets, we performed quality control steps based on the number of genes and
14   UMIs detected per cell, and percent of mitochondrial genes expressed per cell. For the
15   NATGEN dataset, we filtered cells based on similar thresholds from the original manuscript[10].
16   For the PNAS dataset, we filtered cells as previously published[11]. After filtering the datasets, we
17   normalized the expression levels of each cell to compare gene expression between sample
18   cells; gene counts for each cell were normalized by total expression, multiplied by a scale factor
19   of 10,000 and transformed to a log scale. We then performed PCA analysis based on the top
20   2,000 highly variable genes detected and clustered cells based on graph-based methods (SNN
21   and Louvain community detection method)[3] based on the top significant PCs for each data set
22   implemented in Seurat. We used the UMAP algorithm[4] to visualize the clusters of single cells
23   and other known annotations.
24
25   *Subpopulation Identification and Harmonization.* To define the set of consensus immune cell
26   types across regular aging and longevity, we first identified subpopulations of each cell type
27   using immune cell type signatures from literature [8,9]. The average expression score of each
28   signature was calculated for each single cell across datasets as described for the NECS dataset
29   using Seurat. In addition, we compared the expression of canonical gene markers of the
30   immune populations identified for comparison [12]. After identifying subpopulations in each
31   dataset, we then integrated these scRNA-seq datasets of PBMCs by correcting the PCA
32   embeddings using the Harmony algorthim[4] and assessed batch correction using average
33   silhouette width (ASW) score[5,6] for each cell type population based on the top 20 principal
34   components before correction and the top 20 harmony components after batch correction
35   reported with Wilcoxon statistic and p-value significance threshold of 0.05. We determined the
36   four age groups in the integrated datasets by grouping subjects into four approximate quantile
37   groups based on age in decades.
38
39   Statistical Methods.
40   *Heterogeneity of the overall cell type distribution.* We compared the overall cell type composition
41   differences across samples and age groups by calculating the cell type diversity statistic for
42   each sample[13]. The cell type diversity statistic $E_s$ is the normalized Shannon entropy based on
43   the cell type proportions $p_i$ for the sample s. The normalization is based on the total number of
44   cell types k to make the entropy measure independent of the number of cell types.

1

$$E_s = \frac{-\sum_{i=1}^{k} p_{i_s} \log(p_{i_s})}{\log(k)} - 1$$

2

3

4   The normalization is important because we do not have the same number of cell types in the
5   various age groups. For example, if we compare two groups with perfectly uniform cell type
6   proportions, but different number of cell types, the Shannon entropy would show a difference
7   while both distributions are uniformly distributed. We performed ANOVA to assess differences
8   between age groups, and used p-value $< 0.05$ to determine statistical significance.

9

10  *Effect of age and sex on cell type distribution.* To investigate cell type specific differences
11  across age and sex groups, at subject level, we applied a Bayesian multinomial regression
12  model to the cell type abundances. The outcome of this analysis was the vector of counts of the
13  13 cell types in each subject, and the exposure was age group and the confounder was sex.
14  Since high throughput sequencing data is compositional [14–16] and cell type frequencies are
15  constrained to the total cells per sample [17], we used a multinomial model that accounts for this
16  constraints and estimates cell type proportions that add up to 1 for each sample. We
17  implemented a Bayesian analysis in the R-package rjags that allows for the fitting of a
18  multinomial regression model without any additional specification. The model set up is

19

20  $$Y_{i,1:J} \sim Multinomial(p_{i,1:J}, N.total_i)$$
21  $$\log(q_{i,j}) = \alpha_j + \beta_{age.group_i,j} + \gamma_{sex_i,j}$$
22  $$p_{i,j} = \frac{q_{i,j}}{\sum_{k=1}^{N.ct} q_{i,k}}$$

23

24  $$\alpha_j \sim Normal(0, 0.001)$$
25  $$\beta_{age.group_i,j} \sim Normal(0, 0.001)$$
26  $$\gamma_{sex_i,j} \sim Normal(0, 0.001)$$

27

28  where $Y_{i,1:J}$ represents the abundances of cell type $1:J$ for sample $i$ that are modeled using a
29  Multinomial distribution with probabilities $p_{i,1:J}$ $\sum_{j=1}^{J} Y_{i,j} = N.total_i$ for all sample $i$ and
30  $\sum_{j=1}^{J} p_{i,j} = 1$. The probabilities $p_{i,1:J}$ depend on age and sex through the function $\log(q_{i,j})$. We
31  chose a reference for each age group (younger age) and sex (male). The implementation of the
32  analysis in rjags does not require a logistic parameterization since rjags can work with
33  unnormalized probabilities that are internally summed up to 1 [18]. This allows us to calculate
34  explicitly the predicted probabilities of all the cell types for every combination of age and sex,
35  with no need for a reference cell type (Supplementary Table S10). Therefore, for each cell
36  type $1:j$, we can estimate the probability $p_{i,j}$ for each group profile $i$ from the estimates of the
37  parameters $\alpha_j$, $\beta_{age.group_i,j}$, $\gamma_{sex_i,j}$ as shown below:

38  $$\hat{p}_{i,j} = \frac{\exp(\hat{\alpha}_j + \hat{\beta}_{age.group_i,j} + \hat{\gamma}_{sex_i,j})}{\sum_{k=1}^{N.ct} \exp(\hat{\alpha}_j + \hat{\beta}_{age.group_i,j} + \hat{\gamma}_{sex_i,j})}$$

39  where the notation ^ represents the estimated parameters using rjags.

1

2 The model was estimated using Markov Chain Monte Carlo (MCMC) sampling using rjags, the

3 R package for JAGS[19]. We ran parameter inference for all coefficients for group level predicted

4 probabilities of composition ($p_{i,j}$) and the age group and sex effect coefficients ($\beta, \gamma$) using

5 1,000 iterations with 500 iterations for burn-in. We estimated the group level composition

6 predicted probabilities and 95 percent credible interval for males and female subjects for

7 younger, middle, older, and EL age groups across all immune cell types. In addition, to assess

8 the significance of the effect of age and sex in each cell type, we calculated the z-score (Z) for

9 parameters $\beta$ and $\gamma$ based on the mean estimate and standard error of the posterior

10 distribution. We subsequently calculated the two-sided large sample p-value based on the

11 standard normal distribution: $2\Phi(-|Z|)$) where $\Phi$ is the standard normal cumulative distribution

12 function. We calculated the adjusted p-value based on the Benjamin and Hochberg correction

13 for multiple testing across all coefficients tested.

14 In addition, to investigate cell type specific differences in other age group comparisons

15 including older age vs. EL, we applied the multinomial regression analysis as described above

16 with the exception of the age group reference set to EL.

17

18

19 *Analysis of the hierarchy of peripheral immune compartments*. To estimate the hierarchy of

20 peripheral immune compartments, we utilized K2Taxonomer (v1.0.5)[20], which performs top-

21 down partitioning of cell types based on the relative similarity of their transcriptomic profiles.

22 Prior to running K2Taxonomer, we performed several data processing steps. First, for each

23 subject, we removed profiles of lowly represented cell types with fewer than 10 profiles for that

24 subject only. Following this filter, we further removed all plasma cell profiles because they were

25 represented in fewer than 10 subjects. Next, for each subject we aggregated single-cell profiles

26 of each cell type into a "pseudo-bulk" profile by summing the counts of each gene, followed by

27 normalization to log2(counts-per-million). The resulting data set included 515 total profiles

28 across 66 subjects and 12 cell types. The number of subjects for which each of these 12 cell

29 types were identified in each of the four batches is given in Supplementary Table S12. Next, we

30 removed lowly expressed genes, i.e., genes that failed to reach 2 counts-per-million in at least 2

31 profiles across all batches, which left 12,354 genes. Finally, we performed batch correction on

32 these data using ComBat (v3.40.0)[21], parameterized to preserve the cell type-specific signals

33 using "mod" parameter of the ComBat() function. K2Taxonomer was run on these data, using

34 the package's "group-level" workflow, and setting the number of features parameter, "nFeats",

35 to use 5% of the total genes, i.e., 618 genes per partition estimate. Finally, we calculated the

36 cell type diversity statistics for each K2Taxonomer generated cell type subgroup.

37

38 *Cell type specific differential gene expression analysis*. To investigate the cell type specific

39 differences across age groups, we used a Bayesian mixed effects model with the rjags R

40 package[19] to perform differential gene expression analysis across the four age groups. For each

41 gene, we applied the model to the expression values across cells from all subjects with

42 exposure variable for age group comparing middle, older, and EL age groups in reference to the

43 younger age group, and confounding variables including sex, ethnicity, and batch, and a

44 random effect that accounts for within subject correlations of cells. For each cell type, we first

filtered genes to keep genes with expression in at least fifty percent of the smallest cell type population. We then applied the Bayesian mixed effects model where for each cell type, to each gene i:

$$Gene_i \sim Normal(\mu_i, \tau)$$
$$\mu_i = \beta_{S,sample} + \beta_1 Middle_i + \beta_2 Older_i + \beta_3 EL_i + \beta_4 Sex_i + \beta_5 Ethnicity_i + \beta_6 Batch_i$$

$$\beta_{S,sample} \sim Normal(\beta_0, \tau_s)$$
$$\beta_k \sim Normal(0, 0.0001) \; \forall \; k \in [0,\ldots,6]$$
$$\tau \sim Gamma(0.0001, 0.0001)$$
$$\tau_s \sim Gamma(0.00001, 0.00001)$$

where $Gene_i$ is the log-normalized expression of a gene for each cell $i$ in the cell type; $\beta_1, \beta_2, \beta_3, \beta_4, \beta_5, \beta_6, \beta_0$ are model parameters. The subject specific intercept $\beta_{S,sample}$ follows a normal distribution with population mean, $\beta_0$. The model is adjusted by fixed covariates sex, batch, and ethnicity, as well as a random effect based on samples to account for differences in cell abundances between samples within groups. Using this model, we monitored the age-dependent coefficients ($\beta_1$, $\beta_2$, $\beta_3$) across 10,000 MCMC iterations with 2,500 burn-in iterations to obtain the log fold change (logFC) based on age group. Then, we calculated the z-score (Z) for the age-dependent parameter based on the mean estimate and standard error of the posterior distribution. We subsequently calculated the two-sided p-value based on the standard normal distribution: $2*\Phi(-|Z|)$ where $\Phi$ is the standard normal cumulative distribution function. We calculated the FDR based on the Benjamin and Hochberg correction for multiple testing across all genes tested. Significant differential genes were selected based on a significance of FDR < 0.05 and fold change cutoff of 10 percent ($|logFC| > log(1.1)$).

In addition, to investigate cell type specific differences between other age group comparisons including older age vs. EL, we applied the mixed effects model as described above with the exception of the age group reference set to EL.

*Bulk level differential gene expression analysis.* We performed differential gene expression analysis at the bulk level between age groups using DESeq2 [22]. We filtered genes in the single cell data to keep genes expressed in at least 50% of the smallest cell type population. We then aggregated the raw counts per sample and ran DESeq2 to perform normalization and fit a negative binomial generalized linear model with covariates age group, sex, batch, and ethnicity. We used Wald test to identify differentially expressed genes between middle, older, and EL v. younger age with a log fold change greater than log2(1.5) and FDR < 0.05.

**References**

1 Park LM, Lannigan J, Jaimes MC. OMIP-069: Forty-Color Full Spectrum Flow Cytometry Panel for Deep Immunophenotyping of Major Cell Subsets in Human Peripheral Blood. *Cytometry Part A* 2020; **97**: 1044–51.

2 Alpert A, Pickman Y, Leipold M, *et al.* A clinically meaningful metric of immune age derived from high-dimensional longitudinal monitoring. *Nature Medicine* 2019; **25**: 487–95.

3 Stuart T, Butler A, Hoffman P, *et al.* Comprehensive Integration of Single-Cell Data. *Cell* 2019; **177**: 1888-1902.e21.

4 Korsunsky I, Millard N, Fan J, *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nature Methods* 2019; **16**: 1289–96.

5 Tran HTN, Ang KS, Chevrier M, *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biology* 2020; **21**: 12.

6 Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics* 1987; **20**: 53–65.

7 McInnes L, Healy J, Melville J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv:180203426 [cs, stat]* 2020; published online Sept 17. http://arxiv.org/abs/1802.03426 (accessed April 12, 2021).

8 Chen B, Khodadoust MS, Liu CL, Newman AM, Alizadeh AA. Profiling tumor infiltrating immune cells with CIBERSORT. *Methods Mol Biol* 2018; **1711**: 243–59.

9 Popescu D-M, Botting RA, Stephenson E, *et al.* Decoding human fetal liver haematopoiesis. *Nature* 2019; **574**: 365–71.

10 van der Wijst MGP, Brugge H, de Vries DH, Deelen P, Swertz MA, Franke L. Single-cell RNA sequencing identifies celltype-specific cis-eQTLs and co-expression QTLs. *Nature Genetics* 2018; **50**: 493–7.

11 Hashimoto K, Kouno T, Ikawa T, *et al.* Single-cell transcriptomics reveals expansion of cytotoxic CD4 T cells in supercentenarians. *PNAS* 2019; **116**: 24242–51.

12 Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nature Biotechnology* 2018; **36**: 411–20.

13 Karagiannis TT, Monti S, Sebastiani P. Cell Type Diversity Statistic: An Entropy-Based Metric to Compare Overall Cell Type Composition Across Samples. *Frontiers in Genetics* 2022; **13**. https://www.frontiersin.org/article/10.3389/fgene.2022.855076 (accessed April 10, 2022).

14 Gloor GB, Wu JR, Pawlowsky-Glahn V, Egozcue JJ. It's all relative: analyzing microbiome data as compositions. *Annals of Epidemiology* 2016; **26**: 322–9.

15 Gloor GB, Macklaim JM, Pawlowsky-Glahn V, Egozcue JJ. Microbiome Datasets Are Compositional: And This Is Not Optional. *Front Microbiol* 2017; **8**. DOI:10.3389/fmicb.2017.02224.

16 Lin H, Peddada SD. Analysis of microbial compositions: a review of normalization and differential abundance analysis. *npj Biofilms Microbiomes* 2020; **6**: 1–13.

17 Luecken MD, Theis FJ. Current best practices in single-cell RNA-seq analysis: a tutorial. *Mol Syst Biol* 2019; **15**: e8746.

1  18 Lunn D, Jackson C, Best N, Thomas A, Spiegelhalter D. The BUGS Book: A Practical
2      Introduction to Bayesian Analysis. London: Chapman Hall, 2013.

3  19 Plummer M. Penalized loss functions for Bayesian model comparison. *Biostatistics* 2008; **9**:
4      523–39.

5  20 Reed ER, Monti S. Multi-resolution characterization of molecular taxonomies in bulk and
6      single-cell transcriptomics data. *Nucleic Acids Res* 2021; **49**: e98.

7  21 Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
8      empirical Bayes methods. *Biostatistics* 2007; **8**: 118–27.

9  22 Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-
10     seq data with DESeq2. *Genome Biology* 2014; **15**: 550.