

R code and downstream analysis objects for the scRNA-seq atlas of human breast spanning normal, preneoplastic and tumorigenic states

Yunshun Chen^{1,3}, Bhupinder Pal^{2,4}, Geoffrey J Lindeman^{1,5}, Jane E Visvader^{1,3}, and Gordon K Smyth^{1,6*}

¹The Walter and Eliza Hall Institute of Medical Research, Parkville, Victoria 3052, Australia

²Olivia Newton-John Cancer Research Institute, Heidelberg, Vic 3084 Australia

³Department of Medical Biology, The University of Melbourne, Parkville, Victoria 3010, Australia

⁴School of Cancer Medicine, La Trobe University, Bundoora, Vic, Australia

⁵Department of Medicine, The University of Melbourne, Parkville, Victoria 3010, Australia

⁶School of Mathematics and Statistics, The University of Melbourne, Parkville, Victoria 3010, Australia

*corresponding author: Gordon Smyth (smyth@wehi.edu.au)

ABSTRACT

Breast cancer is a common and highly heterogeneous disease. Understanding the cellular diversity in the mammary gland and its surrounding micro-environment across different states can provide insight into the cancer development in human breast. Recently, a large-scale single-cell RNA expression atlas was constructed of the human breast spanning normal, preneoplastic and tumorigenic states. Single-cell expression profiles of nearly 430,000 cells were obtained from 69 distinct surgical tissue specimens from 55 patients. This article extends the study by providing downstream processed R data objects, complete cell annotation and R code to reproduce all the analyses. Details of all the bioinformatic analyses that produced the results described in the study are provided.

Background & Summary

Breast cancer is the most commonly diagnosed cancer and the leading cause of cancer death in women¹. It is a very heterogeneous disease at the molecular level². Different breast cancer subtypes can be characterized on the basis of expression profiles of markers such as estrogen receptor (ER), progesterone receptor (PgR), and human epidermal growth factor receptor 2 (HER2)³. The development of certain cancer subclasses is also known to be associated with mutations such as BRCA1⁴. Recently, we and colleagues constructed a large-scale single-cell RNA expression atlas of the human breast spanning normal, preneoplastic and tumorigenic states (subsequently referred to as the ScBrAtlas)⁵. Single-cell expression profiles of nearly 430,000 cells were obtained from 69 distinct surgical tissue specimens from 55 patients (Figure 1). This article extends the ScBrAtlas by providing downstream processed R data objects, complete cell annotation and R code to reproduce all the analyses.

The ScBrAtlas spanned several stages of breast cancer genesis. First, reduction mammoplasties were obtained from women with no family history of breast cancer to explore cellular diversity in normal breast epithelia as well as complexity within the normal breast ductal micro-environment. Three major epithelial cell populations revealed in literature⁶: basal, luminal progenitor (LP), and mature luminal (ML), were confirmed by the bulk RNA-seq signatures for sorted epithelial populations as well as the cell clustering of the integrated single cell transcriptomic data on normal breast epithelia. Similar cell type composition within the normal epithelium was observed across multiple healthy donors with different hormonal status (pre- and post-menopausal). For the immune and stromal micro-environment of normal breast tissue, integration analysis and the pseudo-bulk differential expression analysis identified different cell clusters including fibroblasts, endothelial cells (vascular and lymphatic), pericytes, myeloid, and lymphoid cells. Differential abundance analysis revealed that fibroblasts are more abundant whereas vascular endothelial cells are less abundant in post-menopausal tissue compared to pre-menopausal tissue⁵.

Next, breast tissue from BRCA1 mutation carriers was obtained for investigating cellular changes in precancerous state. Overall, the differences of stromal and immune subsets between normal and BRCA1+/- preneoplastic tissue were not significant, nor was the proportions of different cell clusters. However, extensive changes in the tissue micro-environment were observed between the preneoplastic and the neoplastic states in BRCA1 mutation carriers⁵.

39 Finally, ER+, HER2+ and triple negative breast cancer (TNBC) tumors were obtained from treatment-naïve patients for
40 exploring the degree of heterogeneity within the cancer cell compartment and its micro-environment across different tumor
41 subtypes. Extensive inter-patient heterogeneity was revealed by single cell integration analyses across all cancer subtypes.
42 Within the tumor populations, a discrete cluster of cycling MKI67+ tumor cells were observed for all three major breast cancer
43 subtypes. Within the tumor micro-environment, different immune landscapes were observed in different cancer subtypes. Both
44 TNBC and HER2 featured a proliferative CD8+ T-cell cluster, whereas ER+ tumors primarily comprised cycling TAMs. In
45 addition, matched pairs of ER+ tumors and involved lymph nodes were profiled for examining the relationship between primary
46 breast tumors and malignant cells that seed lymph nodes. Clonal selection and expansion were observed in some patients,
47 whereas mass migration of cells from the primary tumor to the LN was observed in some other patients⁵.

48 The ScBrAtlas provides a valuable resource for understanding cellular diversity and cancer genesis in human breast. The
49 examination and exploration of the single cell data presented in this study required large-scale bioinformatics analyses for
50 multiple groupings of the original data. While genewise read counts were previously made publicly available for all 421,761
51 individual cells⁷, downstream results after quality filtering, data integration and cell clustering were not provided.

52 In this report we describe the bioinformatics analysis used in the ScBrAtlas in greater detail. We provide a complete
53 description of the quality control filters used to select 341,874 cells for downstream analyses. The technical quality of both the
54 10X single-cell transcriptomic data sets and the bulk RNA-seq reference data set is assessed to demonstrate the reliability of
55 the data. We provide downstream R data objects corresponding to each data integration and cell clustering presented in the
56 ScBrAtlas, together with R code to reproduce the data objects. Crucially, the data objects provided here include cell barcodes
57 by which each individual cell can be tracked through all the analyses. We also provide detailed information allowing the copy
58 number variation analyses to be mapped back to individual samples and cell clustered, thus providing a way to distinguish
59 putative malignant cancer cells from normal epithelial cells in the cancer tumors. All the resources and the detailed information
60 can be easily accessed and utilized by researchers for further exploration and clinical validation, which may lead to discoveries
61 of novel approaches for personalized breast cancer treatment in the future.

62 **Methods**

63 **Read alignment and count quantification**

64 Single-cell RNA-seq expression profiles of 69 samples from 55 patients were generated by the 10x Genomics Chromium
65 platform and an Illumina NextSeq 500 system (Fig. 1a, Supplementary Table 1). The original Illumina BCL output was
66 converted to FASTQ files and then aligned to the human reference genome GRCh38 (cellranger ref v3.0.0) using Cell Ranger
67 v3.0.2 (<https://support.10xgenomics.com>). The outputs for each individual sample contain a count matrix in matrix market
68 mtx.gz format, barcode information and feature information both in tab-delimited tsv.gz format (Supplementary Table 1). Any
69 cell with at least 500 sequence reads assigned to genes was included in this output. All the downstream bioinformatics analyses
70 were performed based on the cellranger outputs.

71 **Quality control and cell filtering**

72 Sequence read counts were obtained for a total of 421,761 cells across the 69 samples (Supplementary Table 2). Quality
73 control (QC) was performed individually for each scRNA-seq sample. Cells with high proportion of mitochondrial reads were
74 considered as of low quality and hence were filtered⁸. A lower bound of 500 was generally applied to the number of detected
75 genes for each cell, although this was reduced to 400 or 300 for a small number of samples with low read coverage. Upper
76 bounds of a combination of number of detected genes and library size were also applied to each sample to remove potential
77 doublets. The threshold values of these QC metrics for each individual sample are shown in Supplementary Table 2 and are
78 also supplied in machine-readable form as part of this data submission⁹. A total of 341,874 cells remained after quality filtering
79 for downstream analysis.

80 **Single-cell RNA-seq integration analysis**

81 The samples included breast tissues from normal healthy donors, BRCA1 mutation carriers and patients diagnosed with different
82 types of breast cancer (triple negative, ER+ and HER2+). Matching pairs of tumor and lymph node (LN) samples, as well as
83 tumor samples from male patients, were also included. The single-cell analysis strategy involved grouping together comparable
84 samples, integrating the profiles, then clustering cells into putative cell types. A total of 16 different sample-groups were
85 integrated (fig. 1b). Some samples were involved in more than one integration, for example the pre-neoplastic samples with
86 BRCA1 mutations were integrated first with the normal samples and later with the BRCA1 triple negative (TN) tumor samples.
87 For some sample-groups analyses, subsets of cells were extracted, re-integrated and re-clustered. The total number of cell
88 cluster analyses is shown in Table 1.

89 Samples were integrated using the Seurat anchor-based integration method¹⁰. To perform dimensionality reduction, the
90 first 30 principal component were computed and used for the cell clustering and t-distributed stochastic neighbor embedding

91 (t-SNE) visualization¹¹. The default Louvain clustering algorithm¹² was used for cell cluster identification. Different resolution
92 parameters were used in different cell clustering analyses to ensure repeatability and the best interpretation of the data⁵.

93 We provide here the Seurat data objects containing each of the cluster analyses as R data files (Table 2). The R data objects
94 contain cell cluster details for each cell. The R code by which each R object was constructed is also provided (Table 2).

95 **Differential expression and pathway analysis**

96 Differential expression analyses were performed to detect marker genes for different cell clusters. In order to account for the
97 biological variation between different patients, a pseudo-bulk approach was used in most cases where read counts from all cells
98 under the same cluster-sample combination were summed together to form pseudo-bulk samples. The edgeR's quasi-likelihood
99 pipeline was used for pseudo-bulk differential expression analysis, where the baseline differences between patients were
100 incorporated into the linear model¹³. The Seurat's FindMarkers function was applied where pseudo-bulk samples were not
101 satisfactory due to low cell numbers or imbalanced cluster-sample combination. KEGG pathway analyses were performed
102 using the kegg function of the limma package¹⁴.

103 **Data visualization**

104 Ternary plot visualization was performed as previously described¹⁵. Ternary plots position cells according to the proportion
105 of basal, LP- or ML-positive signature genes expressed by that cell and were generated using the vcd package¹⁶. The t-SNE
106 visualization for all the integration analyses were generated using the RunTSNE function in Seurat with a random seed of
107 2018 for reproducibility. Diffusion plots were generated using the destiny package¹⁷. MDS plots were created with edgeR's
108 plotMDS function. Log2-CPM values for each gene across cells were calculated using edgeR's cpm function with a prior
109 count of 1. Heat maps were generated using the pheatmap package. Log2-CPM values were standardized to have mean 0 and
110 standard deviation 1 for each gene before producing the heat maps, after which genes and cells were clustered by the Ward's
111 minimum variance method¹⁸.

112 **Bulk RNA-seq data and differential expression analysis**

113 RNA-seq experiments were performed to obtain signature genes of basal, luminal progenitor (LP), mature luminal (ML) and
114 stromal cell populations. Epithelial cells for basal, LP, and ML populations were sorted from eight independent patients and
115 stroma from five patients. For one particular patient, samples were collect from both left and right breast for each of the four
116 cell populations. For another patient, ML cell population was collected twice. The complete RNA-seq data contains 9 basal, 9
117 LP, 10 ML and 6 stroma samples. RNA-seq libraries were prepared using Illumina's TruSeq protocol and were sequenced on
118 an Illumina NextSeq 500.

119 Reads were aligned to the hg38 genome using Rsubread v1.5.3¹⁹. Gene counts were quantified by Entrez Gene IDs using
120 featureCounts and Rsubread's built-in annotation²⁰. Gene symbols were provided by NCBI gene annotation dated 29 September
121 2017. Immunoglobulin genes as well as obsolete Entrez Ids were discarded. Genes with count-per-million above 0.3 in at least
122 3 samples were kept in the analysis. TMM normalization was performed to account for the compositional biases between
123 samples.

124 Differential expression analysis was performed using limma-voom²¹. Patients were treated as random effects and the
125 intra-patient correlation was estimated by the duplicateCorrelation function in limma. Pairwise comparisons between the four
126 cell populations were performed using TREAT with a fold change threshold of 1.5²². An FDR cut-off of 0.05 was applied for
127 each comparison. Genes were considered as signature genes for a particular cell type if they were upregulated in that cell type
128 in all the pairwise comparisons. The analysis yielded 515, 323, 765, and 1094 signature genes for basal, LP, ML, and stroma,
129 respectively. In this submission we provide gene symbols of the signature genes as an R data file and R code to reproduce the
130 bulk RNA-seq analysis^{9,23}.

131 **Differential abundance analysis**

132 Differential abundance analyses were performed to examine the differences in cell cluster frequencies between pre-menopause
133 and post-menopause groups in normal breast tissue micro-environment. Quasi-multinomial and quasi-binomial generalized
134 linear models were used in order to account for the inter-patient variability. The numbers of cells under all the clusters from
135 each individual donor were counted and used as the response variable in the model. The glm function of the stats package was
136 used to fit the cell numbers against cell clusters, donors, plus a cluster-menopausal interaction term. The quasi-Poisson family
137 was used in the glm function.

138 A quasi-multinomial F-test was performed to test for differences in cluster frequencies across all the clusters between pre-
139 and post-menopausal samples, which yielded a p-value of 0.007. To test for cluster frequency differences for each individual
140 cluster, we compared the cell numbers of that cluster with the aggregated cell numbers of all the other clusters across all the
141 donors. Quasi-binomial generalized linear models were fitted and quasi-binomial F-tests were performed for each cluster
142 separately. The p-values are 0.040 and 0.032 for cluster 1 and cluster 2, respectively, indicating these two clusters have

143 significantly different sizes between pre- and post-menopause conditions after accounting for inter-patient variability. Sizes are
144 not significantly different for other clusters. The R code to reproduce the differential abundance analysis is provided in the files
145 NormEpi.R and NormTotal.R (Table 2).

146 Copy number variation analysis

147 Copy number variation (CNV) analysis was performed using inferCNV of the Trinity CTAT Project ([https://github.com/
148 broadinstitute/inferCNV](https://github.com/broadinstitute/inferCNV)), which compares gene expression intensity across genomic locations in the tumor or lymph-node
149 samples with those in a normal reference sample. The single-cell RNA expression profile of a normal breast total cells sample
150 (N-0372-total) was adopted as a reference for all the CNV analyses presented in the ScBrAtlas study. The results of each CNV
151 analysis were visualized in a heatmap, which showed the relative expression intensities of the tumor samples with respect
152 to the normal reference. For ease of visualization, cells from the same patient within the same cluster were grouped into a
153 single column block, and only the blocks containing more than 100 cells were used in the heatmap. All the column blocks
154 were assigned an equal width in each of the heatmap. The column block annotation of all the CNV heatmaps in this study is
155 available as part of the Figshare deposition, indicating which clusters in which samples were classified as normal or tumor⁹.

156 Data Records

157 Cell Ranger genewise read counts for the 69 scRNA-seq profiles, prior to quality filtering, are available as GEO series
158 GSE161529⁷. Quality filtering thresholds, downstream R data objects storing cell cluster identities and associated R code are
159 available from Figshare⁹. Specific files available from Figshare are listed in Table 2.

160 The bulk RNA-seq genewise read counts are available as GEO series GSE161892²⁴. The cell-type signature genes generated
161 from the bulk RNA-seq and associated R code are available from Figshare⁹.

162 Technical Validation

163 Technical quality of the 10X single-cell transcriptomic datasets was assessed by examining the number of mapped reads and
164 the number of detected genes (genes with at least one read count mapped to it) for all cells across all the samples (Fig. 2a-b).

165 Quality control was performed to remove cells of low quality. Cells with a high proportion of mitochondrial reads or a low
166 number of detected genes were removed. For each sample, an upper limit of library size was also used in combination with an
167 upper limit of number of detected genes to remove potential multiplets. The proportion of cells retained after filtering is 82.2%
168 across all 69 samples, indicating good data quality (Fig. 2c).

169 Technical quality of the bulk RNA-seq data was assessed using MDS and BCV plots (Figure 3).

170 Usage Notes

171 The code provided may be run using the free R programming environment with Bioconductor and Seurat R software packages
172 <https://www.r-project.org>. The RDS files may be read using R's readRDS() function. The Seurat objects allow readers to use
173 and extend the results of the major analyses conducted as part of the ScBrAtlas study. Cell barcodes and Seurat cell clustering
174 information are stored in the meta.data component of each Seurat object.

175 Code availability

176 The R code files provided on Figshare contain complete code and input files for reproducing the analyses of the BrScAtlas
177 study⁹ (Table 2). All the bioinformatics analyses were performed in R 3.6.1 on x86_64-pc-linux-gnu (64-bit) platform, running
178 under CentOS Linux 7. The following software packages were used for the analyses: Seurat v3.1.1, limma v3.40.6, edgeR
179 v3.26.8, pheatmap v1.0.12, ggplot2 v3.2.1, org.Hs.eg.db v3.8.2 and vcd v1.4-5.

180 References

- 181 1. Bray, F. *et al.* Global cancer statistics 2018: GLOBOCAN estimates of incidence and mortality worldwide for 36 cancers
182 in 185 countries. *CA: A Cancer J. for Clin.* **68**, 394–424 (2018).
- 183 2. Visvader, J. E. Keeping abreast of the mammary epithelial hierarchy and breast tumorigenesis. *Genes & Dev.* **23**, 2563–2577
184 (2009).
- 185 3. Sotiriou, C. *et al.* Breast cancer classification and prognosis based on gene expression profiles from a population-based
186 study. *Proc. Natl. Acad. Sci.* **100**, 10393–10398 (2003).
- 187 4. Turner, N. C. & Reis-Filho, J. S. Basal-like breast cancer and the BRCA1 phenotype. *Oncogene* **25**, 5846–5853 (2006).

- 188 5. Pal, B. *et al.* A single-cell RNA expression atlas of normal, preneoplastic and tumorigenic states in the human breast.
189 *EMBO J.* **40**, e107333 (2021).
- 190 6. Lim, E. *et al.* Aberrant luminal progenitors as the candidate target population for basal tumor development in BRCA1
191 mutation carriers. *Nat. Medicine* **15**, 907–913 (2009).
- 192 7. Smyth, G. K., Chen, Y. & Visvader, J. E. scRNA-seq profiling of breast cancer tumors, BRCA1 mutant pre-neoplastic
193 mammary gland cells and normal mammary gland cells. GEO series [https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=
194 GSE161529](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE161529) (2021).
- 195 8. Ilicic, T. *et al.* Classification of low quality cells from single-cell RNA-seq data. *Genome Biol.* **17**, 1–15 (2016).
- 196 9. Chen, Y. & Smyth, G. K. Data, R code and output Seurat objects for single cell RNA-seq analysis of human breast tissues.
197 figshare <https://figshare.com/s/c584dda937d346cc9a80> (2021).
- 198 10. Stuart, T. *et al.* Comprehensive integration of single-cell data. *Cell* **177**, 1888–1902 (2019).
- 199 11. Van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. *J. Mach. Learn. Res.* **9** (2008).
- 200 12. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of communities in large networks. *J. Stat.
201 Mech. Theory Exp.* **2008**, P10008 (2008).
- 202 13. Chen, Y., Lun, A. T. L. & Smyth, G. K. From reads to genes to pathways: differential expression analysis of RNA-seq
203 experiments using Rsubread and the edgeR quasi-likelihood pipeline. *F1000Research* **5**, 1438 (2016).
- 204 14. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic
205 Acids Res.* **43**, e47 (2015).
- 206 15. Pal, B. *et al.* Construction of developmental lineage relationships in the mouse mammary gland by single-cell RNA
207 profiling. *Nat. Commun.* **8**, 1–14 (2017).
- 208 16. Meyer, D., Zeileis, A. & Hornik, K. vcd: Visualizing categorical data. R package available from [https://cran.r-project.org/
209 package=vcd](https://cran.r-project.org/package=vcd) (2008).
- 210 17. Angerer, P. *et al.* destiny: diffusion maps for large-scale single-cell data in R. *Bioinformatics* **32**, 1241–1243 (2016).
- 211 18. Ward Jr, J. H. Hierarchical grouping to optimize an objective function. *J. Am. Stat. Assoc.* **58**, 236–244 (1963).
- 212 19. Liao, Y., Smyth, G. K. & Shi, W. The R package Rsubread is easier, faster, cheaper and better for alignment and
213 quantification of RNA sequencing reads. *Nucleic Acids Res.* **47**, e47 (2019).
- 214 20. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general-purpose read summarization program. *Bioinformatics*
215 **30**, 923–930 (2014).
- 216 21. Law, C. W., Chen, Y., Shi, W. & Smyth, G. K. Voom: precision weights unlock linear model analysis tools for RNA-seq
217 read counts. *Genome Biol.* **15**, R29 (2014).
- 218 22. McCarthy, D. J. & Smyth, G. K. Testing significance relative to a fold-change threshold is a TREAT. *Bioinformatics* **25**,
219 765–771 (2009).
- 220 23. Chen, Y. HumanBreast10X. GitHub repository <https://github.com/yunshun/HumanBreast10X> (2021).
- 221 24. Smyth, G. K., Chen, Y., Pal, B. & Visvader, J. E. RNA-seq expression profiling of stromal and epithelial cell subpopulations
222 from human breast tissue. GEO series <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE161892> (2021).

223 Acknowledgements

224 This work was supported by the Chan Zuckerberg Initiative (EOSS4 grant number 2021-237445), the National Breast Cancer
225 Foundation (NBCF, IIRS-20-022), Australian National Health and Medical Research Council (NHMRC) grants (#1054618,
226 1100807,1113133, 1153049); NHMRC IRIISS; the Victorian State Government Operational Infrastructure Support; the
227 Australian Cancer Research Foundation and the Ian Potter Foundation. G.J.L., G.K.S. and J.E.V. were supported by NHMRC
228 Fellowships (G.J.L. #1078730 and 1175960; G.K.S. #1058892; J.E.V. #1037230 and 1102742); Y.C. was supported by Medical
229 Research Future Fund (MRFF) Investigator Grant (#1176199).

230 Author contributions statement

231 Y.C. and G.K.S performed bioinformatic analyses, deposited analysis code and data objects, and wrote the article; B.P., J.E.V
232 and G.J.L designed the human breast study and collected data. All authors reviewed the manuscript.

233 **Competing interests**

234 The authors declare no competing interests.

235 **Figures & Tables**

Label	Tissue Sample Type	Cell Family	Figure
NormEpi	Normal breast	epithelial cells	EV1C
NormEpiSub	Normal breast	epithelial cells without stroma	1E
NormTotal	Normal breast	total cells	2B
NormTotalSub	Normal breast	non-epithelial	2D
NormTotalFib	Normal breast	fibroblast cells	3D
NormB1Total	Normal and BRCA1 preneoplastic	total cells	4B
NormB1TotalSub	Normal and BRCA1 preneoplastic	non-epithelial	4C
BRCA1Tum	BRCA1 preneoplastic and BRCA1 TNBC	total cells	4E
BRCA1TumSub	BRCA1 preneoplastic and BRCA1 TNBC	non-epithelial	5A
TNBC	TNBC	total cells	6A
HER2	HER2+ breast tumor	total cells	6B
ERTotal	ER+ breast tumor	total cells	6C
ERTotalTum	ER+ breast tumor	epithelial cells	6E
PairedER	Two ER+ breast tumors from patient 0029	total cells	6H
TNBCSub	TNBC	non-epithelial	7A
HER2Sub	HER2+ breast tumor	non-epithelial	7B
ERTotalSub	ER+ breast tumor	non-epithelial	7C
TNBCTum	TNBC	epithelial cells	EV3B (top)
HER2Tum	HER2+ breast tumor	epithelial cells	EV3B (bottom)
TNBCTC	TNBC	T-cells	EV4A (left)
HER2TC	HER2+ breast tumor	T-cells	EV4A (middle)
ERTotalTC	ER+ breast tumor	T-cells	EV4A (right)
Male	ER+ breast tumors from male patients	total cells	EV5A
TumLN	ER+ breast tumor & lymph-node pairs from 7 patients	total cells	9A

Table 1. Cell cluster analyses. Each row corresponds to one integration and cell clustering, except for TumLN, where one clustering was done for each of the 7 patients. Columns indicate the group of samples integrated, the cell subset clustered and the figure reference in the original ScBrAtlas study⁵.

Label	Data filename	Code filename
NormEpi	SeuratObject_NormEpi.rds	NormEpi.R
NormEpiSub	SeuratObject_NormEpiSub.rds	NormEpi.R
NormTotal	SeuratObject_NormTotal.rds	NormTotal.R
NormTotalSub	SeuratObject_NormTotalSub.rds	NormTotal.R
NormTotalFib	SeuratObject_NormTotalFib.rds	NormTotal.R
NormB1Total	SeuratObject_NormB1Total.rds	NormBRCA1.R
NormB1TotalSub	SeuratObject_NormB1TotalSub.rds	NormBRCA1.R
BRCA1Tum	SeuratObject_BRCA1Tum.rds	BRCA1Tum.R
BRCA1TumSub	SeuratObject_BRCA1TumSub.rds	BRCA1Tum.R
TNBC	SeuratObject_TNBC.rds	TNBC.R
TNBCSub	SeuratObject_TNBCSub.rds	TNBC.R
TNBCTC	SeuratObject_TNBCTC.rds	TNBC.R
TNBCtum	SeuratObject_TNBCtum.rds	TNBC.R
HER2	SeuratObject_HER2.rds	HER2.R
HER2Sub	SeuratObject_HER2Sub.rds	HER2.R
HER2TC	SeuratObject_HER2TC.rds	HER2.R
HER2tum	SeuratObject_HER2tum.rds	HER2.R
ERTotal	SeuratObject_ERTotal.rds	ER.R
ERTotalSub	SeuratObject_ERTotalSub.rds	ER.R
ERTotalTC	SeuratObject_ERTotalTC.rds	ER.R
ERTotaltum	SeuratObject_ERTotaltum.rds	ER.R
Male	SeuratObject_Male.rds	Male.R
PairedER	SeuratObject_PairedER.rds	PairedER.R
TumLN	SeuratObject_TumLN.rds	TumLN.R

Table 2. Files deposited on Figshare⁹. Data files are in RDS format. Each data file contains one Seurat object except for TumLN, which contains a list of 7 Seurat objects. Each Seurat data object provides cell cluster identities and associated information for the corresponding cell cluster analysis. Code files contain the R code used to produce the corresponding Seurat objects.

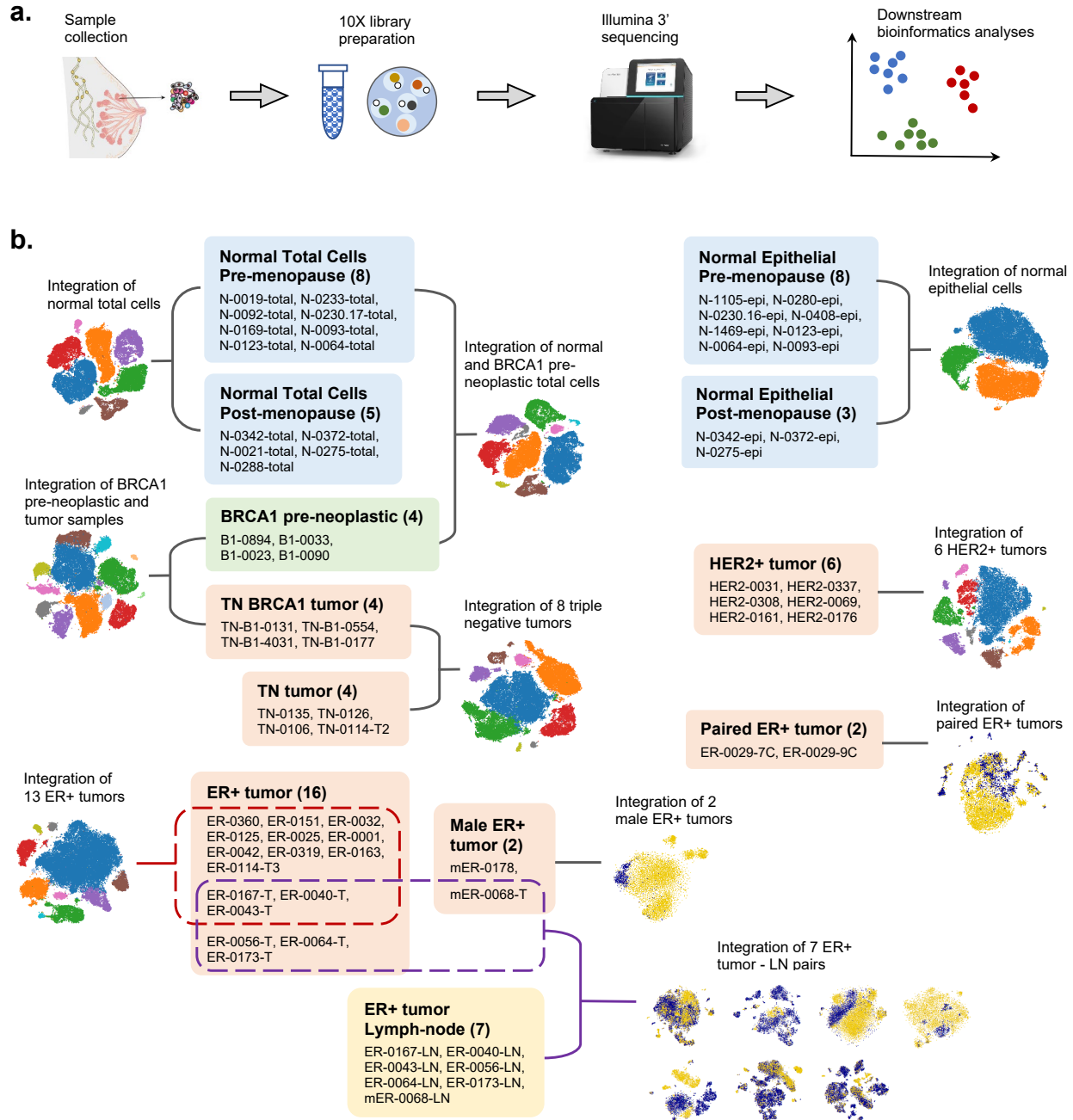


Figure 1. (a) Diagram showing the data processing pipeline from sample collection to downstream bioinformatics analyses. (b) Schematic overview of the all the integration analyses and the samples involved in each integration analysis. Under each category, the names of the samples are listed and the total number of samples is shown in the bracket.

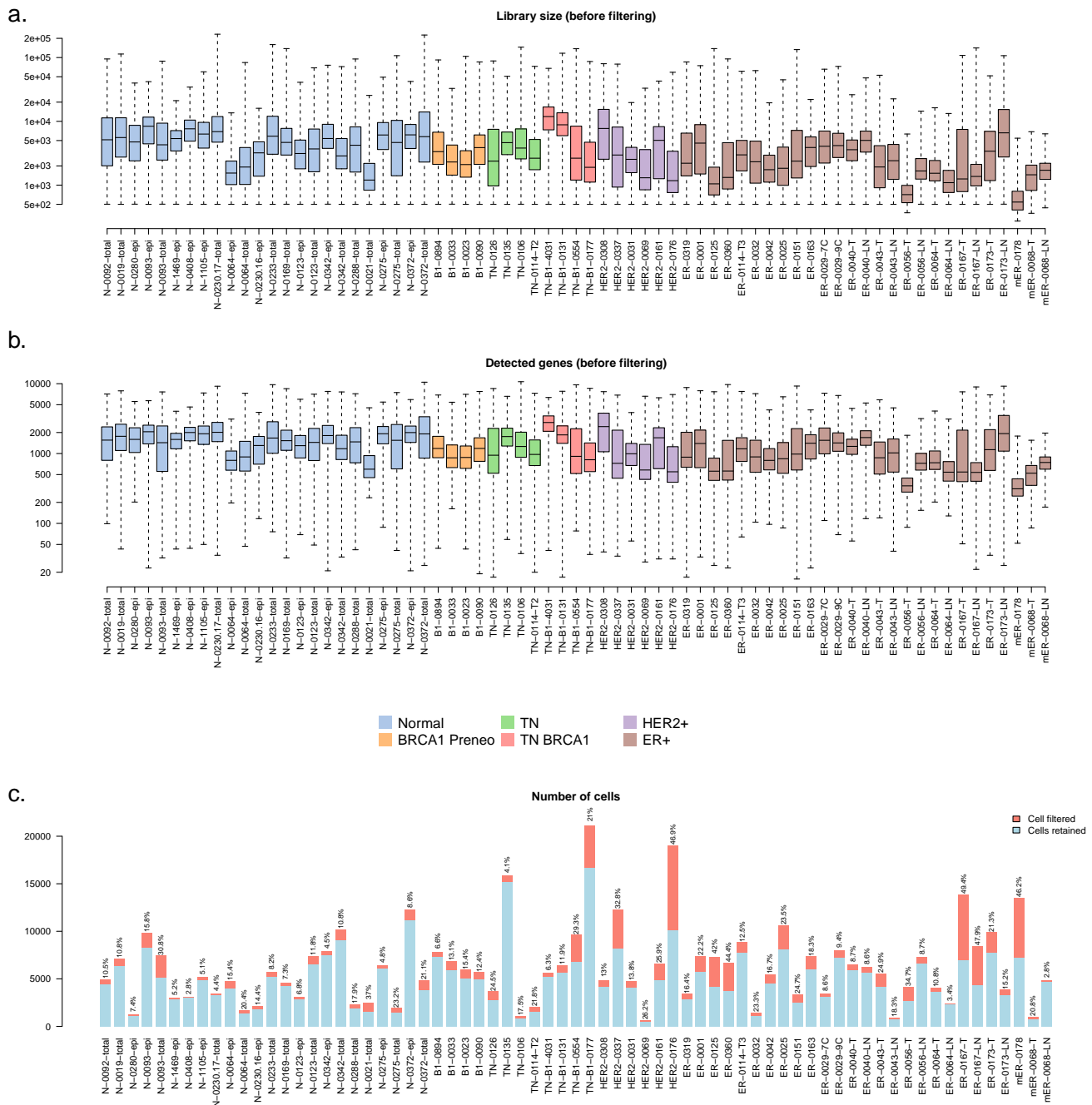


Figure 2. Box plots of (a) the library sizes and (b) the numbers of detected genes for all the cells in each of the 69 samples before filtering. Boxes are coloured by tumor type. (c) Bar plots of the number of cells in each of the 69 samples. The blue segments show the number of cells that are kept after the cell filtering while the red segments show the filtered cells. The proportions of filtered cells are labelled on top of the bars for all 69 samples.

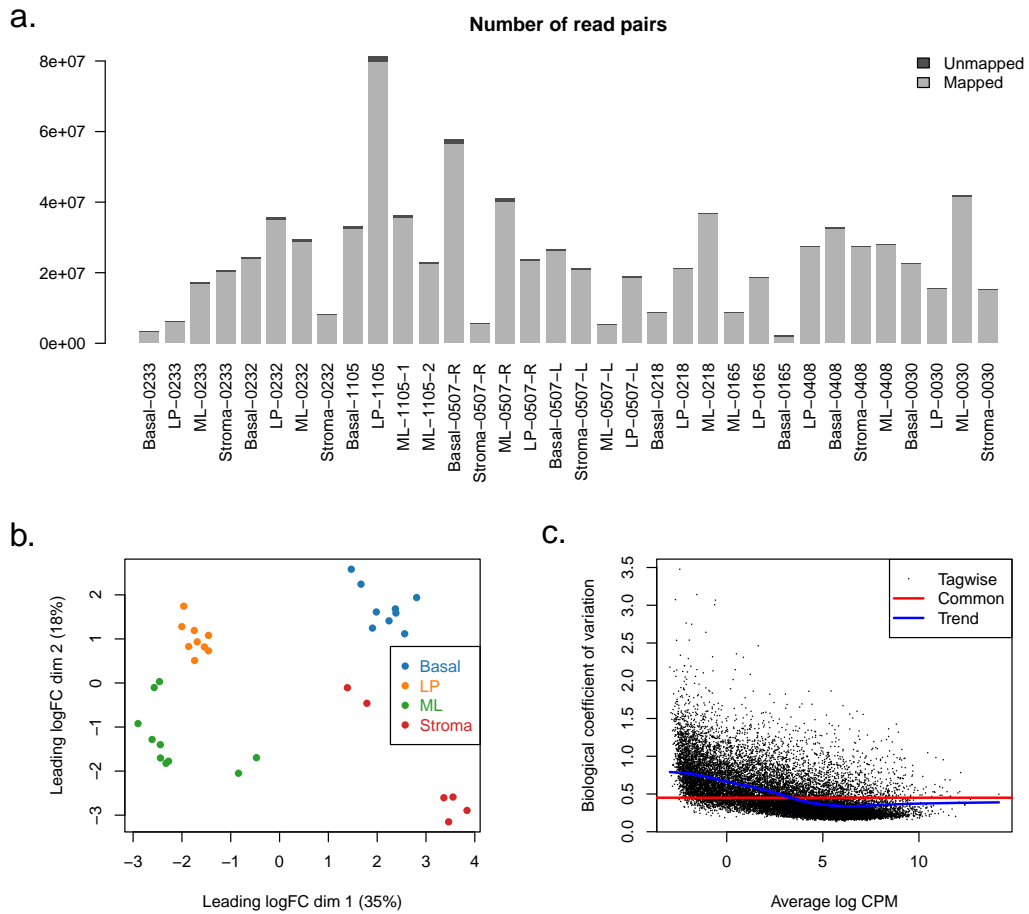


Figure 3. (a) Bar plots of the numbers of read pairs in the human mammary gland bulk RNA-seq samples. The light grey segments represent the mapped read pairs whereas the dark grey segments represent the unmapped ones. (b) MDS plot of all the bulk RNA-seq samples. (c) BCV plot of the bulk RNA-seq data.