# 1    scMuffin: an R package for resolving solid tumor

# 2    heterogeneity from single-cell expression data

3    Valentina Nale[1]*, Noemi Di Nanni[1]*, Alice Chiodi[1], Ingrid Cifola[1], Marco Moscatelli[1], Cinzia Cocola[1],

4    Matteo Gnocchi[1], Eleonora Piscitelli[1], Rolland Reinbold[1], Ileana Zucchi[1], Luciano Milanesi[1],

5    Alessandra Mezzelani[1], Paride Pelucchi[1]*+, and Ettore Mosca[1]*+

6

7    [1]National Research Council, Institute of Biomedical Technologies, Via Fratelli Cervi 93, 20054,

8    Segrate (Milan), Italy

9

10    e-mail addresses:

11    - ettore.mosca@itb.cnr.it

12    - paride.pelucchi@itb.cnr.it

13    *Equal contribution

14    +Corresponding author

15

16

17

18

19

20

21

## Abstract

INTRODUCTION: Single-cell (SC) gene expression analysis is crucial to dissect the complex cellular heterogeneity of solid tumors, which is one of the main obstacles for the development of effective cancer treatments. Such tumors typically contain a mixture of cells with aberrant genomic and transcriptomic profiles affecting specific sub-populations that might have a pivotal role in cancer progression, whose identification eludes bulk RNA-sequencing approaches. We present scMuffin, an R package that enables the characterization of cell identity in solid tumors on the basis of multiple and complementary criteria applied on SC gene expression data.

RESULTS: scMuffin provides a series of functions to calculate several different qualitative and quantitative scores, such as: expression of marker sets for normal and tumor conditions, pathway activity, cell state trajectories, CNVs, chromatin state and proliferation state. Thus, scMuffin facilitates the combination of various evidences that can be used to distinguish normal and tumoral cells, define cell identities, cluster cells in different ways, link genomic aberrations to phenotypes and identify subtle differences between cell subtypes or cell states. As a proof-of-concept, we applied scMuffin to a public SC expression dataset of human high-grade gliomas, where we found that some chromosomal amplifications might underlie the invasive tumor phenotype and identified rare quiescent cells that may deserve further investigations as candidate cancer stem cells.

CONCLUSIONS: The analyses offered by scMuffin and the results achieved in the case study show that our tool helps addressing the main challenges in the bioinformatics analysis of SC expression data from solid tumors.

**Keywords**: single-cell transcriptomics, cancer, tumor heterogeneity, cell identity.

## 1. Background

45   Single-cell (SC) gene expression analysis is crucial to dissect the complex cellular

46 heterogeneity of solid tumors, which is one of the main obstacles for the development of effective

47 cancer treatments (1). A relevant number of software tools has been developed in recent years in

48 the field of SC data analysis (2), a fact that stresses the key opportunities and challenges in this

49 field. A recent study has shown that the development of tools that address common tasks (e.g.

50 clustering of similar cells) and ordering of cells (e.g. definition of cell trajectories) is decreasing,

51 while a greater focus is being paid on data integration and classification (2). These observations

52 reflect the growing availability, scale and complexity of SC datasets (2).

53   SC datasets of solid tumors are typical examples of complex datasets that present a series

54 of computational challenges and whose analysis demands domain-specific and integrative

55 approaches. In fact, solid tumors typically contain a mixture of cells with aberrant genomic and

56 transcriptomic profiles affecting specific sub-populations that might play a pivotal role in cancer

57 progression, whose identification eludes bulk RNA-sequencing approaches. The use of cell type-

58 specific markers (when available) is limited, and the alterations of gene expression that mark

59 cancer cells makes the use of markers for normal cells not completely adequate. Moreover, the

60 molecular heterogeneity of cancer cells (due to both intra-tumor and inter-individual differences)

61 poses intrinsic limits in the definition of such markers. In addition, solid tumor samples typically

62 comprise cells from the surrounding tissue or infiltrating cells that need to be distinguished from

63 tumor populations for an effective analysis. Another challenge is the identification of clinically

64 relevant cell subtypes that may be rare in the tumor mass, such as cancer stem cells or drug

65 resistant subclones: because of their relatively low number, these cells are typically clustered

66 together with many others. Lastly, an intrinsic problem of many SC datasets is the sequencing

67 depth limit at the SC level. These limitations bound the number of detectable genes to the few

68    thousands of the highest expressed genes, which implies, for example, that some established

69    markers may not be used for data analysis.

70        To address these challenges, we developed scMuffin, an R package that implements a

71    series of complementary analyses aimed at shedding light on the complexity of solid tumor SC

72    expression data, including: a fast and customizable gene set scoring system; gene sets from

73    various sources, including pathways, cancer functional states and cell markers; cell cluster

74    association with quantitative (e.g. gene set scores) as well as categorical (e.g. mutation states,

75    proliferation states) features; copy number variation (CNV) analysis; chromatin state analysis;

76    proliferation rate quantification; and marker-based two-sample comparisons **(Figure 1)**. scMuffin

77    facilitates the integrative analysis of these multiple features, thus allowing the identification of cell

78    subtypes that elude more general clustering and classification approaches.

## 2.    Implementation

80        scMuffin is implemented in R and provides a series of functions that allows the user to

81    perform various tasks, which can be combined to obtain various data analysis pipelines. The

82    package includes a vignette that describes the use of the tool, and every function is documented.

83    The results from the various analyses (e.g., gene set scores at SC level and cell chromatin state)

84    can be organized in dedicated (simple) objects in order to enable subsequent analyses (e.g.,

85    assessment of associations between features and cell clusters) that jointly consider multiple cell

86    features and various ways of cell clustering. Computationally intensive tasks (in particular, gene

87    set scoring and CNV inference) are parallelized. In this section, we describe the algorithms used to

88    perform the several tasks offered by scMuffin.

89    **2.1 Quantification of gene set expression scores at cell and cluster levels**

90　　　The quantification of gene set expression scores follows the approach described in (8,9), in

91　　which a gene set is scored on the basis of its average deviation from an empirical null. In scMuffin

92　　the implementation of gene set scoring is parallel and the calculation can be tuned acting on a

93　　series of parameters, such as: the number of bins, the number of minimum genes that must have

94　　non negative values in a cell ("nmark_min"), the minimum number of cells in which at least

95　　nmark_min have to be found ("ncell_min"), the number of permutations ($k$), the minimum

96　　number of permutations required ($k'$). This tuning helps to address the issue of missing values -

97　　typical of SC datasets - and therefore maximizes the number of gene sets for which it is actually

98　　possible to obtain an expression score supported by an empirical null distribution. Briefly, given a

99　　gene set $S$:

100　　　1.　the genes occurring in the genes-by-cells matrix are grouped into a number of bins

101　　　　　according to their average expression across cells;

102　　　2.　a number $k$ of random gene sets $S_i^*$ are created, of the same size of $S$, tossing genes from

103　　　　　the same bins of $S$, in order to match the distribution of gene expression of each $S_i^*$ with

104　　　　　that of $S$;

105　　　3.　the averages $m_c$ and $m_{ic}^*$ are calculated, respectively, over the values of $S$ and $S_i^*$ in every

106　　　　　cell $c$;

107　　　4.　the expression score $Y_c$ is calculated as the average difference between $m_c$ and $m_{ic}^*$;

108　　　5.　the average value $\bar{Y_c}$, calculated over the $Y_c$ of a given cluster, is used as the representative

109　　　　　score of $S$ in that cluster.

**2.2 CNV estimation by adjacent gene windows approach**

111　　　CNV inference in scMuffin is based on the "adjacent gene windows" approach, which has

112　　been validated using both single nucleotide polymorphism arrays (10) and whole-exome

113　　sequencing (8) technologies. The approach is implemented in parallel and offers various

114    parameter tuning and data filtering possibilities, which allows the investigator to optimize the

115    analysis on the characteristics of its dataset. The CNV profile of each cell is calculated as a moving

116    average of scaled gene expression levels ordered by genomic location, with the possibility of

117    subtracting a normal reference profile to identify sample-specific CNVs. The main steps are:

118        1. the reference cells are added to the genes-by-cells matrix (optional);

119        2. the expression of each gene is scaled subtracting its average (optional);

120        3. the gene expression matrix is ordered by chromosome and gene location;

121        4. in each chromosome $h$, the estimated copy number $V_{ic}$ of cell $c$ is calculated for all the

122            ordered genes $i \in \left[\frac{w}{2} + 1, n_h - \frac{w}{2}\right]$:

$$V_{ic} = \sum_{j=i-\frac{w}{2}}^{i+\frac{w}{2}} \frac{e_{jc}}{w+1}$$

123            where $w$ is an even number that defines the window size, that is, the number of genes

124            located before and after gene $i$ which contribute to the estimation of $V_{ic}$, and $e_{jc}$ is the

125            gene expression value;

126        5. $V_{ic}$ values are scaled subtracting their average in a cell (optional);

127        6. cells are clustered by their CNV profile;

128        7. the average CNV profile of the normal reference cells is subtracted from all the CNV

129            profiles (optional).

### 2.3 Chromatin state and proliferation rate

131    The chromatin state $R_c$ of a cell is inferred on the basis of the number of expressed genes over the

132    number of total mapped reads:

$$R_c = \frac{\#\{ic \geq \alpha\}}{\sum_i ic}$$

6

133    where $\alpha$ is a threshold over the gene count $_{ic}$, which defines the gene $i$ as "expressed".  High

134    values of $R_c$ indicate cells that are expressing many genes in relation to the number of mapped

135    reads.

136         The proliferation rate $P_c$ of a cell is quantified as the maximum value between the two

137    gene set scores $Y_c(S_1)$ and $Y_c(S_2)$, calculated on the gene sets $S_1$ and $S_2$ that, respectively,

138    characterize the G1/S and G2/M cell cycle phases:

$$P_c = \max\big(Y_c(S_1), Y_c(S_2)\big)$$

139    where $S_1$ and $S_2$ are defined as in Tirosh *et al.* (8).

140    **2.4 Cluster enrichment analysis for quantitative and categorical features**

141    The assessment of cluster enrichment in high values of quantitative features is computed using a

142    procedure that we name "cell set enrichment analysis" (CSEA), because it is analogous to the gene

143    set enrichment analysis (GSEA) (7), but operates on different input types. In particular, instead of a

144    ranked list of genes, the CSEA considers a list of cells ranked by a feature of interest, and instead

145    of testing a gene set, CSEA tests a cell set (i.e., a cluster of cells). Therefore, CSEA tests whether

146    the cells assigned to a cluster are located at the top (or bottom) of the ranked list of cells. The

147    assessment of a cluster enrichment in a particular value of a categorical feature is computed using

148    the over-representation analysis (ORA) approach (11), which is based on the hypergeometric test.

149         Both CSEA and ORA are implemented in parallel in scMuffin. This is particularly important

150    for CSEA, which uses permutations to build an empirical null distribution. Nonetheless, it is also

151    effective for ORAs that are run over a large number of gene sets.

152    **2.5 Cell clustering**

153    Cell clustering is based on the approaches implemented in the R package Seurat (3). The results

154    from multiple clustering procedures are compared by calculating the overlap coefficients among

155  all-pairs of clusters. Given two partitions $A$ and $B$, defined as sets of cell clusters $A = \{a_i\}$ and

156  $B = \{b_j\}$, the similarity between the cell clusters $a_i$ and $b_j$ is calculated as:

$$o_{ij} = \frac{|a_i \cap b_j|}{\min(|a_i|, |b_j|)}$$

157  Meta-clusters are defined as the union of cell clusters that have high $o_{ij}$ and are found by

158  hierarchical clustering of the matrix $O = \{o_{ij}\}$.

### 2.6 Cluster marker-based two-sample comparison

160  The cluster marker-based two-sample comparison is based on assessing the expression of cluster

161  markers of every cluster of a sample ($A$) in every cluster of the other sample ($B$) and *vice versa*.

162  Given a set of markers $S_{a_i}$, which represents the cell cluster $a_i$ of sample $A$, the gene set score

163  $\bar{Y}_{b_j}(S_{a_i})$ quantifies the expression of $S_{a_i}$ in the cell cluster $b_j$ of sample $B$, while *vice versa* $\bar{Y}_{a_i}(S_{b_j})$

164  quantifies the expression of $S_{b_j}$ in the cell cluster $a_i$ of sample $A$. The average value $\bar{Y}_{a_i b_j} =$

165  $\frac{\bar{Y}_{b_j}(S_{a_i}) + \bar{Y}_{a_i}(S_{b_j})}{2}$ quantifies the similarity between cell clusters $a_i$ and $b_j$ on the basis of the

166  expression of their markers. The procedure is repeated for all-pairs of clusters of sample $A$ and

167  sample $B$.

## 3.  Results and Discussion

169  In this section, we present the user interface **(Table 1)**, and the results obtained using our

170  scMuffin package on the SC dataset generated by Yuan *et al.* (12) from human high-grade glioma

171  (HGG) samples, and available on the Gene Expression Omnibus (GEO) repository (GSE103224) (13).

**Table 1. Main tasks and corresponding functions in scMuffin.**

| Task | Description | User interface |
|---|---|---|
| Gene set expression scoring | Average gene set expression deviation from matched empirical reference; provided gene sets from CellMarker (4), PanglaoDB (5), CancerSEA (6) and MSigDB (7) | • `prepare_gsls`<br>• `calculate_gs_scores`<br>• `calculate_gs_scores_in_clusters` |

| Copy number variations | Estimation of CNVs by means of the "moving window" approach, that is, considering the expression of adjacent genes; calculation of CNV deviation from a normal reference profile; processing of normal tissue-specific expression data from GTEx | • `calculate_CNV`<br>• `cluster_by_features`<br>• `apply_CNV_reference`<br>• `CNV_heatmap`<br>• `process_GTEx_gene_reads` |
|---|---|---|
| Chromatin state | Number of expressed genes in relation to the total reads | • `exp_rate` |
| Proliferation | Maximum between G1/S and G2/M gene set scores | • `proliferation_analysis` |
| Cell state trajectory | Diffusion map computation | • `diff_map` |
| Cell cluster annotation | Assessment of cluster enrichment for quantitative and categorical features | • `assess_cluster_enrichment` |
| Two-sample comparison | Quantification of the expression similarity between all-pairs of clusters between two samples | • `quantify_samples_similarity` |
| Assembling cell features and cell partitions | Objects that host cell-level feature values and cell partitions | • `create_features_obj`<br>• `create_partitions_obj` |
| Visualization | Automated UMAP visualizations for multiple features, heatmaps, box plots and dot plots | • `boxplot_cluster`<br>• `dotplot_cluster`<br>• `quantify_samples_similarity`<br>• `heatmap_CNV`<br>• `plot_umap_colored_features`<br>• `plot_heatmap_features_by_clusters` |

172

## 3.1 Gene set scoring

174 scMuffin provides functions to set up one or more gene set collections and perform SC-

175 level estimation of gene set expression scores in relation to an empirical null model (see

176 Implementation section). This can be applied to any gene set and can therefore be used to

177 estimate several different cell phenotypes, like pathway activities or marker set expression.

178 The function `prepare_gsls` allows the user to collect gene sets of cell types, pathways,

179 cancer functional states, as well as other collections of gene sets (e.g. positional gene sets,

180 hallmarks) from CellMarker (4), PanglaoDB (5), CancerSEA (6) and MSigDB (7) databases. Unlike

181 many existing tools that are used to perform marker-based cell annotation (14), the availability of

182 these gene sets within scMuffin package spares the user the effort of data collection and

183 harmonization. The function, which also accepts any user-given gene sets, applies a series of

184  criteria (e.g., minimum and maximum number of genes in a gene set) to filter the chosen gene

185  sets.

186      The cell-level expression scores for these gene sets can be calculated using

187  `calculate_gs_scores`, which requires the expression matrix, the gene sets, as well as a

188  series of optional parameters to fine-tune the calculation in relation to the characteristics of the

189  SC dataset under analysis. This tuning is important to address the heterogeneity of size and

190  sparsity that characterizes different SC datasets, attributable to both the biological specimen

191  under analysis and the SC platform used. For instance, the following code shows how to quantify

192  cell-level, and then cluster-level, expression score of the cancer functional states from CancerSEA

193  ("SIG_CancerSEA") (6) in a normalized genes-by-cells ("gbc") expression matrix:

```
194  gsc <- prepare_gsls(gs_sources = "SIG_CancerSEA", genes = rownames(gbc))
195  gs_scores_obj   <-   calculate_gs_scores(genes_by_cells   =   gbc,   gs_list   =
196  gsc$SIG_CancerSEA)
197  res_sig_cl   <-   calculate_gs_scores_in_clusters(gs_scores_obj   =   gs_scores_obj,
198  cell_clusters = cell_clusters)
```

199      The cell-level value of any gene set (and more generally of any feature) can be visualized

200  over the UMAP by means of `plot_umap_colored_features`, while cluster-level values of

201  multiple gene sets (features) can be visualized as a heatmap using

202  `plot_heatmap_features_by_clusters`, which relies on the ComplexHeatmap R package

203  (15). For example, in our case study, the analysis of the CancerSEA functional states in the HGG

204  sample PJ016 showed that the two groups of cell clusters (**Figure 2a,** left and right of the UMAP)

205  reflect distinct functional states **(Figure 2b)**: for example, the expression of the CancerSEA

206  "Invasion" markers was particularly high in cell clusters 0 and 9 as compared to all the other

207  clusters **(Figure 2b-c)**.

208  **3.2 CNV estimation and association with CancerSEA functional states**

10

209    CNV inference from SC expression data estimates the presence of relevant genomic

210    aberrations (amplifications and deletions) based on the expression of adjacent genes. This

211    knowledge offers crucial clues to address the difficult task of distinguishing normal from malignant

212    cells, and provides quantitative information to reconstruct the tumor clonal substructure.

213    Moreover, CNV pattern allows the investigator to hypothesize link between genomic alterations

214    and cell phenotypes.

215    The function `calculate_CNV` basically retrieves the genomic locations and performs the

216    CNV estimation; `cluster_by_features` executes the cell clustering based on CNV profiles;

217    `apply_CNV_reference` redefines the CNV values on the basis of normal reference cells; the

218    dedicated plotting function `CNV_heatmap` handles the visualization, where the cell cluster that

219    contains the reference is marked in red. Here is an example that illustrates CNV inference using a

220    100 genes window size and a normal reference profile from The Genotype-Tissue Expression

221    (GTEx) portal (16):

222    cnv_res    <-    calculate_CNV(gbc,    wnd_size    =    100,    reference    =

223    GTEx_mean)

224    cnv_clustering <- cluster_by_features(cnv_res, cnv=TRUE)

225    cnv_res_ref <- apply_CNV_reference(cnv = cnv_res, cnv_clustering =

226    cnv_clustering, reference="reference")

227    cnv_res_ref    <-    CNV_heatmap(cnv    =    cnv_res,    cnv_clustering    =

228    cnv_clustering, reference="reference")

229    To illustrate the use of this workflow, we selected two different HGG samples by Yuan *et al.*

230    (12), that is, PJ030, composed by tumor cells as well as not transformed cells  and PJ016, including

231    only transformed cells. We observed that the reference profile (obtained using the average gene

232    expression values of the normal brain samples available in GTEx portal) falls into cluster 3 of PJ030

233    **(Figure 3a-b)**. With respect to cluster 3 (corresponding to the not transformed cells included in this

234    sample), the clusters 0, 1 and 2 showed large recurrent aneuploidies, some of which are typical of

235    HGG, like the amplification of chromosome 7 **(Figure 3a)**. The CNV pattern here inferred is

236    fundamentally coherent with that reported by Yuan *et al.* (12), even if the authors just quantified a

237    summary value per chromosome, while scMuffin provides multiple CNV estimations per

238    chromosome. For sample PJ016, the CNV inference analysis highlighted two groups of "CNV

239    clusters" that map to two distinct components of the UMAP, while it did not identify a diploid

240    cluster, accordingly to the presence of only transformed cells in this sample **(Figure 3c-d).**

241    Interestingly, clusters 1 and 3 were marked by peculiar amplifications in chromosomes 1p and

242    19p.

243        scMuffin enables the comparison of clusters obtained using different procedures. In

244    particular, the overlap among all-pairs of clusters can be quantified using:

245    `cl_list <- partitions_to_list(clust_obj)`

246    `ov_mat <- overlap_matrix(cl_list)`

247    In our case study, the comparison between expression clusters and CNV clusters of sample PJ016

248    confirmed the presence of two groups of cells: for example, CNV clusters 1 and 3 showed a

249    relevant overlap with expression clusters 0, 6, 8 and 9 **(Figure 3e).**

250        An example of integrative analysis enabled by scMuffin is the functional assessment of CNV

251    patterns. We quantified the expression scores of the CancerSEA functional states throughout the

252    CNV clusters of sample PJ016. As expected, the two aforementioned groups of CNV clusters (0-2-4

253    and 1-3) were characterized by different functional states **(Figure 3f)**, like the corresponding

254    expression clusters. In particular, CNV cluster 3 – which is mainly located in the top-left region of

255    the UMAP visualization **(Figure 3d)** and has a strong overlap with expression clusters 0 and 9

256    **(Figure 3e)** – contains cells that highly express the CancerSEA "Invasion" markers **(Figure 3f** and

257    **Figure 2)**. This finding suggests that the peculiar amplifications of chromosomes 1p and 19p found

258 in this cluster might be linked to the invasive phenotype. This hypothesis is supported by the

259 finding of two CancerSEA invasion markers, Y-Box Binding Protein 1 *(YBX1)* and Heterogeneous

260 Nuclear Ribonucleoprotein M *(HNRNPM)*, located within the amplified regions of chromosome 1p

261 and 19p specifically found in CNV clusters 1 and 3 **(Figure 3c)**. *YBX1* is a DNA/RNA-binding protein

262 and transcription factor which plays a central role in coordinating tumor invasion in glioblastoma

263 (17). *HNRNPM* belongs to a family of spliceosome auxiliary factors and is involved in the regulation

264 of splicing; the upregulation of these factors results in tumor-associated aberrant splicing, which

265 promotes glioma progression and malignancy (18,19). In particular, HNRNPM was identified as an

266 interactor of the DNA/RNA binding protein SON, which drives oncogenic RNA splicing in

267 glioblastoma (20). While it is beyond the scope of this article to further study this hypothesis,

268 these findings clearly highlight the usefulness of the integrative analysis of CNVs and CancerSEA

269 functional states provided by our scMuffin tool.

270 **3.3 Clustering, features and annotation**

271 scMuffin contains functions for assessing the association between cell clusters and

272 quantitative as well as categorical features, by means of CSEA and ORA, respectively. Here is the

273 user interface, where, firstly, the objects containing cell clusters and cell features are set up; then,

274 the enrichment is quantified for all partitions (various ways of clustering cells) and all features:

275 `clust_obj <- create_partitions_obj(cell_clusters)`

276 `feat_obj <- create_features_obj(feature_values)`

277 `cl_enrich   <-   assess_cluster_enrichment(features   =   feat_obj,`

278 `partitions = clust_obj)`

279 The results of CSEA and ORA can be extracted to produce features-by-clusters matrices

280 that contain any score calculated by CSEA or ORA, like, for example, normalized enrichment scores

281 (NES) values and enrichment ratios (er):

13

282  `cl_enrich_table`                                                    `<-`
283  `extract_cluster_enrichment_table(cl_enrich, q_type = "nes", c_type`
284  `= "er")`
285      The results of enrichment analysis can be visualized as box plots (quantitative features) and
286  dot plots (categorical features):
287  `top_feat_lab_CSEA   <-   boxplot_cluster(features   =   feat_obj,`
288  `partitions    =    clust_obj,    clustering_name    =    "global",`
289  `clust_enrich_res = cl_enrich, criterion = "fdr")`
290  `top_feat_lab_ORA   <-   dotplot_cluster(features   =   feat_obj,`
291  `partitions    =    clust_obj,    clustering_name    =    "global",`
292  `clust_enrich_res = cl_enrich, text_val = "p")`
293  These plots show, for each cluster, the distribution of values of the most significant features in the
294  cluster in comparison to all the other clusters, and the related scores (e.g., NES, p-value and FDR).
295  In addition, `boxplot_cluster` and `dotplot_cluster` provide the labels of the most
296  significant features associated with any cluster. These labels can be extracted from the enrichment
297  analysis results also by means of `extract_cluster_enrichment_tags`, according to
298  various criteria (e.g., NES, enrichment ratio, p-value, FDR) that are specific to CSEA or ORA.
299      To illustrate these functions, we assessed the enrichment of the expression clusters of
300  sample PJ016 in terms of both CancerSEA functional states (quantitative features) and three
301  categorical features, namely: cell clusters obtained analyzing ribosomal gene expression, a gene
302  set included in scMuffin because changes in ribosomal gene expression were associated with
303  specific cancer phenotype and can reveal specific malignant subpopulations (21–23); cell clusters
304  obtained using a glioblastoma signature of 500 genes (24) whose expression can be used to
305  classify glioblastoma subtypes; cell cycle phase. Considering as an example the cluster 0 of sample
306  PJ016, the analysis showed that it was significantly enriched in cells that, in comparison with cells
307  of other clusters, highly express the gene markers of CancerSEA "Invasion" state **(Figure 4a)** and

14

308     are in S and G2M phases **(Figure 4b).** The labels of the most significant features of any cluster can

309     be used, by means of `plot_umap`, to plot an annotated UMAP **(Figure 4c)** .

310     **3.4 Chromatin state, proliferation rate and cell state trajectories**

311         Chromatin state and proliferation rate carry two relevant pieces of information for the

312     characterization of a cancer cell.

313         In particular, an open chromatin state is peculiar of stem cells (and cancer stem cells

314     (CSCs)), might indicate de-differentiation processes of tumor progression and might influence cell

315     plasticity, favoring cancer cell adaptability and drug resistance (25,26). In a recent study on

316     glioblastoma, chromatin accessibility was associated to a specific subpopulation of putative tumor-

317     initiating CSCs with invasive phenotype and low survival prediction (27). The global state of the

318     chromatin at SC level can be inferred from SC transcriptomic data and provides a simple and useful

319     score that can be used to distinguish specific cell types, such as CSCs. The chromatin state can be

320     quantified by means of the function `exp_rate` on the genes-by-cells count matrix:

321     `res <- exp_rate(gbc, min_counts = 5)`

322     where 5 is the required threshold above which a gene is considered expressed.

323         The proliferation rate is a relevant indicator for distinguishing cell types in solid tumors and

324     helps to identify cells with potential clinical relevance and interest as candidate therapeutic

325     targets (28,29). In scMuffin, we quantify cell proliferation rate on the basis of the expression of

326     G1/S and G2/M genes:

327     `res <- proliferation_analysis(gbc)`

328         As a proof-of-concept, we show the joint analyses of chromatin state and proliferation rate

329     in sample PJ016 and visualize the results in the state space of cell differentiation trajectories. In

330     scMuffin, cell state trajectories are inferred using the "diffusion maps" approach available in the

331     destiny R package (30), by means of the wrapper function:

332 ```res <- diff_map(gbc)```

333 Interestingly, we observed that cells showing high values of chromatin state score - cells that are

334 expressing a relatively high number of genes (i.e., an open chromatin state) - are located at the

335 root of the trajectory state space **(Figure 5a)**, while the cells that show the highest proliferation

336 rates are located at a corner of the state space **(Figure 5b)**. This pattern suggests that the cells

337 with high values of chromatin state could be quiescent cells, which express a large number of

338 genes but are not actively dividing. Therefore, these cells are interesting candidates for further

339 analysis aimed at studying CSCs in HGG. More generally, this proof-of-concept demonstrates the

340 usefulness of the chromatin state score defined here, especially if used in combination with the

341 proliferation rate for the identification of particular cell types or cell states.

342 **3.5 Comparison of samples**

343 A SC dataset carries an extensive amount of information. The integration of multiple SC

344 datasets is a challenging task and multiple approaches have been proposed to address it (31).

345 Typically, the integrated datasets are computationally demanding due to their huge size. An

346 alternative possibility lies in cross-checking the expression of cluster markers between two

347 samples: the expression of the cluster markers of a sample is assessed in the other sample – and

348 *vice versa* – obtaining the similarities among all pairs of clusters. For example, Nguyen *et al.* (9)

349 used this approach to study the occurrence of the characteristic cell types of normal mammary

350 gland across samples collected from different subjects.

351 scMuffin provides a function to quantify the similarity between all-pairs of clusters of two

352 samples on the basis of cluster-specific markers:

353 ```sim_res <- quantify_samples_similarity(gbc_1, gbc_2, clusters_1,```
354 ```clusters_2, cluster_markers_1, cluster_markers_2)```

355         Concerning our case study, the comparison of samples PJ016 and PJ018 showed a series of

356   similarities between their clusters. For instance, the clusters 0 and 9 of sample PJ1016 are

357   composed of cells highly similar to those grouped into clusters 2 and 5 of sample PJ018 **(Figure**

358   **6a)**. This analysis revealed a pattern of cluster-cluster similarities that is fundamentally coherent

359   with the results obtained by performing the alternative approach of integrating the two datasets

360   and then clustering the cells **(Figure 6b-c)**. For example, both approaches showed that clusters 0

361   and 9 of PJ1016 are similar to clusters 2 and 5 of PJ018, and cluster 4 of PJ016 is close to cluster 7

362   of PJ018. There were also some differences, which, yet again, remark the challenge of this task: for

363   example, cluster 8 of PJ016 is similar to cluster 9 of PJ018 using the marker-based similarity

364   **(Figure 6a)**, while the UMAP obtained by the integrated dataset places cluster 8 of PJ016 close to

365   clusters 6 and 3 of PJ018 **(Figure 6b-c)**.

## 366   **4. Conclusions**

367   Here, we presented scMuffin, an R package that we developed to offer a series of useful functions

368   to perform and integrate multiple types of analyses on SC expression data. As a proof-of-concept,

369   we applied scMuffin on a publicly available SC expression dataset of human HGG. We described

370   two examples of integrative analyses which returned particularly interesting findings that would

371   deserve further investigations. The functional characterization of CNVs highlighted a possible link

372   between amplifications of chromosomes 1p and 19p and invasive tumor phenotype. The joint

373   analysis of chromatin state, proliferation rate and cell state trajectories suggested possible

374   candidates of CSCs in HGG. The analyses offered by scMuffin and the results achieved in this case

375   study show that scMuffin helps addressing the main challenges in the bioinformatics analysis of SC

376   datasets from solid tumors.

# 5. Figure captions

377

378 **Figure 1. Overview of scMuffin package.** scMuffin offers the possibility to perform several

379 different analyses and data integration approaches to address the main challenges of SC gene

380 expression analysis in solid tumors.

381 **Figure 2. Quantification of CancerSEA functional states in the HGG sample PJ016. a)** UMAP

382 visualization where cells are coloured by expression clusters. **b)** Cluster-level expression scores of

383 all the CancerSEA functional states. **c)** UMAP visualization where cells are colored by

384 "CSEA_Invasion" gene set score.

385 **Figure 3. CNV analysis. a-d)** CNV heatmaps (a, c) where cells (columns) are grouped into CNV

386 clusters, and UMAP visualizations (b, d) where cells are colored by CNV clusters, for sample PJ030

387 (a, b) and sample PJ016 (c, d). **e)** Overlap between cell clusters of sample PJ016 obtained by

388 analyzing gene expression (rows, "global_" prefix) and CNV clusters (columns, "cnv_" prefix); *YBX1*

389 and *HNRNPM* are two CancerSEA invasion markers located within the amplified 1p and 19p

390 regions found in CNV clusters 1 and 3. **f)** Expression scores for CancerSEA functional states in CNV

391 clusters of sample PJ016.

392 **Figure 4. Cluster enrichment in HGG sample PJ016. a)** The top five most significant (fdr < 0.05)

393 CancerSEA functional states in cluster 0: distribution of expression scores in cluster 0 (red) in

394 comparison with all the other clusters (grey); normalized enrichment score (NES) and false

395 discovery rate (FDR) values. **b)** Distribution of cells by their values (red labels) in cluster 0 (red

396 dots) in comparison with all the others (grey dots) for three categorical variables, namely, the

397 clusters obtained analyzing ribosomal gene expression ("ribosomes"), the clusters obtained

398 analyzing the expression of a Glioblastoma signature ("GB500"), and cell cycle phase ("Phase",

399 obtained with the Seurat package function "CellCycleScoring"); the numbers over each cell group

400    are ORA p-values. **c)** UMAP visualization with expression clusters annotated with the names of the

401    top two CancerSEA gene sets with the highest enrichment (CSEA) for each cluster.

402    **Figure 5. Chromatin state, proliferation rate and cell state trajectories of HGG sample PJ016. (a-**

403    **b).** Cell state trajectories colored by chromatin state (a) and proliferation rate (b).

404    **Figure 6. Cluster marker-based comparison of HGG samples PJ016 and PJ018. a)** Similarity among

405    all-pairs of clusters. **b-c)** UMAP visualizations obtained by integrating the two samples with the

406    "FindIntegrationAnchors" and "IntegrateData" Seurat functions, showing PJ016 cells (b) and PJ018

407    cells (c) colored by the clusters found by independent analysis of each sample.

# 408   6.    Availability and requirements

409    **Project name:** scMuffin

410    **Project home page:** https://github.com/emosca-cnr/scMuffin

411    **Operating system:** Platform independent

412    **Programming language:** R (>= 4.0.0)

413    **Other requirements:** The R Project for Statistical Computing.

414    **License:** GPL-3

415    **Any restrictions to use by non-academics:** According to GPL-3

# 416   7.    List of abbreviations

417    CNV: Copy Number Variation

418    CSEA: Cell Set Enrichment Analysis

419    GSEA: Gene Set Enrichment Analysis

420    HGG: high grade glioma

421    ORA: over representation analysis

422 SC: single cell

423 CSC: cancer stem cells

424 UMAP: Uniform Manifold Approximation and Projection

425 GTEx: The Genotype-Tissue Expression project

426 NES: normalized enrichment score

427 FDR: False discovery rate

## 428 8. Declarations

429 **Ethics approval and consent to participate.** Not applicable.

430 **Consent for publication.** Not applicable.

431 **Availability of data and materials.** The data used for the analyses described in this manuscript

432 were obtained from: The Gene Expression Omnibus (13), under the accession GSE103224; the

433 GTEx Portal (16) on 04/08/2020.

434 **Competing interests.** The authors declare that they have no competing interests.

435 **Funding.** The work was supported by: the Italian Ministry of Education, University and Research

436 (MIUR) [INTEROMICS PB05].

437 **Author's contributions.** VN implemented the CNV functions, carried out the analyses, interpreted

438 the results and wrote the article. NDN drafted the package, carried out the analyses and wrote the

439 article. AC implemented clustering functions, carried out the analyses and wrote the manuscript.

440 IC curated the biological aspects of CNV analysis and revised the manuscript. MM and MG set up

441 the computational infrastructure for data analysis. CC and EP curated the biological aspects of

442 solid tumor data analysis. RR, IZ, LM and AM contributed to the design of the study. PP

443 contributed to the software design, case study definition, interpretation of the results and wrote

444 the article. EM designed the study, implemented the software, performed the analysis,

20

445     interpreted the results and wrote the manuscript. All authors read and approved the final

446     manuscript.

## 9.    References

448     1.     Baslan T, Hicks J. Unravelling biology and shifting paradigms in cancer with single-cell

449            sequencing. Nat Rev Cancer [Internet]. 2017 Sep 24;17(9):557–69. Available from:

450            https://www.nature.com/articles/nrc.2017.58

451     2.     Zappia L, Theis FJ. Over 1000 tools reveal trends in the single-cell RNA-seq analysis

452            landscape. Genome Biol. 2021 Dec;22(1):301.

453     3.     Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, et al. Integrated analysis of

454            multimodal single-cell data. Cell [Internet]. 2021 Jun;184(13):3573-3587.e29. Available

455            from: https://linkinghub.elsevier.com/retrieve/pii/S0092867421005833

456     4.     Zhang X, Lan Y, Xu J, Quan F, Zhao E, Deng C, et al. CellMarker: a manually curated resource

457            of cell markers in human and mouse. Nucleic Acids Res [Internet]. 2019 Jan 8;47(D1):D721–

458            8. Available from: https://academic.oup.com/nar/article/47/D1/D721/5115823

459     5.     Franzén O, Gan L-M, Björkegren JLM. PanglaoDB: a web server for exploration of mouse and

460            human single-cell RNA sequencing data. Database [Internet]. 2019 Jan 1;2019. Available

461            from: https://academic.oup.com/database/article/doi/10.1093/database/baz046/5427041

462     6.     Yuan H, Yan M, Zhang G, Liu W, Deng C, Liao G, et al. CancerSEA: a cancer single-cell state

463            atlas. Nucleic Acids Res [Internet]. 2019 Jan 8;47(D1):D900–8. Available from:

464       https://academic.oup.com/nar/article/47/D1/D900/5133662

465    7.    Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set

466          enrichment analysis: A knowledge-based approach for interpreting genome-wide

467          expression profiles. Proc Natl Acad Sci [Internet]. 2005 Oct 25;102(43):15545–50. Available

468          from: https://pnas.org/doi/full/10.1073/pnas.0506580102

469    8.    Tirosh I, Izar B, Prakadan SM, Wadsworth MH, Treacy D, Trombetta JJ, et al. Dissecting the

470          multicellular ecosystem of metastatic melanoma by single-cell RNA-seq. Science (80- )

471          [Internet]. 2016 Apr 8;352(6282):189–96. Available from:

472          https://www.science.org/doi/10.1126/science.aad0501

473    9.    Nguyen QH, Pervolarakis N, Blake K, Ma D, Davis RT, James N, et al. Profiling human breast

474          epithelial cells using single cell RNA sequencing identifies cell diversity. Nat Commun

475          [Internet]. 2018 Dec 23;9(1):2028. Available from: http://www.nature.com/articles/s41467-

476          018-04334-1

477    10.   Patel AP, Tirosh I, Trombetta JJ, Shalek AK, Gillespie SM, Wakimoto H, et al. Single-cell RNA-

478          seq highlights intratumoral heterogeneity in primary glioblastoma. Science (80- ) [Internet].

479          2014 Jun 20;344(6190):1396–401. Available from:

480          https://www.science.org/doi/10.1126/science.1254257

481    11.   Khatri P, Sirota M, Butte AJ. Ten Years of Pathway Analysis: Current Approaches and

482          Outstanding Challenges. Ouzounis CA, editor. PLoS Comput Biol [Internet]. 2012 Feb

483          23;8(2):e1002375. Available from: https://dx.plos.org/10.1371/journal.pcbi.1002375

484    12.    Yuan J, Levitin HM, Frattini V, Bush EC, Boyett DM, Samanamud J, et al. Single-cell

485            transcriptome analysis of lineage diversity in high-grade glioma. Genome Med [Internet].

486            2018 Dec 24;10(1):57. Available from:

487            https://genomemedicine.biomedcentral.com/articles/10.1186/s13073-018-0567-9


488    13.    Barrett T, Wilhite SE, Ledoux P, Evangelista C, Kim IF, Tomashevsky M, et al. NCBI GEO:

489            archive for functional genomics data sets—update. Nucleic Acids Res [Internet]. 2012 Nov

490            26;41(D1):D991–5. Available from:

491            http://academic.oup.com/nar/article/41/D1/D991/1067995/NCBI-GEO-archive-for-

492            functional-genomics-data


493    14.    Clarke ZA, Andrews TS, Atif J, Pouyabahar D, Innes BT, MacParland SA, et al. Tutorial:

494            guidelines for annotating single-cell transcriptomic maps using automated and manual

495            methods. Nat Protoc [Internet]. 2021 Jun 24;16(6):2749–64. Available from:

496            http://www.nature.com/articles/s41596-021-00534-0


497    15.    Gu Z, Eils R, Schlesner M. Complex heatmaps reveal patterns and correlations in

498            multidimensional genomic data. Bioinformatics [Internet]. 2016 Sep 15;32(18):2847–9.

499            Available from: https://academic.oup.com/bioinformatics/article-

500            lookup/doi/10.1093/bioinformatics/btw313


501    16.    GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. Available from:

502            https://www.gtexportal.org/home/index.html


503    17.    Gupta MK, Polisetty RV, Sharma R, Ganesh RA, Gowda H, Purohit AK, et al. Altered

504            transcriptional regulatory proteins in glioblastoma and YBX1 as a potential regulator of

505           tumor invasion. Sci Rep [Internet]. 2019 Dec 29;9(1):10986. Available from:

506           http://www.nature.com/articles/s41598-019-47360-9

507   18.   LeFave C V, Squatrito M, Vorlova S, Rocco GL, Brennan CW, Holland EC, et al. Splicing factor

508           hnRNPH drives an oncogenic splicing switch in gliomas. EMBO J [Internet]. 2011 Oct

509           5;30(19):4084–97. Available from:

510           http://emboj.embopress.org/cgi/doi/10.1038/emboj.2011.259

511   19.   Golan-Gerstl R, Cohen M, Shilo A, Suh S-S, Bakàcs A, Coppola L, et al. Splicing Factor hnRNP

512           A2/B1 Regulates Tumor Suppressor Gene Splicing and Is an Oncogenic Driver in

513           Glioblastoma. Cancer Res [Internet]. 2011 Jul 1;71(13):4464–72. Available from:

514           http://cancerres.aacrjournals.org/lookup/doi/10.1158/0008-5472.CAN-10-4410

515   20.   Kim J-H, Jeong K, Li J, Murphy JM, Vukadin L, Stone JK, et al. SON drives oncogenic RNA

516           splicing in glioblastoma by regulating PTBP1/PTBP2 switching and RBFOX2 activity. Nat

517           Commun [Internet]. 2021 Dec 21;12(1):5551. Available from:

518           https://www.nature.com/articles/s41467-021-25892-x

519   21.   Bastide A, David A. The ribosome, (slow) beating heart of cancer (stem) cell. Oncogenesis

520           [Internet]. 2018 Apr 20;7(4):34. Available from: http://www.nature.com/articles/s41389-

521           018-0044-8

522   22.   Guimaraes JC, Zavolan M. Patterns of ribosomal protein expression specify normal and

523           malignant human cells. Genome Biol [Internet]. 2016 Dec 24;17(1):236. Available from:

524           http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-1104-z

525    23.    Dolezal JM, Dash AP, Prochownik E V. Diagnostic and prognostic implications of ribosomal

526            protein transcript expression patterns in human cancers. BMC Cancer [Internet]. 2018 Dec

527            12;18(1):275. Available from:

528            https://bmccancer.biomedcentral.com/articles/10.1186/s12885-018-4178-z


529    24.    Teo W-Y, Sekar K, Seshachalam P, Shen J, Chow W-Y, Lau CC, et al. Relevance of a TCGA-

530            derived Glioblastoma Subtype Gene-Classifier among Patient Populations. Sci Rep

531            [Internet]. 2019 Dec 15;9(1):7442. Available from: http://www.nature.com/articles/s41598-

532            019-43173-y


533    25.    Wainwright EN, Scaffidi P. Epigenetics and Cancer Stem Cells: Unleashing, Hijacking, and

534            Restricting Cellular Plasticity. Trends in Cancer [Internet]. 2017 May;3(5):372–86. Available

535            from: https://linkinghub.elsevier.com/retrieve/pii/S240580331730081X


536    26.    Gaspar-Maia A, Alajem A, Meshorer E, Ramalho-Santos M. Open chromatin in pluripotency

537            and reprogramming. Nat Rev Mol Cell Biol [Internet]. 2011 Jan 22;12(1):36–47. Available

538            from: http://www.nature.com/articles/nrm3036


539    27.    Guilhamon P, Chesnelong C, Kushida MM, Nikolic A, Singhal D, MacLeod G, et al. Single-cell

540            chromatin accessibility profiling of glioblastoma identifies an invasive cancer stem cell

541            population associated with lower survival. Elife [Internet]. 2021 Jan 11;10. Available from:

542            https://elifesciences.org/articles/64090


543    28.    Feitelson MA, Arzumanyan A, Kulathinal RJ, Blain SW, Holcombe RF, Mahajna J, et al.

544            Sustained proliferation in cancer: Mechanisms and novel therapeutic targets. Semin Cancer

545            Biol [Internet]. 2015 Dec;35:S25–54. Available from:

546        https://linkinghub.elsevier.com/retrieve/pii/S1044579X15000140

547    29.    Al-Hajj M, Clarke MF. Self-renewal and solid tumor stem cells. Oncogene [Internet]. 2004

548        Sep 20;23(43):7274–82. Available from: https://www.nature.com/articles/1207947

549    30.    Angerer P, Haghverdi L, Büttner M, Theis FJ, Marr C, Buettner F. destiny: diffusion maps for

550        large-scale single-cell data in R. Bioinformatics [Internet]. 2016 Apr 15;32(8):1241–3.

551        Available from: https://academic.oup.com/bioinformatics/article-

552        lookup/doi/10.1093/bioinformatics/btv715

553    31.    Forcato M, Romano O, Bicciato S. Computational methods for the integrative analysis of

554        single-cell data. Brief Bioinform [Internet]. 2021 May 20;22(3). Available from:

555        https://academic.oup.com/bib/article/doi/10.1093/bib/bbaa042/5828125

556

## Challenges in solid tumor analysis at single-cell level

- Limited availability of markers for definition of cell subtype identity
- Potentially strongly altered and highly heterogeneous gene expression profiles
- Presence of infiltrating cells and cells from the surrounding (healthy) tissue
- Potentially clinically relevant cell subtypes at very low number (e.g., drug resistant subclones)
- Often limited number of detected genes

**scMuffin**

*Single-cell*

*multi-feature*

*integrative*

*analysis*

### Gene sets from various sources

- General: MSigDB
- Cell markers: CellMarker, PanglaoDB
- Functional states: CancerSEA
- Proliferation: G1/S & G2/M
- Ribosomal proteins

### CNV inference

- Adjacent gene windows
- Reference (optional)
- Support of normal tissue expression data from GTEx

### Gene set scoring

- Against empirical null
- Optimized to handle missing data

### Proliferation rate

- Proportional to G1/S and G2/M genes

### Cluster-marker based two-samples comparison

- All-pairs of clusters comparison

### Chromatin state

- Relative number of expressed genes

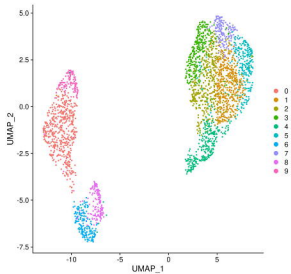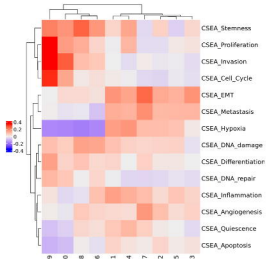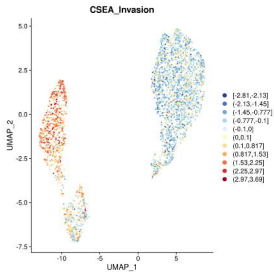### Cell cluster enrichment

- Quantitative features: CSEA
- Categorical features: ORA

### Comparison of multiple partitions

- Overlap matrix

### Cell state trajectory

- Diffusion map

### Implementation

- Computationally intensive tasks are parallelized
- Integration of various results in dedicated objects to enable automated subsequent analyses
- Parametrization of analyses to address dataset-specific characteristics

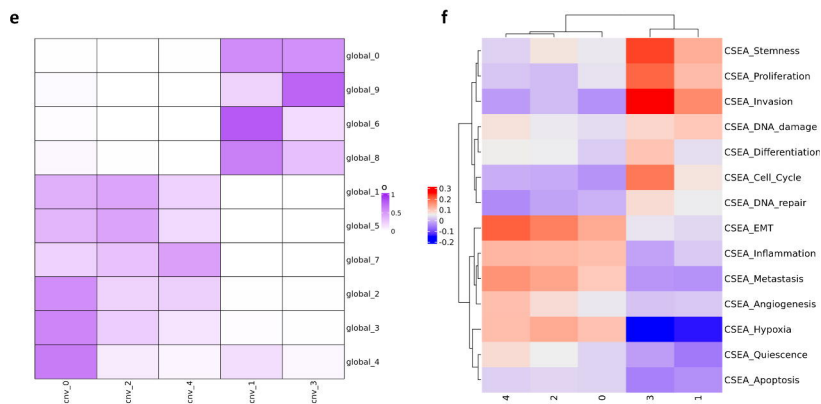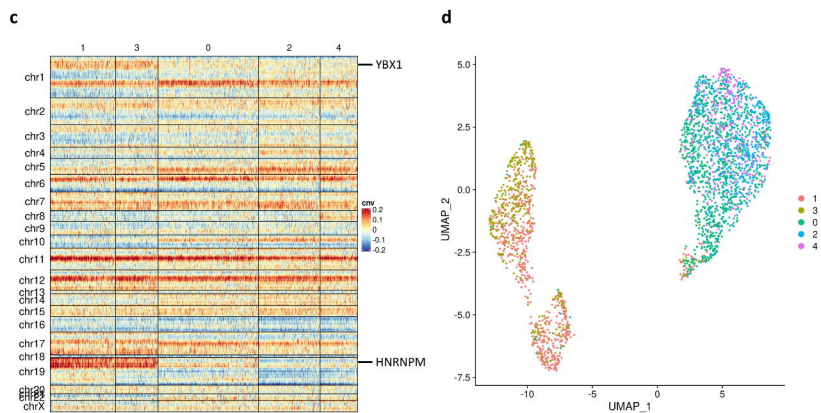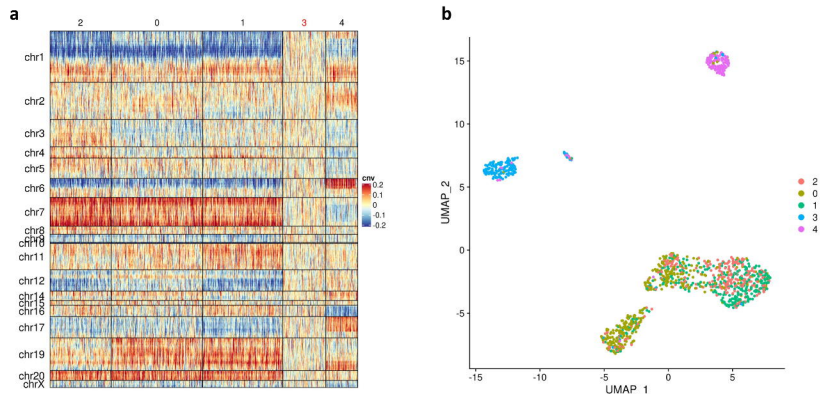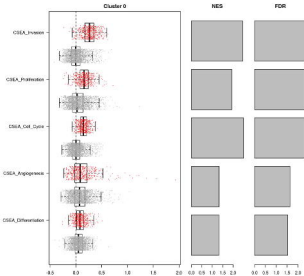### Visualization

- Automated UMAP visualization for multiple quantitative and categorical features
- Clusters-by-cells heatmaps
- CNV heatmap
- Two-samples similarity heatmap
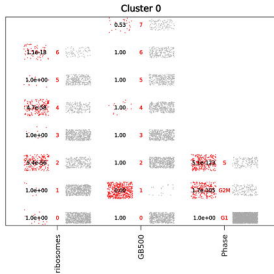- Cluster enrichment boxplots and dotplots

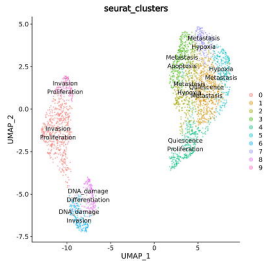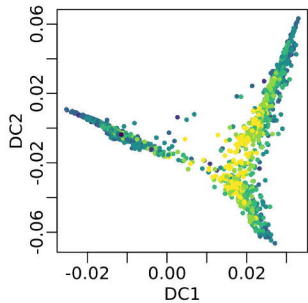### Cell cluster statistics

- mean values
- variability

$1\ (\mu_1, \sigma_1)$
$2\ (\mu_2, \sigma_2)$
$3\ (\mu_3, \sigma_4)$
$4\ (\mu_5, \sigma_5)$

a

b

c

d

e

f

**a**

Cluster 0    NES    FDR

CSEA_Invasion

CSEA_Proliferation

CSEA_Cell_Cycle

CSEA_Angiogenesis

CSEA_Differentiation

**b**

Cluster 0

**c**

seurat_clusters

**a**

**b**

Chromatin state

[0.043,0.056)
[0.056,0.069)
[0.069,0.081)
[0.081,0.094)
[0.094,0.11)
[0.11,0.12)
[0.12,0.13)

Proliferation

[-0.2,-0.1)
[-0.1,0.02)
[0.02,0.2)
[0.2,0.3)
[0.3,0.4)
[0.4,0.5)
[0.5,0.7)

a

b PJI016

c PJI018