

## Description of scripts used for single-cell multiomic analysis of hepatoblastomas

This folder contains all the scripts used in the manuscript entitled 'Single-cell multiomics reveals the interplay of clonal evolution and cellular plasticity in hepatoblastoma' by Roehrig *et al.*

Scripts were run in R version 4. The following **tools and packages** are required:

*cellranger-arc* (v2.0.0)

*Seurat* (v3)

*ArchR* (v1.0.1)

*Signac* (v1.6)

*InferCNV* (v1.6)

*spaceranger* (v1.3.1)

*screadcounts* (v1.1.8)

*samtools view* (v1.14)

Installing these tools should take ~2-3 hours on a desktop computer.

**Input data** are the raw sequencing files available on the European Genome Archive (EGA, accession code: EGAS00001006932).

Below is a description of the **functionalities** of each script. Scripts related to each data type (snRNAseq, snATACseq etc.) are grouped in separate folders.

### snRNAseq

#### 1. Preprocessing

Quality control of the 8 Multiome samples based on number of UMI counts, detected genes and percentage of mitochondrial reads is performed in *1\_Preprocessing.R* and visualized in *2\_Plot\_QC\_all\_samples.R*.

#### 2. Global visualization

Normalization, dimensionality reduction, clustering and 2D plot projection for 1) each sample individually (in *1\_Seurat\_SCTransform\_per\_sample.R*) or 2) merging cells from all 8 samples (in *2\_Seurat\_SCTransform\_merge\_all.R*). Normalization is performed with SCTransform for the main analyses and with NormalizeData when focusing on expression levels between cells. Projection of known signatures to identify the cell subtypes is performed in *3\_Seurat\_signature\_analyses.R* and *4\_Seurat\_H\_LP\_M\_signature\_violin\_plots.R*.

#### 3. Denoising

In-house smoothing of missing values by aggregating the expression values in hundreds of small-size cell clusters. Imputation is performed in *Imputation\_in\_house\_tumor\_cells.R*.

#### 4. Pseudo CNV

Computation of virtual copy-number alteration profiles from cell-level expression data, using normal hepatocytes from one non-tumor sample as reference cells. Input matrices are generated in *1\_InferCNV\_table\_preparation.R* (*Seurat*). *InferCNV* is run on each sample separately or on all cells from the 8 samples in *2\_InferCNV\_pipeline\_individual\_sample.R* and

*2\_InferCNV\_pipeline\_no\_cluster\_by\_sample.R*. Alteration clusters are defined in *3\_InferCNV\_analyses\_no\_cluster\_by\_sample.R*.

## 5. Tumor study

Identification of non-tumor, tumor cells and potential doublets based on PCA of all cells, projection of known signatures, presence of chromosomal aberrations and ATAC doublet enrichment is performed in *1\_Identification\_normal\_tumor\_on\_merged\_UMAP.R*. Analysis pipeline is performed again on tumor cells from all 6 samples (*2\_Seurat\_analyses\_tumor\_cells\_strict.R*) or each sample separately (*4\_Seurat\_analyses\_tumor\_cells\_individual\_strict.R*). scH, scH/LP, scLP and scM Multiome subtypes are identified in *3\_Identification\_H\_LP\_M\_strict.R* and the correlation between PCA components PC1 and PC2, and known HB markers, is realized in *3\_PCA\_component\_identification\_H\_LP\_M\_correlation.R*. Differential expression analysis between subtypes is performed in *5\_Differential\_expression\_subtypes.R*. Signatures related to normal liver development or signatures from other HB single-cell studies are investigated in *Gene\_signatures/1\_Development\_signatures.R* and *Gene\_signatures/2\_Other\_HB\_scRNAseq\_studies.R*.

## 6. Clonal evolution

Definition of genetic subclones from copy-number alteration profiles in *1\_Final\_definition\_CNV\_clusters.R*, and their characterization with diverse annotations in *2\_Characterization\_of\_CNV\_clusters.R* and *2\_Characterization\_of\_CNV\_clusters\_CHC2959T\_2960T.R*. Somatic mutations are investigated in the « *Somatic mutations* » folder.

BAM files are filtered on reads kept for UMI counting (tag xf=25) in *0\_Filter\_BAM\_on\_xf\_tag\_before\_screadcounts.R*.

Lists of somatic mutations from bulk WGS are prepared in *1\_Preparation\_tables\_screadcounts.R*. *scReadCounts* is run on the single-nucleus BAM files of each sample to detect WGS mutations in *2\_scReadCounts\_command\_line.R*. Study of mutations clusters from patient 2959 is performed in *3\_scReadCounts\_analyses\_2959\_patient.R*. Study of tumor-specific sets of mutations is performed in *4\_scReadCounts\_analyses\_global\_mutations.R*. Study of specific mutations like *CTNNB1* is performed in *4\_scReadCounts\_analyses\_individual\_mutations.R*. Investigation of clonal mutation specificity and sensitivity is performed in *5\_scReadCounts\_clonal\_mutations\_all\_samples.R*. Detection of 11p15 alterations is performed in the folder « *11p15* » in *1\_scReadCounts\_analyses\_11p15\_cnLOH.R*.

## 7. Cancer stem cells

Identification of cancer stem cell and liver cancer stem cell markers is performed in *1\_Cancer\_stem\_cells\_study\_all\_csc\_markers.R* and *1\_Cancer\_stem\_cells\_study\_liver\_csc\_markers.R*.

# snATACseq

*ATAC\_tracks\_per\_subtype.R* enables to visualize chromatin accessibility in selected regions aggregated for each Multiome HB subtype.

*Create\_promoter\_gene\_body\_table.R* enables to create a table containing genomic features for peaks.

## 1. Preprocessing

Quality control of the 8 Multiome samples based on number of fragments and TSS enrichment, and creation of the corresponding Arrow files, are performed in

1\_Preprocessing\_quality\_control\_individual\_sample\_QC.R. Arrow files are gathered in an *ArchR* project in 2\_Create\_project\_from\_individual\_QC.R.

## 2. Global visualization

Normalization, dimensionality reduction, clustering and 2D plot projection for 1) merging cells from all 8 samples (in 1\_Normalization\_clustering\_UMAP\_all\_samples.R), 2) merging tumor cells from all 6 tumor samples (in 1\_Normalization\_clustering\_UMAP\_tumor\_cells.R), or 3) tumor cells from each individual patient (in 1\_Normalization\_clustering\_UMAP\_tumor\_cells\_by\_sample.R and 1\_Normalization\_clustering\_UMAP\_tumor\_cells\_by\_sample\_2959.R).

## 3. Peak calling

We used *ArchR* for most snATAC-seq analyses and *Signac* for peak calling. Intersection between cells identified as non-empty droplets by *cellranger* and cells retained after quality control in *ArchR* is performed in 1\_Preprocess\_individual\_samples\_ArchR\_QC.R. Peak calling is then performed on each HB subtype in each tumor sample in 2\_Peak\_calling\_individual\_samples.R, and the resulting peak sets are merged in 3\_Peak\_calling\_combine\_samples.R. The global peak set obtained with *Signac* is integrated in the *ArchR* pipeline and the corresponding peak count matrix is computed in 4\_Integrate\_Signac\_peakset\_in\_ArchR\_peak\_calling.R. Peaks are annotated with their nearest gene in 5\_Compute\_nearest\_gene\_per\_peak.R.

## 4. Differential peaks

Differential peak analysis is performed between tumor subtypes and non-tumor cells on the one hand, and between tumor subtypes on the other hand, in 1\_Differential\_peaks\_tables.R. The resulting differential peaks are annotated by their logFC, FDR, nearest gene, distance to nearest TSS, linked genes and regulating transcription factors in 2\_Annotate\_differential\_peaks.R.

## 5. Peaks study

The number of differential peaks for each comparison is computed in 1\_Differential\_study\_nb\_peaks.R. Differential peaks are characterized by their chromatin state enrichment in 1\_Differential\_study\_chromatin\_states.R.

## 6. Peak to gene linkage

Linkage between peaks and genes based on correlation between peak accessibility and gene expression is performed in 1\_Peak\_to\_gene\_linkage\_multiome.R.

## 7. Motif TF

A custom position weight matrix list to include motif information from version 2 of the cisbp database is created in 0\_Create\_custom\_PWM\_list\_cisbp\_v2.R. Motif enrichment analysis in differential peaks related to each Multiome HB subtype is performed in 1\_Motif\_enrichment\_differential\_peaks\_cisbp\_v1\_v2.R. TF footprints for specific TF are computed in 2\_TF\_footprints\_cisbp\_v2.R. TF motif deviations are computed using cisbp versions 1 and 2 in 2\_TF\_motif\_deviations\_cisbp\_v1.R and 2\_TF\_motif\_deviations\_cisbp\_v2.R.

# Integration\_snRNAseq\_snATACseq

Identification of key TF regulators based on specific criteria relying on bulk RNA-seq, snRNA-seq and snATAC-seq results is performed in *1\_Identify\_TF\_regulators.R*. Gene regulatory networks between the defined key TF and their potential target genes are identified in *2\_Identify\_GRN\_networks\_H\_LP\_M.R*. Aggregated single-nucleus expression of each resulting module (TF and target genes) is projected in boxplots in *3\_Boxplots\_GRN\_modules\_per\_Multiome\_subtype.R*. Single-nucleus expression and motif deviation of the corresponding TF, as well as single-nucleus expression of their targets, are projected on a PC2-ordered heatmap in *3\_Heatmap\_Multiome\_TF\_targets\_expression.R*. Bulk expression of the corresponding TF and their targets is projected on a PC2-ordered heatmap in *3\_Heatmap\_bulk\_RNA\_seq\_TF\_targets\_expression.R*. Recapitulating tables of correlation between TF and their targets, as well as mean expression of the TF and targets in each Multiome HB subtype, can be found in *Tables.R*.

## Visium

Multiome scH and scLP markers are projected in the Visium data of one tumor from the Multioem cohort (#3133T) in *1\_Project\_H\_LP\_snRNAseq\_markers\_on\_visium\_3133T.R*.