

‘FIRM’ Package for flexible integration of heterogeneous scRNA-seq datasets across multiple tissue types, platforms and experimental batches

November 15, 2021

1 Overview

This vignette provides an introduction to the ‘FIRM’ package. R package ‘FIRM’ implements FIRM, an algorithm for flexible integration of heterogeneous scRNA-seq datasets across multiple tissue types, platforms and experimental batches.

The FIRM package can be loaded with the command:

```
R> library("FIRM")
```

Also load required package:

```
R> library(Seurat)
```

```
R> library(RANN)
```

2 Workflow

In this vignette, we use the `ExampleData` in the package which is the aorta data from Tabula Muris generated using Smart-seq2 (SS2) and 10X. The `ExampleData$SS2` and `ExampleData$tenx` provide the gene expression matrices \mathbf{X} for SS2 and 10X respectively, where X_{ij} is the number of reads (for SS2) and unique molecular identified (UMI, for 10X). We only considered the cells with annotations in the original study. The cell type annotations for SS2 and 10X datasets are provided in `ExampleData$meta_SS2` and `ExampleData$meta_tenx`.

```
R> data("ExampleData")
```

```
R> dim(SS2)
```

```
[1] 23341 408
```

```
R> length(meta_SS2)
```

```
[1] 408
```

```
R> dim(tenx)
```

```
[1] 23433 526
```

```
R> length(meta_tenx)
```

```
[1] 526
```

2.1 Data preprocessing

We performed the standard preprocessing workflow to prepare the scaled data for integration.

```
R> prep_SS2 <- prep_data(SS2)
R> Dataset1 <- prep_SS2$Dataset
R> hvg1 <- prep_SS2$hvg
R> prep_tenx <- prep_data(tenx)
R> Dataset2 <- prep_tenx$Dataset
R> hvg2 <- prep_tenx$hvg
```

2.2 Integration using FIRM

The integrated data is provided in 'Dataset'.

```
R> dims <- 15
R> coreNum <- 4
R> Dataset <- FIRM(Dataset1, Dataset2, hvg1, hvg2, dims = dims, coreNum = coreNum)

R> dim(Dataset)

[1] 23433 934
```

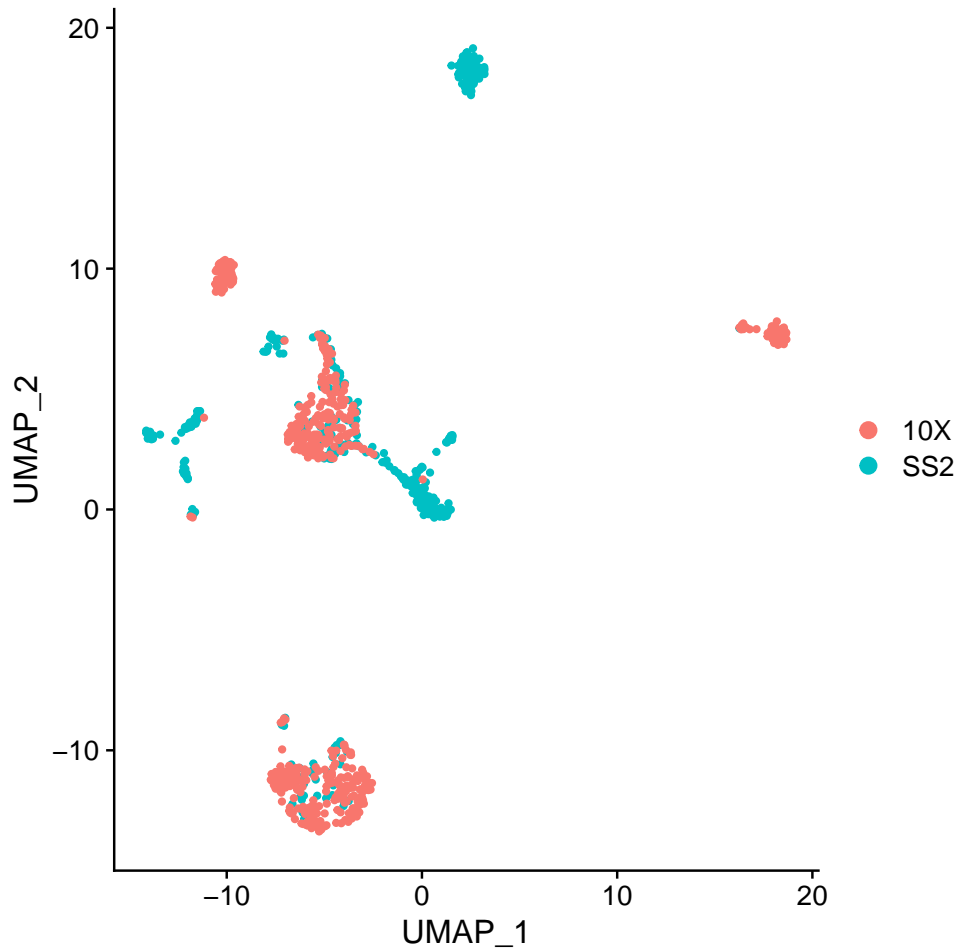
2.3 Downstream analysis

We can use the integrated data for downstream analysis, such as PCA and visaulization. For example, we can create a Seurat object.

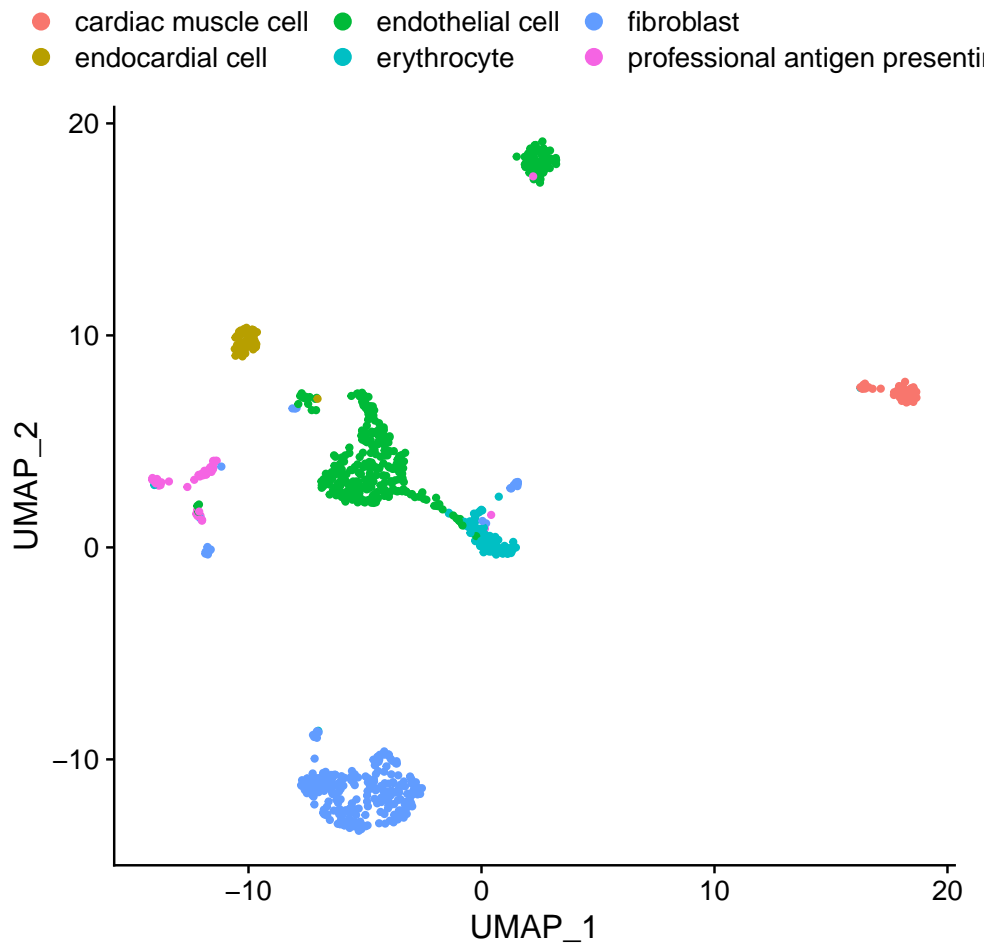
```
R> # counts
R> integrated_counts <- matrix(0, nrow(Dataset), ncol(Dataset))
R> rownames(integrated_counts) <- rownames(Dataset)
R> colnames(integrated_counts) <- colnames(Dataset)
R> integrated_counts[rownames(SS2), colnames(SS2)] <- as.matrix(SS2)
R> integrated_counts[rownames(tenx), colnames(tenx)] <- as.matrix(tenx)
R> # create Seurat object
R> integrated <- CreateSeuratObject(integrated_counts)
R> # normalization
R> integrated <- NormalizeData(integrated)
R> # put the centered integrated data in scaled data
R> integrated[["RNA"]]@scale.data <- Dataset - rowMeans(Dataset)
R> # add meta data
R> integrated <- AddMetaData(integrated, metadata = c(meta_SS2, meta_tenx),
+                           col.name = "annotation")
R> integrated <- AddMetaData(integrated,
+                           metadata = c(rep("SS2", ncol(SS2)), rep("10X", ncol(tenx))),
+                           col.name = "dataset")
R> # hvg for PCA and visualization
R> hvg <- intersect(hvg1, hvg2)
R> # PCA
R> integrated <- RunPCA(integrated, features = hvg, npcs = dims)
R> # UMAP
```

```
R> integrated <- RunUMAP(integrated, reduction = "pca", dims = 1:dims,
+                         umap.method = 'umap-learn', metric = "correlation")

R> DimPlot(integrated, reduction = "umap", group.by = "dataset")
```



```
R> library(ggplot2)
R> DimPlot(integrated, reduction = "umap", group.by = "annotation") +
+   theme(legend.position = "top")
```



References

- [1] Jingsi Ming, Zhixiang Lin, Jia Zhao, Xiang Wan, Can Yang, Angela Ruohao Wu, FIRM: Flexible Integration of single-cell RNA-sequencing data for large-scale Multi-tissue cell atlas datasets. <https://www.biorxiv.org/content/10.1101/2020.06.02.129031v2>
- [2] Schaum, N. et al. Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris. *Nature* 562, 367–372 (2018).
- [3] Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* 177, 1888-1902.e21 (2019).