

# scAlign Tutorial

Nelson Johansen, Gerald Quon

2022-08-24

## Introduction

This tutorial provides a guided alignment for two groups of cells from cellbench RNA mixture experiments. In this tutorial we demonstrate the unsupervised alignment strategy of **scAlign** described in Johansen et al, 2018 along with typical analysis utilizing the aligned dataset, and show how **scAlign** can identify and match cell types across platforms without using the labels as input.

## Alignment goals

The following is a walkthrough of a typical alignment problem for **scAlign** and has been designed to provide an overview of data preprocessing, alignment and finally analysis in our joint embedding space. Here, our primary goals include:

1. Learning a low-dimensional cell state space in which cells group by function and type, regardless of condition (platform).
2. Accurately labeling old cells with cell cycle and cell type information using only the young cell annotations.

## Installation

```
## Install scAlign
install.packages('devtools')
devtools::install_github(repo = 'quon-titative-biology/scAlign')
library(scAlign)

## Install Tensorflow
library(tensorflow)
install_tensorflow(version = "gpu") ## Removing version will install CPU version of Tensorflow
```

Guide to installing python and tensorflow on different operating systems.

On Windows:

Download Python 3.6.8. Note, newer versions of Python (e.g. 3.7) cannot use TensorFlow at this time. Make sure pip is included in the installation.

On Ubuntu:

```
sudo apt update
sudo apt install python3-dev python3-pip
```

On MacOS (homebrew):

```
/usr/bin/ruby -e "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/master/install)"
export PATH="/usr/local/bin:/usr/local/sbin:$PATH"
brew update
```

```
brew install python # Python 3
```

Further details at: <https://www.tensorflow.org/install>

## Data setup

The gene count matrices used for this tutorial are hosted on the cellbench github: [here](#).

First, we load in the normalized cellbench data. The data was normalized following the procedures defined on the cellbench github.

```
library(scAlign)
library(SingleCellExperiment)
library(ggplot2)

## Load in cellbench data
data("cellbench", package = "scAlign", envir = environment())

## Extract RNA mixture cell types
mix.types = unlist(lapply(strsplit(colnames(cellbench), "-"), "[", 2))

## Extract Platform
batch = c(rep("CEL", length(which(!grepl("sortseq", colnames(cellbench)) == TRUE))),
          rep("SORT", length(which(grepl("sortseq", colnames(cellbench)) == TRUE))))
```

## scAlign setup

The general design of scAlign's makes it agnostic to the input RNA-seq data representation. Thus, the input data can either be gene-level counts, transformations of those gene level counts or a preliminary step of dimensionality reduction such as canonical correlates or principal component scores. Here we create the scAlign object from the previously normalized cellbench data and perform CCA on the unaligned data.

```
## Create SCE objects to pass into scAlignCreateObject
youngMouseSCE <- SingleCellExperiment(
  assays = list(scale.data = cellbench[,batch=='CEL'])
)

oldMouseSCE <- SingleCellExperiment(
  assays = list(scale.data = cellbench[,batch=='SORT'])
)

## Build the scAlign class object and compute PCs
scAlignCB = scAlignCreateObject(sce.objects = list("CEL"=youngMouseSCE,
                                                  "SORT"=oldMouseSCE),
                              labels = list(mix.types[batch=='CEL'],
                                             mix.types[batch=='SORT']),
                              data.use="scale.data",
                              pca.reduce = FALSE,
                              cca.reduce = TRUE,
                              ccs.compute = 5,
                              project.name = "scAlign_cellbench")

## [1] "Computing CCA using Seurat."
## Centering and scaling data matrix
## Centering and scaling data matrix
```

```
## Running CCA
## Merging objects
```

## Alignment of cellbench RNAmixture

Now we align the cell populations from both protocols. `scAlign` returns a low-dimensional joint embedding space where the effect of platform is removed allowing us to use the complete dataset for downstream analyses such as clustering or differential expression.

```
## Run scAlign with all_genes
scAlignCB = scAlign(scAlignCB,
                    options=scAlignOptions(steps=1000,
                                           log.every=1000,
                                           norm=TRUE,
                                           early.stop=TRUE),
                    encoder.data="scale.data",
                    supervised='none',
                    run.encoder=TRUE,
                    run.decoder=FALSE,
                    log.dir=file.path(tempdir(), 'gene_input'),
                    device="CPU")
```

```
## Loaded Tensorflow version 2.4.1
```

```
## Found tf version: 2.4.1
```

```
## [1] "===== Step 1/3: Encoder training ====="
```

```
## [1] "Graph construction"
```

```
## [1] "Adding source walker loss"
```

```
## [1] "Error during alignment, returning scAlign class."
```

```
## <Rcpp::exception in py_call_impl(callable, dots$args, dots$keywords): ValueError: Tried to convert 'x' to a vector of type 'double' but the object is of type 'character' and is not coercible.
```

## >

```
# ## Additional run of scAlign with CCA
# scAlignCB = scAlign(scAlignCB,
#                       options=scAlignOptions(steps=1000,
#                                               log.every=1000,
#                                               norm=TRUE,
#                                               early.stop=TRUE),
#                       encoder.data="CCA",
#                       supervised='none',
#                       run.encoder=TRUE,
#                       run.decoder=FALSE,
#                       log.dir=file.path(tempdir(), 'cca_input'),
#                       device="CPU")
```

## Plot aligned data in tSNE space, when the data was processed in three different ways:

```
## 1) either using the original gene inputs,
```

```
## 2) after CCA dimensionality reduction for preprocessing.
```

```
## Cells here are colored by input labels
```

```
set.seed(5678)
```

```
gene_plot = PlotTSNE(scAlignCB,
                     "ALIGNED-GENE",
                     title="scAlign-Gene",
                     perplexity=30)
```

```
## [1] "Error plotting the data provided."
## <simpleError in h(simpleError(msg, call)): error in evaluating the argument 'x' in selecting a method for function 'plot'>
## Show plot
gene_plot

## NULL

# cca_plot = PlotTSNE(scAlignCB,
#                     "ALIGNED-CCA",
#                     title="scAlign-CCA",
#                     perplexity=30)
#
# multi_plot_labels = grid.arrange(gene_plot, cca_plot, nrow = 1)
```

## Session Info

```
sessionInfo()

## R version 4.2.1 (2022-06-23)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 20.04.4 LTS
##
## Matrix products: default
## BLAS:   /home/biocbuild/bbs-3.16-bioc/R/lib/libRblas.so
## LAPACK: /home/biocbuild/bbs-3.16-bioc/R/lib/libRlapack.so
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_GB            LC_COLLATE=C
##  [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=en_US.UTF-8     LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats4      stats      graphics  grDevices  utils      datasets  methods
## [8] base
##
## other attached packages:
##  [1] scAlign_1.11.0      FNN_1.1.3.1
##  [3] ggplot2_3.3.6       Rtsne_0.16
##  [5] irlba_2.3.5         Matrix_1.4-1
##  [7] purrr_0.3.4         tensorflow_2.9.0
##  [9] sp_1.5-0            SeuratObject_4.1.0
## [11] Seurat_4.1.1        SingleCellExperiment_1.19.0
## [13] SummarizedExperiment_1.27.2 Biobase_2.57.1
## [15] GenomicRanges_1.49.1  GenomeInfoDb_1.33.5
## [17] IRanges_2.31.2       S4Vectors_0.35.1
## [19] BiocGenerics_0.43.1   MatrixGenerics_1.9.1
## [21] matrixStats_0.62.0
##
## loaded via a namespace (and not attached):
##  [1] plyr_1.8.7          igraph_1.3.4        lazyeval_0.2.2
```

##	[4]	splines_4.2.1	listenv_0.8.0	scattermore_0.8
##	[7]	tfruns_1.5.0	digest_0.6.29	htmltools_0.5.3
##	[10]	fansi_1.0.3	magrittr_2.0.3	tensor_1.5
##	[13]	cluster_2.1.4	config_0.3.1	ROCR_1.0-11
##	[16]	globals_0.16.0	spatstat.sparse_2.1-1	colorspace_2.0-3
##	[19]	rappdirs_0.3.3	ggrepel_0.9.1	xfun_0.32
##	[22]	dplyr_1.0.9	RCurl_1.98-1.8	jsonlite_1.8.0
##	[25]	progressr_0.10.1	spatstat.data_2.2-0	survival_3.4-0
##	[28]	zoo_1.8-10	glue_1.6.2	polyclip_1.10-0
##	[31]	gtable_0.3.0	zlibbioc_1.43.0	XVector_0.37.0
##	[34]	leiden_0.4.2	DelayedArray_0.23.1	future.apply_1.9.0
##	[37]	abind_1.4-5	scales_1.2.1	DBI_1.1.3
##	[40]	spatstat.random_2.2-0	miniUI_0.1.1.1	Rcpp_1.0.9
##	[43]	viridisLite_0.4.1	xtable_1.8-4	reticulate_1.25
##	[46]	spatstat.core_2.4-4	htmlwidgets_1.5.4	httr_1.4.4
##	[49]	RColorBrewer_1.1-3	ellipsis_0.3.2	ica_1.0-3
##	[52]	pkgconfig_2.0.3	uwot_0.1.14	deldir_1.0-6
##	[55]	utf8_1.2.2	here_1.0.1	tidyselect_1.1.2
##	[58]	rlang_1.0.4	reshape2_1.4.4	later_1.3.0
##	[61]	munsell_0.5.0	tools_4.2.1	cli_3.3.0
##	[64]	generics_0.1.3	ggridges_0.5.3	evaluate_0.16
##	[67]	stringr_1.4.1	fastmap_1.1.0	yaml_2.3.5
##	[70]	goftest_1.2-3	knitr_1.40	fitdistrplus_1.1-8
##	[73]	RANN_2.6.1	pbapply_1.5-0	future_1.27.0
##	[76]	nlme_3.1-159	whisker_0.4	mime_0.12
##	[79]	compiler_4.2.1	plotly_4.10.0	png_0.1-7
##	[82]	spatstat.utils_2.3-1	tibble_3.1.8	stringi_1.7.8
##	[85]	rgeos_0.5-9	lattice_0.20-45	vctrs_0.4.1
##	[88]	pillar_1.8.1	lifecycle_1.0.1	spatstat.geom_2.4-0
##	[91]	lmtest_0.9-40	RcppAnnoy_0.0.19	data.table_1.14.2
##	[94]	cowplot_1.1.1	bitops_1.0-7	httpuv_1.6.5
##	[97]	patchwork_1.1.2	R6_2.5.1	promises_1.2.0.1
##	[100]	KernSmooth_2.23-20	gridExtra_2.3	parallelly_1.32.1
##	[103]	codetools_0.2-18	MASS_7.3-58.1	assertthat_0.2.1
##	[106]	rprojroot_2.0.3	withr_2.5.0	sctransform_0.3.4
##	[109]	GenomeInfoDbData_1.2.8	mgcv_1.8-40	parallel_4.2.1
##	[112]	grid_4.2.1	rpart_4.1.16	tidyr_1.2.0
##	[115]	rmarkdown_2.15	shiny_1.7.2	base64enc_0.1-3