

Welcome to the documentation of singletCode!

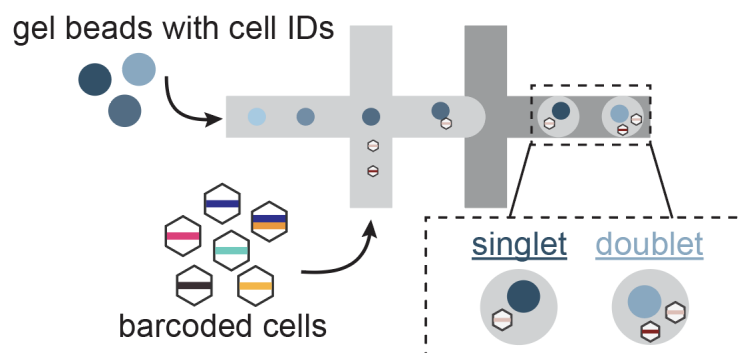
singletCode ([link to preprint](#)) uses cell barcoding technology to identify **true singlets** in scRNAseq data. In our preprint, we used singletCode to benchmark existing doublet detection methods and train a classifier to detect singlets in non-barcoded datasets as well of the same cell type.

doublets in sequencing

- [How do doublets form?](#)
- [Downstream impact of doublets](#)

How do doublets form?

Single-cell RNA sequencing (scRNA-seq) datasets contain true single cells, or singlets, in addition to cells that randomly coalesce during the protocol, or doublets. Sometimes, there are higher rates of doublets which can be attributed to cellular physiology and experimental protocols can lead to cell clumping. Doublet percentage in a sample be as high as 40%. Doublets can be two very transcriptionally different cells captured together (heterotypic) or two transcriptionally similar cells captured together (homotypic). Doublets are difficult to identify because just because two cells are captured together does not mean there is simply more absolute RNA fragments present or sequenced in doublet cases, and cells exist on a transcriptional continuum, making identification of valid singlet cells difficult, especially if cells are in a transitioning or reprogrammed state.



Schematic of how cells with lineage barcodes appear in the single-cell sequencer where droplets add a unique cell ID. The singlet depicted is one reaction droplet with a single cell such that there is a 1:1:1 mapping of a cell to a barcode to a cell ID. Alternatively, an example of a doublet is when one droplet has two or more cells, each having a unique barcode.

Downstream impact of doublets

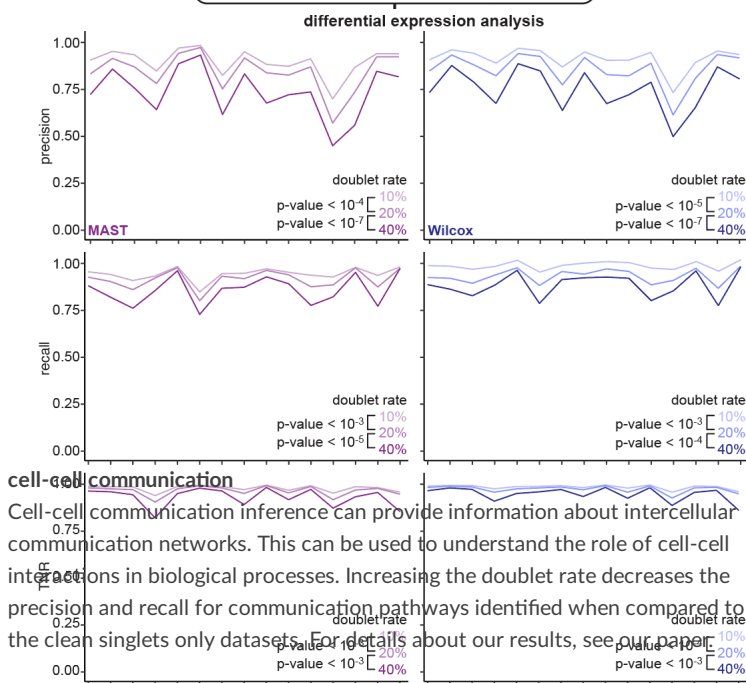
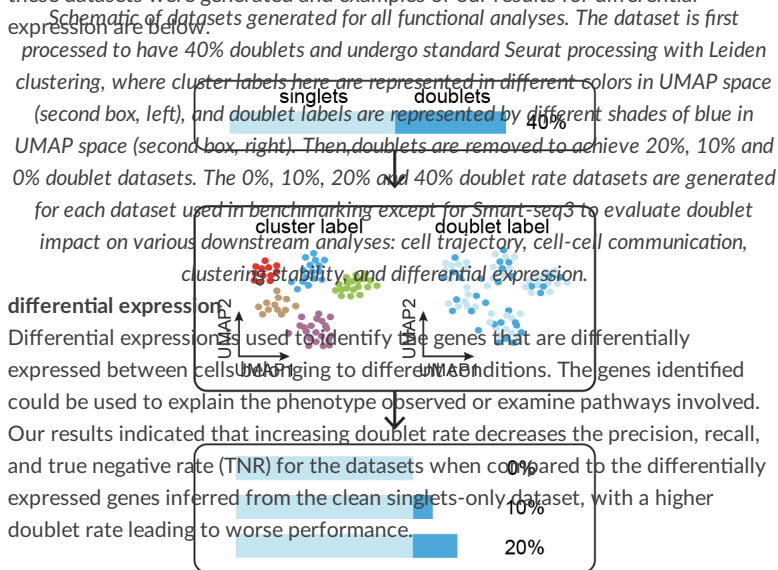
Doublets are problematic because they can impact the conclusions of scRNA-seq downstream functional analyses. We found that doublets confound downstream analyses when we tested common downstream scRNA-seq data analysis protocols (differential expression, cell trajectory, clustering stability, and cell-cell communication) on datasets of various doublet percentages. A schematic for how these datasets were generated and examples of our results for differential expression are below.

Schematic of datasets generated for all functional analyses. The dataset is first processed to have 40% doublets and undergo standard Seurat processing with Leiden clustering, where cluster labels here are represented in different colors in UMAP space (second box, left), and doublet labels are represented by different shades of blue in UMAP space (second box, right). Then, doublets are removed to achieve 20%, 10% and 0% doublet datasets. The 0%, 10%, 20% and 40% doublet rate datasets are generated for each dataset used in benchmarking except for Smart-seq3 to evaluate doublet impact on various downstream analyses: cell trajectory, cell-cell communication, clustering stability, and differential expression.

Schematics of how some common droplet-based methods appear in the single cell sequencing process. Droplets add a unique cell ID. The singlet depicted is one reaction droplet with a single cell such that there is a 1:1:1 mapping of a cell to a barcode to a cell ID. Alternatively, an example of a doublet is when one droplet has two or more cells, each having a unique barcode.

Downstream impact of doublets

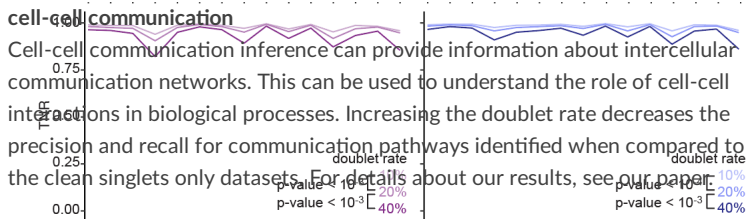
Doublets are problematic because they can impact the conclusions of scRNA-seq downstream functional analyses. We found that doublets confound downstream analyses when we tested common downstream scRNA-seq data analysis protocols (differential expression, cell trajectory, clustering stability, and cell-cell communication) on datasets of various doublet percentages. A schematic for how these datasets were generated and examples of our results for differential expression are below.



cell trajectory
Differential expression analysis results for all datasets using MAST (purple, left) and Wilcoxon (blue, right) tests. The color of the trajectory clusters is shown in the trajectory color bar. Increasing doublet rate caused the trajectory clusters to deviate from the trajectory color bar. For details about our results, see our paper.

what is singletCode?
singletCode is a framework to extract true singlets from barcoded scRNA-seq data. Our pipeline identified barcoded singlets from 10 different publications and 6 original experimental scRNA-seq samples generated in this study, encompassing 7 different barcoding technologies, 94 scRNA-seq samples, 3 unique sequencing technologies, and a total of 384,579 cells. Of the 384,579 barcoded cells, we extracted 293,618 singlets. Read more about the barcoded datasets we used singletCode with [here](#).

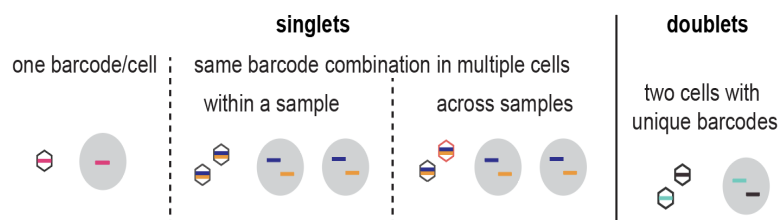
Clustering is one of the most common analyses done to infer the identity of similar cells. When the number of doublets were increased, the probability of getting the correct number of cell clusters decreased. This might lead to spurious clusters or a genuine sub-type of cells not being identified as a distinct group. For details about our results, see our paper.



Cell-cell communication inference can provide information about intercellular communication networks. This can be used to understand the role of cell-cell interactions in biological processes. Increasing the doublet rate decreases the precision and recall for communication pathways identified when compared to the clean singlets only dataset. For details about our results, see our paper.

expression analysis results for all datasets using MAST (purple, left) and increasing doublet rate caused the transition to be more obvious. The transition is inferred from the clean singlet only dataset and there were no significant differences when doublets were introduced. For details about our results, see our paper.

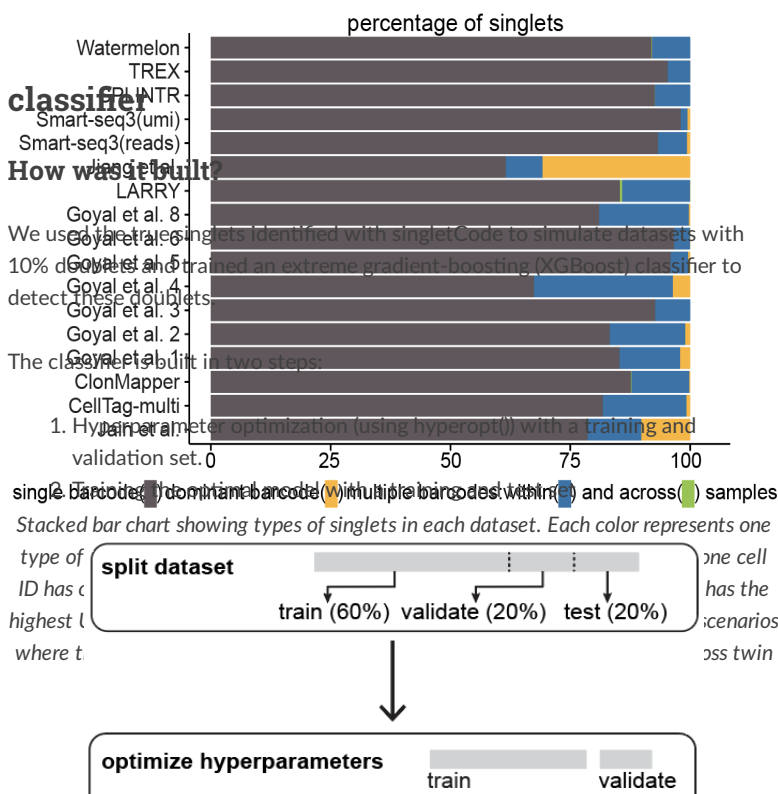
what is singletCode? singletCode is a framework to extract true singlets from barcoded scRNA-seq data. Our pipeline identified barcoded singlets from 10 different publications and 6 original experimental scRNA-seq samples generated in this study, encompassing 7 different barcoding technologies, 94 scRNA-seq samples, 3 unique sequencing technologies, and a total of 354,579 cells. Of the 354,579 barcoded cells, we extracted 293,616 singlets. Read more about the barcoded datasets we used singletCode with [here](#). For details about our results, see our paper.

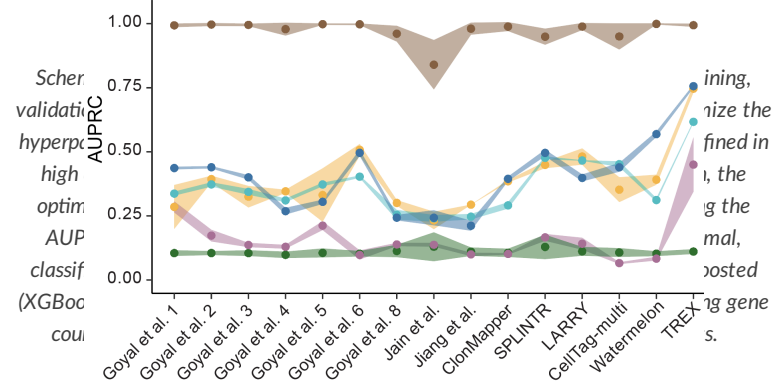
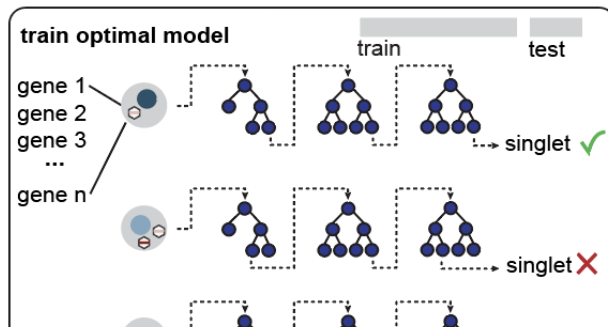
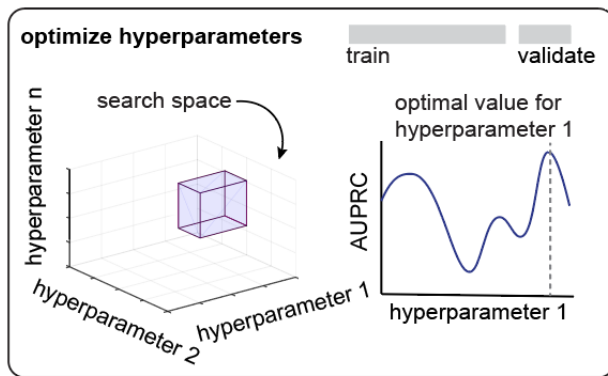
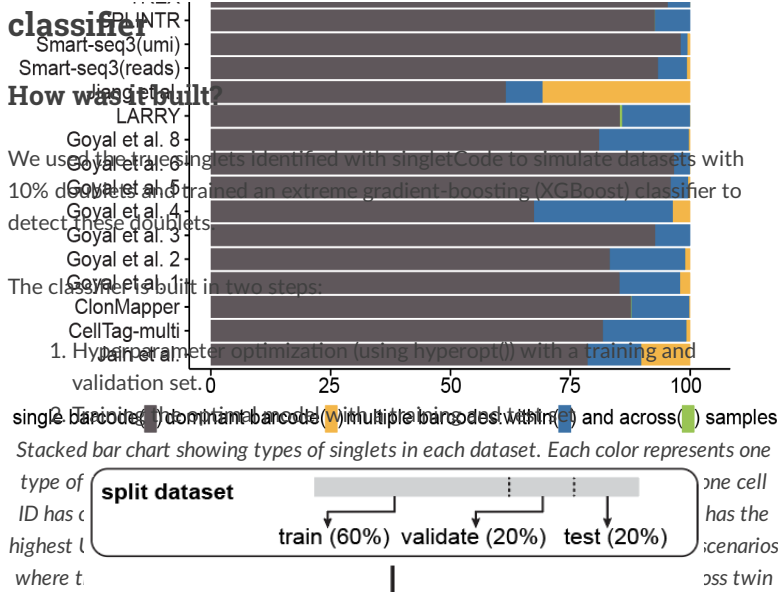


Singlet identification

singletCode leverages the fact that the lineage barcodes are present within the cell before sequencing. As a result of multiplicity of infection, singletCode identifies cells which meet any of the following conditions as a true singlet:

- 1 barcode/1 cell ID
- >1 barcode per 1 cell ID, but 1 barcode has significantly more UMI counts than the other barcodes within the same cell
- M barcodes per 1 cell ID, but the same combination of M barcodes are found in other cells in the same sample
- M barcodes per 1 cell ID, but the same combination of M barcodes are found in other cells across samples within the same experimental design (common in barcoding studies)



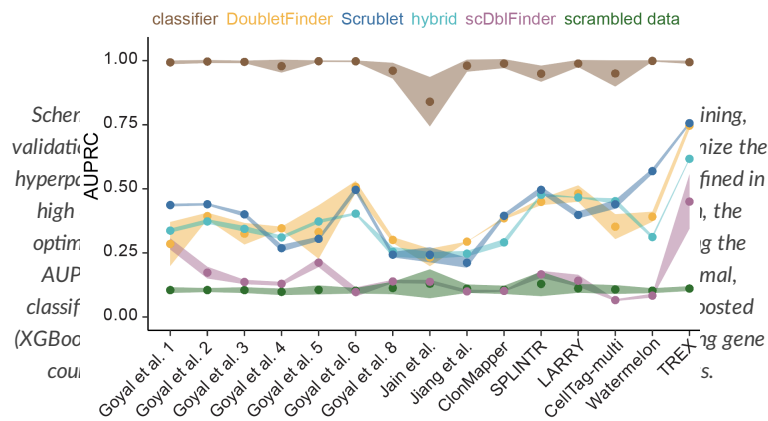


How well does it perform?

measured by AUPRC. Each dot represents the average AUPRC score of a doublet We achieved significantly higher AUPRC and AUPOT scores than using a classifier the compared to the other methods. Each marker for doublet detection method.

Classifying doublets in non-barcoded datasets

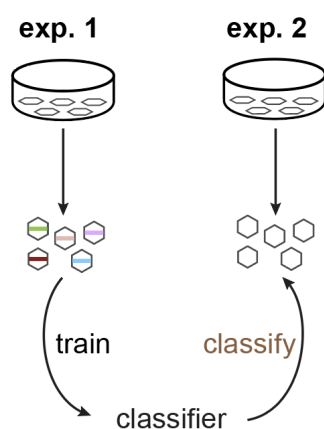
Although barcoding experiments are becoming increasingly prevalent, they are still relatively uncommon. Therefore, we sought to train a doublet classifier on barcoded data that could detect doublets in non-barcoded data. We trained a classifier on cell samples from melanoma mouse brain leukemia and bone



How well does it perform? DoubletFinder is compared to other doublet detection methods as measured by AUPRC. Each dot represents the average AUPRC score of a doublet detection method (with higher AUPRC scores indicating better performance). The lines represent the standard deviation of the AUPRC scores for each method.

Classifying doublets in non-barcoded datasets

Although barcoding experiments are becoming increasingly prevalent, they are still relatively uncommon. Therefore, we sought to train a doublet classifier on barcoded data that could detect doublets in non-barcoded data. We trained a classifier on cell samples from melanoma, mouse brain, leukemia, and bone marrow/leukemia using true singlet labels from singletCode and successfully identified doublets in a similar cell sample without needing barcoding data. Further, we integrated all the data from mouse and human samples together, split it in half, and it to train a classifier that was then tested on individual sample components of the opposite half. This, too, proved to outperform other doublet detection methods.



Schematic of training a doublet classifier on barcoded data from 1 experiment and using that classifier to identify doublets in a biological replicate, experiment 2.

Non-barcoded datasets

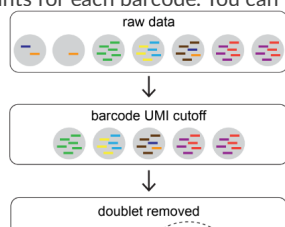
As barcoded scRNA-seq data becomes more abundant, experimenters can train a classifier specific to their cell type using a barcoded subset of data, and, eventually, there could be enough data to train a classifier on multiple barcoded cell types which can be used more generally. Find more information about this in the classifier page [here](#).

implementing singletCode on your data benchmarking doublet detection methods

singletCode is a framework that can be extended to most scRNA-seq or even scATAC-seq data to accurately identify singlets.

Barcoded datasets

To simulate doublets for benchmarking, we randomly selected the gene expression counts data from two cells that were found to be true singlets by singletCode. We averaged the counts from these two cells to generate simulated doublets and we created both a python package called singletCode and a command-line interface. The only input needed is a .csv file with cell ID, barcode, sample information, and UMI counts for each barcode. You can find more information to use it [here](#)



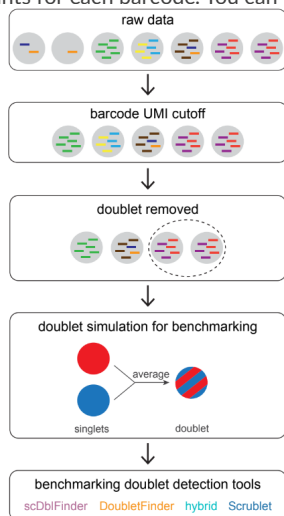
Non-barcoded datasets

As barcoded scRNA-seq data becomes more abundant, experimenters can train a classifier specific to their cell type using a barcoded subset of data, and, eventually, there could be enough data to train a classifier on multiple barcoded cell types which can be used more generally. Find more information about this in the classifier page [here](#).

implementing singletCode on your data benchmarking doublet detection methods

singletCode is a framework that can be extended to most scRNA-seq or even simulation datasets to accurately identify singlets.

To simulate doublets for benchmarking, we randomly selected the gene expression counts data from two cells that were found to be true singlets by singletCode. We averaged the counts from these two cells to generate simulated doublets and we created both a python package called singletCode and a command-line interface. The only input needed is a .csv file with cell ID, barcode, sample information, and UMI counts for each barcode. You can find more information to use it [here](#)



Workflow schematic of how raw data is processed to generate true singlets based on a barcode UMI cutoff and doublet removal based on singletCode specifications, followed by the simulation of doublets for benchmarking doublet detection methods.

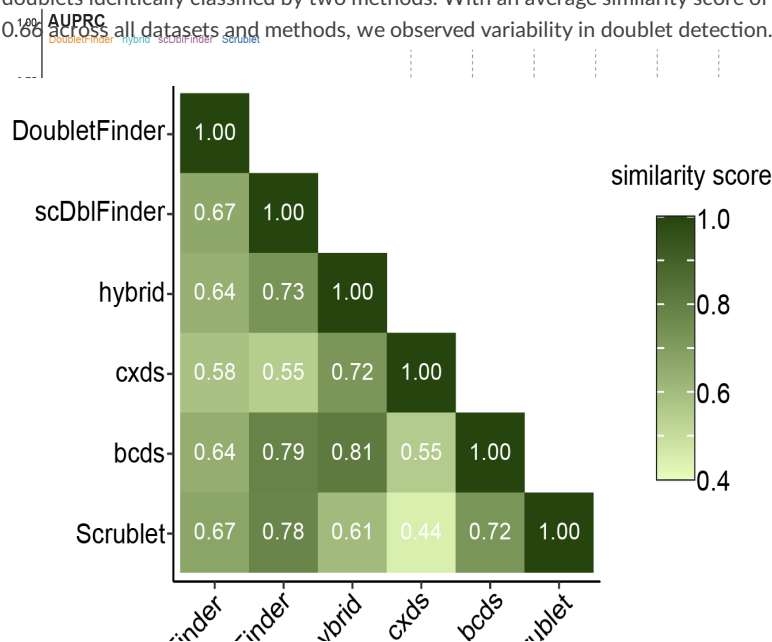
Benchmarking

We used these datasets to benchmark four doublet detection methods:

1. [scDbtFinder](#),
2. [DoubletFinder](#),
3. [Scrublet](#) and
4. [Hybrid](#)

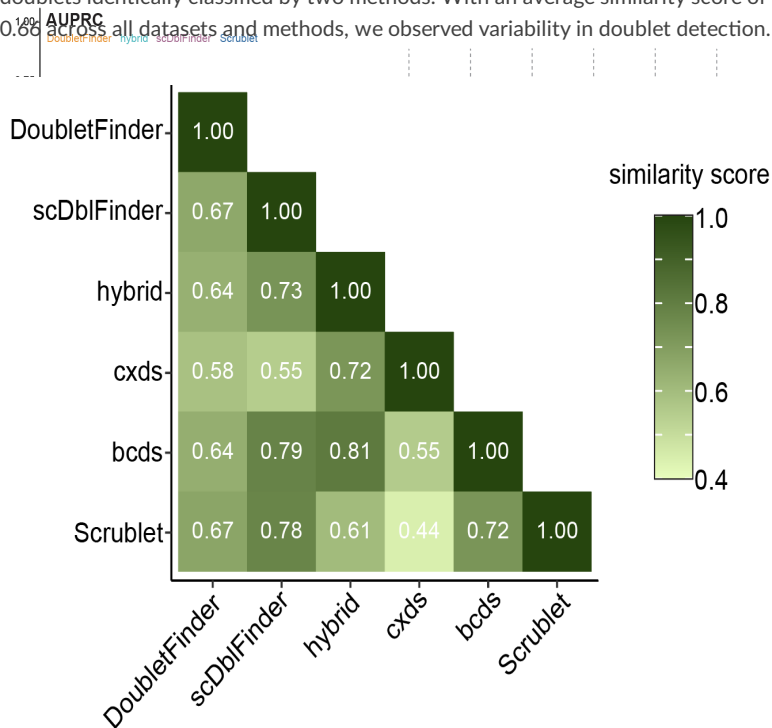
We evaluated the AUPRC, AUROC, TNR, and doublet scores and calls of the four methods and found lower than expected performance for all methods. A plot of our results for AUPRC value is found below.

We examined the consistency of doublet labeling across different doublet detection methods by introducing a similarity score—a measure of the fraction of doublets identically classified by two methods. With an average similarity score of 0.66 across all datasets and methods, we observed variability in doublet detection.



methods and found lower than expected performance for all methods. A plot of our results for AUPRC value is found below.

We examined the consistency of doublet labeling across different doublet detection methods by introducing a 'similarity score'—a measure of the fraction of doublets identically classified by two methods. With an average similarity score of 0.66 across all datasets and methods, we observed variability in doublet detection.

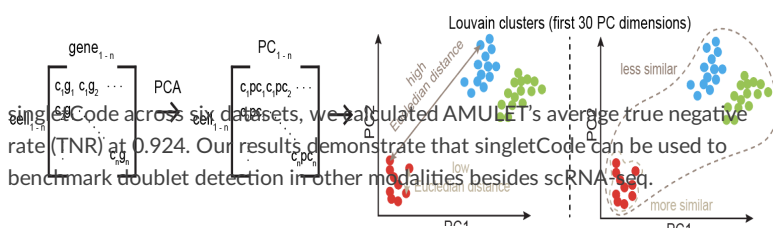


Pairwise similarity score between all benchmarked methods. A score of 1 indicates perfect identity and a 0 indicates complete disagreement. Average similarity score across all methods and datasets is 0.66.

We further evaluated doublet detection on ensemble doublet detection methods (hybrid, Chord) and across sequencing technologies (10X Genomics, Smart-seq3). For a more detailed evaluation of our results, refer to our paper.

Heterogeneity effects

We wanted to know whether heterogeneity of a dataset affects the performance of doublet detection methods. Because heterogeneity can be impacted by many properties of a dataset, such as experimental design and data processing, we made conclusions based on heterogeneity within a sample. We did this by subsampling singlets and doublets within a single PC cluster for a sample (less heterogeneous, low Euclidean distance), and across all clusters for a sample (more heterogeneity, higher Euclidean distance).

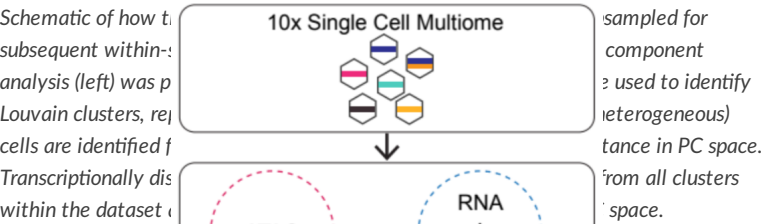


Schematic of how the 10x Single Cell Multiome dataset workflow resulting in TNR calculation. The same cells in the Multiome dataset undergo scATAC-seq, scRNA-seq, and barcode sequencing. Barcodes were processed with singletCode to identify ground truth singlets.

singletCode for

We evaluated our datasets using AMULET and singletCode. Watermelon-barcode scATAC-seq fragments identified true singlets. By comparing the true negative and false positive rates between AMULET and singletCode, we calculated AMULET's average true negative rate (TNR) at 0.924. Our results demonstrate that singletCode can be used to benchmark doublet detection in other modalities besides scRNA-seq.

singleCode across six datasets, we calculated AMULET's average true negative rate (TNR) at 0.924. Our results demonstrate that singleCode can be used to benchmark doublet detection in other modalities besides scRNA-seq.



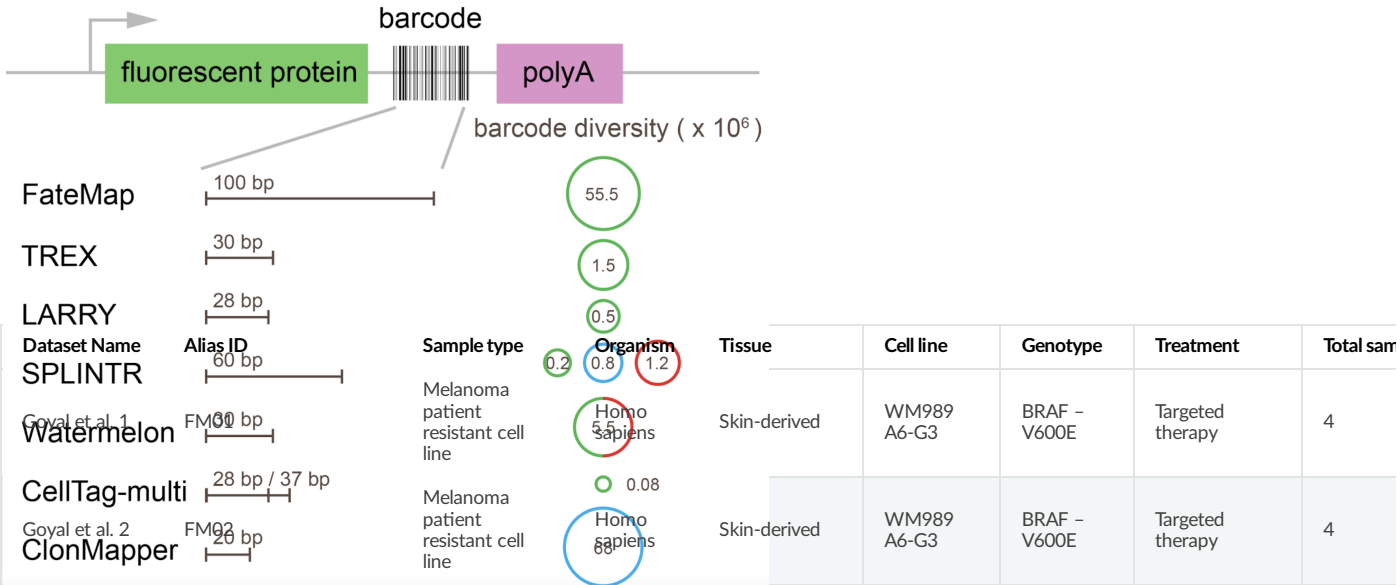
singleCode for

We evaluated our datasets using AMULET and singleCode. Watermelon-barcode scATAC-seq fragments identified true singlets. By comparing the true negative and false positive rates between AMULET and singleCode, we calculated TNR. The same cells in the Multiome dataset undergo scATAC-seq, scRNA-seq, and barcode sequencing. Barcodes were processed with singleCode to identify ground truth singlets while the respective fragments file from ATAC was used to label doublets with AMULET. The AMULET-labeled singlets were compared with ground truth singleCode singlets to calculate TNR.

dataset information

barcoding technologies assessed

7 different barcoding technologies We incorporated datasets using FateMap, ClonMapper, SPLINTR, LARRY, CellTag-multi, Watermelon, and TREX barcodes. Each have their own recovery and analysis specifications that we outline in Box 1 of our paper.



Schematic of the structure of different lineage-tracing barcodes, their length, diversity, and fluorescent color. Each barcode contains an expressed fluorescent protein, a unique sequence, known as the barcode, and a polyA tail for capture. The lengths of barcodes in each method are shown as scale bars, and library diversities are shown as circles, sized proportional to their diversity and colored according to fluorescent proteins used.

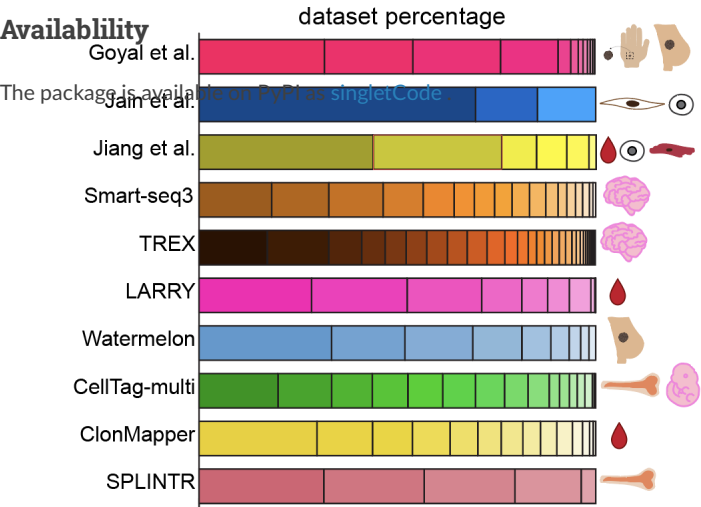
dataset details

Specifications for barcoded and non-barcoded datasets

| Dataset Name | Alias ID | Sample type | Organism | Tissue | Cell line | Genotype | Treatment | Total sam |
|----------------|----------|---|--------------|-----------------------|---------------|--|-------------------------|-----------|
| Goyal et al. 3 | FM03 | Melanoma patient resistant cell line | Homo sapiens | Skin-derived | WM989 A6-G3 | BRAF - V600E | Targeted therapy | 2 |
| Goyal et al. 4 | FM04 | Breast cancer patient resistant cell line | Homo sapiens | Breast; mammary gland | MDA-MB-231-D4 | BRAF - G464V, KRAS - G13D, TERT - c.1-124C>T, TP53 - R280K | Chemotherapy, cytotoxic | 2 |

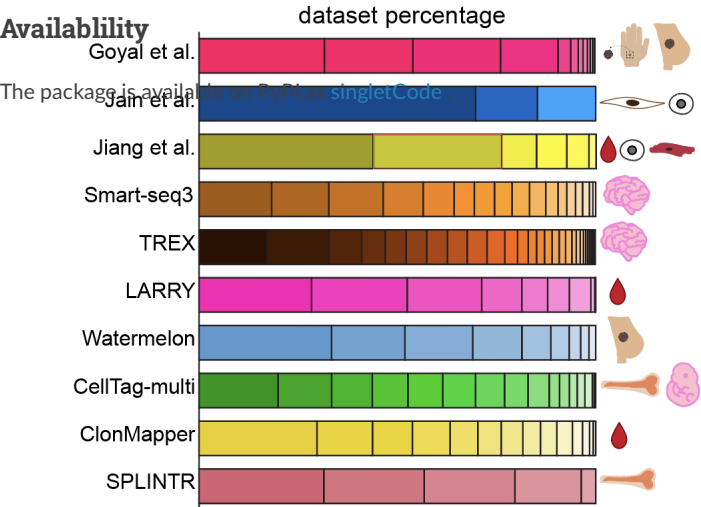
| Dataset Name | Alias ID | Sample type | Organism | Tissue | Cell line | Genotype | Treatment | Total samples |
|---|----------------------------------|---|----------------------------|--|-----------------|---------------------------------------|---|---------------|
| SPLINTR | SPLINTR_clone | Acute myeloid leukemia | Mus musculus | Bone marrow | C57BL/6 | MLL-AF9 + KrasG12D, MLL-AF9 + Flt3ITD | No treatment | 4 |
| Smart-seq3 | Smart-seq3 | Epithelial progenitor cells differentiation | Mus musculus | Neural | CD-1 | | No treatment | 2 |
| SPLINTR | nonbarcoded sc-RNA seq datasets | Acute myeloid leukaemia | Mus musculus | Bone marrow and spleen | C57BL/6 | MLL-AF9 + KrasG12D | No treatment | 1 |
| non-barcoded Watermelon | hm-12k Watermelon | Breast cancer dataset resistant cell line | Homo sapiens, Mus musculus | Human kidney, Murine mammary gland | HEK293T, NIH3T3 | PIK3CA - H1047R; Tp53 - L194F | No treatment targeted therapy | 16 |
| non-barcoded CellTag-multi | hm-6k CellTag-multi_d | Synthetic multi-lineage differentiation (hematopoiesis) | Homo sapiens, Mus musculus | Human kidney, Murine bone marrow | HEK293T, NIH3T3 | WT | No treatment cytokines and growth factors | 13 |
| non-barcoded CellTag-multi | HMEC-orig-MULTI CellTag-multi_B4 | HMECs iEP reprogramming | Homo sapiens musculus | mammary gland embryo | HMEC C57BL/6J | WT | No treatment EGF | 11 |
| non-barcoded | HMEC-rep-MULTI | HMECs | Homo sapiens | mammary gland | HMEC | WT | No treatment | 1 |
| ClonMapper | ClonMapper | CLL resistant cell lines | Homo sapiens | blood | HG3 | | Chemotherapy | 4 |
| non-barcoded | HEK-HMEC-MULTI | Synthetic dataset | Homo sapiens | Human kidney, mammary gland | HEK293T, HMEC | WT | No treatment | 1 |
| non-barcoded | mkidney-ch | Mouse kidney cells | Mus musculus | Mouse kidney | C57BL/6J | WT | No treatment | 1 |
| non-barcoded | pbmc-2ctrl-dm | Patient PBMCs | Homo sapiens | Systemic lupus erythematosus (SLE) PBMCs | Patient-derived | Patient-derived | No treatment | 1 |
| non-barcoded | pbmc-2stim-dm | Patient PBMCs | Homo sapiens | Systemic lupus erythematosus (SLE) PBMCs | Patient-derived | Patient-derived | No treatment | 1 |
| non-barcoded | cline-ch | Human HEK, K562, KG1, and THP1 | Homo sapiens | blood | human - derived | human - derived | No treatment | 1 |
| non-barcoded | pbmc-ch | Human PBMCs | Homo sapiens | blood | human - derived | human - derived | No treatment | 1 |
| non-barcoded | pdx-MULTI | Synthetic dataset | Homo sapiens, Mus musculus | Human breast cancer, mouse immune | PDX mouse model | PDX mouse model | No treatment | 1 |
| non-barcoded | nuc-MULTI | Synthetic dataset | Homo sapiens, Mus musculus | nuclei | | | No treatment | 1 |
| Stacked bar chart of cell types in all datasets. Each distinct color corresponds to one study, with the proportions of different cell types ordered from most to least prevalent. Adjacent to the bar chart, skeletal tissues and cells provide a visualization of the cell types used in the particular study. | | | | | MCF7 | GATA3 - D336G; PIK3CA - E545K | Targeted therapy | 6 |

singleCode PyPI package



| | | | | | | | | |
|---|-------------------|---|--------------|------------------------|------|-------------------------------|------------------|---|
| Stacked bar chart of cell types in all datasets. Each distinct color corresponds to one study, with the proportions of different cell types arranged from most to least prevalent; (deepest to lightest) and separated by vertical lines. Adjacent to the bar chart, skeletal tissues and cells provide a visualization of the cell types used in the particular study. | Multiple datasets | Breast cancer patient resistant cell line | Homo sapiens | Bladder, mammary gland | MCF7 | GATA3 - D336G; PIK3CA - E545K | Targeted therapy | 6 |
|---|-------------------|---|--------------|------------------------|------|-------------------------------|------------------|---|

singletCode PyPI package



How do I use it?

1. Installation

It can be installed from PyPI using the following in the terminal:

```

pip3 install singletCode
from singletCode import check_sample_sheet, get_singlets

```

2. Preparing the input sample sheet.

The input sample sheet is a .csv file that contains the information about cell ID (added while sequencing), lineage barcode, and sample name. Each row should be repeated n times where n is the number of UMIs associated with that barcode and cell ID combination. For creating the input, you can read the .csv file in as a pandas dataframe.

```

import pandas as pd
df = pd.read_csv("path/to/csv/file.csv")

```

You can check if the format and the column names are valid for running

How do I use it?

1. Installation

It can be installed from PyPI using the following in the terminal:

```
pip3 install singletCode
from singletCode import check_sample_sheet, get_singlets
```

2. Preparing the input sample sheet.

The input sample sheet is a .csv file that contains the information about cell ID (added while sequencing), lineage barcode, and sample name. Each row should be repeated n times where n is the number of UMIs associated with that barcode and cell ID combination. For creating the input, you can read the .csv file in as a pandas dataframe.

```
import pandas as pd
df = pd.read_csv("path/to/csv/file.csv")
```

You can check if the format and the column names are valid for running `singletCode` using `check_sample_sheet()` function.

```
check_sample_sheet(df)
```

It will either return an error with information about how to modify the sample sheet to make it a valid input or will print the message, “The sample sheet provided can be used as input to `get_singlets` to get a list of singlets identified”, in which case you can move to the next step.

3. Running `get_singlets()` to get an assignment of singlet status for each cell ID and barcode combination.

This is the step where singlet identification is done using the `singletCode` framework. You can read more about the parameters for `get_singlets()` below. In this example, default values are used.

```
cellLabelList, stats = get_singletons(df, dataset_name = "Sample1")
```

cellLabelList is a pandas dataframe that contains 5 rows: cellID, barcode, sample, nUMI and label. Label is whether the particular cell ID and barcode combination has been called a single or not.

Detailed information about the parameters for the functions in the package

To get a dataframe containing just the singlets, you can run this:

get_singlets

```
singletList = cellLabelList[cellLabelList["label"] == "Singlet"]
```

```
get_singlets(sample_sheet, output_path=None, dataset_name=None,
save_all_singlet_categories=False, save_plot_umi=False, umi_cutoff_method='ratio',
umi_cutoff_ratio=7.5e-06, umi_cutoff_percentile=None, min_umi_cutoff=2,
umi_diff_threshold=50, umi_dominant_threshold=10) \[source\]
```

Function that inputs the sample sheet and other parameters and runs it through count_duplicates to get a list of singlets in the sample. If a row is repeated, it is assumed to reflect the UMI associated with a barcode in this cell (identified by the cellID). Indeterminate cells since singletCode can identify only truly singlet cells

Parameters `sample_sheet` (DataFrame) A data frame that contains 3 columns: cellID, barcode, sample.

You can save the `dataset_name` (str), `file_name` and the list of `file_ids` using each category of singlets by setting `save_all_singlet_categories` to `TRUE`. This will add a new column in the `singlet_stats` sheet returned.

- **output_path** (*str*, *optional*) – The path to store any output files, including plots to show UMI distribution and what the umi_cutoff used is, csv files containing singlets of different categories. If None, then the list of singlets will

and barcode combination has been called a singlet or not.

Detailed information about the parameters for the functions in the package

To get a dataframe containing just the singlets, you can run this:

get_singlets

```
singletList = cellLabelList[cellLabelList["label"] == "Singlet"]
```

get_singlets(*sample_sheet*, *output_path*=None, *dataset_name*=None, *save_all_singlet_categories*=False, *save_plot_umi*=False, *umi_cutoff_method*='ratio', *umi_cutoff_ratio*=7.5e-06, *umi_cutoff_percentile*=None, *min_umi_cutoff*=2, *umi_diff_threshold*=50, *umi_dominant_threshold*=10) [\[source\]](#)

Function that inputs the sample sheet and other parameters and runs it through different categories of singlets (such as single-barcode singlets, multi-barcode singlets, dominant-UMI singlets), number of cells removed due to repeated, it is assumed to reflect the UMI associated with a barcode in this cell (identified by the cellID), to low barcode UMI counts for the barcode and number of indeterminate cells since singletCode can identify only truly singlet cells

Parameters

- sample_sheet** (*DataFrame*) – A dataframe that contains 3 columns: cellID, barcode, sample.
- dataset_name** (*str*) – The name of the dataset being used. This file name is used to save all singlet categories files and the singlet_stats sheet returned.

- output_path** (*str, optional*) – The path to store any output files, including plots to show UMI distribution and what the umi_cutoff used is, csv files containing singlets of different categories. If None, then the list of singlets will be returned but it won't contain information about what category of singlet each cell is. Defaults to None.
- save_all_singlet_categories** (*bool, optional*) – If true, then singlets of each category are saved separately in csv files along with all singlets and all non-singlets. Defaults to False.
- save_plot_umi** (*bool, optional*) – If true, then plots showing UMI distribution indicating the UMI cutoff used will be saved for each sample. Defaults to False.
- umi_cutoff_method** (*str, optional*) – Specify if quality control for barcodes using UMI counts should be based on "ratio" or "percentile". Defaults to 'ratio'.
- umi_cutoff_ratio** (*float, optional*) – The ratio used to determine the umi_cutoff if umi_cutoff_method is "ratio". Defaults to 3/4e5.
- umi_cutoff_percentile** (*float, optional*) – If umi_cutoff_method is "percentile", then the umi_cutoff will be the minimum UMI count required to be in the top umi_cutoff_percentile'th percentile. There is no default and if umi_cutoff_method is set to "percentile", then manually set this parameter.
- min_umi_cutoff** (*int, optional*) – This is the absolute minimum number of UMIs that need to be associated with a barcode for it to be considered a barcode.

Return type tuple
However, the actual umi_cutoff used will be the greater of min_umi_cutoff and the cutoff calculated using umi_cutoff_method. Defaults to 2.

check_sample_sheet

check_sample_sheet(*data_frame*) [\[source\]](#)

Function to check if the data frame can be used for get_singlets function. It checks if the data frame has potential dominant barcode present and in same cell. If this data frame can be used for get_singlets function, then a status of 10 is returned. This is the minimum difference between UMI counts associated with a potential dominant barcode present within a cell and the median UMI count of all barcodes associated with the cell. If a cell has only one dominant barcode, it will be counted. Defaults to 50.

command line interface

Availability

A 2-tuple containing:

The source code for the command line interface is available on [GitHub repository \(repo\)](#). You can clone the repository and the code locally along with singlet assignment to each cell ID.

```
git clone https://github.com/GoYallan/singletCodeTools
```

- pandas DataFrame**: A dataframe which contains the statistics for total singlets, different categories of singlets, and cells removed due to low UMI counts.

Return type tuple

check_sample_sheet with a barcode for it to be considered a barcode. However, the actual umi_cutoff used will be the greater of min_umi_cutoff and the cutoff calculated using umi_cutoff_method. Defaults to 2.

check_sample_sheet(data_frame) [source] old (int, optional) – The minimum UMI count to be associated with a barcode singlet. If a cell has more than one potential dominant barcode present within a cell and the difference between UMI counts associated with a potential dominant barcode present within a cell and the median UMI count of all barcodes associated with the cell is greater than old, it will be counted. Defaults to 50.

Parameters sample_sheet – A dataframe that contains your sample sheet.

Returns None

command line interface

Availability

A 2-tuple containing:

The source code for the commandLineTools folder, you will find 3 files and you will need to run singletCodeCommandLine.py. There are 2 modules available, one to run singletCode and the other to create a sample sheet if you have used Watermelon barcoding technology using the fastq sequenced files from MISEQ.

git clone https://github.com/GoyalLab/singletCodeTools

pandas.DataFrame: A dataframe which contains the statistics for total singlets, different categories of singlets, and cells removed due to low UMI counts.

Using the interface

Navigating to commandLineTools folder, you will find 3 files and you will need to run singletCodeCommandLine.py. There are 2 modules available, one to run singletCode and the other to create a sample sheet if you have used Watermelon barcoding technology using the fastq sequenced files from MISEQ.

Detailed information about the modules in the command line interface

Description:

This script contains two modules: - Count Module: Generates the singlet files for the input data sheet. - Watermelon Module: Uses the MiSeq dial-out files to create the cell ID, barcode, and sample file. These outputs can then be used as input for the singlet code module.

Usage:

For the Count Module:

```
python singletCode.py count -i /path/to/input.txt -o /path/to/output
```

For the Watermelon Module:

```
python fastq_files_to_csv.py -i /path/to/fastq/files -o /path/to/save/csv/file -s path/sample/sheet -use10X False -input10X path/to/barcodes/tsv
```

-o, --outputFolder: Specify the path to save the output CSV file containing the barcode, cell ID information.

-outputName: Specify the name of the output CSV file.

Options for Count Module:

-use10X: Specify if a 10X object is provided which has the same cells as those in fastq. If provided, cells in the fastq file will be filtered out if not present in the 10X object.

-input10X: Specify the path to the input barcode file.

-output10X: Specify the path to the output barcode file.

-f, --force: Force overwrite if the output file already exists.

Authors:

-cutoff: UMI cutoff ratio.

-d, --umi_diff_threshold: Minimum difference in UMI between UMI count for dominant UMI and median UMI count.

-m, --min_umi_good_data_cutoff: Minimum UMI count for a barcode to be considered a barcode.

Vignette to use singletCode Command Line tool to analyse single-cell data with watermelon barcodes

Options for Watermelon Module:

-i, --inputFolder: Specify the path to the folder containing the fastq folders

The barcode region is assumed to be amplified using Illumina MiSeq.

All the data for this vignette and the files output from it can be downloaded [here](#).

For the **sampleSheet**, specify the path to the sample sheet in .csv format that contains sample name and sample number. This should match the names of

```
python3 singletCodeCommandLine.py path/to/inputFiles/fastq/files -o path/to/save/csv/file -s path/sample/sheet -use10X False -input10X path/to/barcodes/tsv
```

- **-outputFolder**: Specify the path to save the output CSV file containing the barcode, cell ID information.

- **-outputName**: Specify the name of the output CSV file.

- **-use10X**: Specify if a 10X object is provided which has the same cells as those in fastq. If provided, cells in the fastq file will be filtered out if not

Options for Count Module:

- **-input10X**: Specify the path to the input barcode file.
- **-input10XPath**: Specify the path to the barcode tsv which contains cell IDs.
- **-f**: Force: Force overwrite if the output file already exists.

-u, **-cutoff**: UMI cutoff ratio.

Authors:

- **-d**, **-umi_diff_threshold**: Minimum difference in UMI between UMI count for dominant UMI and median UMI count.
- **Ziyang Zhang (Charles)**
- **Keerthana M Arun**
- **-m**, **-min_umi_good_data_cutoff**: Minimum UMI count for a barcode to be

Vignette to use singletCode Command Line tool to analyse single cell data with watermelon barcodes

Options for Watermelon Module:

- **-i**, **-inputFolder**: Specify the path to the folder containing the fastq folders

The barcode region is assumed to be amplified using Illumina MiSeq.

All the data for this vignette and the files output from it can be downloaded [here](#). It contains inputFiles and the outputFiles (it contains a test folder which has the expected output files). This vignette can be downloaded as a jupyter notebook from the [singletCode Tools repo](#).

Installing singletCode Command line tool

To use the singletCode command line tool, clone the repository from GitHub. Let the Path to the folder you are running this command be **Path**. We can also install other packages needed to run the tool.

```
!git clone https://github.com/GoyalLab/singletCodeTools
Path = "path/to/singletCodeTools/repo"
%conda install scipy tqdm matplotlib biopython python-levenshtein pandas
```

Next step is to understand the samples present in the FASTQ files.

The sample fastq files are in the inputFolder. We can identify the sample name and number from the FASTQ file. For example, sampleName_S1_L001_R1_001.fastq.gz means that the sample name is sampleName and sample number is 1. Make sure that both read 1 and read 2 for each sample are present in the same folder (R1 and R2).

1. **inputFiles**/ 2. **outputFolder** will be p/outputFiles/ 3. **sampleSheet** will be

p/inputFiles/sampleSheet.csv

Creating sample sheet for these two samples.

```
import subprocess
import pandas as pd
command = [
    'python',
    f'{Path}/commandLine/singletCodeCommandLine.py',
    'watermelon',
    p = ^path/to/download/unzipped/data"
    '-o', f'{p}/outputFiles',
    '-s', f'{p}/inputFiles/sampleSheet.csv',
    '--outputName', 'watermelonBarcodeUmi.csv'
]
sampleSheet = pd.read_csv(f'{p}/inputFiles/sampleSheet.csv')
sampleSheet
```

Arguments: **sampleName**, **sampleNumber**

0 **sampleName**

1 **inputFolder**

2 **otherSampleName**

Now, to run the watermelon module of singletCodeTools, you need to run this

command. If we are going by the folder structure of the zipped file and p is path to

the unzipped folder containing example files, then 1. **inputFolder** will be

```
sampleSheet:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette
outputName: watermelonBarcodeUmi.csv
use10X: False
```

2. **outputFolder** will be p/outputFiles/ 3. **sampleSheet** will be p/inputFiles/sampleSheet.csv

Creating sample sheet for these two samples.

```
import subprocess
import pandas as pd
command = [
    'python',
    f'{Path}/commandLine/singletCodeCommandLine.py',
    'watermelon',
    p = 'path/to/download/unzipped/data'
    '-o', f'{p}/outputFiles',
    '-s', f'{p}/inputFiles/sampleSheet.csv',
    '--outputName', 'watermelonBarcodeUmi.csv'
]

sampleSheet = pd.read_csv(f'{p}/inputFiles/sampleSheet.csv')
subprocess.run(command)
```

Arguments received:
sampleName: sampleNumber
0 sampleName
1 inputFolder:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignettes
Now, to run the watermelon module of singletCodeTools, you need to run this command. If we are going by the folder structure of the zipped file and p is path to the unzipped folder containing example files, then 1. inputFolder will be
sampleSheet:
outputName: watermelonBarcodeUmi.csv
use10X: **False**
input10X: **None**
All the inputs **for** the command are valid **and** will proceed **with** creating the barcode sheet **for** all the samples **in** the sheet.
Filtered rows of dataframe: 940
Filtered rows of dataframe: 718

```
result = subprocess.run([
    'python',
    f'{Path}/commandLine/singletCodeCommandLine.py',
    'watermelon',
    '-i', f'{p}/inputFiles/',
    '-o', f'{p}/outputFiles/',
    '-s', f'{p}/inputFiles/sampleSheet.csv',
    '--outputName', 'watermelonBarcodeUmiWith10X.csv',
    '--use10X',
    '--input10X', f'{p}/inputFiles/barcodes.tsv'
], capture_output=True, text=True)

# Check if the command was successful
if result.returncode == 0:
    print("Command executed successfully")
    print("Output:\n", result.stdout)
else:
    print("Command failed")
    print("Error:\n", result.stderr)
```

Command executed successfully
Output:
Arguments received:
command: count
input_file:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignettes
out_prefix:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignettes
umi_cutoff_ratio: 7.5e-06
umi_diff_threshold: 50
dominant_threshold: 10
sampleSheet:
outputName: watermelonBarcodeUmiWith10X.csv
use10X: **True**
input10X: **True**
All the inputs **for** the command are valid **and** will proceed **with** creating the barcode sheet **for** all the samples **in** the sheet.
Filtered rows of dataframe: 791
Filtered rows of dataframe: 632

Command executed successfully
Output:
Arguments received:
command: count
input_file:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignettes
out_prefix:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignettes
umi_cutoff_ratio: 7.5e-06
umi_diff_threshold: 50
dominant_threshold: 10

```

Command executed successfully
import subprocess
Arguments received:
command: count
input_file:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/outputPrefix:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/sampleSheet/outputFiles/watermelon'
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/outputName: watermelonBarcodeUmiWith10X.csv
# Use X if the command was successful
if subprocess.run([
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/sampleSheet/outputFiles/watermelon'
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/outputName: watermelonBarcodeUmiWith10X.csv
]).returncode == 0:
All phenotypes for the removed are valid and will proceed with creating the
basecode sheet for all the samples in the sheet.
Filtering from data frame 791
Filtering from data frame 639

```

```

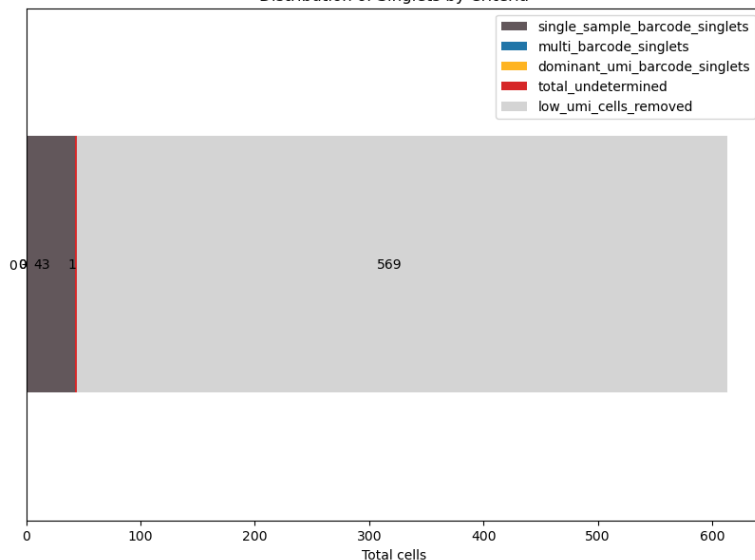
Command executed successfully
Output:
Arguments received:
command: count
input_file:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/outputPrefix:
/home/keerthana/Goyal_Lab/websiteToolData/thingsToAddToWebsite/watermelonVignette/outputPrefix:
umi_cutoff_ratio: 7.5e-06
umi_diff_threshold: 50
dominant_threshold: 10
min_umi_good_data_cutoff: 2
INFO: Raw data counts
sampleNum
sampleName          693
otherSampleName      524
Name: count, dtype: int64
INFO: Using raio based filtering.
Current Sample Adjusted UMI cutoff: 2
Total cells: 45
Sample sampleName singlet: 43
Total Singlets: 43
Total Multiplets: 1
All singlets identified are unique? True
Total Singlets: 43
Total Multiplets: 1
INFO: Using raio based filtering.
Current Sample Adjusted UMI cutoff: 2
Total cells: 22
Sample otherSampleName singlet: 22
Total Singlets: 22
Total Multiplets: 0
All singlets identified are unique? True
Total Singlets: 22
Total Multiplets: 0
All singlets identified are unique? True

```

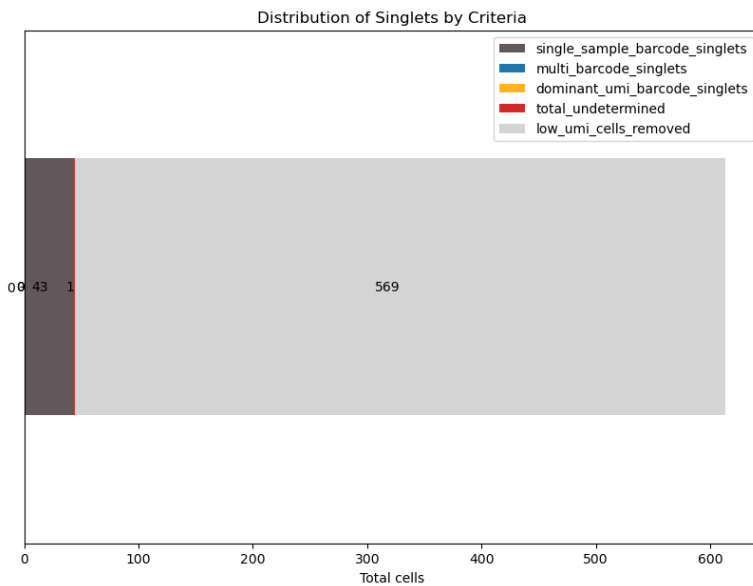
```
import matplotlib.pyplot as plt
```

```
stats = pd.read_csv(f"
```

Distribution of Singlets by Criteria



In the above plot, you see that the original data had 569 cells that were removed due to low barcode UMI count, 43 singlets with a single-barcode associated with them and 1 multiplet (singletCode could not determine if it was a singlet for sure.)



In the above plot, you see that the original data had 569 cells that were removed due to low barcode UMI count, 43 singlets with a single-barcode associated with them and 1 multiplet (singletCode could not determine if it was a singlet for sure.)

Looking at the scRNAseq data associated

Since this data has both scRNAseq and barcodes for the same cells, we can analyse them together

Installing and importing scanpy package to do this

```
#Install scanpy for further single-cell RNAseq analysis
# %conda install -c conda-forge scanpy python-igraph leidenalg
#Import scanpy
import scanpy as sc
```

In case there are version conflicts during this installation or while importing scanpy, we found `%conda update -all` to be a useful command that fixed the version conflict previously. Reading in the 10X h5ad object associated with the same watermelon data

```
adata = sc.read_h5ad(f"{p}/inputFiles/watermelonScRnaSeqData.h5ad")
adata
```

```
AnnData object with n_obs × n_vars = 1093 × 27264
```

Identifying the cells that were below the barcode UMI threshold and were filtered out. **Read in the output files to identify cells as being singlets, multiplets or being removed for low barcode UMI threshold**

```
lowUmiCells = cellidBarcodeUMI[~
First, reading in the cellidBarcodeUMI sheet generated earlier with additional
filter using scRNAseq data
cellidBarcodeUMI['cellID'].isin(otherSampleNameSinglets[0]) |

cellidBarcodeUMI['cellID'].isin(sampleNameMultiplets[0]))
cellidBarcodeUMI =
pd.read_csv(f'{p}/outputFiles/watermelonBarcodeUmiWith10X.csv')
```

Annotating the cells in adata with these labels

Reading in all the singlets and multiplets identified in the two samples. There might not always be no multiplets - check the status file to see if there are any. In this example, there are no multiplets in otherSampleNameSinglets.

```
adata.obs.loc[adata.obs.index.isin(otherSampleNameSinglets[0]),
'singletStatus'] = 'singlet'
adata.obs.loc[adata.obs.index.isin(sampleNameMultiplets[0]),
'sampleNameSinglets'] = pd.read_csv(f"
aparaubstLoc[adata.obs.index.isin(lowUmiCells[cellID])]head[0]['singletStatus']
otherSampleNameSinglets = pd.read_csv(f"
{p}/outputFiles/watermelon_otherSampleName_singlets_all.txt", header = None)
sampleNameMultiplets = pd.read_csv(f"
{p}/outputFiles/watermelon_sampleName_multiplets.txt", header = None)
```

Note that in this dataset we are not doing any actual QC, but in actual analysis

```
adata.obs['cellidBarcodeUMI'] = 1000 * 27207
```

Identifying the cells that were below the barcode UMI threshold and were filtered

Read in the output files to identify cells as being singlets, multiplets or being removed for low barcode UMI threshold

First, loading in the cellidBarcodeUMI sheet generated earlier with additional filter using scRNAseq data

```
lowUmiCells = cellidBarcodeUMI[~
cellidBarcodeUMI['cellID'].isin(otherSampleNameSinglets[0]) |
cellidBarcodeUMI['cellID'].isin(sampleNameMultiplets[0])]
cellidBarcodeUMI =
pd.read_csv(f'{p}/outputFiles/watermelonBarcodeUmiWith10X.csv')
```

Annotating the cells in adata with these labels

Reading in all the singlets and multiplets identified in the two samples. There might not always be multiplets, check the status file to see if there are any. In this example, there are no multiplets in otherSampleNameSinglets[0]),

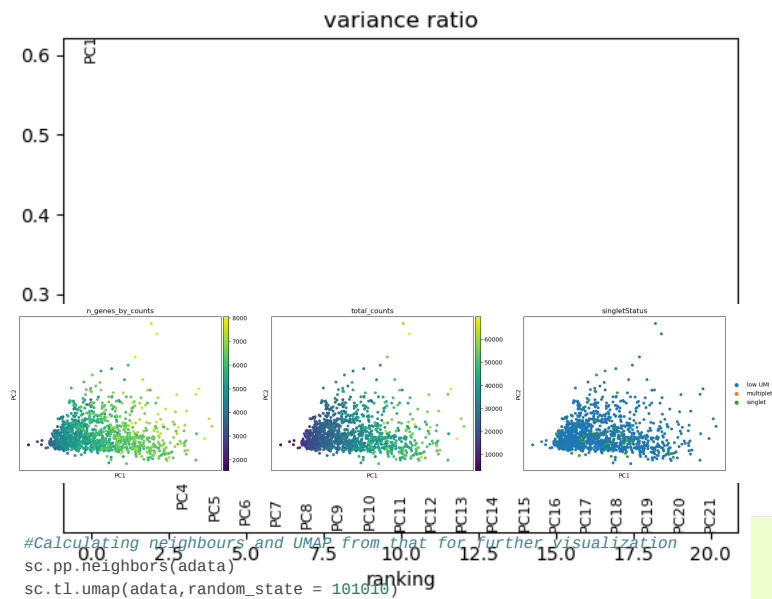
```
adata.obs.loc[adata.obs.index.isin(otherSampleNameSinglets[0]),
'singletStatus'] = 'singlet'
adata.obs.loc[adata.obs.index.isin(sampleNameMultiplets[0]),
'singletStatus'] = 'multiplet'
sampleNameSinglets = pd.read_csv(f"
{p}/outputFiles/watermelon_sampleName_singlets_all.txt", header = None)
sampleNameMultiplets = pd.read_csv(f"
{p}/outputFiles/watermelon_sampleName_multiplets.txt", header = None)
```

Note that in this vignette we are not doing any actual QC - but in actual analysis, it would need to be done.

```
sc.pp.calculate_qc_metrics(adata, inplace=True)
```

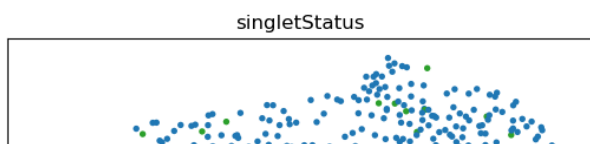
Calculating PCA and UMAP for visualization

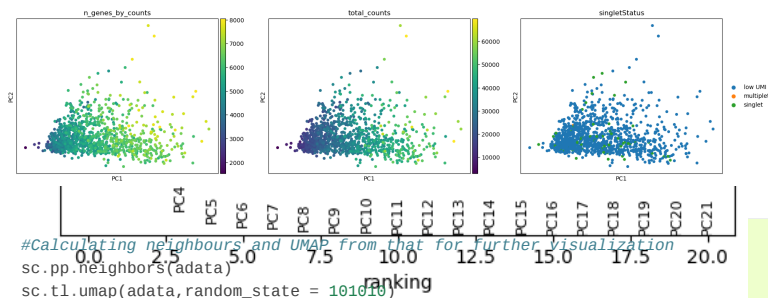
```
#Calculating PCA for the data and plotting variance ratio
sc.tl.pca(adata)
sc.pl.pca_variance_ratio(adata, n_pcs=20)
```



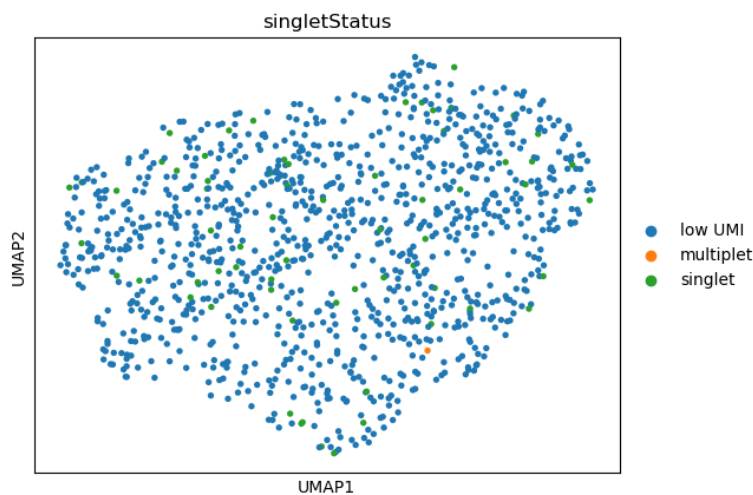
```
#Calculating neighbours and UMAP from that for further visualization
sc.pp.neighbors(adata)
sc.tl.umap(adata, random_state = 101010)
```

```
sc.pl.pca(
adata,
sc.pl.umap(adata, ['n_genes_by_counts', 'total_counts', 'singletStatus'],
size=100,
color=['singletStatus'],
size=60
)
```





```
sc.pl.pca(
    adata,
    sc.pl.umap(
        ['n_genes_by_counts', 'total_counts', 'singletStatus'],
        size=100,
        color=['singletStatus'],
        size=60
    )
)
```



Saving the final adata

```
adata.write(f"{p}/outputFiles/watermelonScRNA.h5ad")
```

Vignette to use singletCode package

The input needed to run singletCode is a .csv file that contains the information about cell ID (added while sequencing), lineage barcode, and sample name. Each row should be repeated n times where n is the number of UMIs associated with that barcode and cell ID combination. You can download a sample input sheet [here](#). It is a subset of data from Jiang Et al and details about it are described in the singletCode paper in detail. The folder also contains expected output files in the test folder within the outputFiles folder. This vignette can be downloaded as a jupyter notebook from the [singletCode Tools repo](#).

```
import pandas as pd
path = "path/to/downloaded/and/unzipped/folder"
pathToInputSheet = f"{path}/inputFiles/JiangEtAlSubset_InputSheet.csv"
df = pd.read_csv(pathToInputSheet)
```

Install singletCode package

```
!pip3 install singletCode
```

Check formatting of input sheet

check sample sheet (df)

Import necessary functions from it

```
from singletCode import check_sample_sheet, get_singlets
```

The sample sheet provided can be used as input to get_singlets to get a list of singlets identified.

The sample sheet provided can be used as input to get_singlets to get a list of singlets identified.

Identify singlets from input sheet

that barcode and cell ID combination. You can download a sample input sheet [here](#). It is a subset of data from Jiang Et al and details about it are described in the singletCode paper in detail. The folder also contains expected output files in the test folder within the outputFiles folder. This vignette can be downloaded as a jupyter notebook from the [singletCode Tools repo](#).

```
path = "path/to/downloaded/and/unzipped/folder"
pathToInputSheet = f"{path}/inputFiles/JiangEtAlSubset_InputSheet.csv"
df = pd.read_csv(pathToInputSheet)
```

Install singletCode package

```
!pip3 install singletCode
```

Check formatting of input sheet

check_sample_sheet(df)

Import necessary functions from it

```
from singletCode import check_sample_sheet, get_singlets
```

The sample sheet provided can be used as input to get_singlets to get a list of singlets identified.

The sample sheet provided can be used as input to get_singlets to get a list of singlets identified.

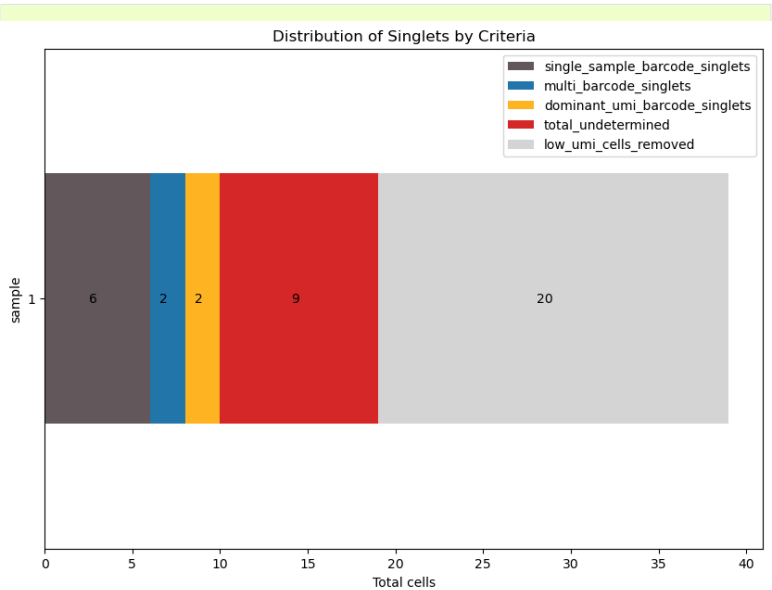
Identify singlets from input sheet

```
outputPath = "path/to/output/folder"
cellLabelList, stats = get_singlets(df, dataset_name= "JiangEtAlSubset",
save_all_singlet_categories = True, output_path=outputPath)
```

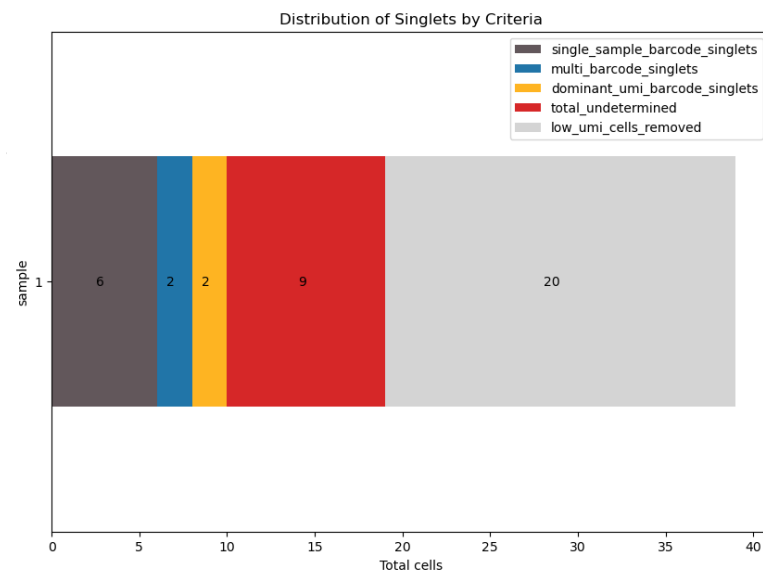
```
INFO: Raw data counts:
sample
1    1306
Name: count, dtype: int64
Total cells for sample 1: 39
INFO: Using ratio based filtering.
Current Sample Adjusted UMI cutoff: 2
```

```
100%|██████████| 122/122 [00:00<00:00, 11475.27it/s]
```

```
All singlets identified with multiple barcodes are unique? True
Total Singlets: 10
Total Multiplets: 9
```



The above plot shows that the data we had contained different kind of singlets: 6 single-barcode cells, 2 cells which had more than one barcode but with same combination being present in more than one cell, 2 cells that had one dominant barcode. The data also contained 9 cells which singletCode could not determine



```
ax.set_xlabel('Total cells')
The above plot shows that the data we had contained different kind of singlets: 6
plt.show()
single-barcode cells, 2 cells which had more than one barcode but with same
combination being present in more than one cell, 2 cells that had one dominant
barcode. The data also contained 9 cells which singletCode could not determine
as being truly singlets and 20 cells whose barcode UMI counts were below the set
threshold.
```

Understanding the output files

To understand some of the files in the output, we can look at cell IDs and their data in the original input sheet

For the dominant_umi_singlets, there are two cell IDs. One of them is TGTAAAGCGTCTCGCGA. If we look at that entry in the input sheet and count the number of UMI associated with each barcode, we see that one barcode has 99 UMI counts while the second highest UMI count is 7. So, the cell most likely has only one barcode associated with it and hence, a singlet.

```
import pandas as pd
df[df['cellID'] == 'TGTAAAGCGTCTCGCGA'].groupby(['cellID', 'barcode',
'sample']).size().reset_index(name='count').sort_values('count',
ascending=False).reset_index(drop=True)
```

| | cellID | barcode | sample | count |
|-----|-------------------|--|--------|-------|
| 0 | TGTAAAGCGTCTCGCGA | ATTGTGTGTTGCAGATGCAGTTGATGCTGATGAAGTTGTACAAGGTC... | 1 | 99 |
| 1 | TGTAAAGCGTCTCGCGA | ATTCGACTTGATCTTCTAGAACATGGTGAAGTAGCAGGTGCTGATC... | 1 | 7 |
| 2 | TGTAAAGCGTCTCGCGA | ATACAGCTCAAGCAGTACTACTTTCGCTTTCATGCAGAACAAAC... | 1 | 6 |
| 3 | TGTAAAGCGTCTCGCGA | ATAGATGCACCTTGGTGGTCGAGTTCTAGTTGTAGCTGATCGTCCAG... | 1 | 6 |
| 4 | TGTAAAGCGTCTCGCGA | ATTCGACCAGAACACATGCAGTTCAACGTGTTTCGAGGTGTAGATG... | 1 | 6 |
| ... | ... | ... | ... | ... |
| | cellID | barcode | sample | count |
| 92 | TGTAAAGCGTCTCGCGA | ATAGGAGTAGTCTGATGTTGACCAATACCTGAGCTGAGAG... | 1 | 13 |
| 93 | TGTAAAGCGTCTCGCGA | ATAGGAGTAGTCTGATGTTGACCAATACCTGAGCTGAGAG... | 1 | 13 |
| 94 | TGTAAAGCGTCTCGCGA | ATAGGAGTAGTCTGATGTTGACCAATACCTGAGCTGAGAG... | 1 | 21 |
| 84 | TGTAAAGCGTCTCGCGA | ATAGTACATGGTGGACCTGGACTTCGAGATGGAGCTCTTGTTCCTG... | 1 | 1 |
| 85 | TGTAAAGCGTCTCGCGA | ATAGGAGTAGTCTGATGTTGACCAATACCTGAGCTGAGAG... | 1 | 1 |
| 86 | TGTAAAGCGTCTCGCGA | ATAGGAGTAGTCTGATGTTGACCAATACCTGAGCTGAGAG... | 1 | 1 |

Further single-cell RNAseq analysis with both scRNAseq data and singlet
information from singletCode output

Next, we can look at multi-barcode singlets. There are two cell IDs:

AGGCTGCTCTTTCCGG and GAGGGATGTAACATCC. If we look at the barcodes with greater than 2 UMI counts, we see that they have the same combination. The only way this can occur is if a cell receives multiple barcode initially and then divides.

```
#Read the scRNAseq data in h5ad format
import pandas as pd
df[df['cellID'] == 'AGGCTGCTCTTTCCGG']
import scanpy as sc
sc.read_h5ad('h5ad')
df = df[df['cellID'] == 'AGGCTGCTCTTTCCGG']
df.groupby(['cellID', 'barcode', 'sample']).size().reset_index(name='count')
df.sort_values('count', ascending=False).reset_index(drop=True)
#Reading the scRNAseq data in h5ad format
```

| | | | | |
|-----|------------------|--|--------|-------|
| 3 | TGTAAGCGTCTCGCGA | ATAGATGCACCTTGGTGGTCGAGTTCTAGTTGTAGCTGATCGTCCAG... | 1 | 6 |
| 4 | TGTAAGCGTCTCGCGA | ATTCGACCAGAACCACATGCAGTTCAACGTGTTGAGGGTGTAGATG... | 1 | 6 |
| ... | cellID | barcode | sample | count |
| 84 | TGTAAGCGTCTCGCGA | AAAGAGTAGTTCGTTGATGTTCAACAGAGCTGAGCTGTCAG... | 11 | 13 |
| 85 | TGTAAGCGTCTCGCGA | AATAGAGAGCACTGCACTTACATATGAGTACCACTAGCACTC... | 11 | 21 |
| 86 | RNAseq data | | | |
| 85 | TGTAAGCGTCTCGCGA | ATAGTACATGGTGGACCTGGACTTCGAGATGGAGCTCTTGTTCCTG... | 1 | 1 |
| 85 | TGTAAGCGTCTCGCGA | ATAGGAGTAGTTGGTGATGGTCTACCAGAAGGTGAAGGTGGAGAAG... | 1 | 1 |
| 86 | TGTAAGCGTCTCGCGA | GCTGCTGCACTTCTGTCTACTTCTAGTTGATGTTGGACGTCATC... | 1 | 1 |

Next, we can look at multi-barcode singlets. There are two cell IDs: AGGCTGCTCTTTCCGG and GAGGGATGTAACATCC. If we look at the barcodes with greater than 2 UMI counts, we see that they have the same combination. The only way this can occur is if a cell receives multiple barcode initially and then divides.

```

{df[df['cellID'] == 'AGGCTGCTCTTTCCGG']}
import scanpy as sc
sc.pp.subsample(df, min_count=2)
df.reset_index(name='count')
df.sort_values('count', ascending=False)
df[df['count'] > 2]
df.reset_index(drop=True)

#Reading the scRNAseq data in h5ad format
adata = sc.read_h5ad(f"{path}/inputFiles/JiangEtAlSubset_scRNAseqData.h5ad")
adata

```

AnnData object with n_obs x n_vars = 39 x 36601
var: 'gene_ids', 'feature_types'

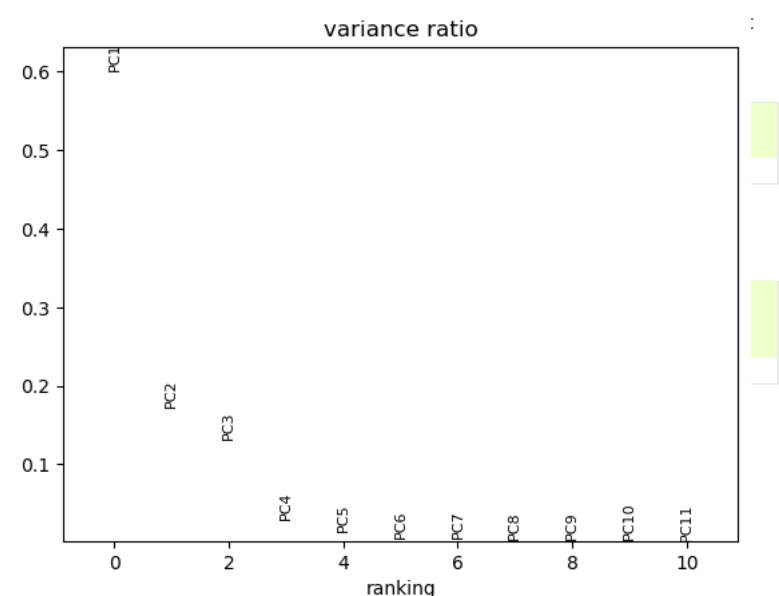
Making copies of singletCode input/output to use them along with scRNAseq data. The -1 is added to cell IDs to match the cell IDs seen in 10x format data.
NOTE: It may not be needed for your actual data.

```

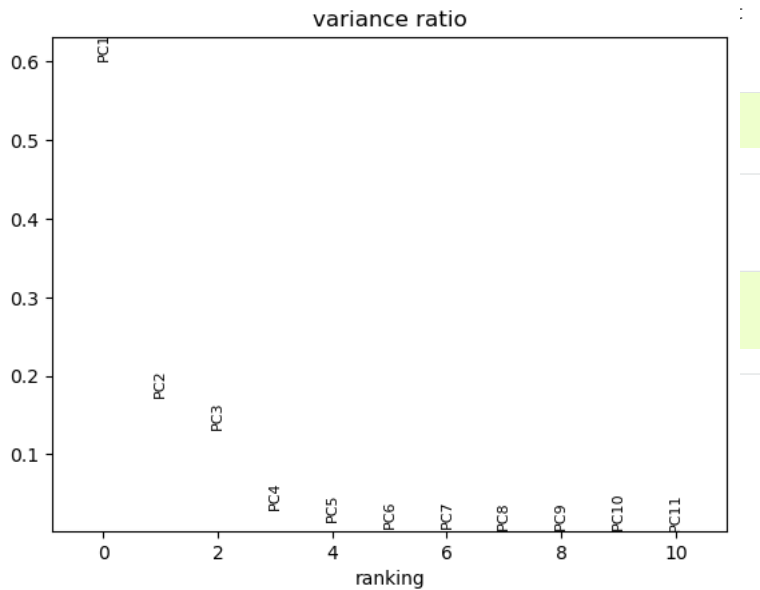
singleCellDf = df.copy()
singleCellDf['cellID'] = singleCellDf['cellID'] + "-1"
singleCellDf = singleCellDf.drop_duplicates(subset = 'cellID')
cellLabelListSingleCell = cellLabelList.copy()
cellLabelListSingleCell['cellID'] = cellLabelListSingleCell['cellID'] + "-1"
cellLabelListSingleCell = cellLabelListSingleCell.drop_duplicates(subset='cellID').reset_index(drop = True)

```

Calculating total counts and genes identified per cell.



Identifying cells that were thresholded by singletCode as low UMI by identifying cells that were in the original list provided to singletCode but not labeled as either singlet or undetermined. Then creating a list of annotations of



Identifying cells that were thresholded by singletCode as low UMI by identifying cells that were in the original list provided to singletCode but not labeled as either singlet or undetermined. Then creating a list of annotations of singletStatus(singlet, multiplet, low UMI) for all cells

```
umiCutoff = pd.DataFrame(

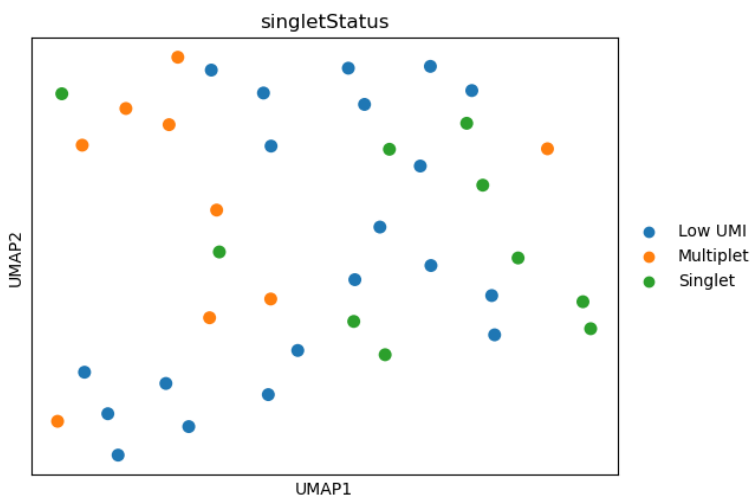
singleCellDf.loc[~singleCellDf['cellID'].isin(cellLabelListSingleCell['cellID']
'cellID')]
    .drop_duplicates()
    .reset_index(drop=True),
    columns=['cellID']
)
umiCutoff['label'] = "Low UMI"
```

```
cellIDLabels = cellLabelListSingleCell.drop(columns = ['barcode', 'sample',
'nUMI']).drop_duplicates().reset_index(drop = True)
```

```
#Creating a list of cell IDs with annotation of whether singlet, multiplet
or low UMI.
labelID = pd.concat([umiCutoff, cellIDLabels]).reset_index(drop=True)
labelID = labelID.set_index(labelID['cellID']).drop(columns = ['cellID'])
#Adding the labels to cells in the adata to visualise it
adata.obs["singletStatus"] = labelID
```

Visualising the cells in PCA space

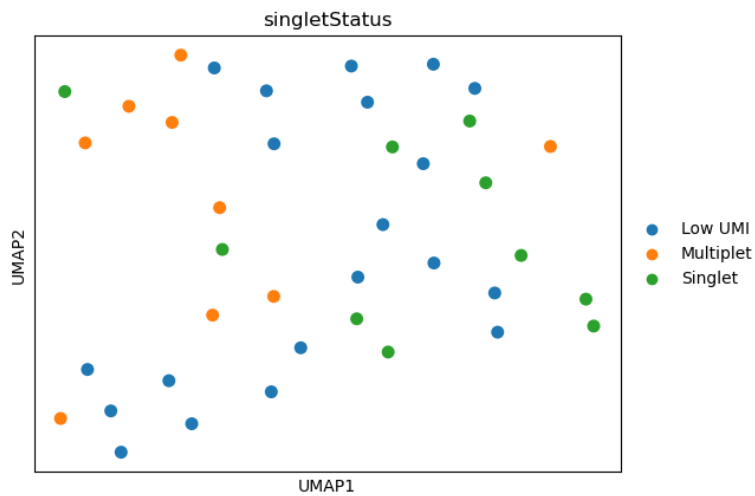
```
sc.pl.umap(
    adata,
    sc.pl.pca(
        adata,
        color=['singletStatus'],
        # Setting a smaller point size to get prevent overlap
        color=['n_genes_by_counts', 'total_counts', 'singletStatus'],
        size=250,
        size = 250
    )
)
```



```

sc.pl.umap(
    adata,
    sc.pl.pca(
        adata,
        color=['singletStatus'],
        # Setting a smaller point size to get prevent overlap
        color = ['n_genes_by_counts', 'total_counts', 'singletStatus'],
        size=250
    )
)

```



Saving the AnnData

```
adata.write(f"{outputPath}/JiangEtAlSubset.h5ad")
```

questions? contact us!

Do you have a doublet detection method that you would like use with singletCode? Do you have a question about how to use singletCode? Do you have a suggestion for how to improve singletCode? Please contact us!

For more information or questions, please contact Yogesh Goyal at yogesh.goyal@northwestern.edu

repositories

GitHub repository for analyses performed in our paper [here](#).

GitHub repository for singletCode tools (command line interface and python package) [here](#).