

## Supplemental Protocol 4

### iPS2-seq design and analysis with *catcheR*

This Supplemental Protocol describes the following procedures:

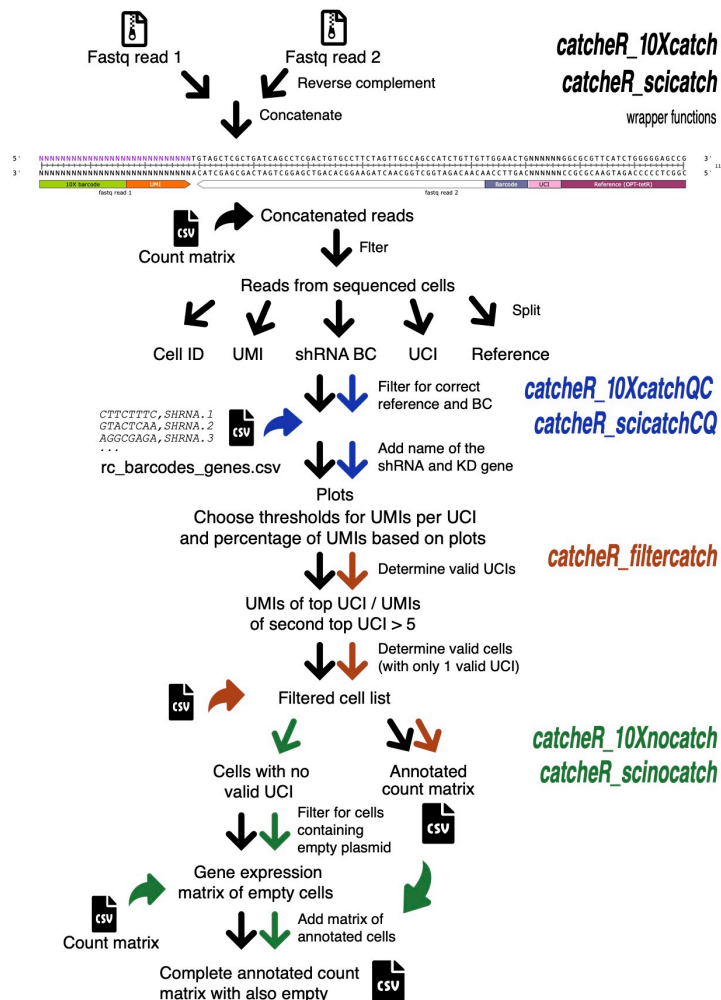
1. [Overview](#)
2. [Installation](#)
3. [Oligonucleotides design](#)
4. [Pooled cloning step 1 plasmid QC](#)
5. [Pooled cloning step 2 plasmid QC](#)
6. [iPS2-10X-seq perturbation deconvolution](#)
7. [iPS2-sci-seq perturbation deconvolution](#)
8. [Barcode reassignment](#)
9. [Perturbation effect analysis](#)

#### Overview

*catcheR* is a comprehensive bioinformatic package for designing and analyzing iPS2-seq experiments. It comprises the following functions ([Figure SP4.1](#)):

1. *catcheR\_design*, which designs oligonucleotides for [Supplemental Protocol 1 - Design shRNA oligonucleotides](#), facilitating shRNA library cloning
2. *catcheR\_step1QC*, which analyzes the results of [Supplemental Protocol 1 - Intermediate plasmid pool QC](#), assessing pooled cloning step 1 plasmids for barcode swaps
3. *catcheR\_step2QC*, which analyzes the results of [Supplemental Protocol 1 - Final plasmid pool QC](#), assessing pooled cloning step 2 for shRNA representation
4. *catcheR\_scicount*, which analyzes 2-level indexing sci-RNA-seq data, facilitating the generation of gene expression matrix for iPS2-sci-seq experiments
5. *catcheR\_scicatch*, which assigns shRNA perturbations to single nuclei transcriptomes obtained by [Supplemental Protocol 2](#), enabling the primary analysis of iPS2-sci-seq
6. *catcheR\_10Xcatch*, which assigns shRNA perturbations to single cell transcriptomes obtained by [Supplemental Protocol 3](#), enabling the primary analysis of iPS2-10X-seq
7. *catcheR\_scicatchQC* and *catcheR\_10XcatchQC*, which use the outputs of *catcheR\_scicatch* and *catcheR\_10Xcatch*, respectively, to fine-tune shRNA assignment thresholds
8. *catcheR\_filtercatch*, which leverages on the output of *catcheR\_scicatchQC* and *catcheR\_10XcatchQC* to filter single nuclei/cell transcriptomes expressing a single shRNA
9. *catcheR\_sortcatch*, which quality controls the cell-by-gene matrix based on the results of *catcheR\_step1QC*, reassigning hPSC clones with barcode swaps to the correct shRNA
10. *catcheR\_scinocatch* and *catcheR\_10Xnocatch*, which identify cells expressing no shRNA in iPS2-sci-seq and iPS2-10X-seq experiments, respectively, adding them to the cell-by-gene matrix to be used as additional controls

- 2689 11. *catcheR\_load* to load the gene expression matrices generated after *iPS2-10X-seq* per-  
 2690 *turbation deconvolution* and *iPS2-sci-seq* *perturbation deconvolution* and pack the data  
 2691 in a *monocle* object, ready to be analyzed
- 2692 12. *catcheR\_pseudotime*, *catcheR\_modules* and *catcheR\_enrichment* allow follow up analy-  
 2693 ses on the annotated perturbed cells



### Figure SP4.1. Overview of *catcheR*

Logical relationship and main inputs/outputs for the main functions used for the *catcheR* analytical pipeline of iPS2-seq experiments. Functions used for oligo design and plasmid QC are not illustrated.

## 2694 Installation

2695 *catcheR* is available at <https://github.com/alessandro-bertero/catcheR>. The GitHub repository  
 2696 folder "*scripts*" contains all the bash and R scripts that can be run independently. However,  
 2697 in order to ensure reproducible analyses, we strongly recommend to install *catcheR* package  
 2698 from GitHub, since its functions run all the analysis inside Docker containers.

2699 1. Install Docker engine, following the instructions at <https://docs.docker.com/engine/install/>

2700 2. Install *catcheR*

2701 (a) In R ( $\geq 3.0.2$ ) install *devtools*, if not already present

```

2702     install.packages("devtools")
2703
2704 (b) Install catcheR from GitHub
2705     install_github("alessandro-bertero/catcheR")
2706
2707 (c) Install rrundocker from GitHub
2708
2709     install_github("Reproducible-Bioinformatics/rrundocker")
2710
2711
2712 (d) Load catcheR and rrundocker in your R environment
2713     library(catcheR)
2714     library(rrundocker)
2715

```

## 2715 **Oligonucleotides design**

2716 This step complements [Supplemental Protocol 1 - Design shRNA oligonucleotides](#)

2717 1. In a new working folder, prepare the following files:

2718 (a) A comma-separated values (CSV) file with three columns listing: (1) the forward  
2719 oligos from the TRC shRNA library; (2) the corresponding barcodes (BC); (3) the  
2720 shRNA names. Below is an example for the shRNA described in [Figure SP1.2](#):

```

2721 CCGGCAAGTACTCCTTGCTGGATTGCTCGAGCAATCCAGCAAGGAGTACTTGT TTTTG , CAGTTCCA , SMAD2 . 1
2722 ...
2723

```

2724 (b) (Optional) - A txt file with a newline-separated list of 5'-3' restriction sites, or other  
2725 sequences, to be avoided in the shRNAs. By default these are Sall, Swal and Ascl:

```

2726 GTCGAC
2727 ATTTAAAT
2728 GGCGCGCC
2729

```

2730 2. Run *catcheR\_design*:

```

2731 catcheR_design(
2732   group=c("docker","sudo"),
2733   folder,
2734   sequences,
2735   gibson.five = "AGTTCCTATCAGTGATAGAGATCCC",
2736   gibson.three = "GTAGCTCGCTGATCAGC",
2737   fixed = "GTCGACATTTAAATGGCGCGCC",
2738   restriction.sites = NULL)
2739

```

2740 *catcheR\_design* arguments:

2741 (a) group: string with two options: sudo or docker, depending on the user group. For a  
2742 detailed explanation of Docker user groups, see [this page](#)

2743 (b) folder: string with the working folder path

2744 (c) sequences: string with the CSV file name from step [1a](#)

2745 (d) gibson.five: (optional) string with the 5' Gibson homology

2746 (e) gibson.three: (optional) string with the 3' Gibson homology

2747 (f) fixed: (optional) character string with the multicloning site

2748 (g) restriction.sites: (optional) string with the txt file name from [1b](#)

## Example usage:

```
catcheR_design(  
  group = "docker",  
  folder = "path/to/folder",  
  sequences = "filename.csv",  
  restriction.sites = "filename.txt")
```

## *catcheR\_design* outputs:

- (a) "output.txt", with the oligo sequences to be used for synthesis ([Figure SP1.2](#))
- (b) "bad\_oligos.txt", with the shRNAs with forbidden restriction enzyme sites (highlighted in lowercase characters); this information can be used to refine the shRNA list

## Example output for "bad\_oligos.txt":

```
AGTTCCTATCAGTGATAGAGATCCCGGACATAATCACTGCGTAATCCTCagatctTACGCAGTGATTATGTCCTTTTTTGT -  
CGACATTTAAATGGCGCGCCNNNNNGCTGAAGAGTAGCTCGCTGATCAGC , GATA4  
...
```

## Pooled cloning step 1 plasmid QC

This step complements [Supplemental Protocol 1 - Intermediate plasmid pool QC](#)

1. In a new working folder, prepare the following files:

- (a) Fastq/fq or fastq.gz files with demultiplexed read 1 from the NGS run
- (b) A CSV file with the shRNA names and their full sequences

```
SMAD2.1 , GCAAGTACTCCTTGCTGGATTGCTCGAGCAATCCAGCAAGGAGTACTTG  
...
```

- (c) (Optional) - A txt file with a newline-separated list of clones of interest as "BC\_UCI"

```
CAAGAGCC_CATCGT  
...
```

2. Run *catcheR\_step1QC*:

```
catcheR_step1QC(  
  group=c("docker", "sudo"),  
  folder,  
  fastq.read1,  
  DIs = 100,  
  ratio = 10,  
  plot.threshold = 2000,  
  clones = NULL)
```

## *catcheR\_step1QC* arguments:

- (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- (b) folder: string with the working folder path
- (c) fastq.read1: string with the read 1 filename from step [1a](#)
- (d) DIs: integer of the minimum number of diversity indexes (DIs, pseudo-unique reads) of the most represented shRNA matching to a given UCI-BC; in combination with "ratio", it selects UCI-BCs for which it is possible to reliably assign an shRNA
- (e) ratio: integer of the minimum ratio between the number of DIs of the most represented and second most represented shRNAs matching to a given UCI-BC

2796 (f) plot.threshold: integer of the minimum number of DIs per UCI-BC for output [2b](#)

2797 (g) clones: (optional) a string with the txt file from step [1b](#)

2798 Example usage:

```
2799 catcher_step1QC(  
2800     group = "docker",  
2801     folder = "path/to/folder",  
2802     fastq.read1 = "filename.fq",  
2803     clones = "filename.txt")
```

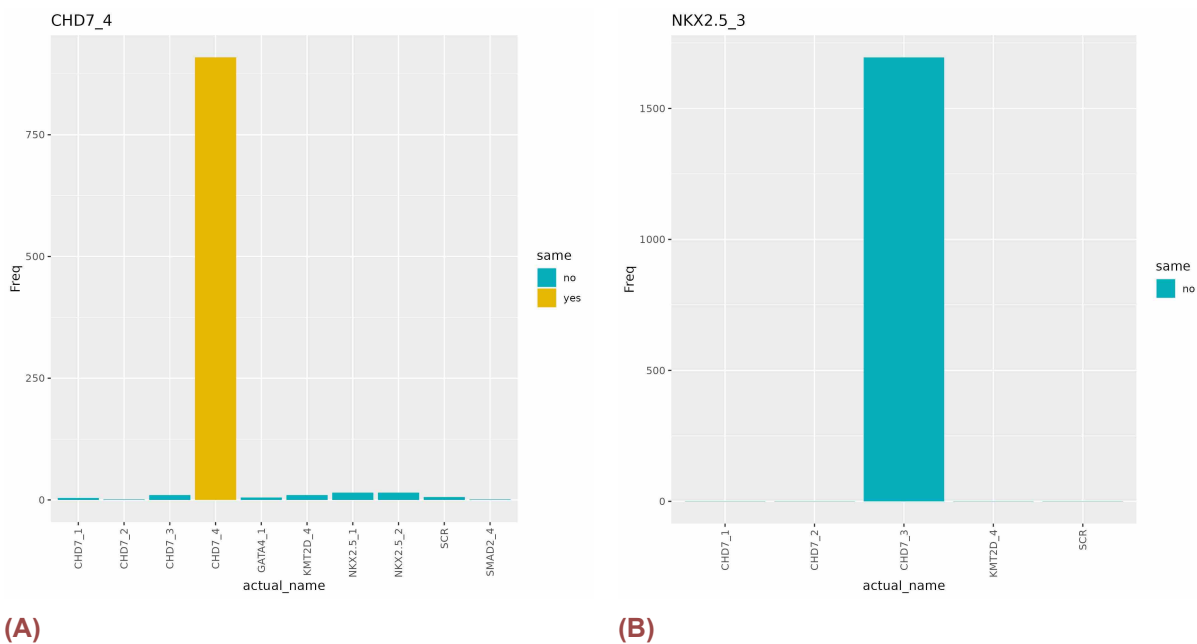
2804 *catcher\_step1QC* key outputs:

2805 (a) "reliable\_clones\_swaps.csv", which lists UCI-BCs with reliable evidence of shRNA-  
2806 barcode swap, and is used as input for [Barcode reassignment](#)

2807 (b) Bar chart of the number of DIs for each clone above "plot.threshold"

2808 (c) Bar chart of the number of DIs associated to each BCs, shRNAs, and reliable swaps

2809 (d) (If "clones" argument is provided) - csv files and bar charts with the number of DIs  
2810 for each shRNA matching to each clone of interest ([Figure SP4.2](#))



**Figure SP4.2.** Example of quantification of DIs for bacterial clones carrying the expected shRNA (A) or with robust evidence of a barcode swap that can be reassigned (B; [Figure 1B](#)). These graphs allow visual evaluation of appropriate thresholds for "DI" and "ratio".

## 2811 Pooled cloning step 2 plasmid QC

2812 This step complements [Supplemental Protocol 1 - Final plasmid pool QC](#)

2813 1. In a new working folder, prepare the following files:

2814 (a) Fastq/fq or fastq.gz files with demultiplexed read 1 from the NGS run

2815 (b) "rc\_barcodes\_genes.csv", a CSV file with two columns: (1) the shRNA BCs; (2) the  
2816 matching shRNA names in the format "GENE.shRNAID"

```

2817         CAAGAGCC , SMAD2 . 1
2818         ...

```

2819 (c) (Optional) - A txt file with clones of interest (step 1a of the previous section)

## 2821 2. Run *catcheR\_step2QC*:

```

2822 catcheR_step2QC (
2823     group=c("docker","sudo"),
2824     folder,
2825     fastq.read1,
2826     DIs = 1000,
2827     clones = NULL)

```

2828 *catcheR\_step2QC* arguments:

- 2830 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- 2831 (b) folder: string with the working folder path
- 2832 (c) fastq.read1: string with the read 1 file from step 1a
- 2833 (d) DIs: integer of the minimum number of DIs for a given UCI-BC; it selects reliably-
- 2834 measured UCI-BCs
- 2835 (e) clones: (optional) a string with the txt file from step 1c

2836 Example usage:

```

2837 catcheR_step2QC (
2838     group= "docker",
2839     folder = "path/to/working/folder",
2840     fastq.read1 = "filename.fq",
2841     clones = "filename.txt")

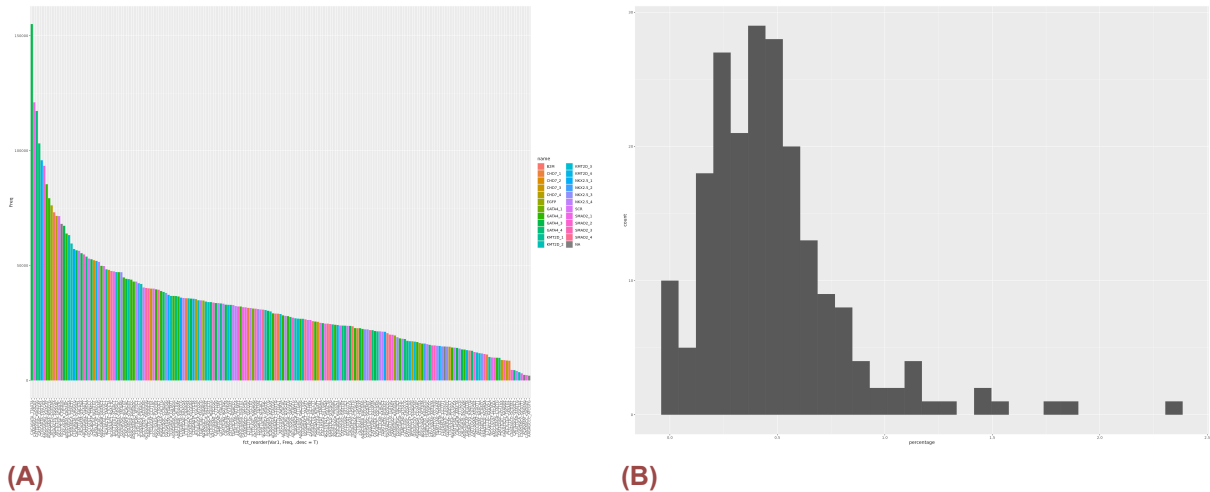
```

2842 *catcheR\_step2QC* key outputs:

- 2844 (a) Pie charts of DIs per shRNA and per gene targets associated to each clone
- 2845 (b) Text file listing all clones above the DI threshold
- 2846 (c) Bar chart of the number of DIs for each clone above the DI threshold ([Figure SP4.3A](#))
- 2847 (d) Frequency histogram of the percentage of clones above the DI threshold for each
- 2848 shRNA ([Figure SP4.3B](#))

## 2849 Obtain iPS2-10X-seq count matrices

- 2850 1. Download and install *cellranger* from the [10X Genomics download resource page](#); al-
- 2851 ternatively, a Docker container with [cellranger 7](#) or [cellranger 9](#) can be pulled from our
- 2852 Docker repository. Preprocess iPS2-multi-seq with cellranger ARC, which can be either
- 2853 downloaded from the [10X Genomics download resource page](#) or pulled as a [Docker](#) con-
- 2854 tainer from our repository
- 2855 2. Demultiplex Illumina base calls with *cellranger mkfastq*, following 10X Genomics' manual.
- 2856 In the CSV sample sheet, include the 10X Genomics indexes used to generate the GEX,
- 2857 CMO, and UCI-BC libraries in [Supplemental Protocol 3](#)
- 2858 3. Obtain cell-by-gene count matrices using *cellranger count* (single sample) or *cellranger*
- 2859 *multi* (multiplexed experiment).
- 2860 **Note:** count matrices from multiple samples from the same experiment can be aggre-
- 2861 gated with *cellranger aggr*, then can be analyzed in a single run of *catcheR\_10Xcatch* by
- 2862 providing the number of samples to the "samples" argument.
- 2863 4. Run *cell ranger matrix2csv* to transform the sparse count matrix in a dense CSV matrix
- 2864 After obtaining the count matrix, proceed to [iPS2-10X-seq perturbation deconvolution](#).



**Figure SP4.3.** Example of quantification of DIs for bacterial clones represented by more than 1,000 DIs (A), and of the frequency of bacterial clones for each shRNA. These graphs confirm whether shRNAs are normally distributed in the final plasmid pool.

## Obtain iPS2-sci-seq count matrices

1. *catcherR\_scicount* is a wrapper of the "bbi-sci" pipeline developed by the Brotman Baty Institute for Precision Medicine, which was implemented and dockerized in *catcherR* to be used on any operating system.

2. Demultiplex Illumina base calls to fastq files

- (a) Create a "SampleSheet.csv" file with as many sample rows as the PCR wells from [Indexing PCR and NGS of Supplemental Protocol 2](#), where "Sample\_ID" is the well identified with format [A-H][01-12], and index and index2 are the i7 and i5 indexes used in the corresponding row and column (refer to [Table SP2.1](#))

- (b) Run Illumina *bcl2fastq* following Illumina's manual

- (c) Run fastQC to confirm the quality of the fastq files

3. In a new working folder:

- (a) Create the subfolder "fastq", and copy all the demultiplexed fastq.gz files. Ensure that all file names begin with the well coordinate (e.g. A01) followed by an underscore.

- (b) Create a tab separated txt file called "sci-RNA-seq-8.RT.oligos" with the association between RT wells and RT barcode sequences (refer to [Table SP2.1](#))

```
A01 TTCTCGCATG
...
```

- (c) Create the subfolder "GENOMES", and copy the annotated genome (i.e., [GRCh38](#))

4. Run *catcherR\_scicount*:

```
catcherR_scicount(
  group=c("docker", "sudo"),
  folder,
  sample.name,
  UMI.cutoff)
```

*catcherR\_scicount* arguments:

- (a) group: string with two options: sudo or docker, depending on the user group ([info](#))



- 2894 (b) folder: string with the working folder path
- 2895 (c) sample.name: string with the name of the experiment
- 2896 (d) UMI.cutoff: integer of the minimum number of UMI per nucleus needed to consider
- 2897 the single cell transcriptome valid

2898 Example usage:

```
2899 catcheR_scicount(  
2900   group = "docker",  
2901   folder = "path/to/file",  
2902   sample.name = "experiment",  
2903   UMI.cutoff = 500)
```

2904 *catcheR\_scicount* outputs, found in the final-output folder: knee-plot of UMI per cell ([Figure S2F](#)), statistics files, sparse cell-by-gene count matrix, dense gene matrix "exp\_mat.csv" and "exp\_mat\_no0.csv" (filtered to exclude genes with 0 counts in all cells), also in Rdata format.

2909 After obtaining the count matrix, proceed to [iPS2-sci-seq perturbation deconvolution](#).

## 2910 iPS2-10X-seq perturbation deconvolution

2911 shRNA perturbations can be assigned to single cells using *catcheR\_10Xcatch*, which identifies NGS reads containing UCI-BCs, matches them to the corresponding transcriptome based on their shared cellular barcodes, applies filters for background noise, and selects cells with robust evidence of a single integration. UCI-BC filtering involves: (1) filtering out UCI-BCs supported by less than a certain number of UMIs in a given transcriptome, to account for PCR and sequencing artifacts; (2) filtering out UCI-BCs that represent less than a certain fraction of all UMIs for UCI-BCs in a given transcriptome, to reduce noise arising from free mRNA; and (3) filtering out UCI-BCs whose UMI fraction in a given transcriptome is not several fold greater than the second most common UCI-BC, to eliminate cells that contain multiple UCI-BCs when the second most common one falls just below one of the first two thresholds. The resulting list of *bona fide* shRNA integrations per cell is then used to filter those expressing a single shRNA.

2922 This section describes a typical analysis pipeline starting with *catcheR\_10Xcatch* to perform a full analysis with automatic thresholds, after which *catcheR\_10XcatchQC* can fine-tune the thresholds to re-filter cells with *catcheR\_filtercatch*. Lastly, *catcheR\_10Xnocatch* can add cells expressing no shRNA as additional controls to the final annotated cell-by-gene count matrix.

2926 1. In a new working folder:

- 2927 (a) Copy the fastq/fq or fastq.gz files with demultiplexed read 1 and read 2 of the UCI-BC library (step 2 of [Obtain iPS2-10X-seq count matrices](#))
- 2928
- 2929 (b) Copy the cell-by-gene count matrix in CSV format (step 4 of [Obtain iPS2-10X-seq count matrices](#))
- 2930
- 2931 (c) Create a CSV file named "rc\_barcodes\_genes.csv" with two columns: (1) the shRNA BCs; (2) the matching shRNA names in the format "GENE.shRNAID"
- 2932

```
2933 CAAGAGCC , SMAD2 . 1  
2934 ...  
2935
```

2936 2. Run *catcheR\_10Xcatch* to execute the complete analysis with automatic thresholding:

```
2937 catcheR_10Xcatch(  
2938   group=c("docker", "sudo"),  
2939   folder ,  
2940   fastq.read1 ,
```



```

2941     fastq.read2,
2942     expression.matrix,
2943     reference = "GGCGCGTTCATCTGGGGGAGCCG",
2944     UCI.length = 6,
2945     threads = 2,
2946     percentage = 15,
2947     mode = "bimodal",
2948     ratio = 5,
2949     samples = 1,
2950     x = 100,
2951     y = 400)

```

2952 *catcherR\_10Xcatch* arguments:

- 2954 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- 2955 (b) folder: string with the working folder path
- 2956 (c) fastq.read1: string with the filename of read 1 fastq/fq or fastq.gz (step [1a](#))
- 2957 (d) fastq.read2: string with the filename of read 2 fastq/fq or fastq.gz (step [1a](#))
- 2958 (e) expression.matrix: string with the filename of the count matrix file CSV (step [1b](#))
- 2959 (f) reference: string with the reverse complement of the sequence before the shRNA BC
- 2960 at the start of read 2 (default is the 3' end of the OPTtetR cDNA, optional argument;
- 2961 [Figure S3F](#))
- 2962 (g) UCI.length: integer of the UCI length (default is 6, optional argument)
- 2963 (h) threads: integer of the threads to be used for parallelization (default is 2, optional
- 2964 argument)
- 2965 (i) percentage: integer of the percentage of UMIs supporting a given UCI over the total
- 2966 UMIs supporting all UCIs in a given cell; it is used as threshold to consider the UCI
- 2967 valid ([Figure SP4.4D](#)); the recommended default is 15 (optional argument)
- 2968 (j) mode: string with two options: bimodal or noise, defining the mode for automatic
- 2969 thresholding the minimum number of UMIs per UCI to consider the UCI valid (optional
- 2970 argument)
  - 2971 i. "bimodal" (default) sets the threshold at the valley of the bimodal UMIxUCI dis-
  - 2972 tribution ([Figure SP4.4C](#))
  - 2973 ii. "noise" sets the threshold at  $1.35 \times$  the number of UCIs supported by a single
  - 2974 UMI; used if the UMIxUCI distribution is not bimodal
- 2975 (k) ratio: the ratio between the number of UMIs supporting the most represented UCI
- 2976 and the number of UMIs supporting the second most represented UCI. The parame-
- 2977 ter is needed to identify the cell as a single integration cell. The default is 5 (optional
- 2978 argument)
- 2979 (l) samples: the number of samples present in the same experiment (cells from dif-
- 2980 ferent experiments will have the corresponding number after the cell ID in the gene
- 2981 expression matrix, so that multiple reactions of the same experiment can be unified).
- 2982 The default is 1 (optional argument)
- 2983 (m) x: an integer indicating the upper limit on the x-axis of the cropped version of the
- 2984 plot "UMIxUCI". Default is 100 (optional argument)
- 2985 (n) y: an integer indicating the upper limit on the y axis of the cropped version of the plot
- 2986 "UMIxUCI". Default is 400 (optional argument)

2987 Example usage:

```

2988 catcheR_10Xcatch(group = "docker",
2989                   folder = "path/to/folder",
2990                   fastq.read1 = "R1.fastq.gz",
2991                   fastq.read2 = "R2.fastq.gz",
2992                   expression.matrix = "filename.csv",
2993                   threads = 12)

```

2994 *catcheR\_10Xcatch* key outputs (will be stored inside the "Result" folder of each sample):

- 2995 (a) "log.txt" with the starting number of reads
- 2996 (b) "log2.txt" with the number of cells, UMIs, UCIs, and the calculated thresholds of UMI  
2997 per UCI for both "bimodal" and "noise" (only one is chosen based on the "mode"  
2998 argument; the other one is provided for a potential reanalysis)
- 2999 (c) Bar charts of UMI counts per shRNA and target gene ([Figure SP4.4A](#) and [Fig-  
3000 ure SP4.4B](#))
- 3001 (d) Frequency histograms of UCIs supported by a certain number of UMIs (UMIxUCI),  
3002 to interpret the signal/noise of the experiment and possibly set a custom UMIxUCI  
3003 threshold for subsequent reanalysis ([Figure SP4.4C](#)).
- 3004 **Note:** here and below, copies of the same UCI expressed by different cells are  
3005 plotted and analyzed separately
- 3006 (e) Frequency histogram of UCIs supported by a certain fraction of UMIs over the total  
3007 number of UMIs supporting all UCIs in a given cell (UMIpercentagexUCI), to further  
3008 assess signal/noise and possibly adjust the default threshold ([Figure SP4.4D](#))
- 3009 (f) Dot plots that combine the UMIpercentagexUCI and UMIxUCI data on the x and y  
3010 axes, respectively, with each dot representing one or more UCI (quantified by the  
3011 dot size). Dot colors indicate either the number of valid UCIs (i.e., shRNAs) in the  
3012 cell containing a given UCI ([Figure SP4.4E](#)), or whether said UCI is the only valid  
3013 one in the relevant cell and is thus assigned to such cell ([Figure SP4.4F](#))
- 3014 (g) "log\_part3.txt" with how many cells were assigned to a single shRNA or were filtered  
3015 out due to zero or multiple shRNAs
- 3016 (h) "silencing\_matrix.csv" is the key output used for the secondary analyses: the cell-by-  
3017 gene count matrix provided as input, filtered and annotated with the shRNA encoded  
3018 in each cell. It is also provided in RDS format for easy loading into R. Cell names  
3019 are modified as follows:

```

3020 cellID_UMIxUCI_BC_GENE_UCI

```

3021 Where:

- 3023 i. cellID is the original cell name (i.e., the 16 bp of the 10X RT barcode)
- 3024 ii. UMIxUCI is the number of UMI associated with that perturbation
- 3025 iii. BC is the shRNA barcode
- 3026 iv. GENE is the shRNA target
- 3027 v. UCI is the Unique Clonal Identifier (shared by all cells originating from the same  
3028 hPSC clone)

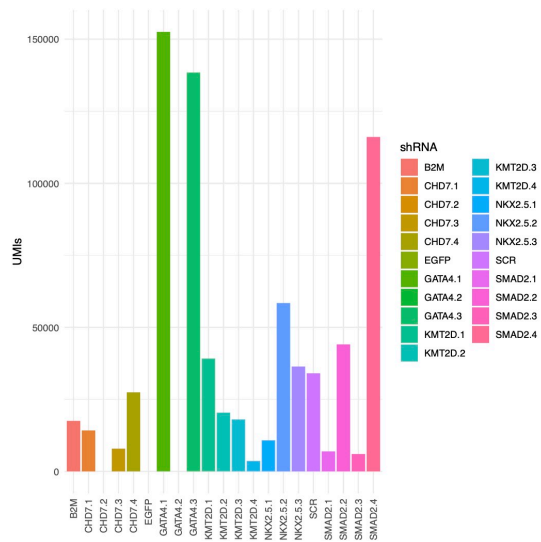
3029 Example of a cell name after the analysis:

```

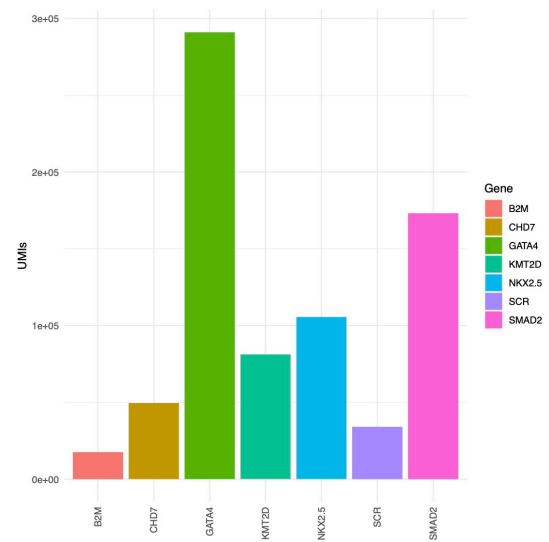
3030 TTCTAACCACAGTCGC_180_CGTGATGC_NKX2.5_ACAGTG

```

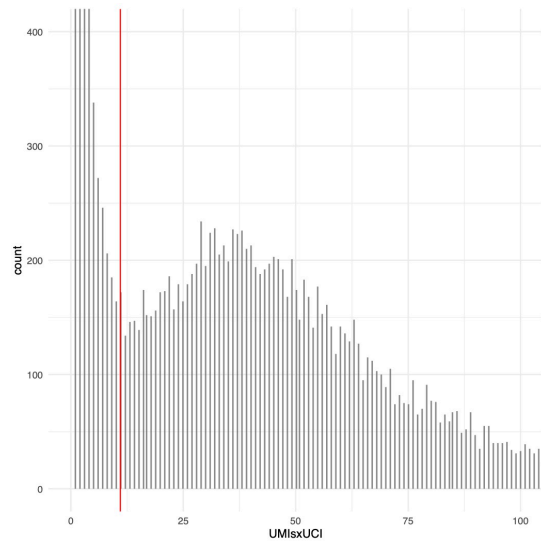
3031 This annotation can be leveraged with the [scripts](#) available on our GitHub repository  
3032 to enable a variety of secondary analyses described in the manuscript, or using cus-  
3033 tom scripts for other types of secondary analyses or with the functions implemented  
3034 within *catcheR* and described in the [Perturbation effect analysis](#).  
3035



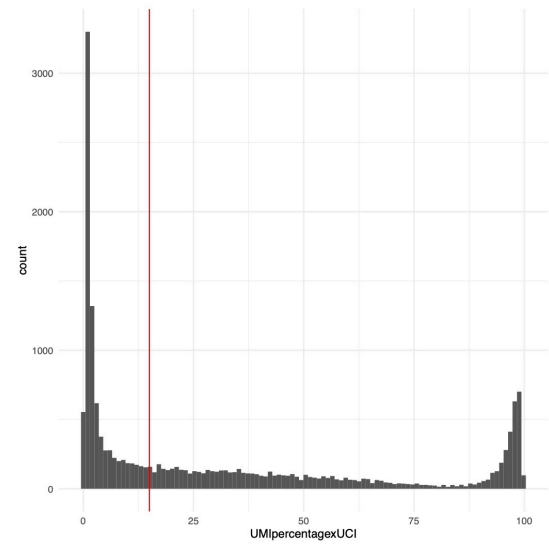
(A)



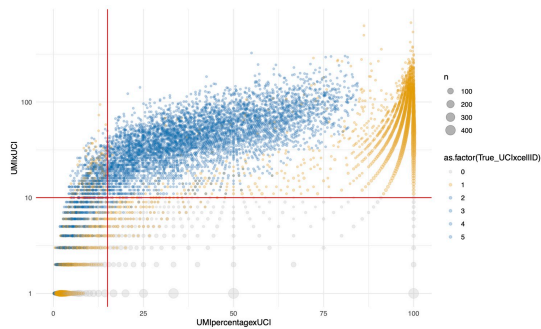
(B)



(C)



(D)



(E)



(F)

**Figure SP4.4.** Exemplary key outputs of a *catcherR\_10Xcatch* analysis: "barcode\_distribution.pdf" (A); "gene\_distribution.pdf" (B); "UMIxUCI\_400\_100.pdf" with the automatic UMIxUCI threshold based on "bimodal" (C); "percentage\_of\_UMIxUCI\_dist.pdf" with the default 15% "percentage" threshold (D); "2D\_percentage\_of\_UMIxUCI\_UMI\_count\_trueorfalse.pdf" (E); and "2D\_percentage\_of\_UMIxUCI\_UMI\_ValidCells.pdf" (F)

3036 **OPTIONAL: fine tune iPS2-10X-seq perturbation assignment**

- 3037 1. In the same working folder used to run *catcheR\_10Xcatch*, run *catcheR\_10XcatchQC*,  
3038 which uses the txt file outputs to regenerate the quality control plots described in [Figure SP4.4](#) and suggest new thresholds to enable subsequent filtering  
3039

```
3040 catcheR_10XcatchQC(  
3041     group=c("docker","sudo"),  
3042     folder,  
3043     reference = "GGCGCGTTCATCTGGGGGAGCCG",  
3044     mode = "bimodal",  
3045     sample = 1  
3046     x = 100,  
3047     y = 400)
```

3048  
3049 Example usage:

```
3050 catcheR_10XcatchQC(  
3051     group = "docker",  
3052     folder = "path/to/working/folder",  
3053     mode = "noise")
```

3054 The function is to be run separately for each sample (specified by argument sample with  
3055 default 1). *catcheR\_10XcatchQC* arguments and outputs are the same as *catcheR\_10Xcatch*

- 3056 2. In the same working folder used to run *catcheR\_10Xcatch*, run *catcheR\_filtercatch*, which  
3057 uses the output of *catcheR\_10Xcatch* to filter cells based on custom thresholds

```
3058 catcheR_filtercatch(  
3059     group=c("docker","sudo"),  
3060     folder,  
3061     expression.matrix,  
3062     UMI.count,  
3063     percentage = 15,  
3064     ratio = 5,  
3065     sample = 1)
```

3066  
3067 Example usage:

```
3068 catcheR_filtercatch(group = "docker",  
3069     folder = "path/to/working/folder",  
3070     expression.matrix = "filename.csv",  
3071     UMI.count = 5,  
3072     sample = 1)
```

3073  
3074 The function is to be run separately for each sample (specified by argument sample with  
3075 default 1). *catcheR\_filtercatch* argument *percentage* is the minimum percentage of UMI  
3076 for a given UCI over the total UMIs of UCIs in that cell, to consider the UCI valid. *UMI.count*  
3077 is an integer of the custom UMI threshold. *ratio* is the minimum ration between UMIs of  
3078 top UCI and UMIs of 2nd top UCI. The other arguments and the outputs are the same as  
3079 *catcheR\_10Xcatch*.

3080 **OPTIONAL: identify cells expressing no shRNA in iPS2-10X-seq**

- 3081 1. In the same working folder used to run *catcheR\_10Xcatch*, run *catcheR\_10Xnocatch*:

```
3082 catcheR_10Xnocatch(  
3083     group=c("docker","sudo"),  
3084     folder,  
3085     expression.matrix,
```

```

3086     threshold,
3087     sample = 1,
3088     reference = "TACGCGTTCATCTGGGGGAGCCG")

```

3089 **Example usage:**

```

3090     catcherR_10Xnocatch(group = "docker",
3091     folder = "path/to/folder",
3092     expression.matrix = "filename.csv",
3093     threshold = 5)

```

3094 *catcherR\_10Xnocatch* argument *threshold* is an integer of the minimal UMI count for the  
3095 "empty" reference [Figure SP4.3B](#) needed to confidently identify it as having integrated  
3096 said plasmid; all other arguments are the same as *catcherR\_10Xcatch*

3097 *catcherR\_10Xnocatch* output is an updated "silencing\_matrix.csv", also available in RDS  
3098 format, in which names of empty cells are modified as "cellID\_?\_empty\_NA\_empty".

### 3099 **OPTIONAL: merge multiple samples**

3100 1. In case fine-tuning is applied after *catcherR\_10Xcatch* and multiple samples are present,  
3101 run *catcherR\_merge* to aggregate the results in a single matrix. This step is done auto-  
3102 matically when running *catcherR\_10Xcatch*.

```

3103     catcherR_merge(
3104     group=c("docker","sudo"),
3105     folder,
3106     samples = 2,
3107     empty = T
3108     )

```

3109 The argument "empty" determines whether the cells identified by *catcherR\_nocatch* should  
3110 be added to the matrix.

3111 **Example usage:**

```

3112     catcherR_merge(group = "docker",
3113     folder = "/20tb/ratto/catcher/test_CM5/",
3114     samples = 4)

```

### 3115 **iPS2-sci-seq perturbation deconvolution**

3116 This section describes a typical analysis pipeline similar to the one described in the previous  
3117 section, but leveraging on the *catcherR\_scicatch*, *catcherR\_scicatchQC*, *catcherR\_filtercatch*,  
3118 and *catcherR\_scinocatch* functions.

3119 1. In a new working folder:

3120 (a) Create a subfolder called "fastq", and copy all the demultiplexed fastq.gz files (step  
3121 [2 of Obtain iPS2-sci-seq count matrices](#)). Ensure that all file names begin with the  
3122 well coordinate (e.g. A01)

3123 (b) Copy the cell-by-gene expression matrix CSV file obtained with *catcherR\_scicount*  
3124 function (step [4 of Obtain iPS2-sci-seq count matrices](#)).

3125 (c) Create a CSV file called "rc\_barcodes\_genes.csv" with two columns: (1) the shRNA  
3126 BCs; (2) the matching shRNA names in the format "GENE.shRNAID"

```

3127     CAAGAGCC , SMAD2 . 1
3128     . . .

```

3129

3130 (d) Copy the text file "sci-RNA-seq-8.RT.oligos" used by *catcheR\_scicount* (step 3b of  
3131 [Obtain iPS2-sci-seq count matrices](#)).

3132 2. Run *catcheR\_scicatch* to execute the complete analysis with automatic thresholding; the  
3133 arguments are the same as for *catcheR\_10Xcatch* (step 2 of [iPS2-10X-seq perturbation](#)  
3134 [deconvolution](#)), except that no filenames are provided since fastq files must be in the  
3135 "fastq" subfolder in the working directory:

```
3136 catcheR_scicatch(  
3137   group=c("docker","sudo"),  
3138   folder,  
3139   expression.matrix,  
3140   reference = "GGCGCGTTCATCTGGGGGAGCCG",  
3141   UCI.length = 6,  
3142   threads = 2,  
3143   percentage = 15,  
3144   ratio = 5,  
3145   mode = "bimodal",  
3146   x = 100,  
3147   y = 400)
```

3148 Example usage:  
3149

```
3150 catcheR_scicatch(group = "docker",  
3151   folder = "path/to/working/folder",  
3152   expression.matrix = "filename.csv",  
3153   threads = 12)
```

3154 *catcheR\_scicatch* key outputs are the same as for *catcheR\_10Xcatch* (step 2 of ?? and  
3155 [Figure SP4.3](#)), except that in "silencing\_matrix.csv" "cell ID" indicates the PCR well and  
3156 RT barcode ID (which can be leveraged to identify nuclei from different samples pooled  
3157 on the same RT plate)  
3158

3159 Example of cell name after the analysis:

```
3160 P24__RT_27_7_GCCTGTGT_SCR_ACGGTC
```

3161 In addition, *catcheR\_scicatch* provides two quality control outputs to evaluate experimen-  
3162 tal biases during sci-RNA-seq library preparation: *demux* and *RT* detail about how many  
3163 cells were identified from each row and column of the PCR and RT plates, respectively.  
3164

## 3165 **OPTIONAL: fine tune iPS2-sci-seq perturbation assignment**

3166 1. Run *catcheR\_scicatchQC* following the steps described for *catcheR\_10XcatchQC* in [OP-](#)  
3167 [TIONAL: fine tune iPS2-10X-seq perturbation assignment](#): the functions share the same  
3168 arguments and outputs, but *catcheR\_scicatchQC* used the output of *catcheR\_scicatch*

3169 Example usage:

```
3170 catcheR_scicatchQC(  
3171   group = "docker",  
3172   folder = "path/to/working/folder",  
3173   mode = "noise")  
3174
```

3175 2. Run *catcheR\_filtercatch* function as described in in [OPTIONAL: fine tune iPS2-10X-seq](#)  
3176 [perturbation assignment](#): this function also works on the output of *catcheR\_scicatch* or  
3177 *catcheR\_scicatchQC*

## OPTIONAL: identify cells expressing no shRNA in iPS2-sci-seq

1. Run *catcheR\_scinocatch* following the steps described for *catcheR\_10Xnocatch* in [OPTIONAL: identify cells expressing no shRNA in iPS2-10X-seq](#): the functions share the same arguments and outputs, but *catcheR\_scinocatch* uses the output of *catcheR\_scicatch*

Example usage:

```
catcheR_scinocatch(group = "docker",
  folder = "path/to/working/folder",
  expression.matrix = "filename.csv",
  threshold = 5,
  reference = "TACGCGTTCATCTGGGGGAG")
```

## Barcode reassignment

*catcheR\_sortcatch* is an optional function that corrects the annotated cell-by-gene count matrix obtained with *catcheR\_10Xcatch* or *catcheR\_scicatch* reassigning the perturbation of any cell belonging to a hPSC clone with reliable evidence of a shRNA-barcode swap, based on the results of *catcheR\_step1QC*.

1. In a new working folder:

- (a) Copy the annotated count matrix (i.e., "silencing\_matrix.csv")
- (b) Copy the CSV file with the list of UCI-BCs with reliable evidence of a shRNA-barcode swap (i.e., "reliable\_clones\_swaps.csv"; step 2a of [Pooled cloning step 1 plasmid QC](#))

2. Run *catcheR\_sortcatch*:

```
catcheR_sortcatch(
  group=c("sudo","docker"),
  folder,
  expression.matrix,
  swaps)
```

Example usage:

```
catcheR_sortcatch(
  group = "docker",
  folder = "path/to/working/folder",
  expression.matrix = "silencing_matrix.csv",
  swaps = "reliable_clones_swaps.csv")
```

*catcheR\_sortcatch* arguments:

- (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- (b) folder: string with the working folder path
- (c) expression.matrix: string with the filename of the annotated count matrix CSV
- (d) swaps: a character string with the filename of the txt file (step 1b)

*catcheR\_sortcatch* output is an updated annotated gene expression matrix CSV file called "silencing\_matrix\_updated.csv", in which "BC" and "GENE" have been modified to reflect the actual shRNA encoded in each cell



## Perturbation effect analysis

Once each cell is assigned to its corresponding perturbation, the impact of each perturbation on gene expression can be assessed by comparing it to control conditions.

The second part of *catcheR* offers a comprehensive exploratory analysis, featuring visualizations and statistical summaries that highlight the biological effects of the perturbations of interest.

### Annotation

Before proceeding, the gene expression matrix should be annotated with gene symbols with *scannobyGtf* function from the R package *rCASC*. Optionally, before that, ribosomal and mitochondrial genes can be assessed (e.g. with the function *mitoRiboUmi* from the same package).

**Note:** After this step, row names of the matrix (the genes) will have the following format:

GeneSymbol:EnsemblID

E.g.:

ENSG000000000003:TSPAN6

### Data loading

This step enables the loading of single-cell data generated by either *catcheR\_10Xcatch* or *catcheR\_scicatch*, following GTF-based annotation. The data is imported into a Monocle object, where experimental design information is added, followed by normalization and clustering.

1. In a new working folder:

- (a) Copy the annotated count matrix (filtered\_annotated\_silencing\_matrix\_complete\_all\_samples.csv).
- (b) Copy the file rc\_barcodes\_genes.csv described above.
- (c) Create a new-line separated file listing the control genes (e.g. SCR, B2M).
- (d) Create a new-line separated file listing the control samples, if any (e.g. 1,3). These are the sample names also used by *aggr* (see CSV file used as input for *aggr*).
- (e) Create a newline-separated file listing the sample replicates (e.g., batch1, batch1, batch2, batch2). When multiple samples from different experiments are included, batch correction is recommended. The replicate names defined in this file must correspond to the samples in the same order.
- (f) Create a CSV file listing each sample along with its corresponding annotation name, which will be used for display instead of the sample number (this file is mandatory).  
Example of CSV file to download on [GitHub](#)
- (g) Create a newline-separated file listing the genes of interest whose expression you wish to visualize on the UMAP (optional).

2. Run *catcheR\_load*:

```
catcheR_load(  
  group="docker",  
  folder,  
  expression.matrix,  
  control_genes,  
  control_samples = NULL,  
  replicates = NULL,  
  sample_names,  
  resolution = 8e-4,  
  genes = NULL)
```

3263 Example usage:

```
3264 catcheR_load(  
3265   group="docker",  
3266   folder="/path/to/working/folder/",  
3267   expression.matrix =  
3268     "annotated_silencing_matrix_complete_all_samples.csv",  
3269   control_genes = "controls.txt",  
3270   control_samples = "noTET.txt",  
3271   replicates = "replicates.txt",  
3272   sample_names = "samples.csv",  
3273   resolution = 8e-4,  
3274   genes = "genelist.txt")
```

3275 The argument *resolution* refers to the resolution parameter used by Monocle's *cluster\_cells*  
3276 function for clustering.

3277 *catcheR\_load* outputs:

- 3278 1. *expression\_data.csv* and *cell\_metadata.csv* : These files can be used to create a Monocle  
3279 Cell Data Set (CDS) and are also included in the ready-to-load R object *starting\_cds.Rdata*.
- 3280 2. *UMAP.pdf* plots the dimensionality reduction and *UMAP\_gene\_expression.pdf* shows the  
3281 gene expression on the UMAP of the genes provided by the argument "genes".
- 3282 3. *UMAP\_clustering.pdf* show the clustering obtained on the UMAP with the provided reso-  
3283 lution.
- 3284 4. *processed\_cds.RData* is the Cell Data Set after normalization, dimensionality reduction,  
3285 clustering and calculation of trajectories.

3286 At the end of this step, the data are structured in a format compatible with the Monocle3  
3287 package. However, switching to other platforms such as Seurat or Scanpy is possible at any  
3288 point in the analysis. **Note:** When performing the standard iPS2-seq perturbation analysis,  
3289 always use the CDS object generated by *catcheR* after each step.

```
3290 library(SeuratWrappers)  
3291 library(Seurat)  
3292 seurat <- as.Seurat(cds, assay = NULL)  
3293 scanpy_sce <- as.SingleCellExperiment(seurat)
```

3294 The follow up analysis are: *catcheR\_pseudotime*, *catcheR\_modules* and *catcheR\_enrichment*.

- 3295 1. With the [Pseudotime](#) function, *catcheR* calculates the cumulative frequency of cells shar-  
3296 ing the same gene, shRNA, or clone, compared to each control, along the pseudotime  
3297 trajectory. *CatcheR* can compare the cumulative frequency curves by computing direc-  
3298 tional Kolmogorov-Smirnov statistics between the target and control.
- 3299 2. With the [Genes modules](#) function, *catcheR* identifies gene modules within the dataset  
3300 and computes the cumulative frequency of cells with the same gene, shRNA, or clone—  
3301 relative to each control—based on the expression of each gene module. Also in this case,  
3302 cumulative frequency curves are compared using a directional Kolmogorov-Smirnov test.
- 3303 3. With the [Enrichment / depletion analysis](#) function, *catcheR* calculates the enrichment of  
3304 cells with the same gene, shRNA, or clone within clusters, in comparison to each control.  
3305 *CatcheR* can compare each target and control quantity and obtain statistics using the  
3306 Fisher exact test.

## 3307 Pseudotime

3308 Since monocle pseudotime calculation requires the user to select the starting point interactively,  
3309 the first step of the pseudotime analysis needs to be done separately within monocle3.

3310 1. After loading the processed\_cds.RData in R, calculate the pseudotime trajectory and plot  
3311 it with the R script below:

```
3312 library(monocle3)
3313 cds <- order_cells(cds)
3314 pt = as.data.frame(pseudotime(cds))
3315 names(pt) = c("pseudotime")
3316 write.csv(pt, paste0(dir, "/pseudotime.csv"))
3317
3318 plot_cells(cds,
3319            color_cells_by = "pseudotime",
3320            label_cell_groups=FALSE,
3321            label_leaves=FALSE,
3322            label_branch_points=FALSE,
3323            graph_label_size=1.5)
```

3324 2. The working folder will contain the output of *catcherR\_load* and the CSV file obtained after  
3325 running the monocle pseudotime function (e.g. "pseudotime.csv").

3326 3. Run *catcherR\_pseudotime*:

```
3327 catcherR_pseudotime(
3328   group=c("docker", "sudo"),
3329   folder,
3330   cds,
3331   pseudotime,
3332   all = FALSE)
```

3333 Example usage:

```
3334 catcherR_pseudotime(
3335   group="docker",
3336   folder="/path/to/working/folder/",
3337   cds = "processed_cds.RData",
3338   pseudotime = "pseudotime.csv")
```

3339 Argument "all" is a logical operator indicating whether to perform the Kolmogorov-Smirnov  
3340 test against the controls together (all = F) or against each control separately (all = T). The  
3341 default argument is false.

3342 Below is a list of the *catcherR\_pseudotime* outputs:

- 3343 (a) The "cumulative\_frequency\_pseudotime" plots show the number of cells at each  
3344 given point of the pseudotime for different groups of cells at different comparison  
3345 levels - gene, shRNA, and clone.
- 3346 (b) The "ks\_statistics" CSV files that include the Kolmogorov-Smirnov test results com-  
3347 paring cumulative the frequency based on the pseudotime between KD and different  
3348 controls
- 3349 (c) The Volcano plots of the results of the Kolmogorov-Smirnov test show the signifi-  
3350 cance and fold change between the pseudotime cumulative curves. Different con-  
3351 trols and levels are used.
- 3352 (d) The "correlated\_pseudotime\_gene\_exp.pdf" showing the expression on the UMAP  
3353 of the genes most correlated with the pseudotime.

## 3354 **Genes modules**

3355 This function uses Monocle to find gene modules which expression varies in different pertur-  
3356 bation groups, i.e., perturbed gene or shRNA or clones.

### 3357 1. Run *catcheR\_modules*:

```
3358     catcheR_modules(group=c("docker", "sudo"),  
3359                     folder,  
3360                     cds,  
3361                     resolution=1e-2)  
3362
```

### 3363 Example usage:

```
3364     catcheR_modules(group="docker",  
3365                     folder="/30tb/3tb/data/ratto/testing/",  
3366                     cds = "processed_cds.RData")  
3367
```

3368 The resolution parameter refers to the value used in Monocle's *find\_gene\_modules* func-  
3369 tion, which influences the number of gene modules identified and the number of genes  
3370 included in each module—higher resolution typically results in more, smaller modules.

3371 Below are listed the outputs of *catcheR\_modules*:

- 3372 (a) gene\_modules.csv, listing the genes present in each module.
- 3373 (b) Heatmap plots display the Z-scores of each module—either all modules or the top  
3374 10 most variable ones—across perturbation groups.
- 3375 (c) The "modules\_cells" folder contains tables with aggregated module expression val-  
3376 ues for each cell. These can be used as input for *catcheR\_pseudotime* in place of  
3377 the pseudotime CSV file, allowing the analysis of cumulative frequency based on  
3378 module expression.

### 3379 2. Run *catcheR\_pseudotime* on the CSV files stored in the "modules\_cells" folder to use 3380 the data generated by *catcheR\_modules* as described in the previous example in the 3381 [Pseudotime](#) section

## 3382 **Enrichment / depletion analysis**

3383 This function evaluates whether perturbation groups are enriched or depleted in cells (number)  
3384 or in specific cell subpopulations (clusters).

### 3385 1. Run *catcheR\_enrichment*:

```
3386     catcheR_enrichment(  
3387         group=c("docker", "sudo"),  
3388         folder,  
3389         file,  
3390         meta,  
3391         timepoint = "PSC",  
3392         control_gene = "SCR",  
3393         min_cells_cluster = 70,  
3394         min_cells_shRNA = 40)  
3395
```

### 3396 Example usage:

```

3397  catcheR_enrichment(group = "docker",
3398  folder = "/3tb/data/ratto/aggr/test /"
3399  file = "processed_cds.RData",
3400  meta = "cell_metadata.csv",
3401  timepoint = "PSC",
3402  control_gene = "SCR")
3403

```

3404 The required input is the cds file generated from *catcheR\_load*.

3405 *Timepoint* refers to the baseline time point used as a control for statistical analysis. It is  
3406 required when experiments include multiple time points and aim to assess enrichment or  
3407 depletion across these time points. *control\_gene* specifies the control gene used as a  
3408 reference for statistical comparisons.

3409 Below are the outputs of *catcheR\_enrichment*:

- 3410 (a) Plots of cells in each perturbation group.
- 3411 (b) Volcano plot showing enrichment or depletion of cell numbers in perturbation groups  
3412 compared to the control, based on fold change (log2ratio and FC) and statistical sig-  
3413 nificance. A corresponding bar plot displays the log2FC values for each perturbation.
- 3414 (c) Barplots showing the distribution of cells from different perturbation groups in the  
3415 clusters.
- 3416 (d) Volcano plot showing the results of Fisher's exact test, comparing the distribution of  
3417 cells from perturbation groups across Monocle-derived clusters. The plot displays  
3418  $-\log_{10}$  adjusted p-values versus log2 fold changes relative to the control group.
- 3419 (e) Table with Fisher's statistics.