

2209 Supplemental Protocol 4

2210 iPS2-seq design and analysis with *catcheR*

2211 This Supplemental Protocol describes the following procedures:

- 2212 1. [Overview](#)
- 2213 2. [Installation](#)
- 2214 3. [Oligonucleotides design](#)
- 2215 4. [Pooled cloning step 1 plasmid QC](#)
- 2216 5. [Pooled cloning step 2 plasmid QC](#)
- 2217 6. [Generation of cell-by-gene matrices](#)
- 2218 7. [iPS2-10X-seq perturbation deconvolution](#)
- 2219 8. [iPS2-sci-seq perturbation deconvolution](#)
- 2220 9. [Barcode reassignment](#)

2221 Overview

2222 *catcheR* is a comprehensive bioinformatic package for designing and analyzing iPS2-seq ex-
2223 periments. It comprises the following functions ([Figure SP4.1](#)):

- 2224 1. *catcheR_design*, which designs oligonucleotides for [Supplemental Protocol 1 - Design](#)
2225 [shRNA oligonucleotides](#), facilitating shRNA library cloning
- 2226 2. *catcheR_step1QC*, which analyses the results of [Supplemental Protocol 1 - Intermediate](#)
2227 [plasmid pool QC](#), assessing pooled cloning step 1 plasmids for barcode swaps
- 2228 3. *catcheR_step2QC*, which analyses the results of [Supplemental Protocol 1 - Final plasmid](#)
2229 [pool QC](#), assessing pooled cloning step 2 for shRNA representation
- 2230 4. *catcheR_scicount*, which analyses 2-level indexing sci-RNA-seq data, facilitating the gen-
2231 eration of gene expression matrix for iPS2-sci-seq experiments
- 2232 5. *catcheR_scicatch*, which assigns shRNA perturbations to single nuclei transcriptomes
2233 obtained by [Supplemental Protocol 2](#), enabling the primary analysis of iPS2-sci-seq
- 2234 6. *catcheR_10Xcatch*, which assigns shRNA perturbations to single cell transcriptomes ob-
2235 tained by [Supplemental Protocol 3](#), enabling the primary analysis of iPS2-10X-seq
- 2236 7. *catcheR_scicatchQC* and *catcheR_10XcatchQC*, which use the outputs of *catcheR_sci-*
2237 *catch* and *catcheR_10Xcatch*, respectively, to fine-tune shRNA assignment thresholds
- 2238 8. *catcheR_filtercatch*, which leverages on the output of *catcheR_scicatchQC* and *catcheR_*
2239 *10XcatchQC* to filter single nuclei/cell transcriptomes expressing a single shRNA
- 2240 9. *catcheR_sortcatch*, which quality controls the cell-by-gene matrix based on the results of
2241 *catcheR_step1QC*, reassigning hPSC clones with barcode swaps to the correct shRNA
- 2242 10. *catcheR_scinocatch* and *catcheR_10Xnocatch*, which identify cells expressing no shRNA
2243 in iPS2-sci-seq and iPS2-10X-seq experiments, respectively, adding them to the cell-by-
2244 gene matrix to be used as additional controls

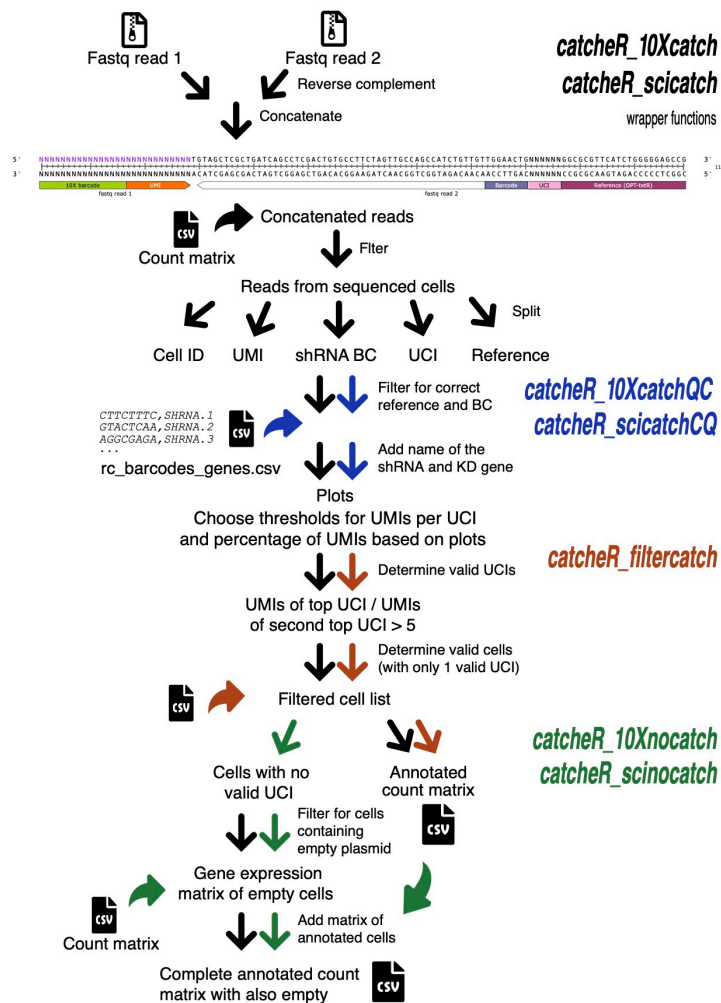


Figure SP4.1. Overview of catcheR

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.

Installation

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

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

2. Install catcheR

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

```
install.packages("devtools")
```

(b) Install catcheR from GitHub

```
install_github("alessandro-bertero/catcheR")
```

(c) Load catcheR in your R environment

```
library(catcheR)
```

Oligonucleotides design

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

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

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

```
CCGGCAAGTACTCCTTGCTGGATTGCTCGAGCAATCCAGCAAGGAGTACTTGTGTTTTC , CAGTTCCA , SMAD2.1
...
```

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

```
GTCGAC
ATTTAAAT
GGCGCGCC
```

2. Run *catcheR_design*:

```
catcheR_design(
  group=c("docker","sudo"),
  folder,
  sequences,
  gibson.five = "AGTTCCTATCAGTGATAGAGATCCC",
  gibson.three = "GTAGCTCGCTGATCAGC",
  fixed = "GTCGACATTTAAATGGCGCGCC",
  restriction.sites = NULL)
```

catcheR_design arguments:

- (a) group: string with two options: sudo or docker, depending on the user group. For a detailed explanation of Docker user groups, see [this page](#)
- (b) folder: string with the working folder path
- (c) sequences: string with the csv file name from step 1a
- (d) gibson.five: (optional) string with the 5' Gibson homology
- (e) gibson.three: (optional) string with the 3' Gibson homology
- (f) fixed: (optional) character string with the multicloning site
- (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 SP1.1](#))
- (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":

```
AGTTCCTATCAGTGATAGAGATCCCGACATAATCACTGCGTAATCCTCagatctTACGCAGTGATTATGTCCTTTTTTGT -
CGACATTTAAATGGCGCGCCNNNNNGCTGAAGAGTAGCTCGCTGATCAGC , GATA4
...
```

2311 Pooled cloning step 1 plasmid QC

2312 This steps complements [Supplemental Protocol 1 - Intermediate plasmid pool QC](#)

2313 1. In a new working folder, prepare the following files:

2314 (a) Fastq/fq or fastq.gz files with demultiplexed read 1 from the NGS run

2315 (b) A tab-separated txt file with the shRNA names and their full sequences

2316 SMAD2.1 GCAAGTACTCCTTGCTGGATTGCTCGAGCAATCCAGCAAGGAGTACTTG

2317 ...

2318 (c) (Optional) - A txt file with a newline-separated list of clones of interest as "BC_UCI"

2319 CAAGAGCC_CATCGT

2320 ...

2321

2322

2323

2. Run *catcheR_step1QC*:

```
2324 catcheR_step1QC(  
2325     group=c("docker","sudo"),  
2326     folder,  
2327     fastq.read1,  
2328     DIs = 100,  
2329     ratio = 10,  
2330     plot.threshold = 2000,  
2331     clones = NULL)
```

2332

2333

catcheR_step1QC arguments:

2334 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))

2335 (b) folder: string with the working folder path

2336 (c) fastq.read1: string with the read 1 filename from step [1a](#)

2337 (d) DIs: integer of the minimum number of diversity indexes (DIs, pseudo-unique reads)
2338 of the most represented shRNA matching to a given UCI-BC; in combination with
2339 "ratio", it selects UCI-BCs for which it is possible to reliably assign an shRNA

2340 (e) ratio: integer of the minimum ratio between the number of DIs of the most repre-
2341 sented and second most represented shRNAs matching to a given UCI-BC

2342 (f) plot.threshold: integer of the minimum number of DIs per UCI-BC for output [2b](#)

2343 (g) clones: (optional) a string with the txt file from step [1b](#)

2344 Example usage:

```
2345 catcheR_step1QC(  
2346     group = "docker",  
2347     folder = "path/to/folder",  
2348     fastq.read1 = "filename.fq",  
2349     clones = "filename.txt")
```

2350 *catcheR_step1QC* key outputs:

2351 (a) "reliable_clones_swaps.txt", which lists UCI-BCs with reliable evidence of shRNA-
2352 barcode swap, and is used as input for [Barcode reassignment](#)

2353 (b) Bar chart of the number of DIs for each clone above "plot.threshold"

2354 (c) Bar chart of the number of DIs associated to each BCs, shRNAs, and reliable swaps

2355 (d) (Optional) - Text files and bar charts with the number of DIs for each shRNA matching
2356 to each clone of interest ([Figure SP4.2](#))

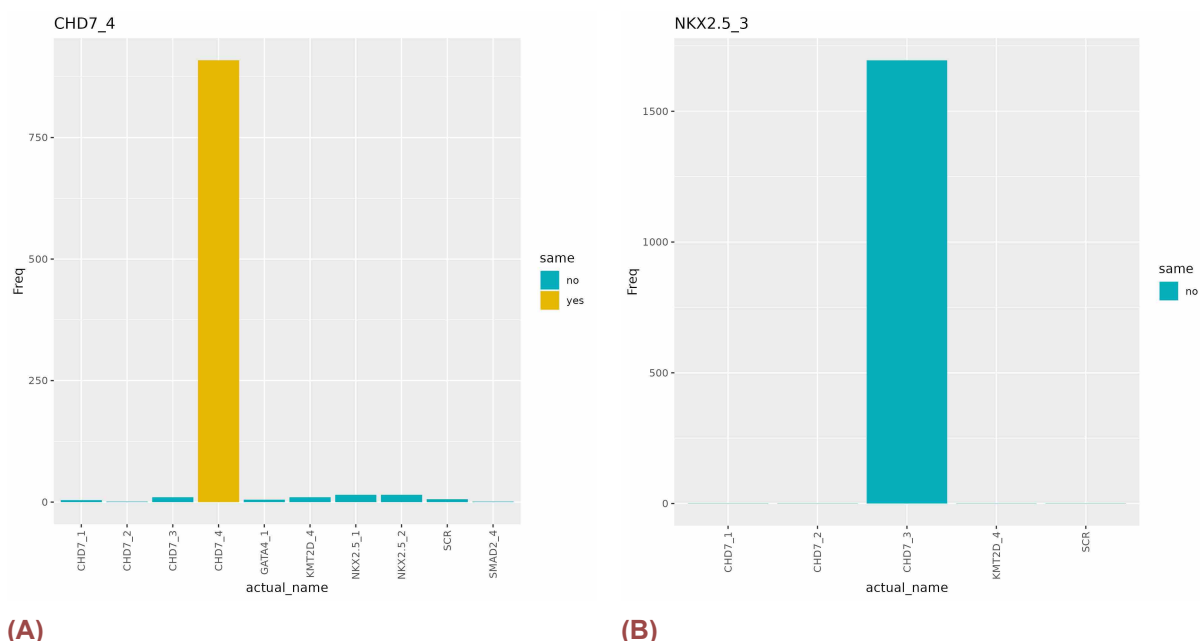


Figure SP4.2. Exemplary quantification of DIs for bacterial clones carrying the expected shRNA (A) or with robust evidence of a barcode swap that can be reassigned (B; [Figure 1E](#)). These graphs allow visual evaluation of appropriate thresholds for "DI" and "ratio".

2357 Pooled cloning step 2 plasmid QC

2358 This steps complements [Supplemental Protocol 1 - Final plasmid pool QC](#)

2359 1. In a new working folder, prepare the following files:

- 2360 (a) Fastq/fq or fastq.gz files with demultiplexed read 1 from the NGS run
- 2361 (b) "rc_barcode_genes.csv", a csv file with two columns: (1) the shRNA BCs; (2) the
- 2362 matching shRNA names in the format "GENE.shRNAID"

2363 CAAGAGCC , SMAD2 . 1

2364 . . .

- 2365 (c) (Optional) - A txt file with clones of interest (step [1a](#) of the previous section)
- 2366

2367 2. Run *catcherR_step2QC*:

```
2368 catcherR_step2QC (
2369     group=c("docker", "sudo"),
2370     folder,
2371     fastq.read1,
2372     DIs = 1000,
2373     clones = NULL)
```

2374 *catcherR_step2QC* arguments:

- 2376 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- 2377 (b) folder: string with the working folder path
- 2378 (c) fastq.read1: string with the read 1 file from step [1a](#)
- 2379 (d) DIs: integer of the minimum number of DIs for a given UCI-BC; it selects reliably-
- 2380 measured UCI-BCs
- 2381 (e) clones: (optional) a string with the txt file from step [1c](#)

Example usage:

```
catcher_step2QC(  
  group= "docker",  
  folder = "path/to/working/folder",  
  fastq.read1 = "filename.fq",  
  clones = "filename.txt")
```

catcher_step2QC key outputs:

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

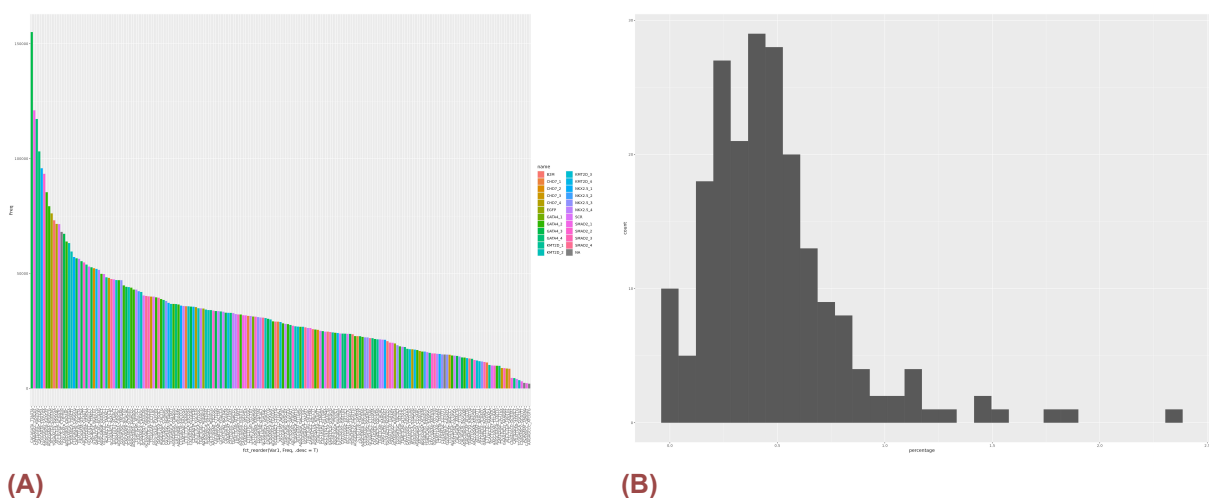


Figure SP4.3. Exemplary quantification of DIs for bacterial clones represented by more than 1.000 DIs (A), and of the frequency of bacterial clones for each shRNA. These graphs allow to confirm whether shRNAs are normally distributed in the final plasmid pool

Generation of cell-by-gene matrices

Obtain iPS2-10X-seq count matrices

- Download and install cellranger from the [10X Genomics download resource page](#); alternatively, a [Docker](#) container with cellranger 7 can be pulled from our Docker repository
- Demultiplex Illumina base calls with *cellranger mkfastq*, following 10X Genomics' manual. In the simple CSV sample sheet, include the 10X Genomics indexes used to generate the GEX, CMO, and UCI-BC libraries in [Supplemental Protocol 3](#)
- Obtain cell-by-gene count matrices using *cellranger count* (single sample) or *cellranger multi* (multiplexed experiment).

Note: count matrices from multiple samples from the same experiment can be aggregated with *cellranger aggr*, but keep the individual count matrices as these are required to run *catcherR_10Xcatch* (which requires GEX and UCI-BC fastq files to have the same library ID to match UCI-BC and cells). Our GitHub repository contains R [scripts](#) that can be used to subsequently integrate the outputs of *catcherR_10Xcatch* and *cellranger aggr*

- Run *cell ranger matrix2csv* to transform the sparse count matrix in a dense csv matrix

After obtaining the count matrix, proceed to [iPS2-10X-seq perturbation deconvolution](#).

2411 Obtain iPS2-sci-seq count matrices

2412 1. *catcheR_scicount* is a wrapper of the "bbi-sci" pipeline developed by the Brotman Baty
2413 Institute for Precision Medicine, which was implemented and dockerized in *catcheR* to be
2414 used on any operating system.

2415 2. Demultiplex Illumina base calls to fastq files

2416 (a) Create a "SampleSheet.csv" file with as many sample rows as the PCR wells from
2417 [Indexing PCR and NGS of Supplemental Protocol 2](#), where "Sample_ID" is the well
2418 identified with format [A-H][01-12], and index and index2 are the i7 and i5 indexes
2419 used in the corresponding row and column (refer to [Table SP2.1](#))

2420 (b) Run Illumina *bcl2fastq* following Illumina's manual

2421 (c) Run fastQC to confirm the quality of the fastq files

2422 3. In a new working folder:

2423 (a) Create the subfolder "fastq", and copy all the demultiplexed fastq.gz files. Ensure
2424 that all file names begin with the well coordinate (e.g. A01)

2425 (b) Create a tab separated txt file called "sci-RNA-seq-8.RT.oligos" with the association
2426 between RT wells and RT barcode sequences (refer to [Table SP2.1](#))

```
2427 A01 TTCTCGCATG  
2428 ...
```

2429 (c) Create the subfolder "GENOMES", and copy the annotated genome (i.e., [GRCh38](#))
2430

2431 4. Run *catcheR_scicount*:

```
2432 catcheR_scicount(  
2433   group=c("docker", "sudo"),  
2434   folder,  
2435   sample.name,  
2436   UMI.cutoff)
```

2437 *catcheR_scicount* arguments:

2438 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))

2439 (b) folder: string with the working folder path

2440 (c) sample.name: string with the name of the experiment

2441 (d) UMI.cutoff: integer of the minimum number of UMI per nucleus needed to consider
2442 the single cell transcriptome valid

2443 Example usage:

```
2444 catcheR_scicount(  
2445   group = "docker",  
2446   folder = "path/to/file",  
2447   sample.name = "experiment",  
2448   UMI.cutoff = 500)
```

2449 *catcheR_scicount* outputs: cell-by-gene count matrix (sparse), knee-plot of UMI per cell
2450 ([Figure S2F](#)), and statistics files, which are all placed in the "final-output" folder
2451

2452 5. Run *catcheR_scicsv* using the "final-output" folder from *catcheR_scicount* as its only input
2453 (a folder containing the files "gene.annotations", "cell.annotations" and "UMI.count.matrix"):

```
2454 catcheR_scicsv(  
2455   group=c("docker", "sudo"),  
2456   folder)  
2457 .
```

2458 Output files are "exp_mat.csv" and "exp_mat_no0.csv" (filtered to exclude genes with 0
2459 counts in all cells)

2460 After obtaining the count matrix, proceed to [iPS2-sci-seq perturbation deconvolution](#).

2461 iPS2-10X-seq perturbation deconvolution

2462 This section describes a typical analysis pipeline starting with *catcheR_10Xcatch* to perform
2463 a full analysis with automatic thresholds, after which *catcheR_10XcatchQC* can fine-tune the
2464 thresholds to re-filter cells with *catcheR_filtercatch*. Lastly, *catcheR_10Xnocatch* can add cells
2465 expressing no shRNA as additional controls to the final annotated cell-by-gene count matrix.

2466 1. In a new working folder:

- 2467 (a) Copy the fastq/fq or fastq.gz files with demultiplexed read 1 and read 2 of the UCI-BC
2468 library (step 2 of [Obtain iPS2-10X-seq count matrices](#))
- 2469 (b) Copy the cell-by-gene count matrix in csv format (step 4 of [Obtain iPS2-10X-seq](#)
2470 [count matrices](#))
- 2471 (c) Create a csv file named "rc_barcodes_genes.csv" with two columns: (1) the shRNA
2472 BCs; (2) the matching shRNA names in the format "GENE.shRNAID"
2473 CAAGAGCC , SMAD2 . 1
2474 ...
2475

2476 2. Run *catcheR_10Xcatch* to execute the complete analysis with automatic thresholding:

```
2477 catcheR_10Xcatch(  
2478   group=c("docker","sudo"),  
2479   folder,  
2480   fastq.read1,  
2481   fastq.read2,  
2482   expression.matrix,  
2483   reference = "GGCGCGTTCATCTGGGGGAGCCG",  
2484   UCI.length = 6,  
2485   threads = 2,  
2486   percentage = 15,  
2487   mode = "bimodal")
```

2488 *catcheR_10Xcatch* arguments:
2489

- 2490 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))
- 2491 (b) folder: string with the working folder path
- 2492 (c) fastq.read1: string with the filename of read 1 fastq/fq or fastq.gz (step [1a](#))
- 2493 (d) fastq.read2: string with the filename of read 2 fastq/fq or fastq.gz (step [1a](#))
- 2494 (e) expression.matrix: string with the filename of the count matrix file csv (step [1b](#))
- 2495 (f) reference: string with the reverse complement of the sequence before the shRNA
2496 BC at the start of read 2 (default is the 3' end of the OPTtetR cDNA; [Figure S3F](#))
- 2497 (g) UCI.length: integer of the UCI length (default is 6)
- 2498 (h) threads: integer of the threads to be used for parallelization (default is 2)
- 2499 (i) percentage: integer of the percentage of UMIs supporting a given UCI over the total
2500 UMIs supporting all UCIs in a given cell; it is used as threshold to consider the UCI
2501 valid ([Figure SP4.4D](#)); the recommended default is 15
- 2502 (j) mode: string with two options: bimodal or noise, defining the mode for automatic
2503 thresholding the minimum number of UMIs per UCI to consider the UCI valid.
 - 2504 i. "bimodal" (default) sets the threshold at the valley of the bimodal UMIxUCI dis-
2505 tribution ([Figure SP4.4C](#))
 - 2506 ii. "noise" sets the threshold at 1.35 * the number of UCIs supported by a single
2507 UMI; used if the UMIxUCI distribution is not bimodal

Example usage:

```
catcheR_10Xcatch(group = "docker",  
  folder = "path/to/folder",  
  fastq.read1 = "R1.fastq.gz",  
  fastq.read2 = "R2.fastq.gz",  
  expression.matrix = "filename.csv",  
  threads = 12)
```

catcheR_10Xcatch key outputs (all inside the working folder):

- (a) "log.txt" with the starting number of reads
- (b) "log2.txt" with the number of cells, UMIs, UCIs, and the calculated thresholds of UMI per UCI for both "bimodal" and "noise" (only one is chosen based on the "mode" argument; the other one is provided for a potential reanalysis)
- (c) Bar charts of UMI counts per shRNA and target gene ([Figure SP4.4A](#) and [Figure SP4.4B](#))
- (d) Frequency histograms of UCIs supported by a certain number of UMIs (UMIxUCI), to interpret the signal/noise of the experiment and possibly set a custom UMIxUCI threshold for subsequent reanalysis ([Figure SP4.4C](#)).
Note: here and below, copies of the same UCI expressed by different cells are plotted and analyzed separately
- (e) Frequency histogram of UCIs supported by a certain fraction of UMIs over the total number of UMIs supporting all UCIs in a given cell (UMIpercentagexUCI), to further assess signal/noise and possibly adjust the default threshold ([Figure SP4.4D](#))
- (f) Dot plots that combine the UMIpercentagexUCI and UMIxUCI data on the x and y axes, respectively, with each dot representing one or more UCI (quantified by the dot size). Dot colors indicate either the number of valid UCIs (i.e., shRNAs) in the cell containing a given UCI ([Figure SP4.4E](#)), or whether said UCI is the only valid one in the relevant cell and is thus assigned to such cell ([Figure SP4.4F](#))
- (g) "log_part3.txt" with how many cells were assigned to a single shRNA or were filtered out due to zero or multiple shRNAs
- (h) "silencing_matrix.csv" is the key output used for the secondary analyses: the cell-by-gene count matrix provided as input, filtered and annotated with the shRNA encoded in each cell. Cells names are modified as follows:

```
cellID_UMIxUCI_BC_GENE_UCI
```

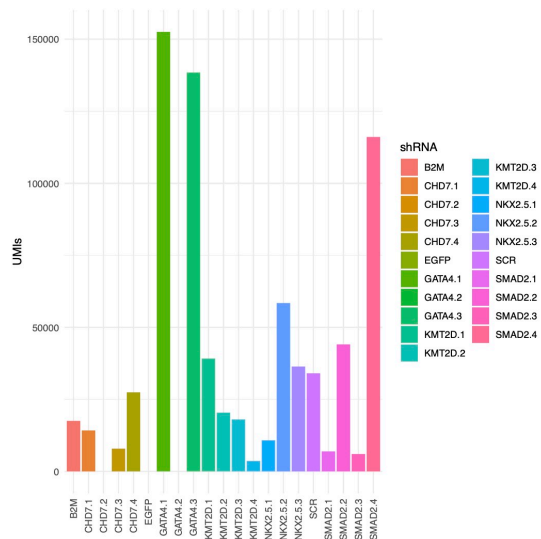
Where:

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

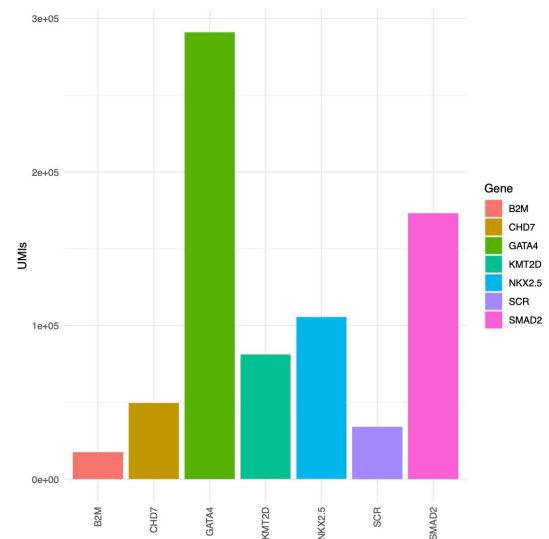
Example of cell name after the analysis:

```
TTCTAACCACAGTCGC_180_CGTGATGC_NKX2.5_ACAGTG
```

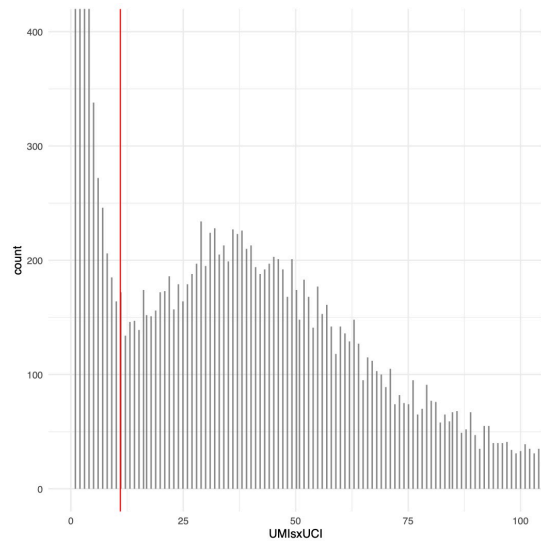
This annotation can be leveraged with the [scripts](#) available on our GitHub repository to enable a variety of secondary analyses described in the manuscript, or using custom scripts for other types of secondary analyses



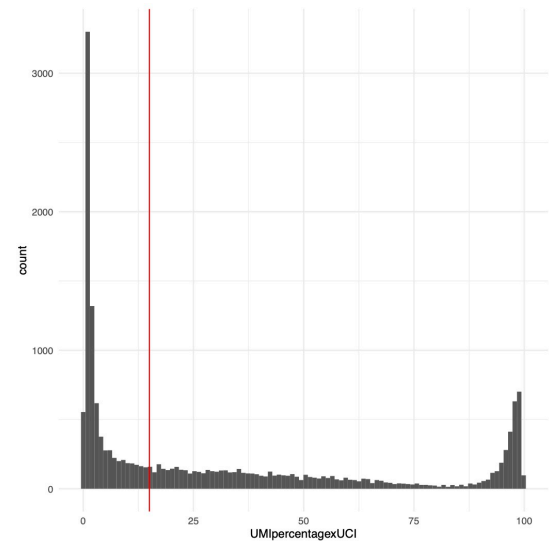
(A)



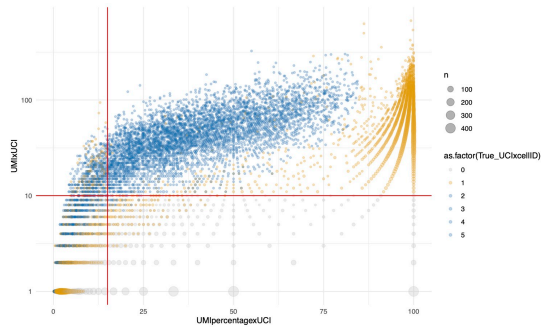
(B)



(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)

2555 **OPTIONAL: fine tune iPS2-10X-seq perturbation assignment**

- 2556 1. In the same working folder used to run *catcheR_10Xcatch*, run *catcheR_10XcatchQC*,
2557 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
2558

```
2559 catcheR_10XcatchQC(  
2560     group=c("docker","sudo"),  
2561     folder,  
2562     reference = "GGCGCGTTCATCTGGGGGAGCCG",  
2563     mode = "bimodal")
```

2564 Example usage:
2565

```
2566 catcheR_10XcatchQC(  
2567     group = "docker",  
2568     folder = "path/to/working/folder",  
2569     mode = "noise")
```

2570 *catcheR_10XcatchQC* arguments and outputs are the same as *catcheR_10Xcatch*

- 2571 2. In the same working folder used to run *catcheR_10Xcatch*, run *catcheR_filtercatch*, which
2572 uses the output of *catcheR_10Xcatch* to filter cells based on custom thresholds

```
2573 catcheR_filtercatch(  
2574     group=c("docker","sudo"),  
2575     folder,  
2576     expression.matrix,  
2577     UMI.count,  
2578     percentage = 15)
```

2579 Example usage:
2580

```
2581 catcheR_filtercatch(group = "docker",  
2582     folder = "path/to/working/folder",  
2583     expression.matrix = "filename.csv",  
2584     UMI.count = 5)
```

2585 *catcheR_10XcatchQC* argument *UMI.count* is an integer of the custom UMI threshold;
2586 the other arguments and the outputs are the same as *catcheR_10Xcatch*. ;
2587

2588 **OPTIONAL: identify cells expressing no shRNA in iPS2-10X-seq**

- 2589 1. In the same working folder used to run *catcheR_10Xcatch*, run *catcheR_10Xnocatch*:

```
2590 catcheR_10Xnocatch(  
2591     group=c("docker","sudo"),  
2592     folder,  
2593     expression.matrix,  
2594     threshold)
```

2595 Example usage:

```
2596 catcheR_10Xnocatch(group = "docker",  
2597     folder = "path/to/folder",  
2598     expression.matrix = "filename.csv",  
2599     threshold = 5)
```

2600 *catcheR_10Xnocatch* argument *threshold* is an integer of the minimal UMI count for the
2601 "empty" reference [Figure SP4.3B](#) needed to confidently identify it as having integrated
2602 said plasmid; all other arguments are the same as *catcheR_10Xcatch*

2603 *catcheR_10Xnocatch* output is an updated "silencing_matrix.csv" in which names of empty
2604 cells are modified as "cellID_?_empty_NA_empty".

2605 iPS2-sci-seq perturbation deconvolution

2606 This section describes a typical analysis pipeline similar to the one described in the previous
2607 section, but leveraging on the *catcheR_scicatch*, *catcheR_scicatchQC*, *catcheR_filtercatch*,
2608 and *catcheR_scinocatch* functions.

2609 1. In a new working folder:

2610 (a) Create a subfolder called "fastq", and copy all the demultiplexed fastq.gz files (step
2611 2 of [Obtain iPS2-sci-seq count matrices](#)). Ensure that all file names begin with the
2612 well coordinate (e.g. A01)

2613 (b) Copy the cell-by-gene expression matrix csv file obtained with *catcheR_scicount*
2614 function (step 4 of [Obtain iPS2-sci-seq count matrices](#)).

2615 (c) Create a csv file called "rc_barcodes_genes.csv" with two columns: (1) the shRNA
2616 BCs; (2) the matching shRNA names in the format "GENE.shRNAID"

2617 CAAGAGCC , SMAD2 .1

2618 ...

2619 (d) Copy the text file "sci-RNA-seq-8.RT.oligos" used by *catcheR_scicount* (step 3b of
2620 [Obtain iPS2-sci-seq count matrices](#)).
2621

2622 2. Run *catcheR_scicatch* to execute the complete analysis with automatic tresholding; the
2623 arguments are the same as for *catcheR_10Xcatch* (step 2 of [iPS2-10X-seq perturbation](#)
2624 [deconvolution](#)), except that no filenames are provided since fastq files must be in the
2625 "fastq" subfolder in the working directory:

```
2626 catcheR_scicatch(  
2627   group=c("docker","sudo"),  
2628   folder,  
2629   expression.matrix,  
2630   reference = "GGCGCGTTCATCTGGGGGAGCCG",  
2631   UCI.length = 6,  
2632   threads = 2,  
2633   percentage = 15,  
2634   mode = "bimodal")
```

2635
2636 Example usage:

```
2637 catcheR_scicatch(group = "docker",  
2638   folder = "path/to/working/folder",  
2639   expression.matrix = "filename.csv",  
2640   threads = 12)
```

2641
2642 *catcheR_scicatch* key outputs are the same as for *catcheR_10Xcatch* (step 2 of [Genera-](#)
2643 [tion of cell-by-gene matrices](#) and [Figure SP4.3](#)), except that in "silencing_matrix.csv" "cell
2644 ID" indicates the PCR well and RT barcode ID (which can be leveraged to identify nuclei
2645 from different samples pooled on the same RT plate)

2646 Example of cell name after the analysis:

2647 P24__RT_27_7_GCCTGTGT_SCR_ACGGTC

2648
2649 In addition, *catcheR_scicatch* provides two quality control outputs to evaluate experimen-
2650 tal biases during sci-RNAseq library preparation: *demux* and *RT* detail about how many
2651 cells were identified from each row and column of the PCR and RT plates, respectively.

2652 **OPTIONAL: fine tune iPS2-sci-seq perturbation assignment**

- 2653 1. Run *catcheR_scicatchQC* following the steps described for *catcheR_10XcatchQC* in [OPTIONAL: fine tune iPS2-10X-seq perturbation assignment](#): the functions share the same
2654 arguments and outputs, but *catcheR_scicatchQC* used the output of *catcheR_scicatch*
2655

2656 Example usage:

```
2657 catcheR_scicatchQC(  
2658     group = "docker",  
2659     folder = "path/to/working/folder",  
2660     mode = "noise")  
2661
```

- 2662 2. Run *catcheR_filtercatch* function as described in in [OPTIONAL: fine tune iPS2-10X-seq](#)
2663 [perturbation assignment](#): this function also works on the output of *catcheR_scicatch* or
2664 *catcheR_scicatchQC*

2665 **OPTIONAL: identify cells expressing no shRNA in iPS2-sci-seq**

- 2666 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
2667 same arguments and outputs, but *catcheR_scinocatch* uses the output of *catcheR_scicatch*
2668

2669 Example usage:

```
2670 catcheR_scinocatch(group = "docker",  
2671     folder = "path/to/working/folder",  
2672     expression.matrix = "filename.csv",  
2673     threshold = 5)  
2674
```

2675 **Barcode reassignment**

2676 *catcheR_sortcatch* is an optional function that corrects the annotated cell-by-gene count matrix
2677 obtained with *catcheR_10Xcatch* or *catcheR_scicatch* reassigning the perturbation of any cell
2678 belonging to a hPSC clone with reliable evidence of a shRNA-barcode swap, based on the
2679 results of *catcheR_step1QC*.

- 2680 1. In a new working folder:

- 2681 (a) Copy the annotated count matrix (i.e., "silencing_matrix.csv")
2682 (b) Copy the txt file with the list of UCI-BCs with reliable evidence of a shRNA-barcode
2683 swap (i.e. "reliable_clones_swaps.txt"; step 2a of [Pooled cloning step 1 plasmid QC](#))

- 2684 2. Run *catcheR_sortcatch*:

```
2685 catcheR_sortcatch(  
2686     group=c("sudo","docker"),  
2687     folder,  
2688     expression.matrix,  
2689     swaps)
```

2690 Example usage:

```
2691 catcheR_sortcatch(  
2692     group = "docker",  
2693     folder = "path/to/working/folder",  
2694     expression.matrix = "silencing_matrix.csv",  
2695     swaps = "reliable_clones_swaps.txt")
```

2696 *catcheR_sortcatch* arguments:

2697 (a) group: string with two options: sudo or docker, depending on the user group ([info](#))

2698 (b) folder: string with the working folder path

2699 (c) expression.matrix: string with the filename of the annotated count matrix csv

2700 (d) swaps: a character string with the filename of the txt file (step [1b](#))

2701 *catcheR_sortcatch* output is an updated annotated gene expression matrix csv file called

2702 "silencing_matrix_updated.csv", in which "BC" and "GENE" have been modified to reflect

2703 the actual shRNA encoded in each cell