

Process scRNA-Seq reads in *scruff*

Zhe Wang

2019-10-29

Package

scruff 1.4.0

Contents

1	Introduction	2
2	Quick Start	2
3	Stepwise Tutorial For CEL-Seq Samples	4
3.1	Load Example Dataset	4
3.2	Demultiplex and Assign Cell Specific Reads	4
3.3	Alignment	5
3.4	UMI correction and Generation of Count Matrix	6
3.5	Visualization of QC metrics	6
3.6	Visualization of Read mapping locations.	13
4	10X BAM File Quality Assessment	14
5	Session Information	16

1 Introduction

scruff is a toolkit for processing single cell RNA-seq FASTQ reads generated by CEL-Seq and CEL-Seq2 protocols. It does demultiplexing, alignment, Unique Molecular Identifier (UMI) filtering, and transcript counting in an automated fashion and generates the gene count matrix, QC metrics and provides visualizations of data quality. This vignette provides a brief introduction to the *scruff* package by walking through the demultiplexing, alignment, and UMI-counting of a built-in publicly available example dataset ([van den Brink, et al. 2017](#)).

2 Quick Start

```
# Run scruff on example dataset
# NOTE: Requires Rsubread index and TxDb objects for the reference genome.
# For generation of these files, please refer to the Stepwise Tutorial.

library(scruff)
## Registered S3 method overwritten by 'GGally':
##   method from
##   +.gg      ggplot2

# Get the paths to example FASTQ, FASTA, and GTF files.
# Please note that because the following files are included in
# scruff R package, we use system.file() function to extract the paths
# to these files. If the user is running scruff on real data, the
# paths to the input FASTQ, FASTA, and GTF files should be provided,
# and there is no need to call system.file() function. For example,
# vlh1R1 <- "/PATH/TO/vandenBrink_1h1_L001_R1_001.fastq.gz"
# fasta <- "/Path/TO/GRCm38_MT.fa"
vlh1R1 <- system.file("extdata",
  "vandenBrink_1h1_L001_R1_001.fastq.gz",
  package = "scruff")
vlh1R2 <- system.file("extdata",
  "vandenBrink_1h1_L001_R2_001.fastq.gz",
  package = "scruff")
vb1R1 <- system.file("extdata",
  "vandenBrink_b1_L001_R1_001.fastq.gz",
  package = "scruff")
vb1R2 <- system.file("extdata",
  "vandenBrink_b1_L001_R2_001.fastq.gz",
  package = "scruff")
fasta <- system.file("extdata", "GRCm38_MT.fa", package = "scruff")
gtf <- system.file("extdata", "GRCm38_MT.gtf", package = "scruff")
```

Build Rsubread alignment index. This is for the alignment step. For test purpose, here we are aligning the example FASTQ files to the genes on mitochondrial chromosome only.

```
# NOTE: Rsubread package does not support Windows environment.
if (!requireNamespace("Rsubread", quietly = TRUE)) {
```

Process scRNA-Seq reads in *scruff*

```
message("Package \"Rsubread\" needed for \"alignRsubread\"",
       " function to work. ",
       "Please install it if you are using Linux or macOS systems. ",
       "The function is not available in Windows environment.")
} else {
  # Create index files for GRCm38.MT.
  # For details, please refer to Rsubread user manual.
  # Specify the basename for Rsubread index
  indexBase <- "GRCm38.MT"
  Rsubread::buildindex(basename = indexBase,
                      reference = fasta,
                      indexSplit = FALSE)
}
```

Now that everything is ready, we can run *scruff*. In experiment 1h1, cell barcodes 95 and 96 are empty well controls. In experiment b1, cell barcode 95 is bulk sample containing 300 cells. These information can be set by the `cellPerWell` argument. We apply cell barcode and UMI correction by setting `bcEdit` to 1 and `umiEdit` to 1. *scruff* makes use of the [SingleCellExperiment](#) package. The following command returns a `SingleCellExperiment` object containing UMI filtered count matrix as well as gene and sample annotations and QC metrics.

```
data(barcodeExample, package = "scruff")

if (!requireNamespace("Rsubread", quietly = TRUE)) {
  message("Package \"Rsubread\" needed. ",
        "Please install it if you are using Linux or macOS systems. ",
        "The function is not available in Windows environment.")
} else {
  sce <- scruff(project = "example",
               experiment = c("1h1", "b1"),
               lane = c("L001", "L001"),
               read1Path = c(v1h1R1, vb1R1),
               read2Path = c(v1h1R2, vb1R2),
               bc = barcodeExample,
               index = indexBase,
               unique = FALSE,
               nBestLocations = 1,
               reference = gtf,
               bcStart = 1,
               bcStop = 8,
               bcEdit = 1,
               umiStart = 9,
               umiStop = 12,
               umiEdit = 1,
               keep = 75,
               cellPerWell = c(rep(1, 94), 0, 0, rep(1, 94), 300, 1),
               cores = 2,
               verbose = TRUE)
}
```

Process scRNA-Seq reads in *scruff*

Visualize data quality.

```
data(sceExample, package = "scruff")
qc <- qcplots(sceExample)
```

3 Stepwise Tutorial For CEL-Seq Samples

3.1 Load Example Dataset

The *scruff* package contains 4 single cell RNA-seq FASTQ example files. Each file has 10,000 sequenced reads.

```
library(scruff)

# Get the paths to example FASTQ files.
# Please note that because the following files are included in
# scruff R package, we use system.file() function to extract the paths
# to these files. If the user is running scruff on real data, the
# paths to the input FASTQ, FASTA, and GTF files should be provided,
# and there is no need to call system.file() function. For example,
# vlh1R1 <- "/PATH/TO/vandenBrink_1h1_L001_R1_001.fastq.gz"
# fasta <- "/Path/TO/GRCm38_MT.fa"
vlh1R1 <- system.file("extdata",
  "vandenBrink_1h1_L001_R1_001.fastq.gz",
  package = "scruff")
vlh1R2 <- system.file("extdata",
  "vandenBrink_1h1_L001_R2_001.fastq.gz",
  package = "scruff")
vb1R1 <- system.file("extdata",
  "vandenBrink_b1_L001_R1_001.fastq.gz",
  package = "scruff")
vb1R2 <- system.file("extdata",
  "vandenBrink_b1_L001_R2_001.fastq.gz",
  package = "scruff")
```

3.2 Demultiplex and Assign Cell Specific Reads

Now the FASTQ files are ready to be demultiplexed. *scruff* package provides built-in predefined cell barcodes `barcodeExample` for demultiplexing the example dataset. In the example FASTQ files, read 1 contains cell barcode and UMI sequence information. Read 2 contains transcript sequences. The barcode sequence of each read starts at base 1 and ends at base 8. The UMI sequence starts at base 9 and ends at base 12. They can be set via `bcStart`, `bcStop`, and `umiStart`, `umiStop` arguments. Cell barcode correction can be set by `bcEdit` parameter. By default, reads with any nucleotide in the barcode and UMI sequences with sequencing quality lower than 10 (Phred score) will be excluded. The following command demultiplexes the example FASTQ reads and trims reads longer than 75 nucleotides. The command returns a `SingleCellExperiment` object.

Process scRNA-Seq reads in *scruff*

`gleCellExperiment` object whose `colData` contains the cell index, barcode, reads, percentage of reads assigned, sample, and FASTQ file path information for each cell. By default, the cell specific demultiplexed fastq.gz files are stored in `./Demultiplex` folder.

```
data(barcodeExample, package = "scruff")
de <- demultiplex(project = "example",
  experiment = c("1h1", "b1"),
  lane = c("L001", "L001"),
  read1Path = c(v1h1R1, vb1R1),
  read2Path = c(v1h1R2, vb1R2),
  barcodeExample,
  bcStart = 1,
  bcStop = 8,
  bcEdit = 1,
  umiStart = 9,
  umiStop = 12,
  keep = 75,
  minQual = 10,
  yieldReads = 1e+06,
  verbose = TRUE,
  overwrite = TRUE,
  cores = 2)
```

3.3 Alignment

scruff provides an alignment function `alignRsubread` which is a wrapper function to `align` in `Rsubread` package. It aligns the reads to reference sequence index and outputs sequence alignment map files in "BAM" or "SAM" format. For demonstration purpose, the built-in mitochondrial DNA sequence from GRCm38 reference assembly `GRCm38MitochondrialFasta` will be used to map the reads. First, a [Rsubread](#) index for the reference sequence needs to be generated.

```
# Create index files for GRCm38-MT. For details, please refer to Rsubread
# user manual.
fasta <- system.file("extdata", "GRCm38_MT.fa", package = "scruff")
# NOTE: Rsubread package does not support Windows environment.
if (!requireNamespace("Rsubread", quietly = TRUE)) {
  message("Package \"Rsubread\" needed.",
    " Please install it if you are using Linux or macOS systems.",
    " The function is not available in Windows environment.\n")
} else {
  # Create index files for GRCm38-MT.
  # For details, please refer to Rsubread user manual.
  # Specify the basename for Rsubread index
  indexBase <- "GRCm38-MT"
  Rsubread::buildindex(basename = indexBase,
    reference = fasta,
    indexSplit = FALSE)
}
```

Process scRNA-Seq reads in *scruff*

The following command maps the FASTQ files to GRCm38 mitochondrial reference sequence `GRCm38_MT.fa` and returns a `SingleCellExperiment` object. By default, the files are stored in BAM format in `./Alignment` folder.

```
# Align the reads using Rsubread
if (requireNamespace("Rsubread", quietly = TRUE)) {
  al <- alignRsubread(de,
    indexBase,
    unique = FALSE,
    nBestLocations = 1,
    format = "BAM",
    overwrite = TRUE,
    verbose = TRUE,
    cores = 2)
}
```

3.4 UMI correction and Generation of Count Matrix

Example GTF file `GRCm38_MT.gtf` will be used for feature counting. Currently, *scruff* applies the union counting mode of the HTSeq Python package. The following command generates the UMI corrected count matrix for the example dataset by allowing correction of UMIs with Hamming distances smaller than 1 for each gene in each cell.

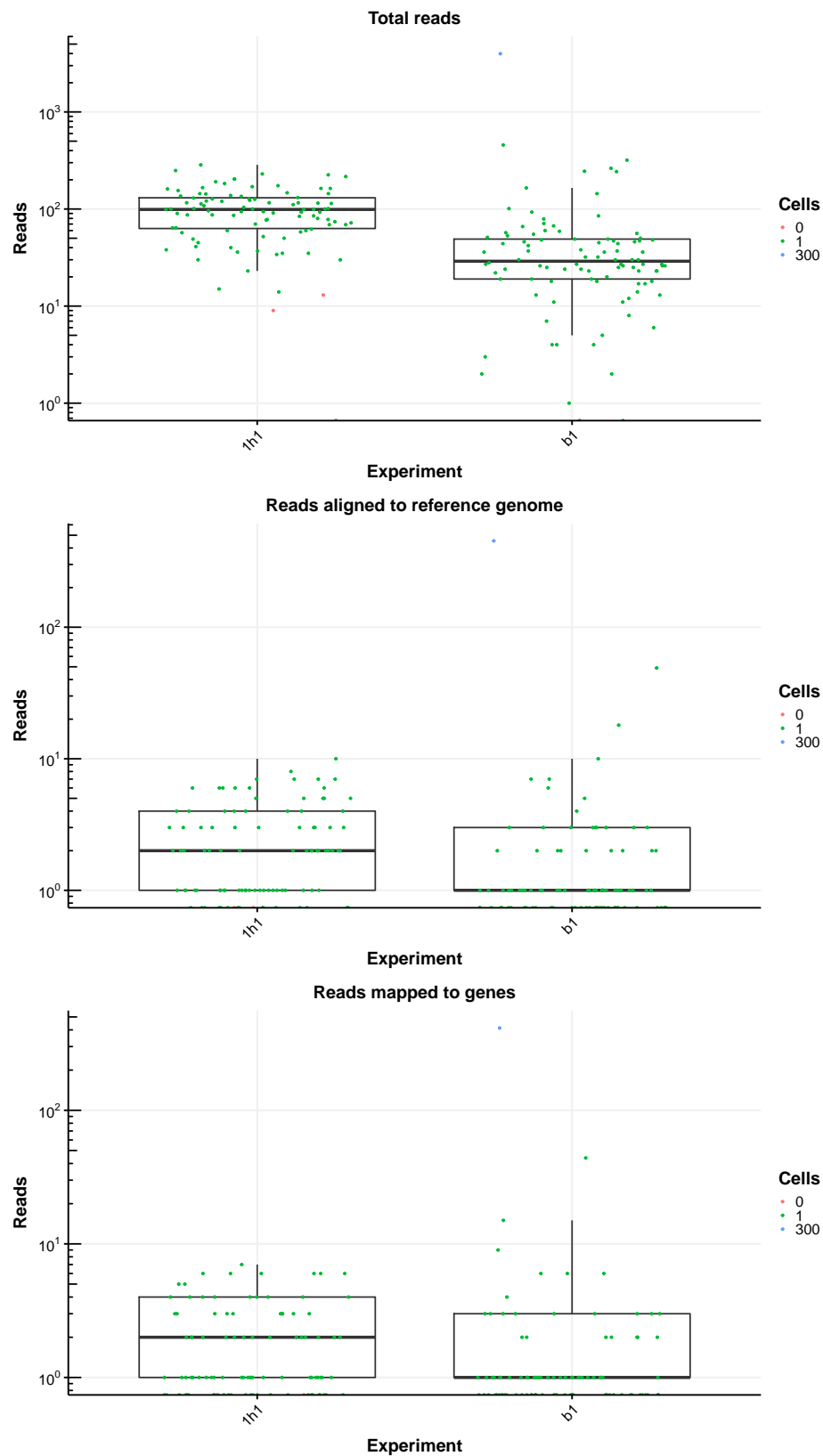
```
gtf <- system.file("extdata", "GRCm38_MT.gtf", package = "scruff")
# get the molecular counts of transcripts for each cell
# In experiment 1h1, cell barcodes 95 and 96 are empty well controls.
# In experiment b1, cell barcode 95 is bulk sample containing 300 cells.
if (requireNamespace("Rsubread", quietly = TRUE)) {
  sce = countUMI(al,
    gtf,
    umiEdit = 1,
    format = "BAM",
    cellPerWell = c(rep(1, 94), 0, 0, rep(1, 94), 300, 1),
    verbose = TRUE,
    cores = 2)
}
```

3.5 Visualization of QC metrics

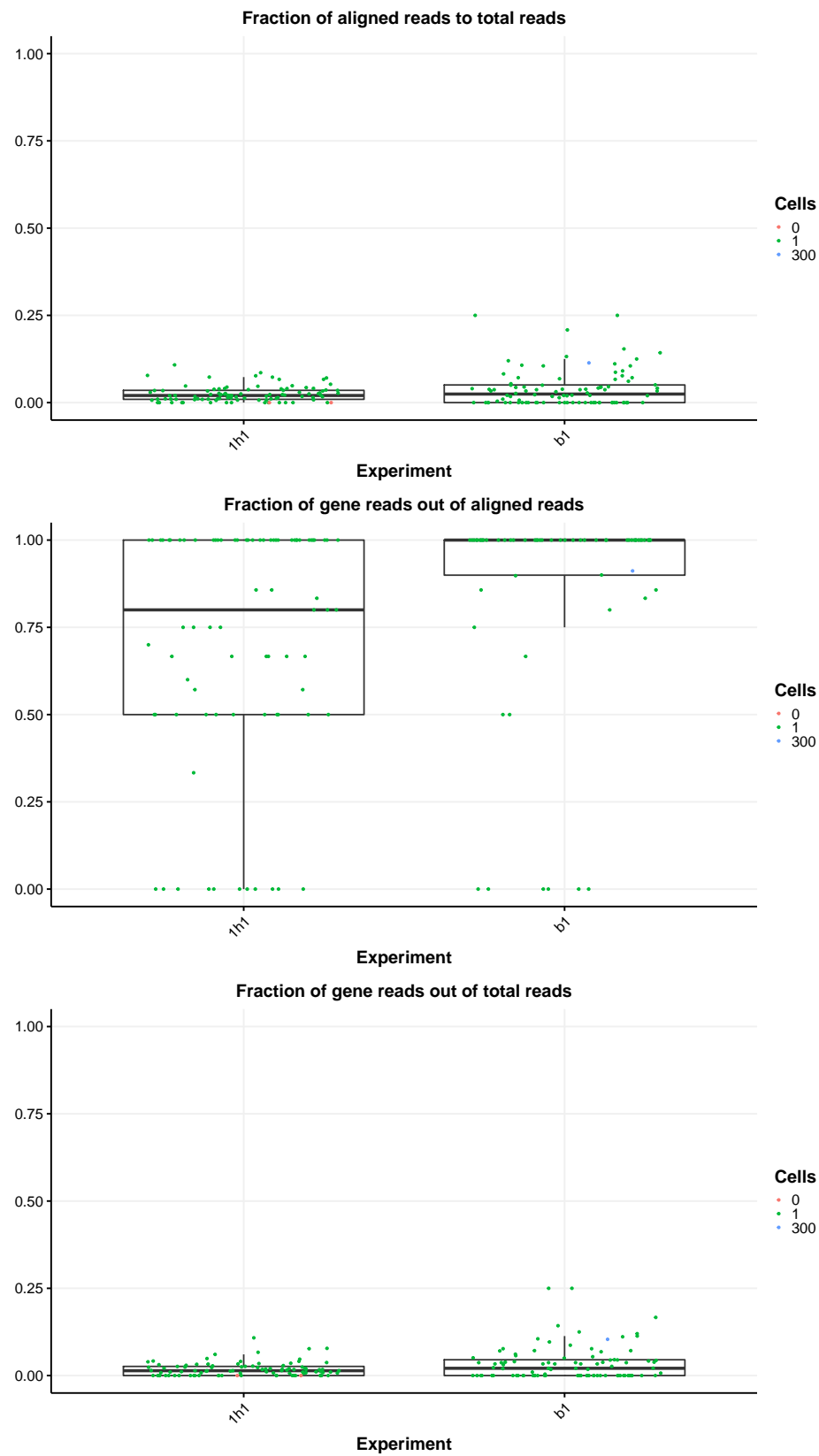
The data quality diagnostic information are contained in the `colData` of the returned `SingleCellExperiment` object `sce`. They can be visualized using the `qcplots` function.

```
data(sceExample, package = "scruff")
qc <- qcplots(sceExample)
qc
```

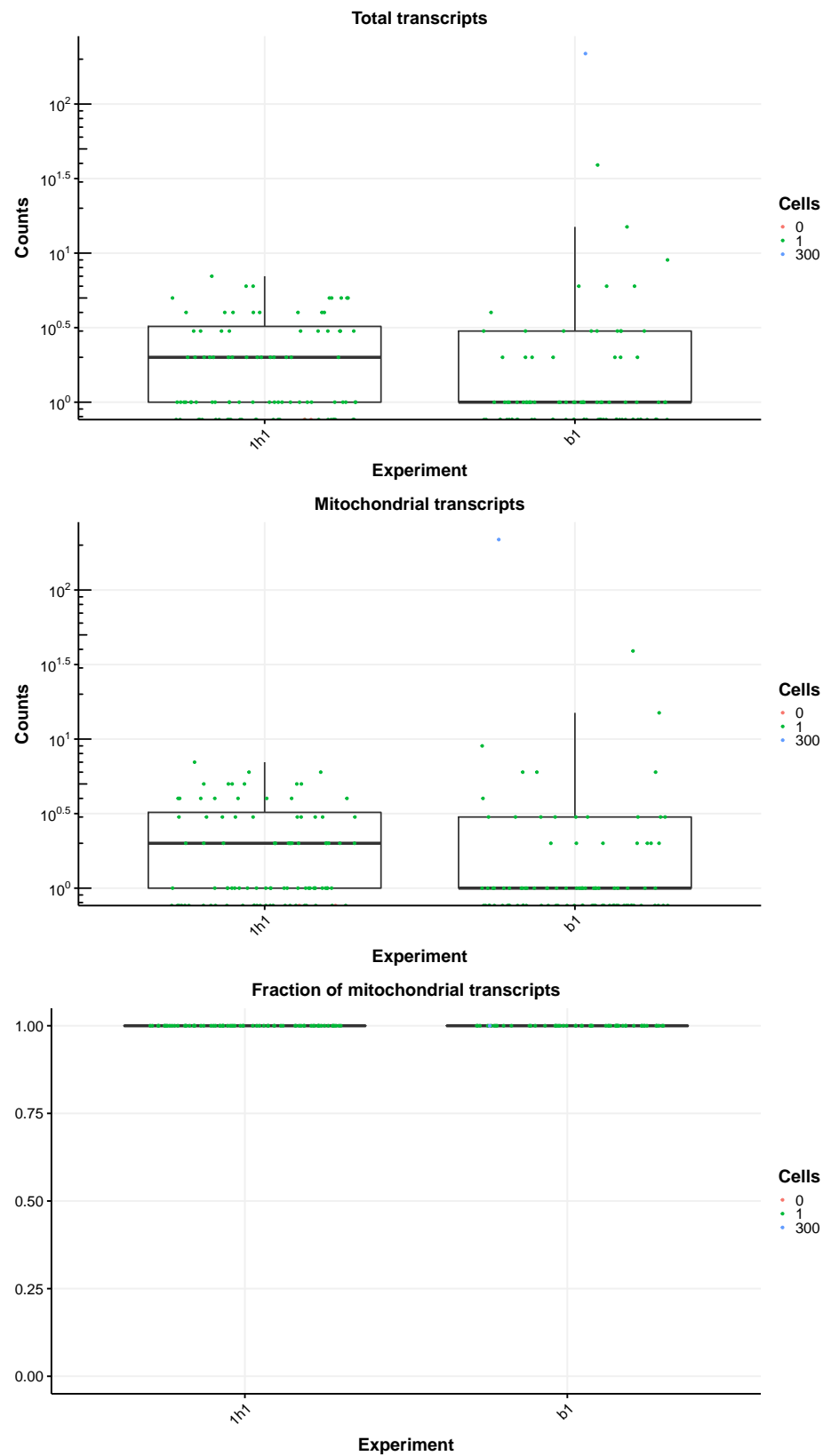
Process scRNA-Seq reads in *scruff*



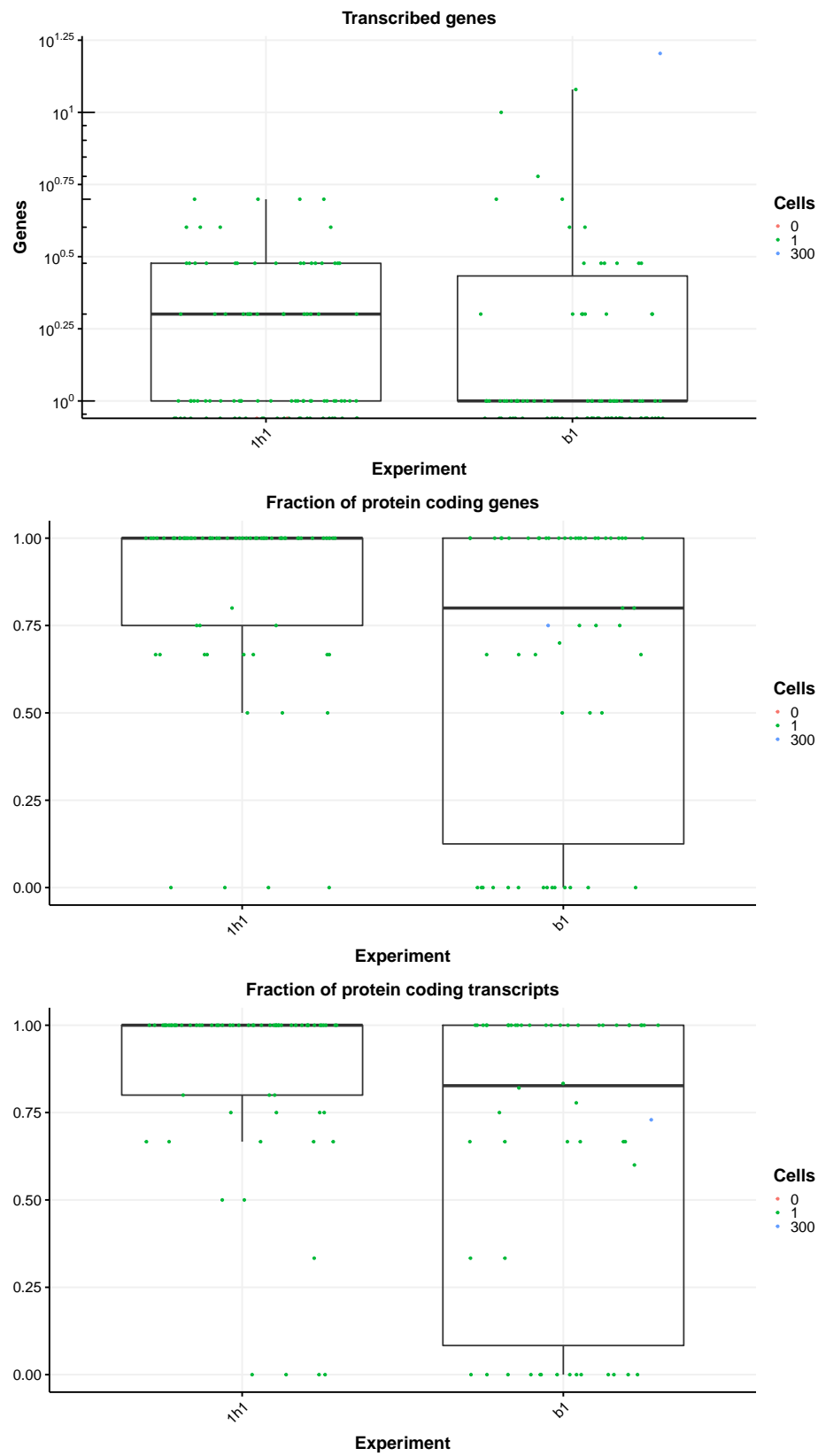
Process scRNA-Seq reads in *scruff*



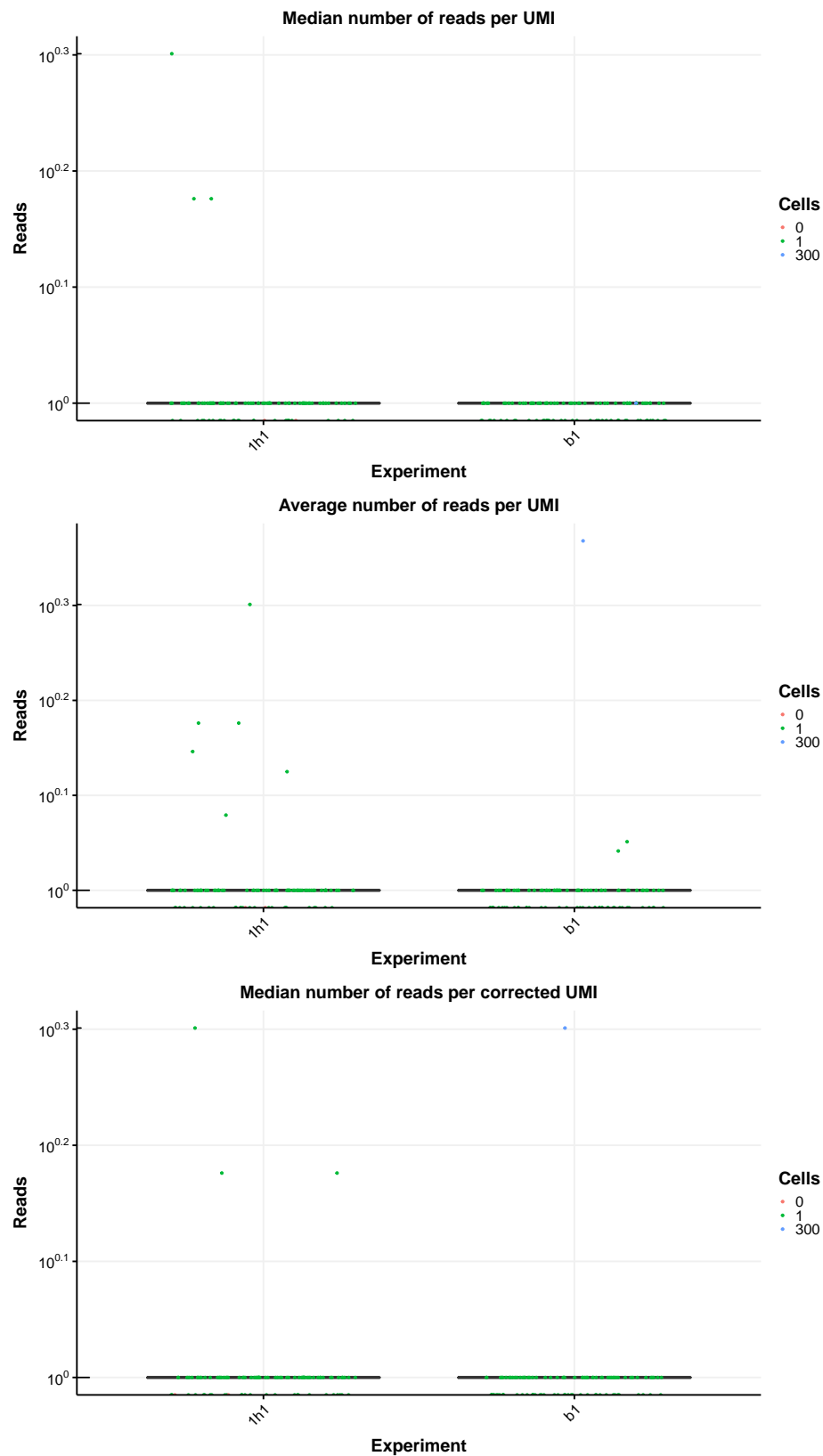
Process scRNA-Seq reads in *scruff*



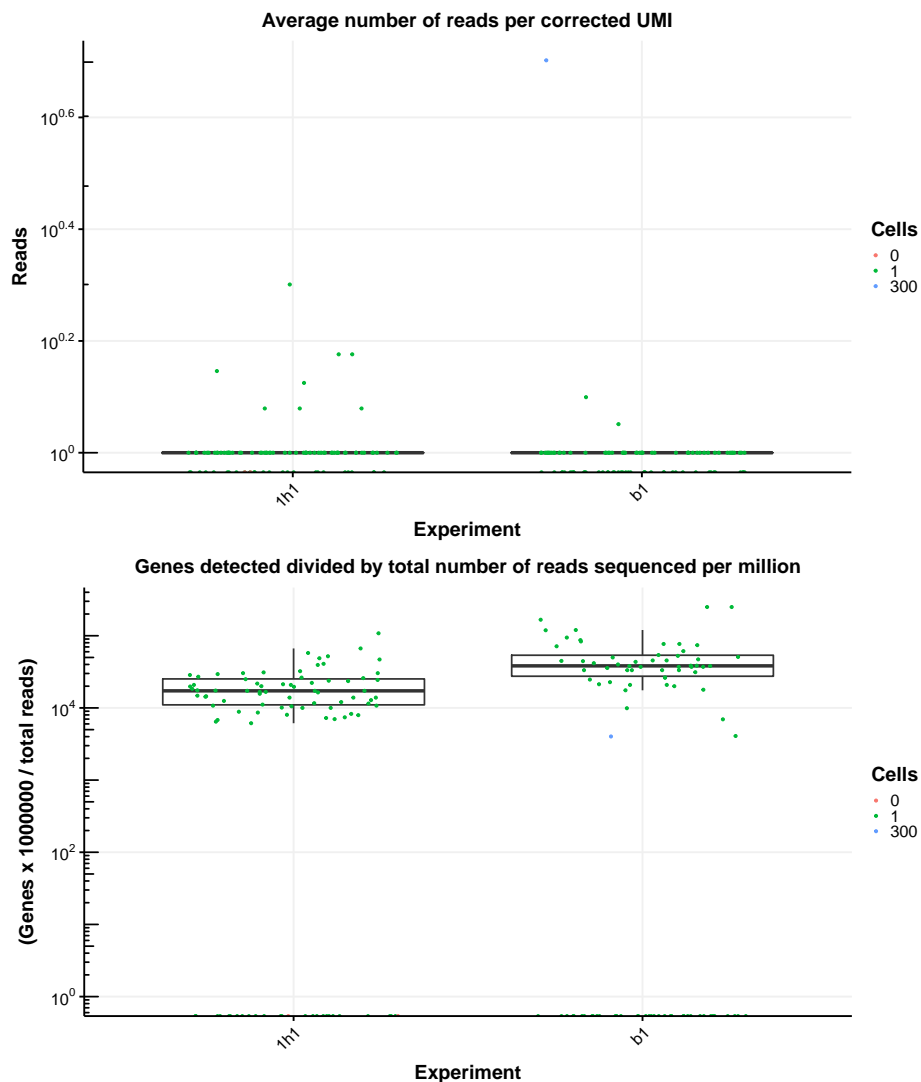
Process scRNA-Seq reads in *scruff*



Process scRNA-Seq reads in *scruff*



Process scRNA-Seq reads in *scruff*



Each of these boxplots shows the distribution of a quality metric within and across experiments. Y axis shows the number or fraction of these metrics. Each point represents a unique cell barcode associated with the cells or bulk samples in the wells of the plate. These points are colored by the `cellPerWell` parameter which is the number of sorted cells in each well. For example, in experiment 1h1, cell barcodes 95 and 96 are empty well controls. In experiment b1, cell barcode 95 is a bulk sample containing 300 cells.

These quality metrics include total reads, aligned reads, reads mapped to genes, fraction of aligned reads, fraction of gene reads out of aligned reads, fraction of gene reads out of total reads, total transcripts, mitochondrial transcripts, fraction of mitochondrial transcripts, transcribed genes, fraction of protein coding genes and transcripts, median and average number of reads per corrected or uncorrected UMI, and gene detection rate.

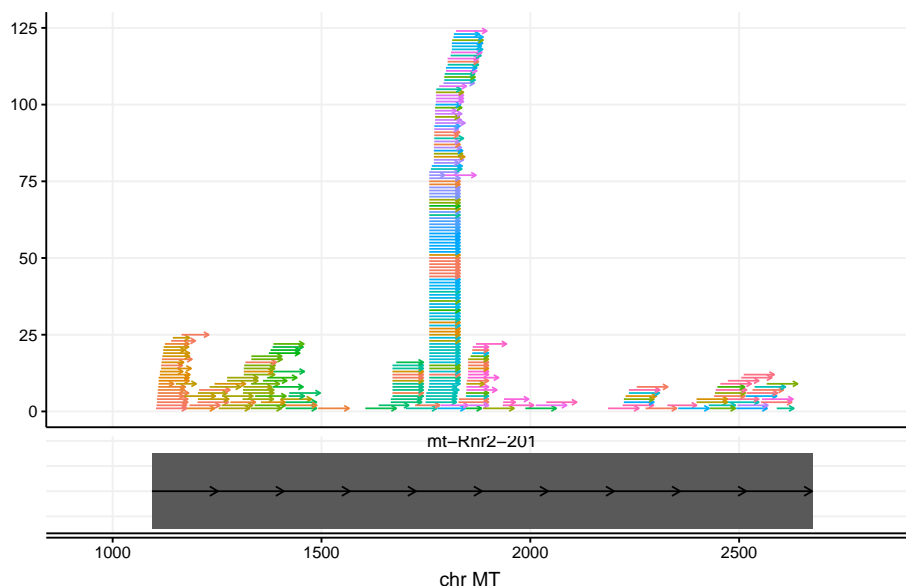
3.6 Visualization of Read mapping locations

scruff package provides function to visualize read alignment locations at specified genomic coordinates on the reference genome. The following command visualize the reads mapped to gene *mt-Rnr2* for the bulk sample *vandenBrink_b1_cell_0095*. Reads are colored by their uncorrected UMI tags. *bamExample* in the following example is a *GAlignments* object generated by the *readGAlignments* function in the *GenomicAlignments* package. Please refer to the documentation of *readGAlignments* and *GenomicAlignments* for details about how to read BAM files into R.

```
# Visualize the reads mapped to gene "mt-Rnr2" in
# cell "vandenBrink_b1_cell_0095".
# bamExample is generated by GenomicAlignments::GAlignments function
data(bamExample, package = "scruff")
```

```
# gene mt-Rnr2 starts at 1094 and ends at 2675
start <- 1094
end <- 2675

g1 <- rview(bamExample, chr = "MT", start = start, end = end)
g2 <- gview(gtf, chr = "MT", start = start, end = end)
g <- ggbio::tracks(g1, g2, heights = c(4,1), xlab = "chr MT")
g
```



This plot shows the read alignment information for gene *mt-RNA2*. The top panel shows the alignment and orientation of each read aligned to this gene. Y axis indicates the number of aligned reads at specific locations. Each arrow represents a read. The length of the arrow indicates the length of the read. The reads are colored by their UMI tags. In the bottom panel, each arrow represents an isoform of the gene. The isoforms are named by their transcript names. The grey boxes indicate the exonic regions of the transcript. X axis is the genomic coordinates on the corresponding chromosome (e.g. chr MT). It is consistent between the top and bottom panels if *legend* is set to FALSE for *rview* function.

4 10X BAM File Quality Assessment

The function `tenxBamqc` collects QC metrics from BAM files generated by 10X Cell Ranger pipeline. The collected QC metrics can be visualized by the function `qcplots`. The collected measurements include number of aligned reads and number of reads aligned to a gene for each valid cell barcode. Users can see the alignment quality for the filtered cells after plugging in the filtered barcode file from the Cell Ranger results. Here we show the visualization of an example BAM file from Cell Ranger output.

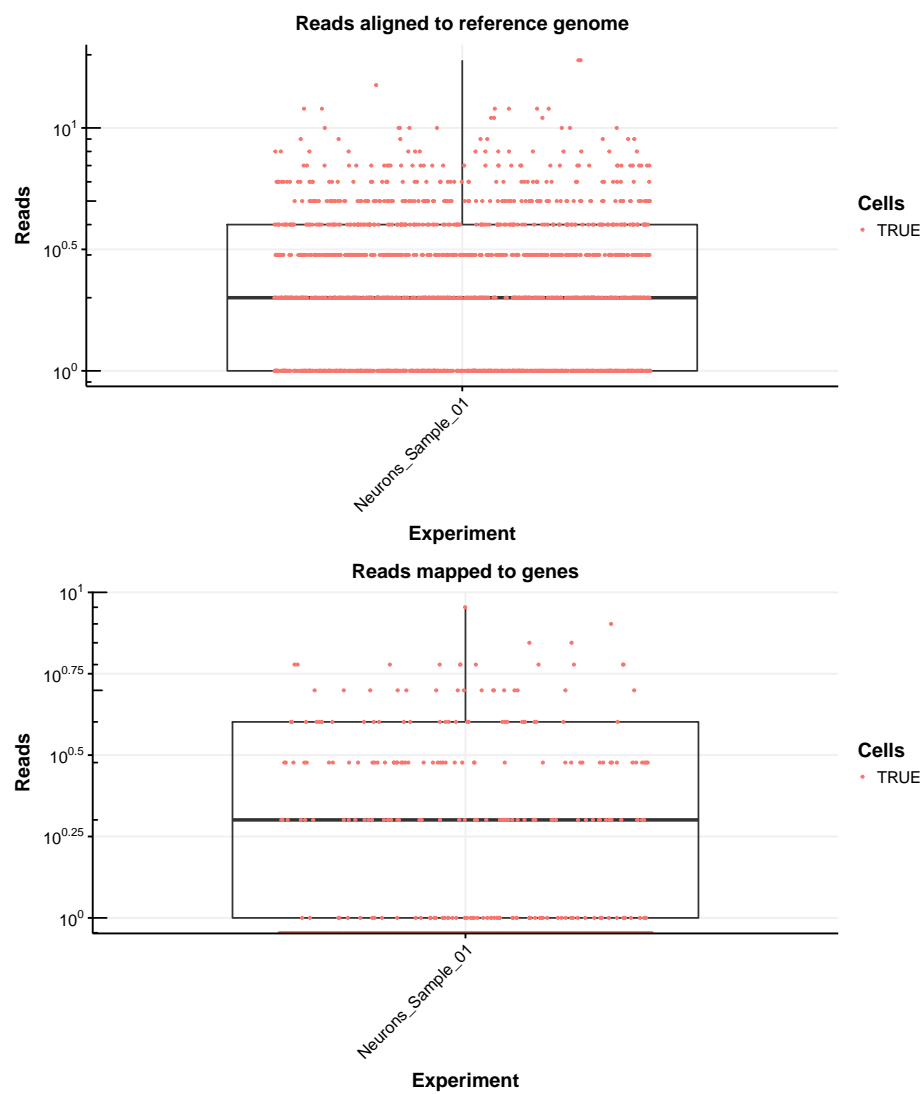
```
# The following example BAM file is the first 5000 BAM file records extracted
# from sample 01 of the 1.3 Million Brain Cells dataset from E18 Mice.
# (https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.3.0/
# 1M_neurons)
# The BAM file for sample 01 is downloaded from here:
# http://sra-download.ncbi.nlm.nih.gov/srapub\_files/
# SRR5167880_E18_20160930_Neurons_Sample_01.bam
# see details here:
# https://trace.ncbi.nlm.nih.gov/Traces/sra/?study=SRP096558
# and here:
# https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE93421
bamfile10x <- system.file("extdata",
  "SRR5167880_E18_20160930_Neurons_Sample_01_5000.bam",
  package = "scruff")

# The filtered cell barcodes are generated using the following script:
# library(TENxBrainData)
# library(data.table)
# tenx <- TENxBrainData()
# # get filtered barcodes for sample 01
# filteredBcIndex <- tstrsplit(colData(tenx)[, "Barcode"], "-")[[2]] == 1
# filteredBc <- colData(tenx)[filteredBcIndex, ]["Barcode"]
filteredBc <- system.file("extdata",
  "SRR5167880_E18_20160930_Neurons_Sample_01_filtered_barcode.tsv",
  package = "scruff")
# QC results are saved to current working directory
qcDt <- tenxBamqc(bam = bamfile10x,
  experiment = "Neurons_Sample_01",
  filter = filteredBc)
## tenxBamqc(bam = bamfile10x, experiment = "Neurons_Sample_01",
##   filter = filteredBc)
qcDt
## class: SingleCellExperiment
## dim: 1 1576
## metadata(0):
## assays(1): ''
## rownames: NULL
## rowData names(0):
## colnames(1576): AACCTGAGGCCCGTT-Neurons_Sample_01
##   AACCTGAGTCCGGTC-Neurons_Sample_01 ...
##   TTTGTCATCGGTGTTA-Neurons_Sample_01
##   TTTGTCATCTTGATC-Neurons_Sample_01
```

Process scRNA-Seq reads in *scruff*

```
## colData names(5): cell_barcode reads_mapped_to_genome  
##   reads_mapped_to_genes experiment number_of_cells  
## reducedDimNames(0):  
## spikeNames(0):  
## altExpNames(0):
```

```
g <- qcplots(qcDt)  
g
```



Process scRNA-Seq reads in *scruff*



5 Session Information

```
## R version 3.6.1 (2019-07-05)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 18.04.3 LTS
##
## Matrix products: default
## BLAS: /home/biocbuild/bbs-3.10-bioc/R/lib/libRblas.so
## LAPACK: /home/biocbuild/bbs-3.10-bioc/R/lib/libRlapack.so
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_US.UTF-8       LC_COLLATE=C
##  [5] LC_MONETARY=en_US.UTF-8   LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=en_US.UTF-8      LC_NAME=C
##  [9] LC_ADDRESS=C              LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] stats4      parallel  stats      graphics  grDevices  utils      datasets
## [8] methods    base
##
## other attached packages:
##  [1] GenomicFeatures_1.38.0 AnnotationDbi_1.48.0 Biobase_2.46.0
##  [4] GenomicRanges_1.38.0  GenomeInfoDb_1.22.0 IRanges_2.20.0
##  [7] S4Vectors_0.24.0      BiocGenerics_0.32.0  scruff_1.4.0
## [10] BiocStyle_2.14.0
##
## loaded via a namespace (and not attached):
##  [1] ProtGenerics_1.18.0      bitops_1.0-6
##  [3] matrixStats_0.55.0      bit64_0.9-7
```


Process scRNA-Seq reads in *scruff*

```
## [5] RColorBrewer_1.1-2      progress_1.2.2
## [7] httr_1.4.1              tools_3.6.1
## [9] backports_1.1.5         R6_2.4.0
## [11] rpart_4.1-15            Hmisc_4.2-0
## [13] DBI_1.0.0               lazyeval_0.2.2
## [15] colorspace_1.4-1       nnet_7.3-12
## [17] withr_2.1.2             tidyselect_0.2.5
## [19] gridExtra_2.3          prettyunits_1.0.2
## [21] GGally_1.4.0           bit_1.1-14
## [23] curl_4.2               compiler_3.6.1
## [25] graph_1.64.0           htmlTable_1.13.2
## [27] DelayedArray_0.12.0    labeling_0.3
## [29] rtracklayer_1.46.0     bookdown_0.14
## [31] ggbio_1.34.0           checkmate_1.9.4
## [33] scales_1.0.0           RBGL_1.62.0
## [35] askpass_1.1            rappdirs_0.3.1
## [37] stringr_1.4.0          digest_0.6.22
## [39] Rsamtools_2.2.0        foreign_0.8-72
## [41] rmarkdown_1.16         XVector_0.26.0
## [43] Rsubread_2.0.0         dichromat_2.0-0
## [45] base64enc_0.1-3       pkgconfig_2.0.3
## [47] htmltools_0.4.0       ensemblDb_2.10.0
## [49] BSgenome_1.54.0       dbplyr_1.4.2
## [51] ggthemes_4.2.0         htmlwidgets_1.5.1
## [53] rlang_0.4.1           rstudioapi_0.10
## [55] RSQLite_2.1.2         hwriter_1.3.2
## [57] BiocParallel_1.20.0   acepack_1.4.1
## [59] dplyr_0.8.3           VariantAnnotation_1.32.0
## [61] RCurl_1.95-4.12       magrittr_1.5
## [63] GenomeInfoDbData_1.2.2 Formula_1.2-3
## [65] Matrix_1.2-17         Rcpp_1.0.2
## [67] munsell_0.5.0         stringi_1.4.3
## [69] yaml_2.2.0            SummarizedExperiment_1.16.0
## [71] zlibbioc_1.32.0       plyr_1.8.4
## [73] BiocFileCache_1.10.0  grid_3.6.1
## [75] blob_1.2.0            crayon_1.3.4
## [77] lattice_0.20-38       Biostrings_2.54.0
## [79] splines_3.6.1         hms_0.5.1
## [81] zeallot_0.1.0        knitr_1.25
## [83] pillar_1.4.2          reshape2_1.4.3
## [85] biomaRt_2.42.0        XML_3.98-1.20
## [87] glue_1.3.1           evaluate_0.14
## [89] ShortRead_1.44.0      biovizBase_1.34.0
## [91] latticeExtra_0.6-28   data.table_1.12.6
## [93] BiocManager_1.30.9    vctrs_0.2.0
## [95] gtable_0.3.0          openssl_1.4.1
## [97] purrr_0.3.3           reshape_0.8.8
## [99] assertthat_0.2.1     ggplot2_3.2.1
## [101] xfun_0.10            AnnotationFilter_1.10.0
## [103] survival_2.44-1.1    SingleCellExperiment_1.8.0
## [105] OrganismDbi_1.28.0   tibble_2.1.3
```

Process scRNA-Seq reads in *scruff*

```
## [107] GenomicAlignments_1.22.0    memoise_1.1.0  
## [109] cluster_2.1.0
```