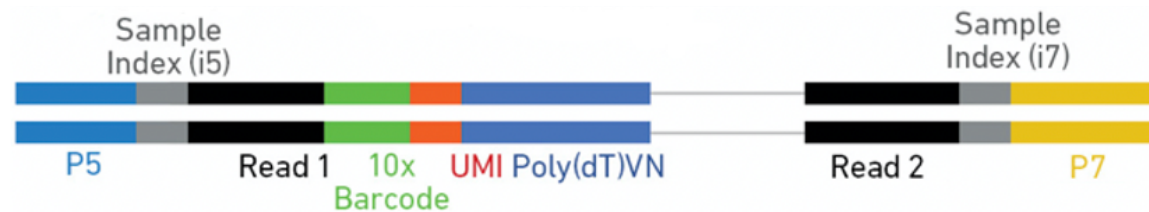# Alignment and feature counting

September 2022

# Single Cell RNAseq Analysis Workflow

# 10x library file structure

The 10x library contains four pieces of information, in the form of DNA sequences, for each "read".
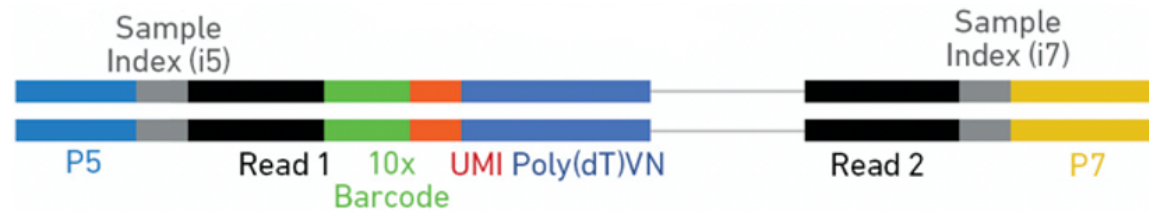
- **sample index** - identifies the library, with one or two indexes per sample
- **10x barcode** - identifies the droplet in the library
- **UMI** - identifies the transcript molecule within a cell and gene
- **insert** - the transcript molecule

# Raw fastq files

The sequences for any given fragment will generally be delivered in 3 or 4 files:

- **I1**: I7 sample index
- **I2**: I5 sample index if present (dual indexing only)
- **R1**: 10x barcode + UMI
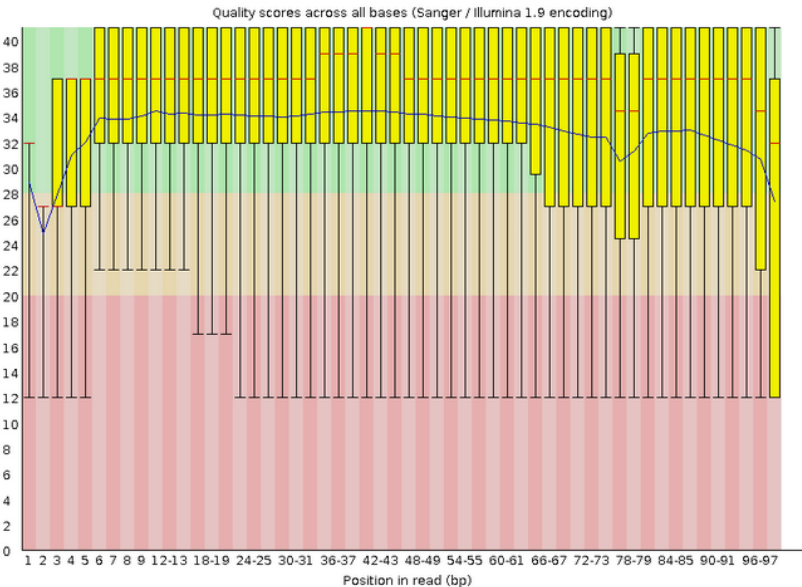- **R2**: insert sequence

# QC of Raw Reads - FASTQC

# QC of Raw Reads - MultiQC - General Statistics

# QC of Raw Reads - MultiQC - Sequence Quality Histograms

# Alignment and counting

The first steps in the analysis of single cell RNAseq data:

- Align reads to genome
- Annotate reads with feature (gene)
- Quantify gene expression

# Cell Ranger

- 10x Cell Ranger - This not only carries out the alignment and feature counting, but will also:
  - Call cells
  - Generate a summary report in html format
  - Generate a "cloupe" file

Alternative methods include:

- STAR solo:
  - Generates outputs very similar to CellRanger minus the cloupe file and the QC report
  - Will run with lower memory requirements in a shorter time than Cell Ranger
- Alevin:
  - Based on the popular Salmon tool for bulk RNAseq feature counting
  - Alevin supports both 10x-Chromium and Drop-seq derived data

# Obtaining Cell Ranger

# Cell Ranger tools

Cell Ranger includes a number of different tools for analysing scRNAseq data, including:

- `cellranger mkref` - for making custom references
- `cellranger count` - for aligning reads and generating a count matrix
- `cellranger aggr` - for combining multiple samples and normalising the counts

# Preparing the raw fastq files

Cell Ranger requires the fastq file names to follow a convention:

`<SampleName>_S<SampleNumber>_L00<Lane>_<Read>_001.fastq.gz`

e.g. for a single sample in the Caron data set we have:

```
SRR9264343_S0_L001_I1_001.fastq.gz
SRR9264343_S0_L001_R1_001.fastq.gz
SRR9264343_S0_L001_R2_001.fastq.gz
```

# Genome/Transcriptome Reference

As with other aligners Cell Ranger requires the information about the genome and transcriptome of interest to be provided in a specific format.

- Obtain from the 10x website for human or mouse (or both - PDX)
- Build a custom reference with `cellranger mkref`

# Running `cellranger count`

- Computationally very intensive
- High memory requirements

```
File  Edit  View  Search  Terminal  Help
%h%-$
%h%-$
%h%-$ cellranger count --id=SRR9264343 \
              --transcriptome=refdata-gex-mm10-2020-A \
              --fastqs=fastq \
              --sample=SRR9264343 \
              --localcores=8 \
              --localmem=64
```

# Cell Ranger outputs

- One directory per sample



File Edit View Search Terminal Help
```
%h%-$ ..
%h%-$ ls SRR9264343/
_cmdline
_filelist
_finalstate
_invocation
_jobmode
_log
_mrosource
outs
_perf
SC_RNA_COUNTER_CS
_sitecheck
SRR9264343.mri.tgz
_tags
_timestamp
_uuid
_vdrkill
_versions
%h%-$
```

# Cell Ranger outputs

# Cell Ranger outputs

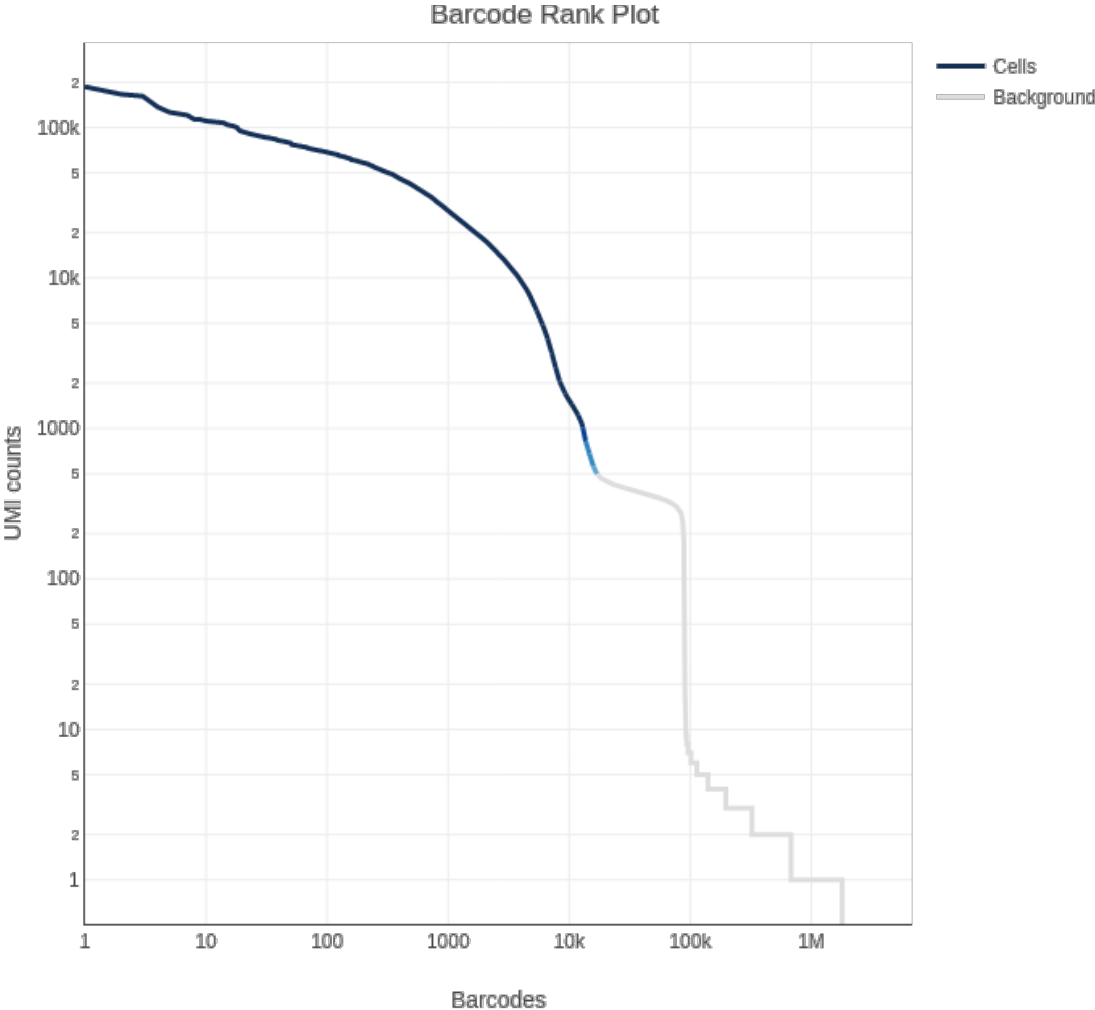# Cell Ranger report

# Cell Ranger outputs

```
File  Edit  View  Search  Terminal  Help
_versions
%h%-$
%h%-$ ls SRR9264343/outs/
analysis
cloupe.cloupe
filtered_feature_bc_matrix
filtered_feature_bc_matrix.h5
metrics_summary.csv
molecule_info.h5
possorted_genome_bam.bam
possorted_genome_bam.bam.bai
raw_feature_bc_matrix
raw_feature_bc_matrix.h5
web_summary.html
%h%-$
```

# Loupe Browser

# Cell Ranger outputs

# Cell Ranger outputs

# Cell Ranger outputs

# Cell Ranger outputs

# Cell Ranger cell calling

# Single Cell RNAseq Analysis Workflow