# KB & Seqspec for share-seq

scRNA-SEQ

scRNA-SEQ (TBD – with the newer version?)

# !kb ref - differences

- Adding multiple references files Vs. cDNA only
  - !kb ref -i ref_cDNA/transcriptome.idx -g ref_cDNA/transcripts_to_genes.txt \ -f1 ref_cDNA/Homo_sapiens.GRCh38.cdna.all.fa.gz \ reference/Homo_sapiens.GRCh38.dna.primary_assembly.fa.gz \ reference/Homo_sapiens.GRCh38.109.gtf.gz
- CDNA Reference file is generated from DNA and GTF Vs. Bulk - cDNA only (Downloaded from ENSEMBL)
  - One time effort

# Fastq files

- **ATAC:**
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC.barcodes_annotated.end1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC.barcodes_annotated.end2.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC_L001_R1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC_L001_R2.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC_L002_R1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_ATAC_L002_R2.fastq.gz
- **RNA:**
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA.barcodes_annotated.UMI.end1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA.barcodes_annotated.UMI.end2.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA_L001_R1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA_L001_R2.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA_L002_R1.fastq.gz
- /oak/stanford/groups/akundaje/marinovg/ENCODE4/single-cell/2023-06-02-scJamboree/BMMC_single_donor_RNA_L002_R2.fastq.gz

# Issues with YAML

- Can not have tabs (->)

- File structure is important. Assays in the root (RNA, ATAC)

- YAML debugging: https://www.yamllint.com/

- !seqspec print shows the YAML structure

# High level comparison

**Broad**

```
#·Assay·region
!Assay
name:·SHARE-Seq
doi_url:·https://doi.org/10.1016/j.cell.2020.09.056
publication_date:·23·October·2020
description:·Simultaneous·high-throughput·ATAC·and·RNA·expression·in·the·sam
lib_struct:·https://teichlab.github.io/scg_lib_structs/methods_html/SHARE-se

modalities:
-·ATAC
-·RNA

assay_spec:
····#·ATAC
····#·Read·1·Fastq
····-·!Region
·······region_id:·ATAC-R1.fastq.gz
······region_type:·fastq
······name:·ATAC·Read·1·FASTQ
······sequence_type:·joined
······sequence:·x
······min_len:·50
······max_len:·50
······onlist:·null
······regions:
·········-·!Region
············region_id:·ATAC-read1
············region_type:·gDNA
············name:·Genomic·DNA·read·1
············sequence_type:·random
············sequence:·X
············min_len:·50
············max_len:·50
············onlist:·null
············parent_id:·ATAC-R1.fastq.gz
····#·Read·2·Fastq
```
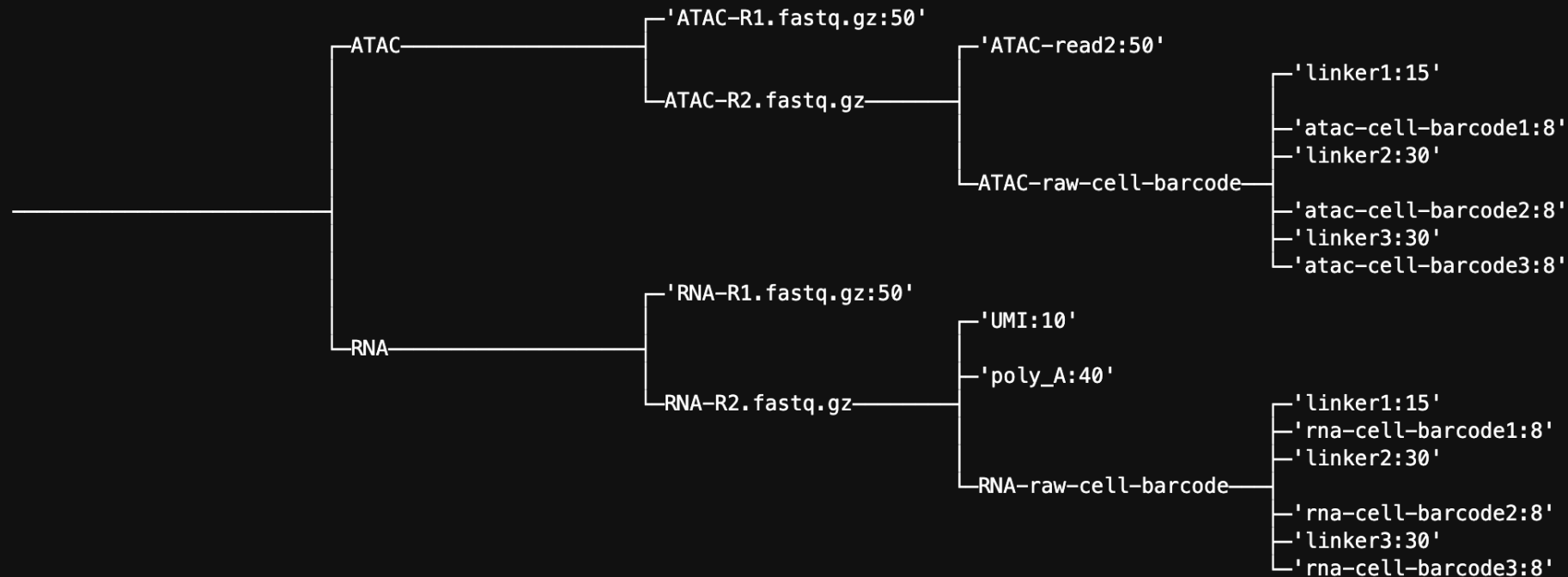
**Pachter**

```
!Assay
seqspec_version:·0.0.0
assay:·null
sequencer:·null
name:·SHARE-seq
doi:·https://doi.org/10.1016/j.cell.2020.09.056
publication_date:·23·October·2020
description:·The·SHARE-seq·method·is·developed·based·on·the·idea·of·combinatorial
··indexing·stratgy·that·is·used·in·sci-RNA-seq·and·SPLiT-seq
modalities:
-·RNA
-·ATAC
lib_struct:·https://teichlab.github.io/scg_lib_structs/methods_html/SHARE-seq.html
assay_spec:
-·!Region
··region_id:·RNA
··region_type:·RNA
··name:·RNA
··sequence_type:·joined
··sequence:·
AATGATACGGCGACCACCGAGATCTACACNNNNNNNNNNTCGTCGGCAGCGTCAGATGTGTATAAGAGACAGXXNNNNNNNNNNNNCTGTCTC
AAGTATGCAGCGCGCTCAAGCACGTGGATNNNNNNNNNAGTCGTACGCCGATGCGAAACATCGGCCACNNNNNNNNNATCTCGTATGCCGT
··min_len:·239
··max_len:·366
··onlist:·null
··regions:
··-·!Region
····region_id:·illumina_p5
····region_type:·illumina_p5
····name:·illumina_p5
····sequence_type:·fixed
····sequence:·AATGATACGGCGACCACCGAGATCTACAC
····min_len:·29
····max_len:·29
```

# After fix (collaborative effort with the Broad)

**>Seqspec print**

# >Kb --list – supported by name

```
[7]:  !kb --list                                          ▸ list                    Aᵃ ab .*  ⊤   5/5         ∧

      List of supported single-cell technologies

      Positions syntax: `input file index, start position, end position`
      When start & end positions are None, refers to the entire file
      Custom technologies may be defined by providing a kallisto-supported technology string
      (see https://pachterlab.github.io/kallisto/manual)
```

| name | description | whitelist | barcode | umi | cDNA |
|------|-------------|-----------|---------|-----|------|
| 10XV1 | 10x version 1 | yes | 0,0,14 | 1,0,10 | 2,None,None |
| 10XV2 | 10x version 2 | yes | 0,0,16 | 0,16,26 | 1,None,None |
| 10XV3 | 10x version 3 | yes | 0,0,16 | 0,16,28 | 1,None,None |
| 10XV3_ULTIMA | 10x version 3 sequenced with Ultima | yes | 0,22,38 | 0,38,50 | 0,62,None |
| BDWTA | BD Rhapsody | yes | 0,0,9 0,21,30 0,43,52 | 0,52,60 | 1,None,None |
| BULK | Bulk (single or paired) | | | | 0,None,None 1,None,None |
| CELSEQ | CEL-Seq | | 0,0,8 | 0,8,12 | 1,None,None |
| CELSEQ2 | CEL-SEQ version 2 | | 0,6,12 | 0,0,6 | 1,None,None |
| DROPSEQ | DropSeq | | 0,0,12 | 0,12,20 | 1,None,None |
| INDROPSV1 | inDrops version 1 | | 0,0,11 0,30,38 | 0,42,48 | 1,None,None |
| INDROPSV2 | inDrops version 2 | | 1,0,11 1,30,38 | 1,42,48 | 0,None,None |
| INDROPSV3 | inDrops version 3 | yes | 0,0,8 1,0,8 | 1,8,14 | 2,None,None |
| SCRUBSEQ | SCRB-Seq | | 0,0,6 | 0,6,16 | 1,None,None |
| SMARTSEQ2 | Smart-seq2  (single or paired) | | | | 0,None,None 1,None,None |
| SMARTSEQ3 | Smart-seq3 | | | 0,11,19 | 0,11,None 1,None,None |
| SPLIT-SEQ | SPLiT-seq | | 1,10,18 1,48,56 1,78,86 | 1,0,10 | 0,None,None |
| SURECELL | SureCell for ddSEQ | | 0,0,6 0,21,27 0,42,48 | 0,51,59 | 1,None,None |
| Visium | 10x Visium | yes | 0,0,16 | 0,16,28 | 1,None,None |

Does not include Share-seq. But it is supported. You will get the relevant string for share-seq

# !kb count

- Require whitelist that matches the barcode lengths on the YAML file
- Few minutes execution (without limiting threads on server machine)
  - 3+ min for one lane (Kali, I didn't set up the threads)
  - 6+ min on 2 lanes (Kali, I didn't set up the threads)

# One lane execution and output

# Two lanes execution and output