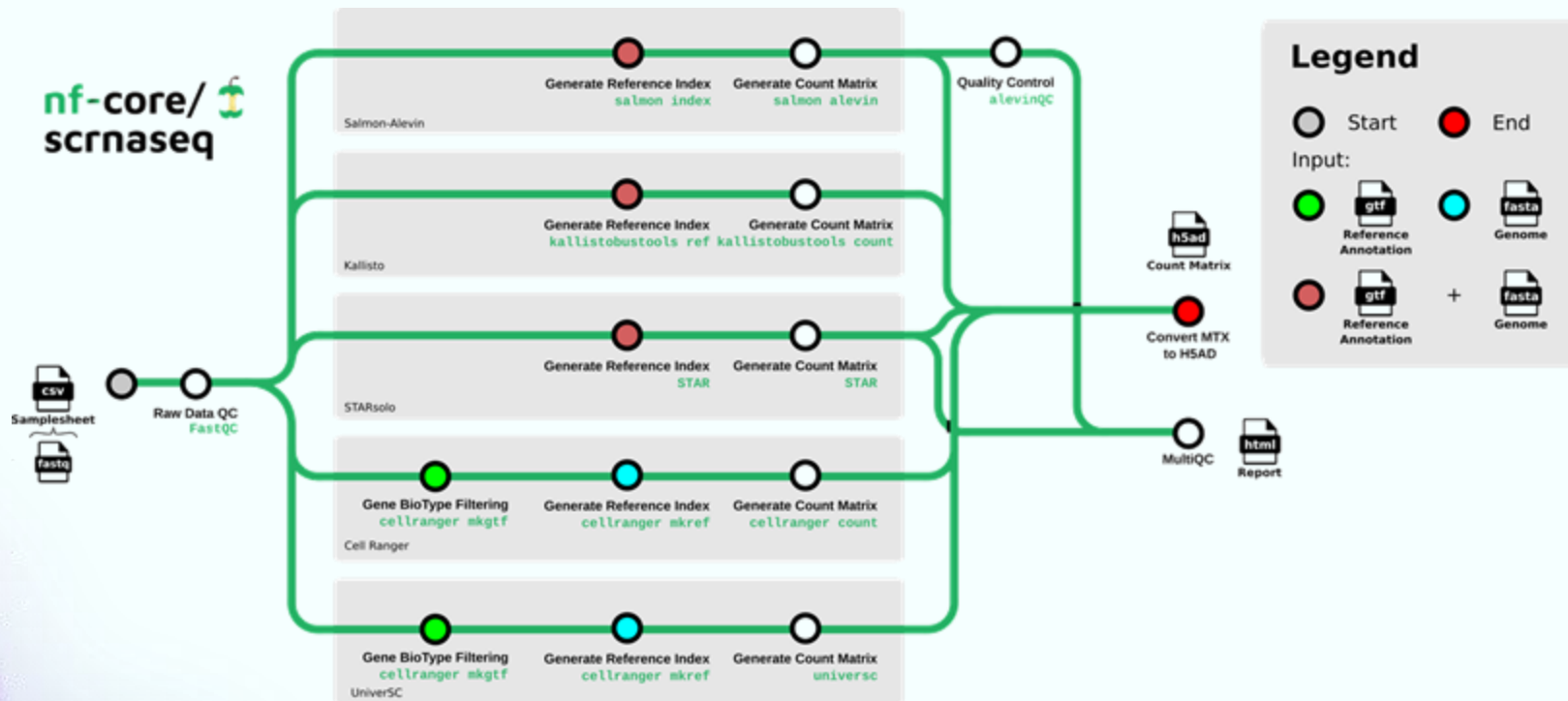# ScRNAseq Pipeline overview
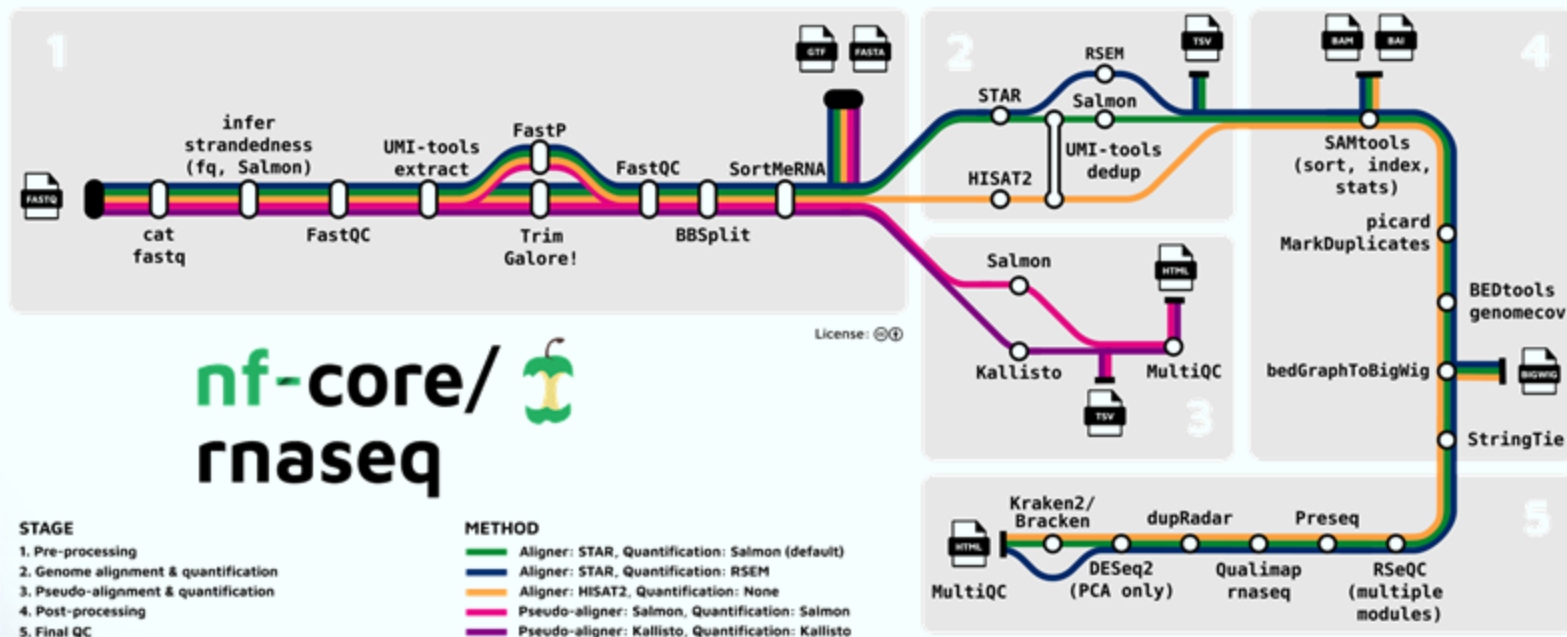
## ScRNAseq in the Cloud

MDIBL Comparative Genomics and Data Science Core
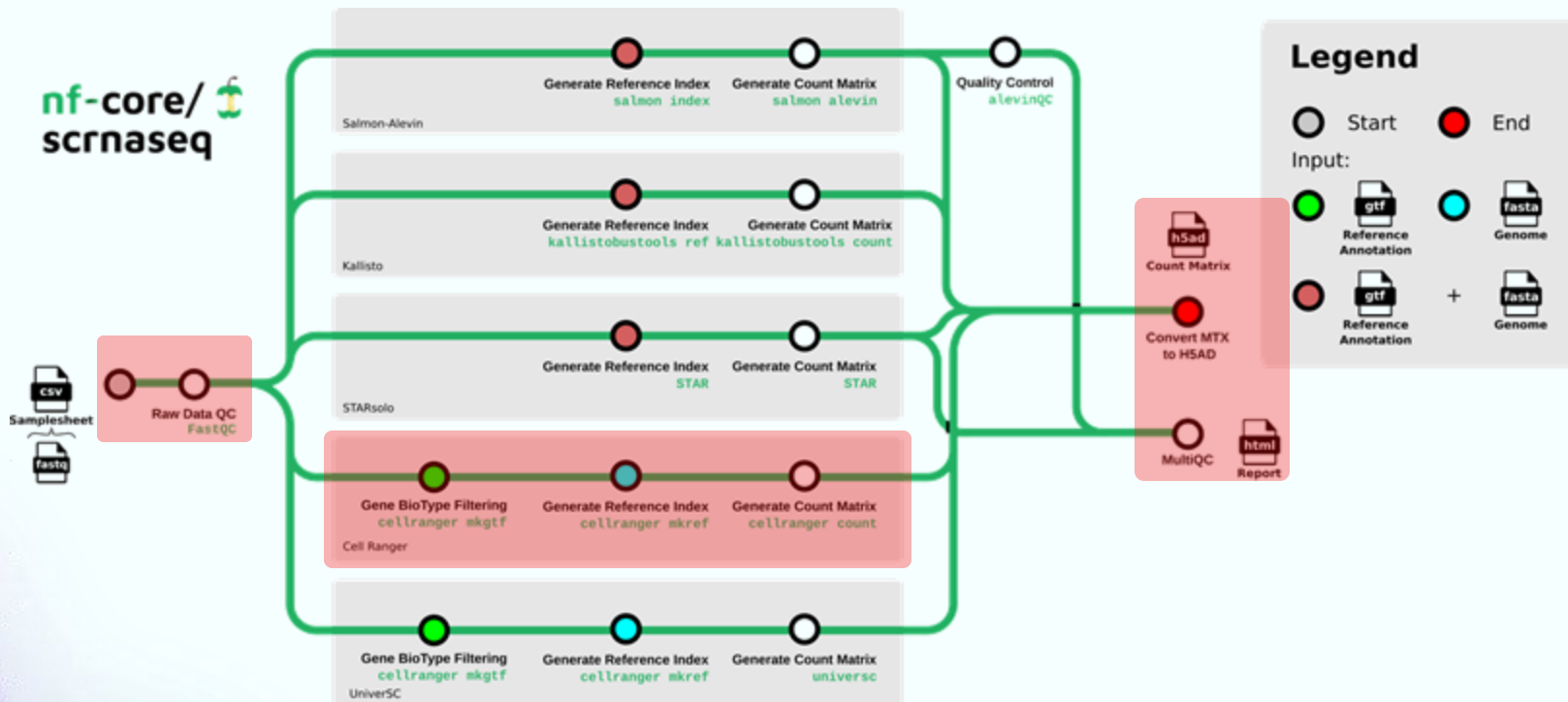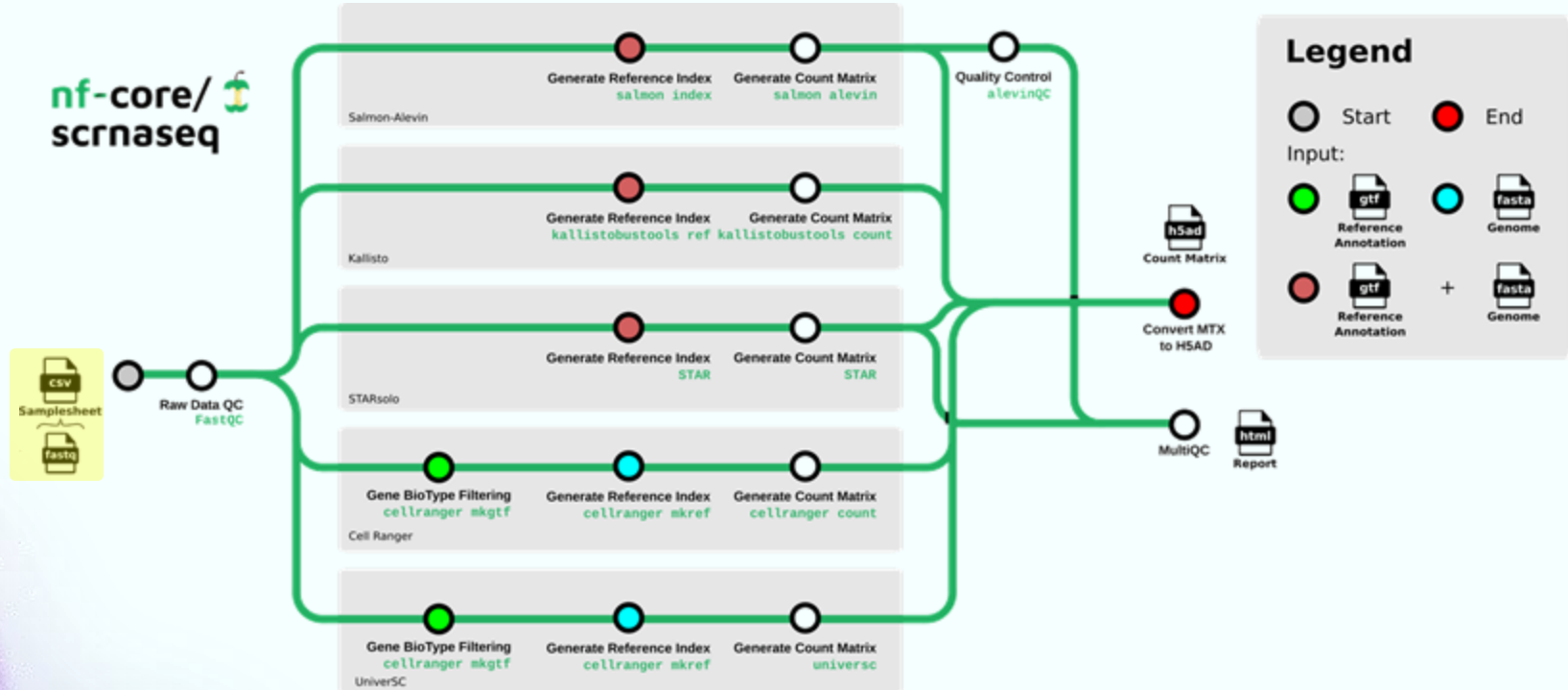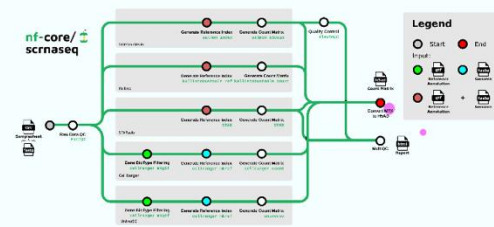
# Metro Map

# Metro Map (rnaseq)

# Metro Map

# Inputs

# Inputs



## Samples

samplesheet.csv
sampleID,fastq_1,fastq_2

- This defines the samples that will be processed
- Each entry needs a Read 1 and Read 2
- Samples need to be g-zipped.

## References

<organism>.fasta.gz
- Text-based file representing nucleotide (or protein) sequences. In our case organized by chromosome.

<organism>.gtf.gz
- Gene Transfer Format file describes gene structure information specifically location of genes.
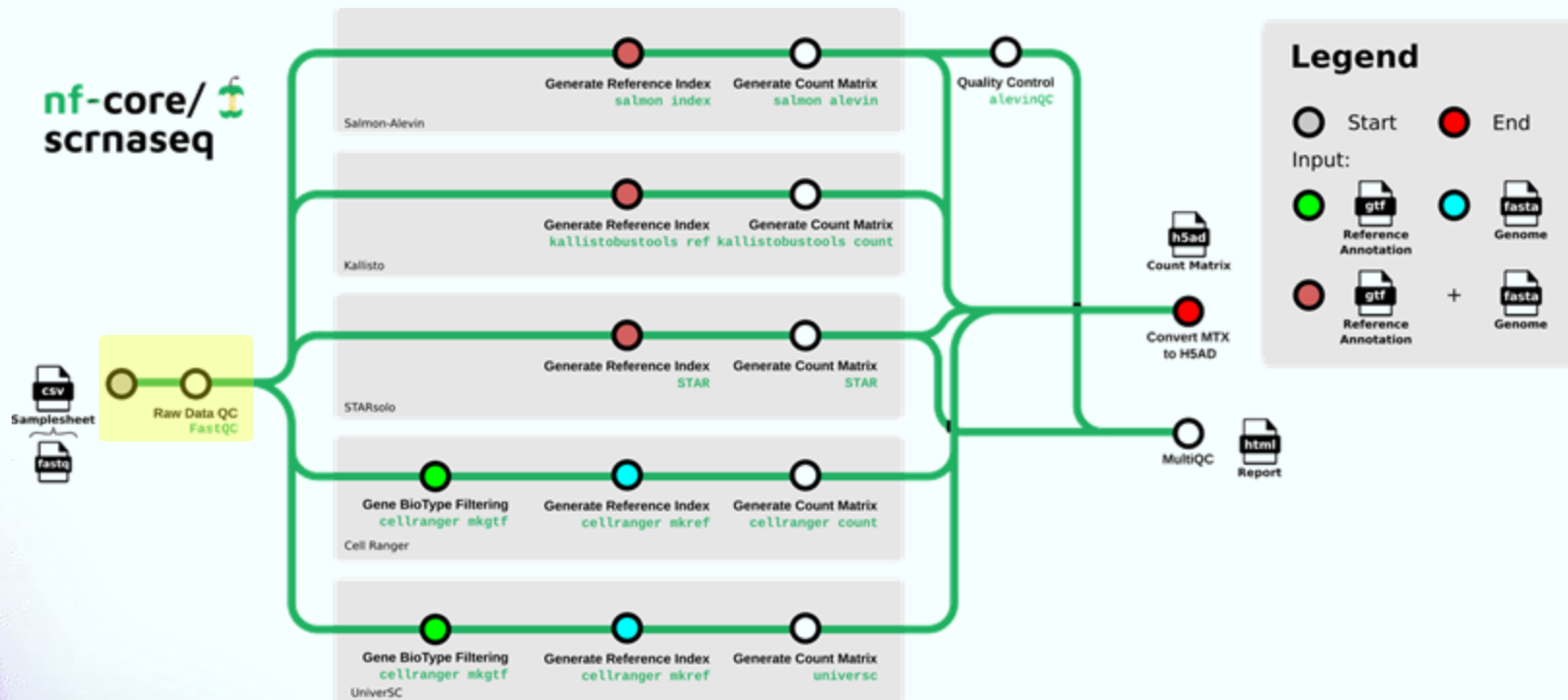
Together, they will be used to make our index.

## Options

- ALIGNER:
  - We will be choosing **cellranger** – specifically designed for 10x generated data.
  - **Cellranger** uses **starSOLO** as the alignment algorithm but make setup simple and easy.
  - Produces QC reports per sample.

nf-core

MDI Biological Laboratory

# FASTQC

# FASTQC

## What

- Tool for assessing the quality of raw sequencing.
- Commonly used for high-throughput sequencing such as ScRNAseq.

## Output

- HTML file
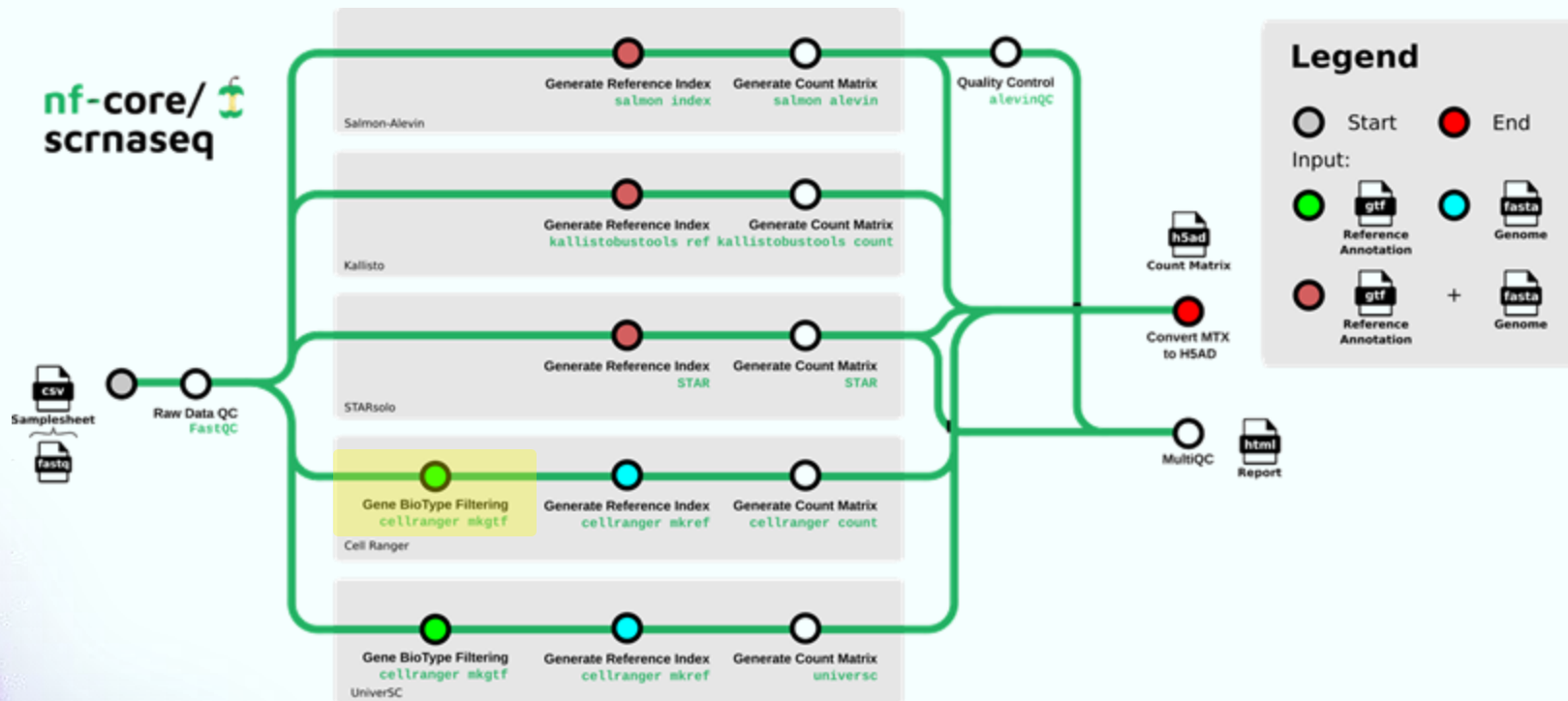- Each sample has 2 files:
    - Read 1
    - Read 2

## Readouts

- Basic Statistics
- Per base sequence quality
- Per sequence quality score
- Per base sequence content
- Per sequence GC content
- Per base N content
- Sequence Length Distribution
- Sequence Duplication Levels
- Overrepresented sequences
- Adapter Content

⚠️ The .fastq files do not get modified in any way during this step.
FASTQC only interprets the .fastq files.

# cellranger mkgtf
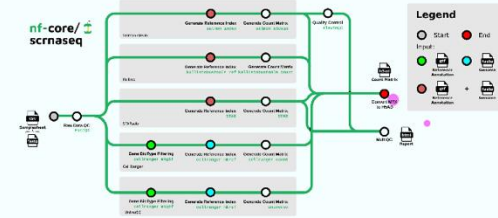
# cellranger mkgtf



## What

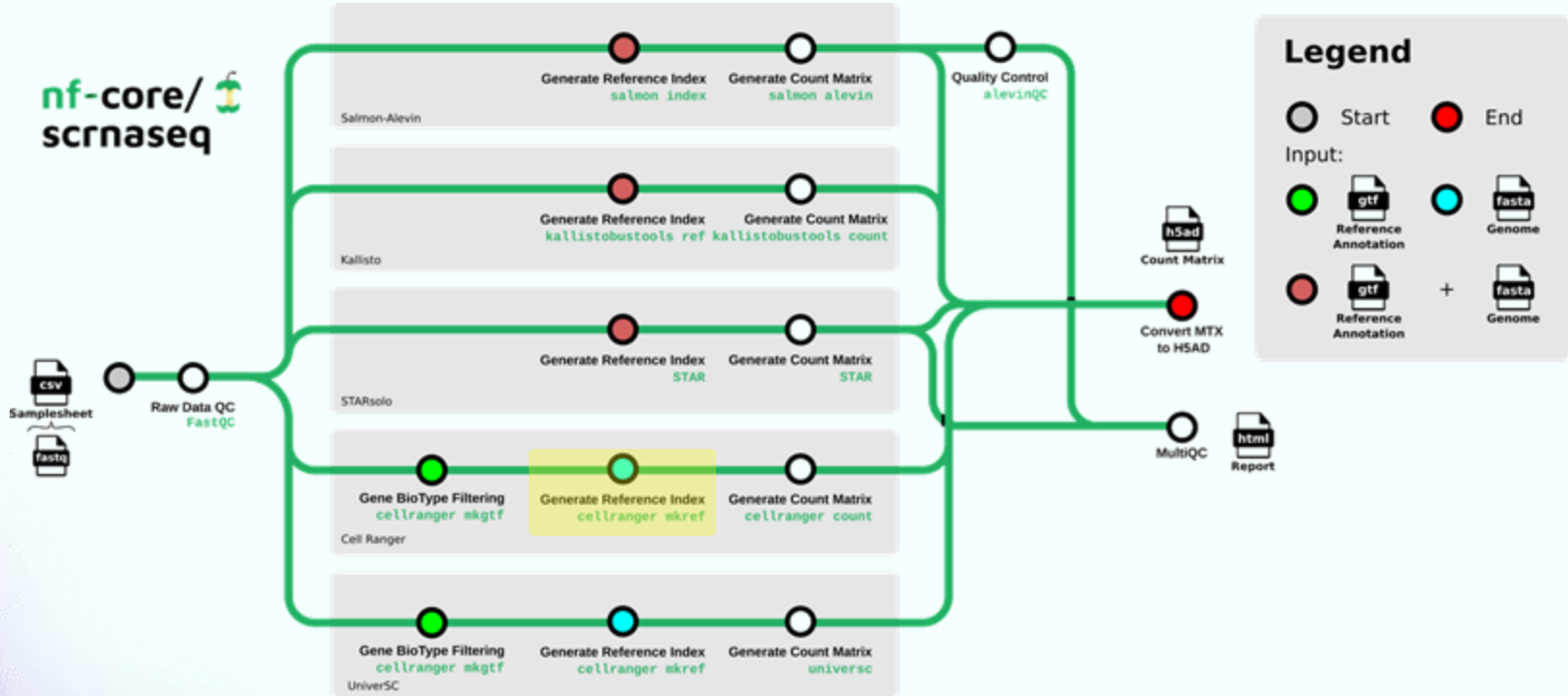- Pre-processing step run by **cellranger** to prepare the .gtf file.

## Why

- .gtf files do not always conform to a strict organization.
- Often, there is additional information in the .gtf file that is not needed for **cellranger**.

## Output

- A slimmed down .gtf that is properly formatted to ensure that the subsequent **cellranger** processes are run correctly and efficiently.

# cellranger mkref
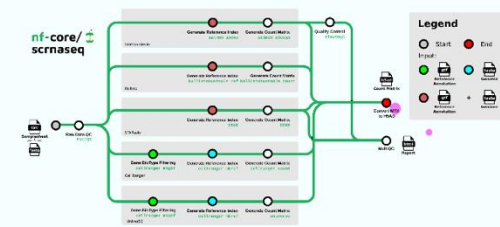
# cellranger mkref



## What

- Builds an index or *map* to be used for aligning your reads.
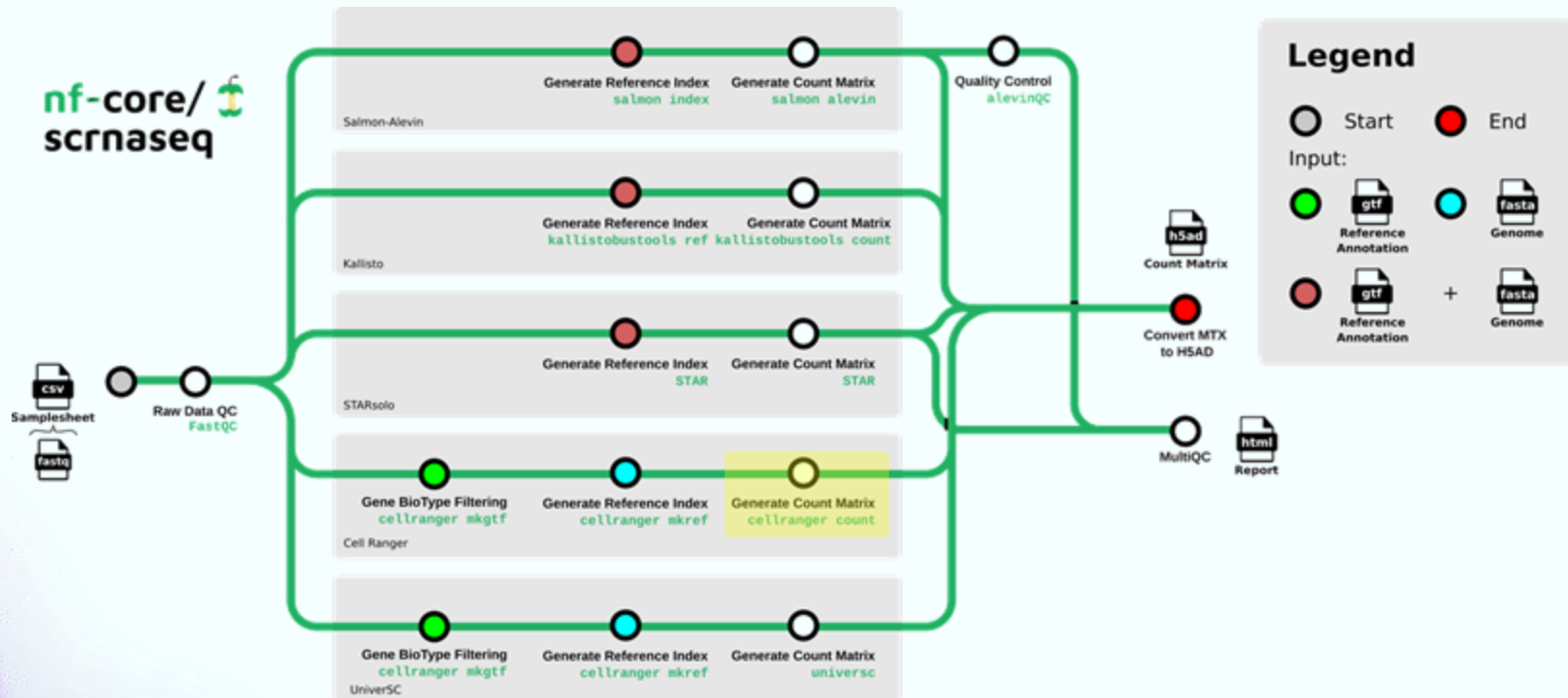- Takes in both the .fasta (sequence data) and .gtf (gene annotations).

## Why

- Having an index ensures the alignment step runs efficiently and accurately.

## Output

- A structured directory containing the index and auxiliary files to be used in the **cellranger count** step.
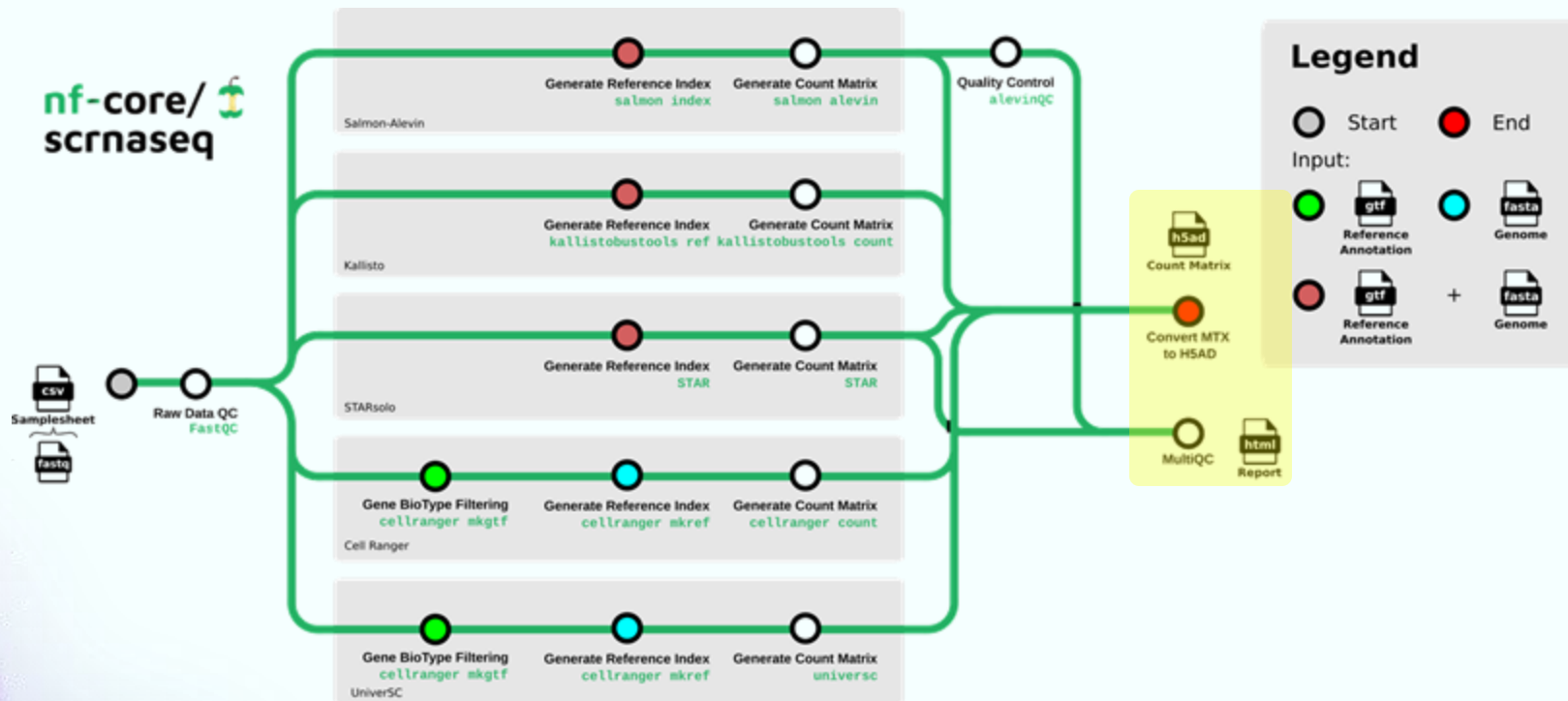
# cellranger count

# cellranger count



## What

- This is the most computationally intensive task in the pipeline.
- Takes the reads (from the .fasta files) and uses the index created in **cellranger mkref** to map what gene each read matches.

## How

- Two pass solution:
- Pass one: reads are aligned to the reference genome. Spice junction (both known and novel) and logged.
- Junctions are then filtered based on number of reads supporting the junction.
- Pass two: reads are once again aligned, but with this filtered set of junctions from the first pass.

## Output

- Key Output: raw and filtered counts
- Structured in a directory that contains 3 files (MEX format):
  - features.tsv.gz
  - barcodes.tsv.gz
  - matrix.mtx.gz

# post-alignment
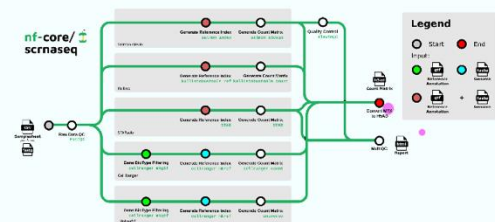
# post-alignment



## Extra Outputs

- This pipeline produces extra outputs mostly consisting of additional file conversions.
    - MEX → H5ad
    - MEX → rds

## pipeline_info

- Overview of the execution of the pipeline.
    - Report
    - Timeline
    - Dag

## multiqc

- A full breakdown of the workflow reporting results and statistics on the steps run.

# Questions?