# WASP - A versatile web-accessible single cell RNA-Seq processing platform
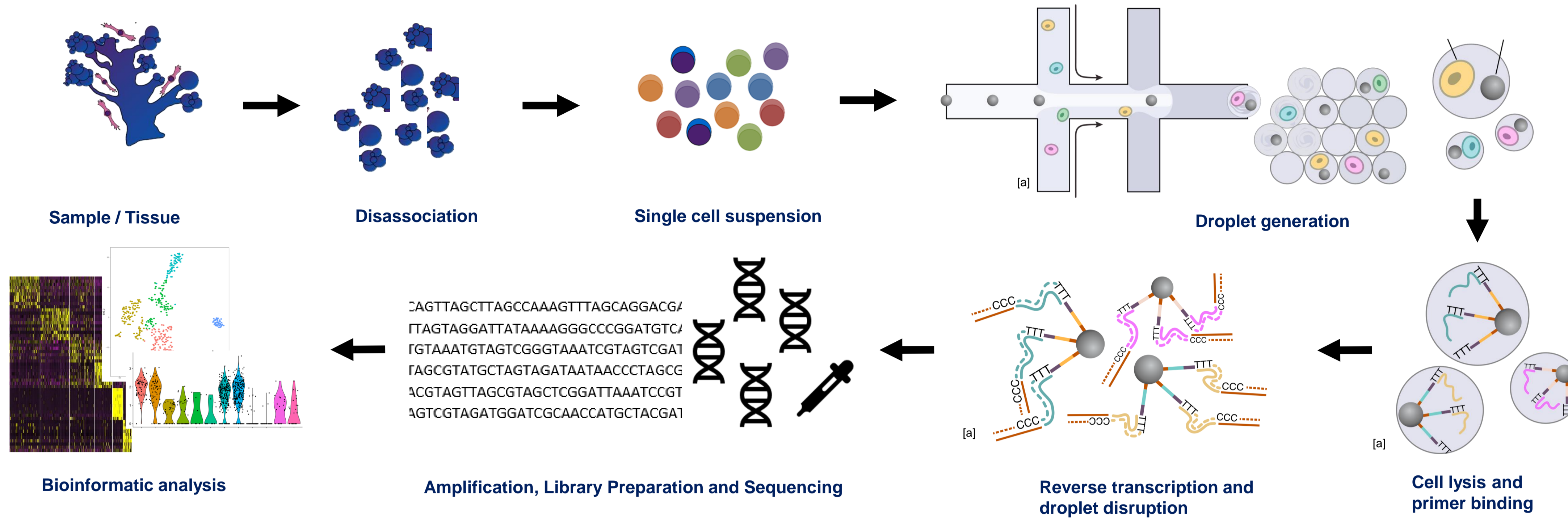
Andreas Hoek[1], Katharina Maibach[1], Ebru Özmen[1], Torsten Hain[2], Susanne Herold[3], Alexander Goesmann[1]

[1]Bioinformatics and Systems Biology, Justus Liebig University Giessen, [2]Institute of Medical Microbiology, Justus Liebig University Giessen, [3]Department of Medicine II, Justus Liebig University Giessen

## Processing of single cells



Sample / Tissue → Disassociation → Single cell suspension → Droplet generation → Cell lysis and primer binding → Reverse transcription and droplet disruption → Amplification, Library Preparation and Sequencing → Bioinformatic analysis
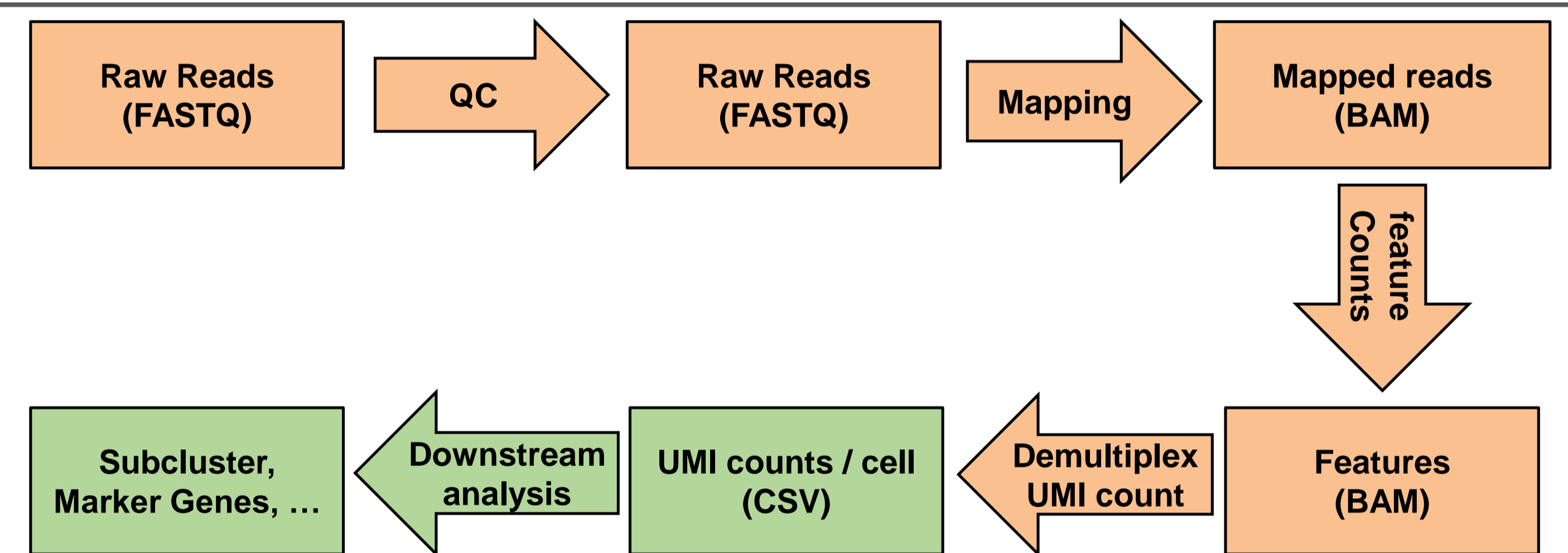
### Single cell RNA sequencing:

- Single cell RNA sequencing (scRNA-Seq) allows analysis of samples, e.g. tissues at high resolutions
- In contrast to bulk RNA-Seq, scRNA-Seq does allow the detection of cellular heterogeneity, cell types and gene expression of specific cells that are masked in bulk RNA-Seq
- Protocols developed in the early stages (2009 Tang et. al) only allowed analysis of a few cells
- Newer protocols such as Drop-Seq lowered costs and increased throughput allowing analysis of up to thousands of cells per run

## Bioinformatic analysis of single cell RNA sequencing data

### Workflow

- Analysis starts with the FASTQ file(s) containing the raw reads obtained from sequencing and ends with detection of cellular clusters, differentially expressed genes and marker genes including visualizations
- This workflow can be separated into two major steps:
  - Pre-Processing – Processing of raw reads to a gene-expression matrix
  - Post-Processing – Processing of gene-expression matrix to cellular clusters, characterization of subclusters, marker-genes, visualization of results



Raw Reads (FASTQ) → QC → Raw Reads (FASTQ) → Mapping → Mapped reads (BAM) → feature Counts → Features (BAM) → Demultiplex UMI count → UMI counts / cell (CSV) → Downstream analysis → Subcluster, Marker Genes, …
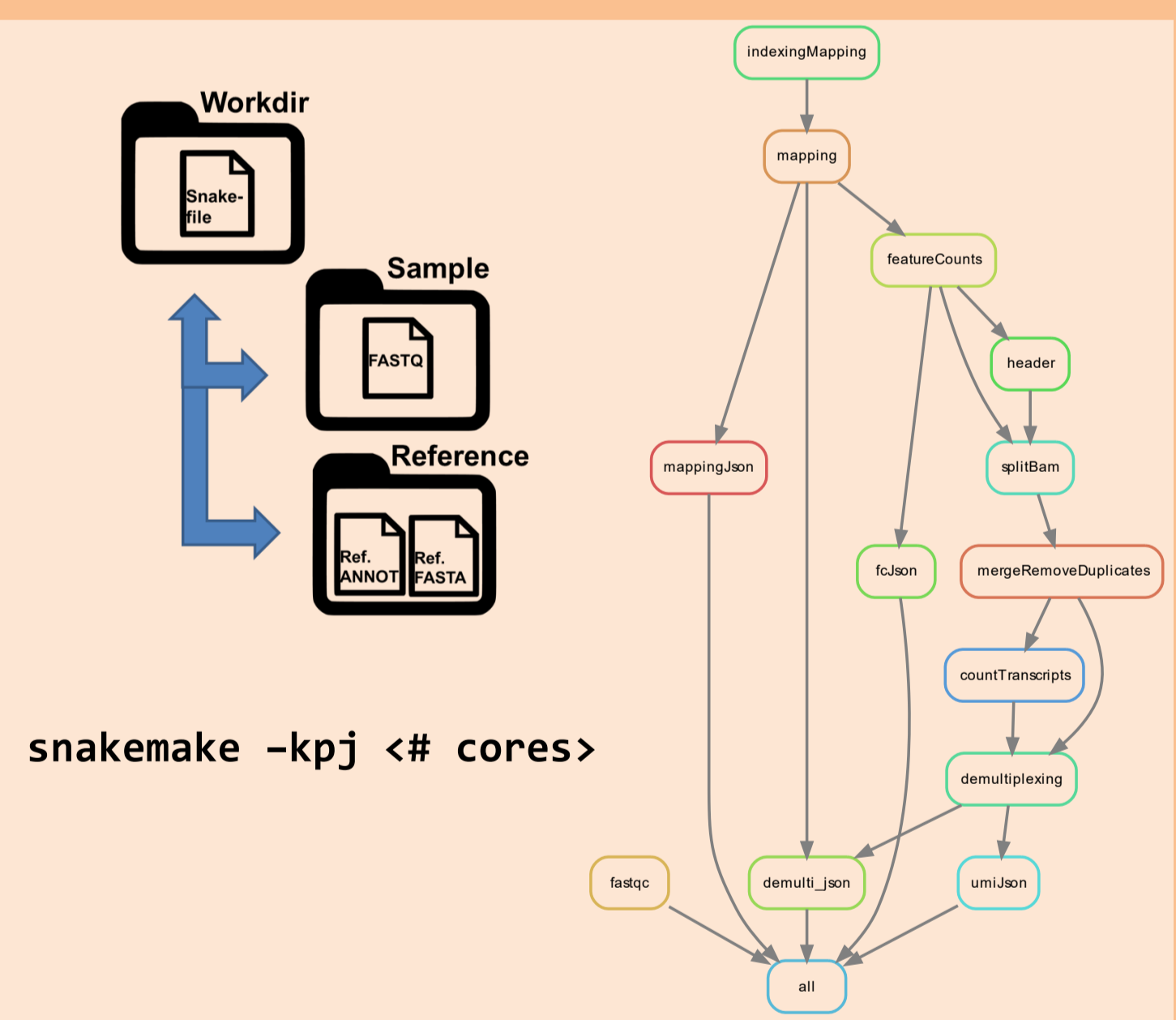
### Pre-Processing:

- Quality control of reads using FASTQC
- Reads are mapped to reference genome using STAR
- Counting of reads mapping to genomic features with featureCounts
- Demultiplexing reads to their cell of origin
- Counting unique mRNA fragments (UMI) per genomic feature
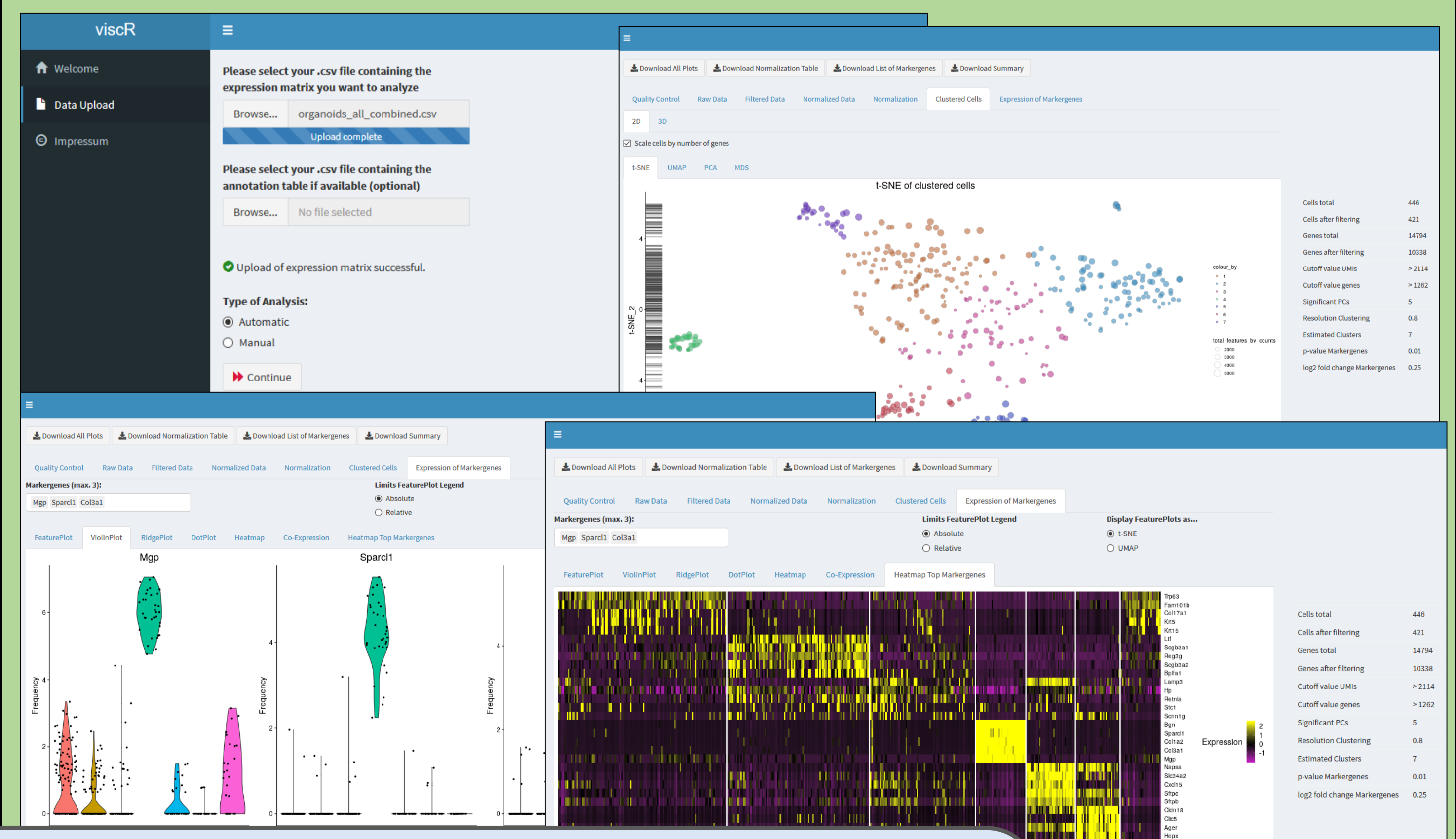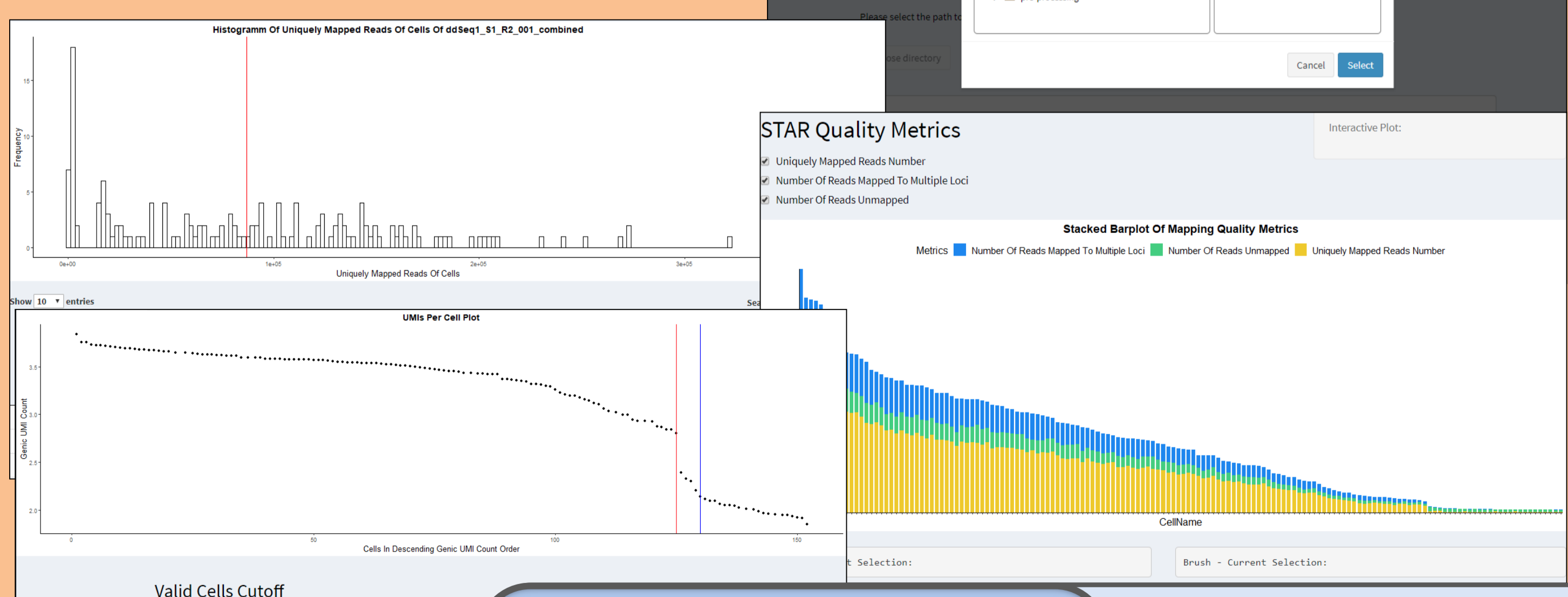
### Post-Processing:

- Steps performed using R packages
- Quality control & normalization
- Detection of highly variable genes
- Dimension reduction (PCA, t-SNE, UMAP)
- Clustering of cells and characterization of clusters, detection of differential gene expression, …
- Visualization of results

## Processing scRNA-seq data with WASP

- Pre-processing is performed using a snakemake workflow
- Users only have to provide reference genome, annotation and the FASTQ file(s)
- Metrics for each step are presented as an interactive R-Shiny web page
- After processing, users can select the number of cells to proceed
- As final result the pre-processing pipeline generates a gene expression matrix

- All dependencies (Tools and packages) can simply be installed via Conda
- Furthermore, a Docker container running the pre-processing is available

- Pre-Processing of ~100M reads running on 8 cores takes less than 3 hours

```
snakemake –kpj <# cores>
```



- Post-processing is performed using a variety of designated R packages
- Visualizations of all analysis steps are presented as an interactive R-Shiny web page
- WASP provides an automated workflow for unexperienced users and a manual workflow
- The manual workflow allows the user to choose parameters for each step and directly shows the impact in R-Shiny

- All dependencies (Python, R and packages) can be installed for Windows and Linux
- Furthermore, a Docker container running the post-processing is available

- Post-Processing of 1,000 cells running on a standard PC/Laptop takes ~5min



## WASP key features

- Automated processing of high-throughput scRNA-seq data
- Interactive browsing of results with high-quality graphics
- Easy usability, reproducible results and a fast analysis
- User-provided reference genome, annotation and FASTQ file(s)
- … or directly start with a gene expression matrix
- Local installation possible, Docker container planned

### References
[a] Macosko EZ., et al., Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets, Cell 2015

### Contact
www.computational.bio/software/wasp
andreas.hoek-2@computational.bio.uni-giessen.de