







Single-cell RNA-seq data analysis using Chipster

9.2.2018

chipster@csc.fi



CSC – Suomalainen tutkimuksen, koulutuksen, kulttuurin ja julkishallinnon ICT-osaamiskeskus



Understanding data analysis - why?

• Bioinformaticians might not always be available when needed

- Biologists know their own experiments best
 - Biology involved (e.g. genes, pathways, etc)
 - Potential batch effects etc
- Allows you to design experiments better > less money wasted

Allows you to discuss more easily with bioinformaticians

What will I learn?



- How to operate the Chipster software
- Analysis of single cell RNA-seq data
 - Central concepts
 - Analysis steps
 - File formats
- Exercises
 - DropSeq data preprocessing (from raw reads to expression values)
 - Clustering analysis of 10X Genomics data with Seurat tools



Introduction to Chipster

CSC

Chipster

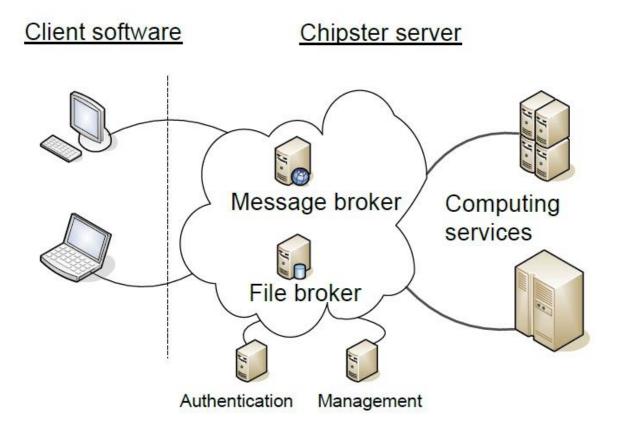
- Provides an easy access to over 360 analysis tools
 - Command line tools
 - R/Bioconductor packages
- Free, open source software

- What can I do with Chipster?
 - o analyze and integrate high-throughput data
 - visualize data efficiently
 - share analysis sessions
 - save and share automatic workflows

Technical aspects



- Client-server system
 - Enough CPU and memory for large analysis jobs
 - Centralized maintenance
- Easy to install
 - Client uses Java Web Start
 - Server available as a virtual machine







- Home
- Getting access
- Analysis tool content
- Screenshots
- Manual
- Tutorial videos
- Cite
- FAQ
- Contact
- For developers:
- o Open source project
- Tool editor

Welcome to Chipster

Chipster is a user-friendly software for analyzing high-throughput data such as NGS and microarrays. It contains over 360 analysis tools and a large collection of reference genomes. Users can save and share automatic analysis workflows, and visualize data interactively using for example the built-in genome browser. Chipster's client software uses Java Web Start to install itself automatically, and it connects to computing servers for the actual analysis. Chipster is open source, and the server environment is available as a virtual machine image free of charge. If you would like to use Chipster running on CSC's server, you need a user account.



Launch Chipster v3.12

...or launch with more memory: 3 GB or 6 GB

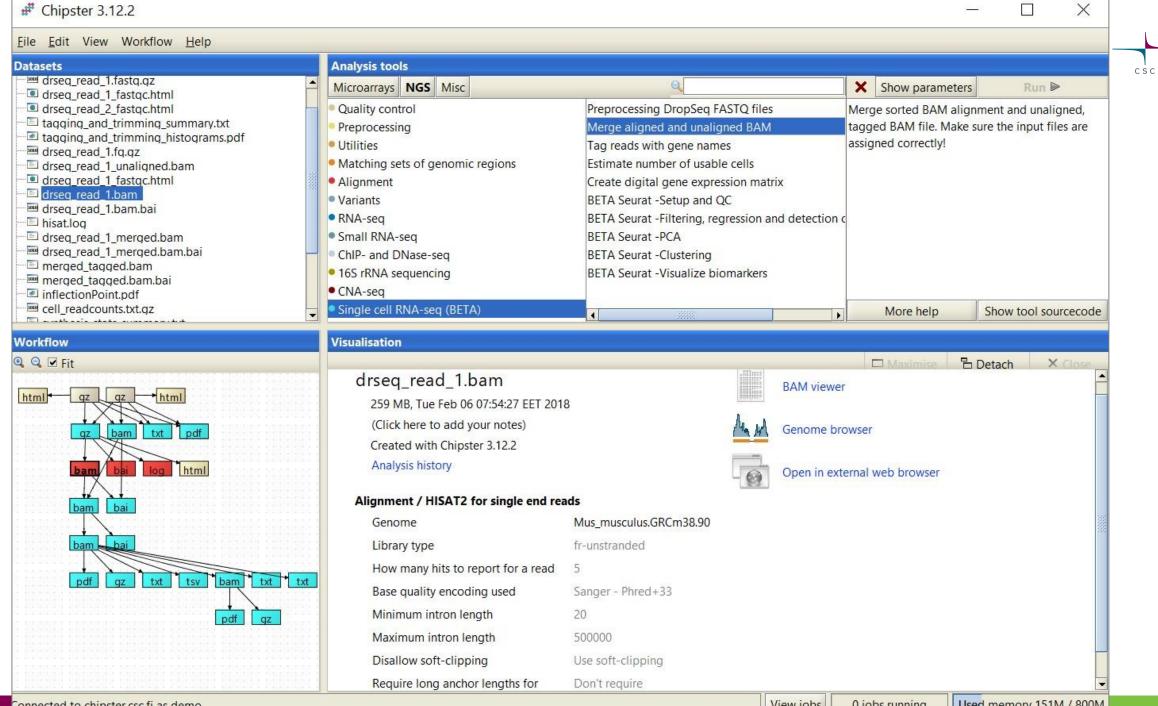
If you have trouble launching Chipster, read this

News and resources:

- 25.10.2017 Version 3.12 released
- 11.10.2016 Webinar: Beginner's introduction to RNA-seq data analysis using Chipster
- 10.10.2016 Instructions for stranded RNA-seq data
- 19.8.2014 RNA-seq data analysis guidebook with Chipster instructions
- News archive

Training:

- 9.2.2018 Single cell RNA-seg data analysis using Chipster, CSC
- 16.1.2018 Webinar: VirusDetect pipeline

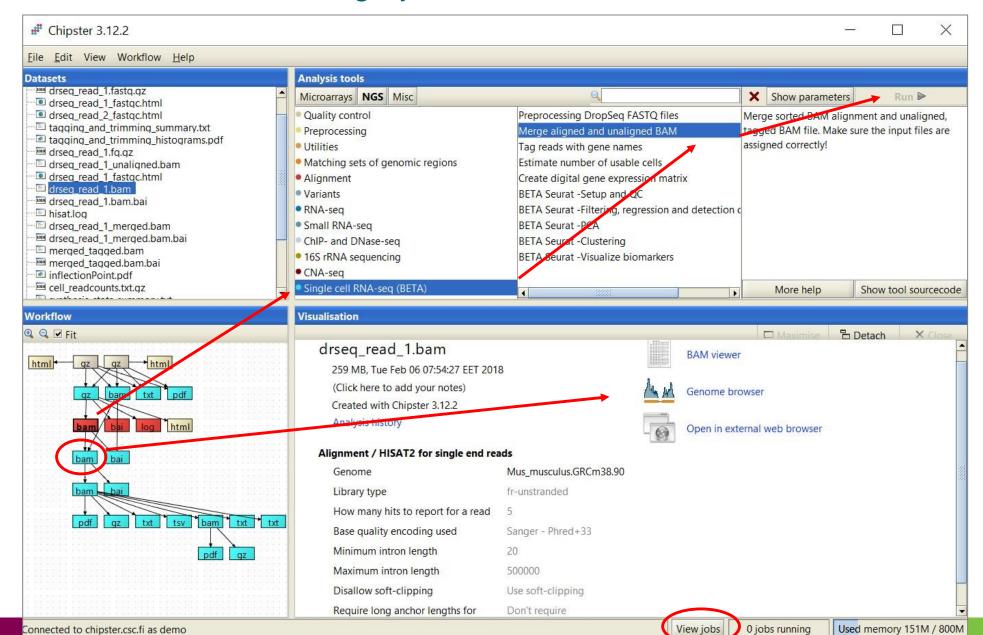


Used memory 151M / 800M View iobs 0 jobs running Connected to chipster.csc.fi as demo

Mode of operation

Select: data → tool category → tool → run → visualize

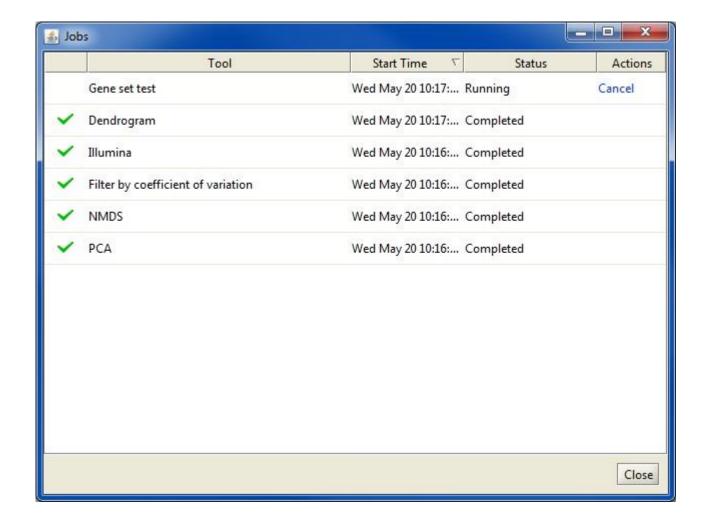




Job manager



- You can run many analysis jobs at the same time
- Use Job manager to:
 - view status
 - o cancel jobs
 - o view time
 - view parameters



Analysis history is saved automatically -you can add tool source code to reports if needed



Show for Datasets: Step title Applied analysis tool User notes	<u>≰</u> History	_
Dataset name: htsq-counts.tsv Created with operation: Alignment / Bowlie2 for single end reads Parameter Genome or transcriptome: hg19 Parameter Alignment strategy to use:sensitive Parameter Quality value format used:phred33 Parameter How many valid alignments are reported per read: 0 Parameter Put unaligned reads to a separate file: no Parameter Match bonus: 2 Parameter Maximum penalty for mismatch: 6 Parameter Maximum penalty for the reads: 5 Parameter Gap opening penalty for the reads: 5 Parameter Gap opening penalty for the reads: 3 Parameter Gap opening penalty for the reference: 5 Parameter Gap extension penalty for the reference: 3 Step 5	✓ Step title ✓ Applied analysis tool ☐ User notes ✓ Dataset name ✓ Parameters	
Parameter Gap extension penalty for the reference: 3 Step 5 Dataset name: htseq-counts.tsv Created with operation: RNA-seq / Count aligned reads per genes with HTSeq Parameter Organism: Homo_sapiens.GRCh37.68 Parameter Chromosome names in my BAM file look like: yes Parameter Does the alignment file contain paired-end data: no Parameter Was the data produced with a strand-specific RNA-seq protocol: no Parameter Mode to handle reads overlapping more than one feature: union Parameter Minimum alignment quality: 1 Parameter Feature type to count: exon Parameter Feature ID to use: gene_id	Dataset name: hESC.bam Created with operation: Alignment / Bowtie2 for single end reads Parameter Genome or transcriptome: hg19 Parameter Alignment strategy to use:sensitive Parameter Quality value format used:phred33 Parameter How many valid alignments are reported per read: 0 Parameter Put unaligned reads to a separate file: no Parameter Match bonus: 2 Parameter Maximum penalty for mismatch: 6 Parameter Penalty for non-ACGTs: 1 Parameter Gap opening penalty for the reads: 5 Parameter Gap extension penalty for the reads: 3	
Save Close	Step 5 Dataset name: htseq-counts.tsv Created with operation: RNA-seq / Count aligned reads per genes with HTSeq Parameter Organism: Homo_sapiens.GRCh37.68 Parameter Chromosome names in my BAM file look like: yes Parameter Does the alignment file contain paired-end data: no Parameter Was the data produced with a strand-specific RNA-seq protocol: no Parameter Mode to handle reads overlapping more than one feature: union Parameter Minimum alignment quality: 1 Parameter Feature type to count: exon Parameter Feature ID to use: gene_id Parameter Add chromosomal coordinates to the count table: yes	

Analysis sessions



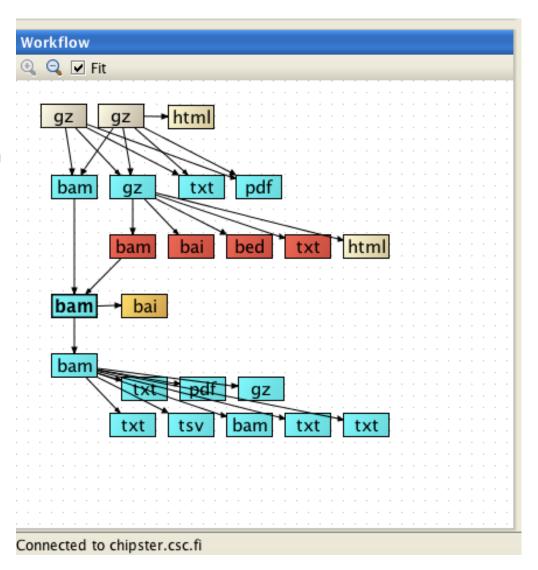
- Remember to save the analysis session.
 - Session includes all the files, their relationships and metadata (what tool and parameters were used to produce each file).
 - Session is a single .zip file.

- You can save a session locally (on your computer)
- and in the cloud
 - obut note that the cloud sessions are not stored forever!

Workflow panel



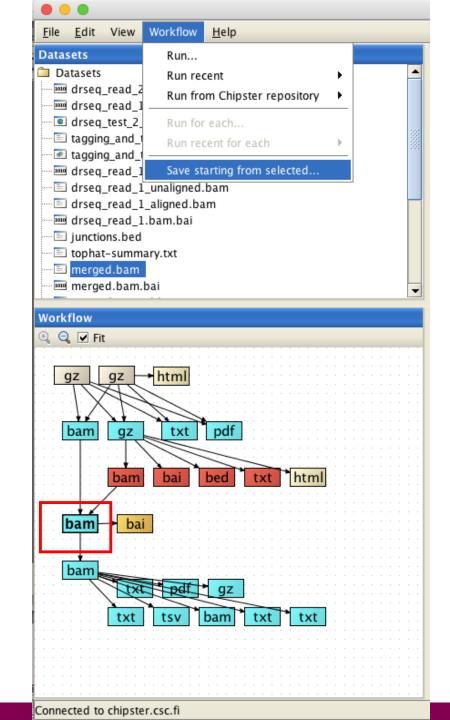
- Shows the relationships of the files
- You can move the boxes around, and zoom in and out.
- Several files can be selected by keeping the Ctrl key down
- Right clicking on the data file allows you to
 - Save an individual result file ("Export")
 - o Delete
 - o Link to another data file
 - Save workflow





Workflow – reusing and sharing your analysis pipeline

- You can save your analysis steps as a reusable automatic "macro"
 - o all the analysis steps and their parameters are saved as a script file
 - you can apply workflow to another dataset
 - you can share workflows with other users



Saving and using workflows



Select the <u>starting point</u> for your workflow

Select "Workflow / Save starting from selected"

- Save the workflow file on your computer with a meaningful name
 - Opon't change the ending (.bsh)!

- To run a workflow, select
 - Workflow -> Open and run
 - Workflow -> Run recent (if you saved the workflow recently).

Analysis tool overview



160 NGS tools for

RNA-seq

single cell RNA-seq

small RNA-seq

exome/genome-seq

ChIP-seq

FAIRE/DNase-seq

CNA-seq

Metagenomics (16S rRNA)

• 140 microarray tools for

gene expression

miRNA expression

protein expression

aCGH

SNP

integration of different data

60 tools for sequence analysis

BLAST, EMBOSS, MAFFT

Phylip



Visualizing the data

Data visualization panel

- Maximize and redraw for better viewing
- o Detach = open in a separate window, allows you to view several images at the same time

Two types of visualizations

1. <u>Interactive visualizations</u> produced by the client program

- Select the visualization method from the pulldown menu
- Save by right clicking on the image

2. Static images produced by analysis tools

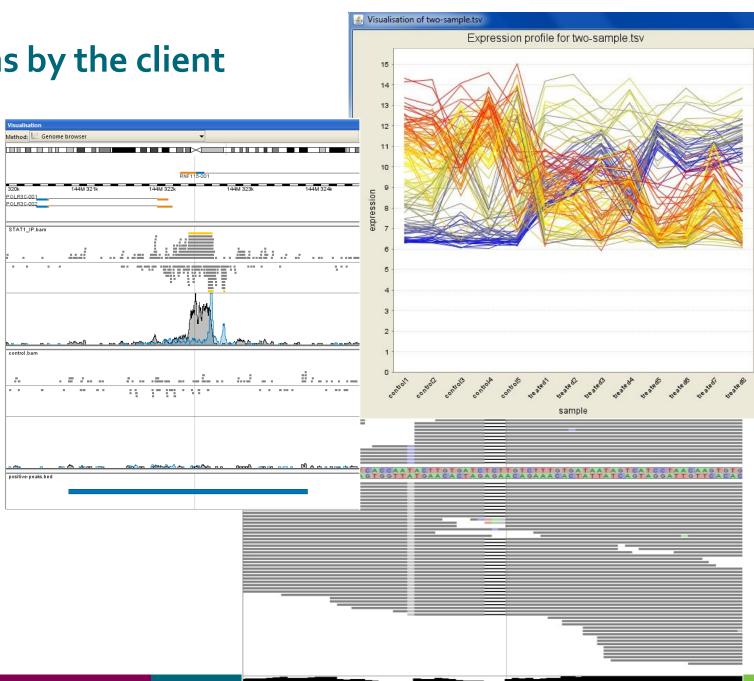
- Select from Analysis tools / Visualisation
- View by double clicking on the image file
- Save by right clicking on the file name and choosing "Export"

Interactive visualizations by the client

- Genome browser
- Spreadsheet
- Histogram
- Venn diagram
- Scatterplot
- 3D scatterplot
- Volcano plot
- Expression profiles
- Clustered profiles
- Hierarchical clustering
- SOM clustering

Available actions:

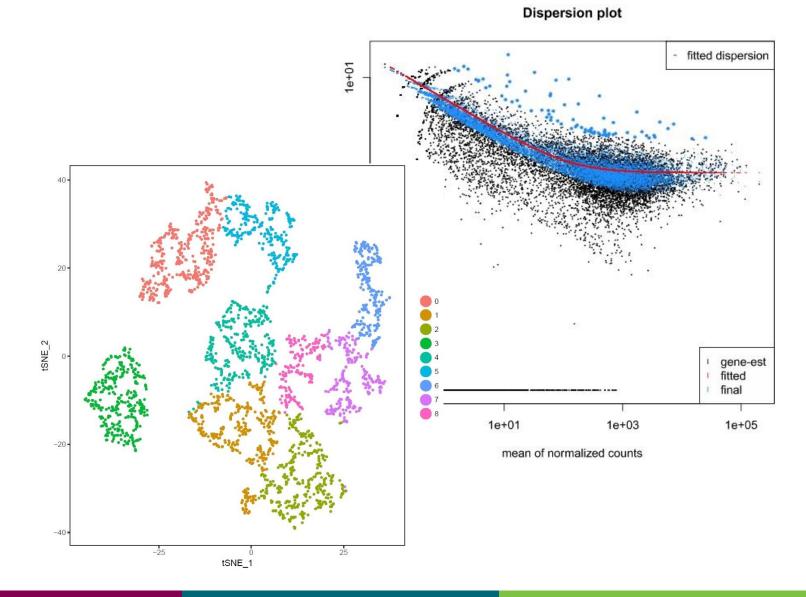
- Select genes and create a gene list
- Change titles, colors etc
- Zoom in/out





Static images produced by R/Bioconductor

- Dispersion plot
- MA plot
- MDS plot
- Box plot
- Histogram
- Heatmap
- tSNE plot
- Violin plot
- PCA plot
- Dendrogram
- K-means clustering
- SOM-clustering
- Etc...



Options for importing data to Chipster



- Import files/ Import folder
- Import from URL
 - Utilities / Download file from URL directly to server
- Open an analysis session
 - Files / Open <u>session</u>
- Import from SRA database
 - Utilities / Retrieve FASTQ or BAM files from SRA
- Import from Ensembl database
 - Utilities / Retrieve data for a given organism in Ensemble
- What kind of RNA-seq data files can I use in Chipster?
 - o Compressed files (.gz) and tar packages (.tar) are ok
 - FASTQ, BAM, read count files (.tsv), GTF

Problems? Send us a support request



-request includes the error message and link to analysis session (optional)

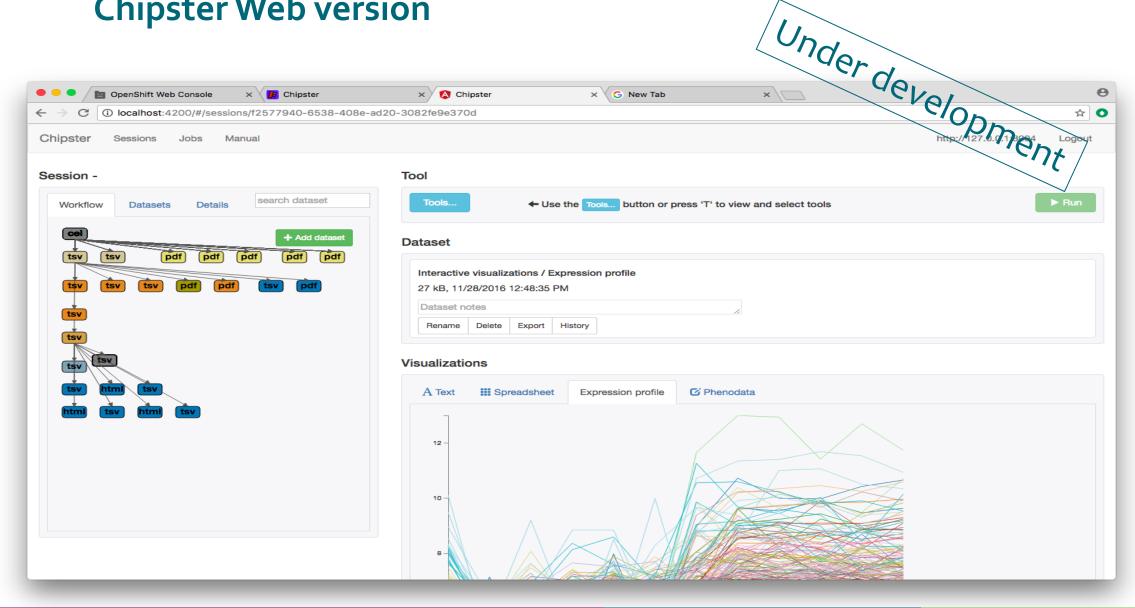
```
Hi.
I'm trying to normalise my Illumina microarray data (obtained with the Illumina HT-12 v4.0)
For that purpose I have selected the Normalisation option "Illumina - lumi pipeline"
However, the normalisation did not complete successfully.
Any advice to solve this problem ?
                                                                                             Contact support
Thank you in advance for your precious help.
                                                                                              Message
Best regards
Error message:
in library(chiptype, character.only = T) :
  there is no package called 'Illumina.db'
> chipster.common.path = '/opt/chipster/comp/modules/common/R-2.12'
> chipster.module.path = '/opt/chipster/comp/modules/microarray'
> setwd("271661a6-946c-450f-bb21-5d5b5a2837aa")
                                                                                              Your email
> probe.identifier <- "Probe ID"
> transformation <- "log2"
> background.correction <- "none"
> normalize.chips <- "quantile"
                                                                                              Attach data and workflow information
> chiptype <- "empty"
> # TOOL norm-illumina-lumi.R: "Illumina - lumi pipeline" (Illumina normalization using

✓ Attach log files

BeadSummaryData files, and using lumi methodology. If you have a BeadSummaryData that re
                                                                                                                  Cancel
```



Chipster Web version



More info

- chipster@csc.fi
- http://chipster.csc.fi
- Chipster tutorials in YouTube

Chipster Tutorials

Uploads





Using Chipster in Taito-shell

Using Chipster on Taito-shell 146 views · 9 months ago

Basic sequence analysis tasks in Chipster

11:43

NGS

BLAST tools in Chipster

Sequence analysis

223 views · 1 year ago

EdgeR for multiva experiments

924 views . 1 year ag

EdgeR for mu

(differential e

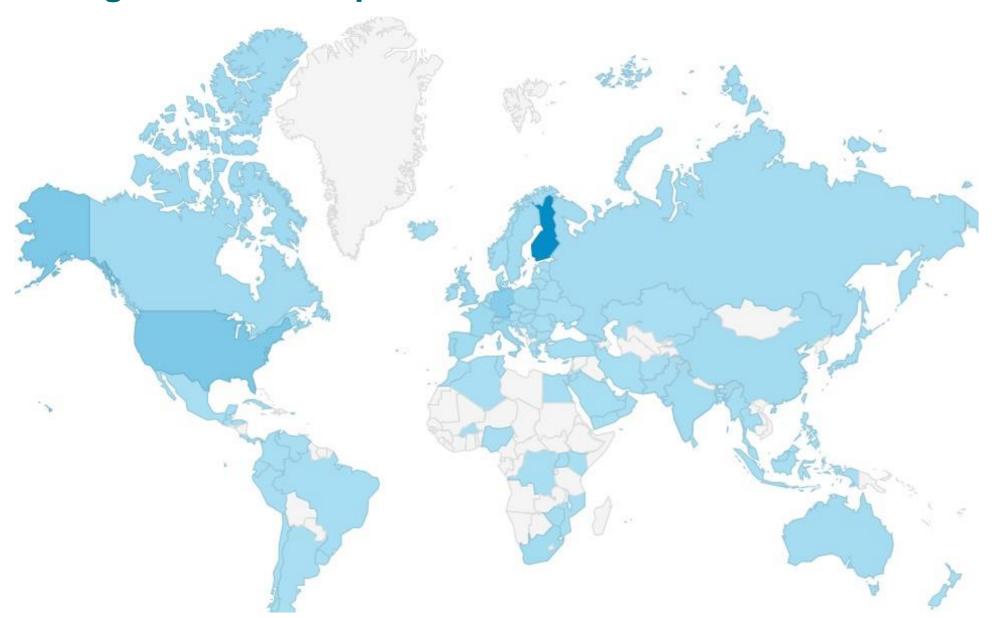
analysis for R

experime



Acknowledgements to Chipster users and contibutors







Introduction to single cell RNA-seq

Single cell RNA-seq



- New technology, data analysis methods are actively developed
- Measures distribution of expression levels for each gene across a population of cells
- Allows to study cell-specific changes in transciptome
- Applications
 - o Identifying cell sub-populations within a biological condition
 - Studying dynamic processes like differentiation using pseudotime ordering
 - Swithces
 - Branch points
 - Transcriptional regulatory networks
- Datasets range from 10² to 10⁵ cells

Single cell RNA-seq technology

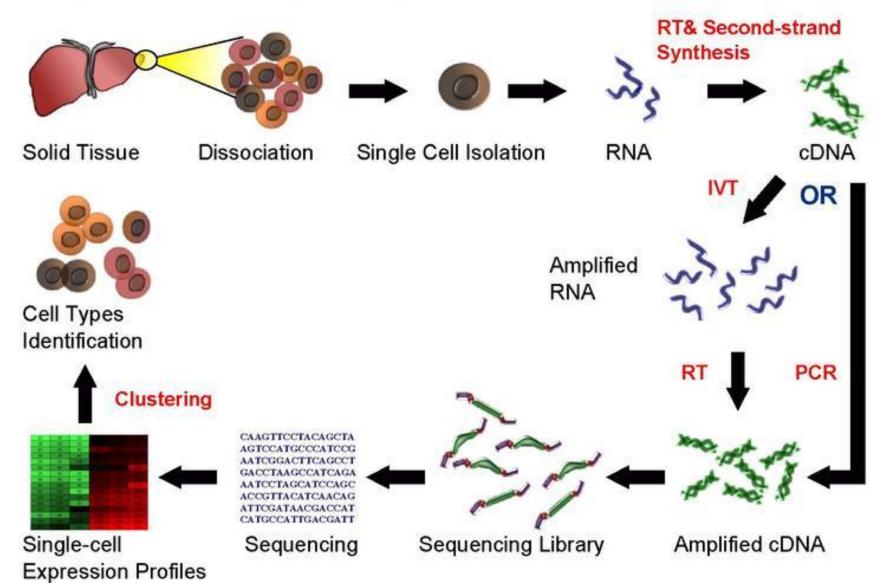


- Protocols:
 - DropSeq
 - \circ InDrop
 - CELL-seq
 - SMART-seq
 - SMARTer
 - MARS-seq
 - SCRB-seq
 - ∘ Seq-well
 - STRT-seq

- Commercial platforms:
 - Fluidigm C1 (FuGU)
 - 10x Genomics Chromium (FIMM, BTK)

Single Cell RNA Sequencing Workflow



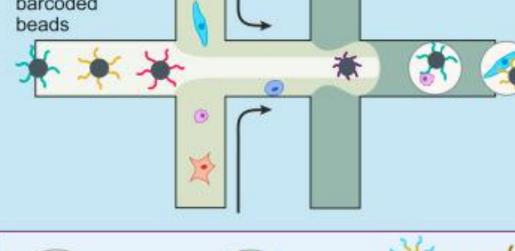


https://hemberg-lab.github.io/scRNA.seq.course/introduction-to-single-cell-rna-seq.html

8.2.2018



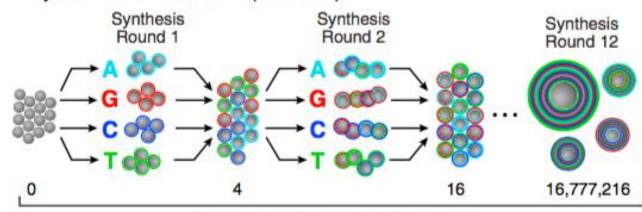
Drop-seq single cell analysis Cells Distinctly barcoded beads



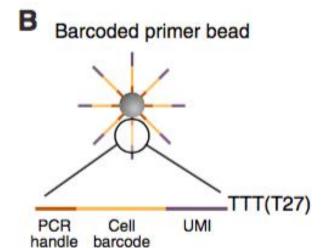


1000s of DNA-barcoded single-cell transcriptomes

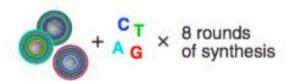
C Synthesis of cell barcode (12 bases)



Number of unique barcodes in pool



Synthesis of UMI (8 bases)

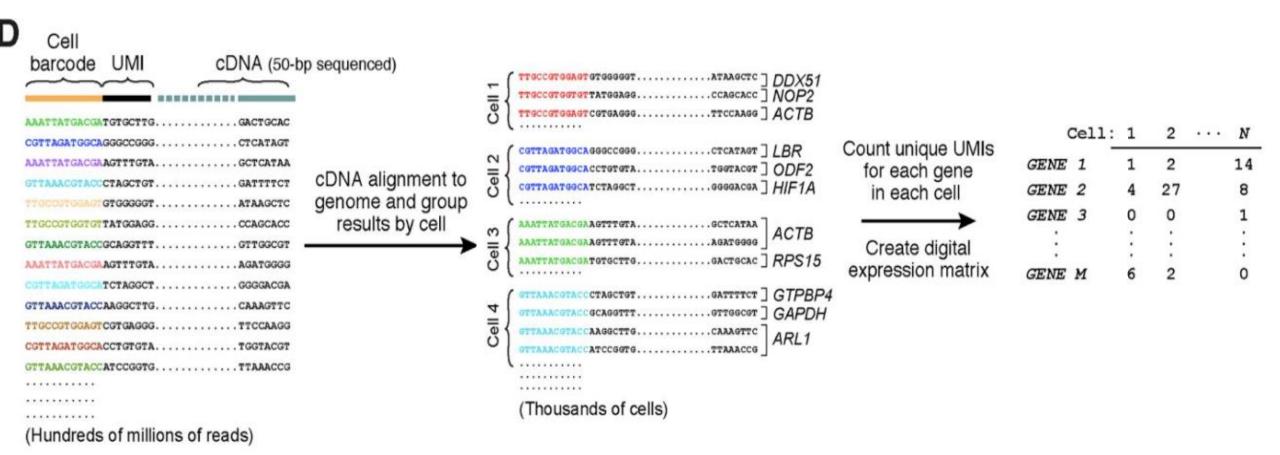


- Millions of the <u>same</u> cell barcode per bead
- 4⁸ different molecular barcodes (UMIs) per bead

DropSeq data preprocessing overview



Overview of DGE extraction



8.2.2018

CSC

What can go wrong in Drop-seq?

- 1. Ideally there is one healthy cell in the droplet so that we get a STAMP (single-cell transcriptome attached to a microparticle). However, sometimes
- There is no cell in the droplet, just ambient RNA
 - → Detect "empties" based on the small number of genes expressed and remove
- There are two (or more) cells in a droplet
 - → Detect duplets (and multiplets) based on the large number of genes expressed and remove
- The cell in the droplet is broken/dead
 - > Detect based on high proportion of reads mapping to mitochondrial genome and remove
- 2. Sometimes barcodes have synthesis errors in them, e.g. one base is missing
 - > Detect by checking the distribution of bases at each position and fix the code or remove the cell



Single cell RNA-seq data analysis

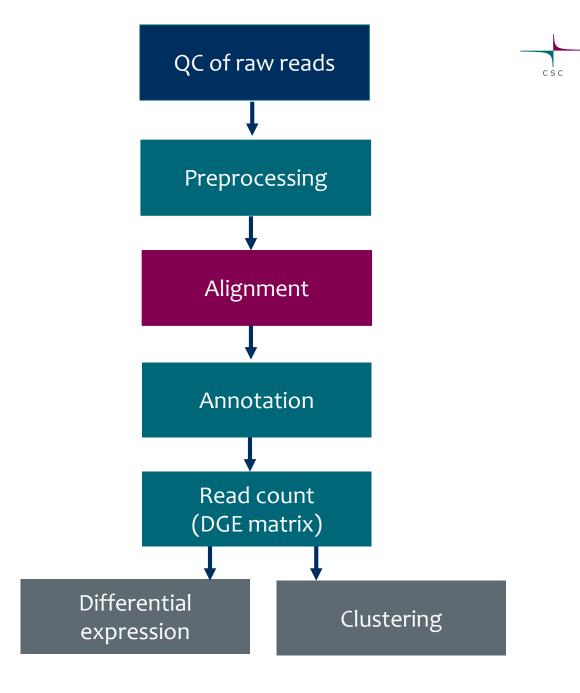
CSC

Single cell RNA-seq data is challenging

- The detected expression level for many genes is zero
- Data is noisy. High level of variation due to
 - Capture efficiency (percentage of mRNAs captured)
 - Amplification bias (non-uniform amplification of transcripts)
 - Cells differ in terms of cell-cycle stage and size
- Complex distribution of expression values
 - o Cell heterogeneity and the abundance of zeros give rise to multimodal distributions

→ Many methods used for bulk RNA-seq data won't work

scRNA-seq data analysis steps



36

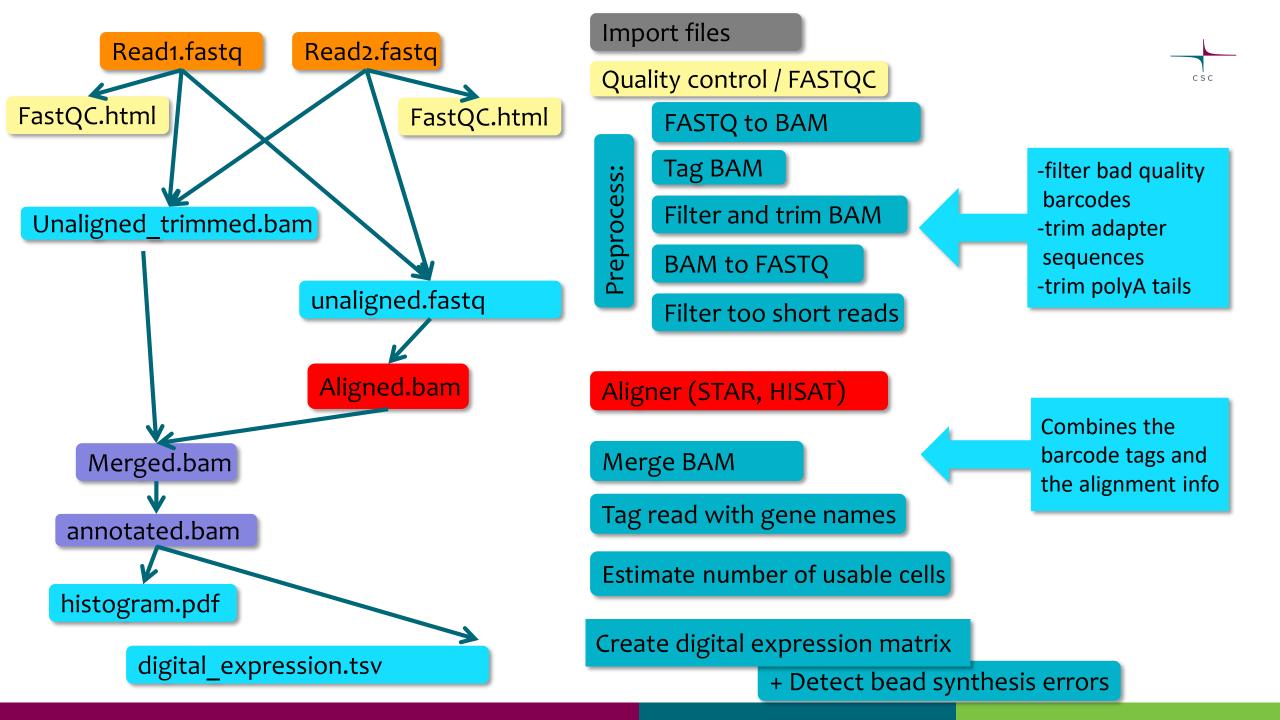


DropSeq data preprocessing



From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix





From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

What and why?



Potential problems

- o low confidence bases, Ns
- o sequence specific bias, GC bias
- \circ adapters
- o sequence contamination
- ounexpected length of reads

O ...

Knowing about potential problems in your data allows you to

- o correct for them before you spend a lot of time on analysis
- o take them into account when interpreting results



Software packages for quality control

- FastQC
- FastX
- PRINSEQ
- TagCleaner
- ...

Raw reads: FASTQ file format



• Four lines per read:

@read name

GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT

+ read name

!"*((((***+))%%%++)(%%%%).1***-+*"))**55CCF>>>>>CCCCCC65

- http://en.wikipedia.org/wiki/FASTQ_format
- Attention: Do not unzip FASTQ files!
 - Chipster's analysis tools can cope with zipped files (.gz)

Base qualities



- If the quality of a base is 20, the probability that it is wrong is 0.01.
 - \circ **Phred quality score** \circ = -10 * \log_{10} (probability that the base is wrong)

T C A G T A C T C G
40 40 40 40 40 40 40 37 35

Probability: 0.0001
Phred score: 40
ASCII coding: I (capital i)

- "Sanger" encoding: numbers are shown as ASCII characters so that 33 is added to the Phred score
 - \circ E.g. 39 is encoded as "H", the 72nd ASCII character (39+33 = 72)
 - Note that older Illumina data uses different encoding
 - o Illumina1.3: add 64 to Phred
 - o Illumina 1.5-1.7: add 64 to Phred, ASCII 66 "B" means that the whole read segment has low quality

Base quality encoding systems

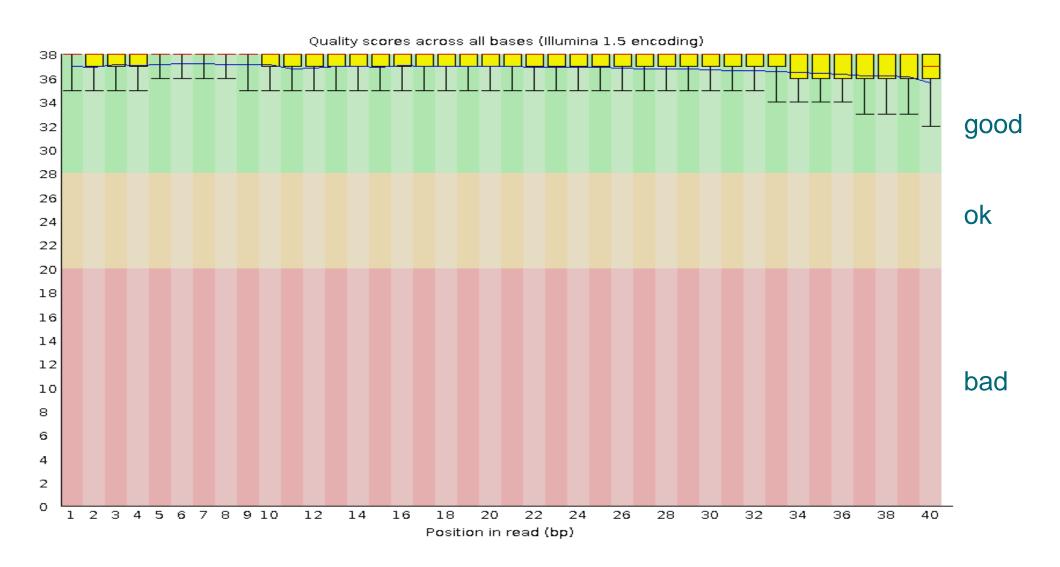


```
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^ `abcdefghijklmr
33
                                            104
0.2.....41
S - Sanger
          Phred+33, raw reads typically (0, 40)
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)
```

http://en.wikipedia.org/wiki/FASTQ_format

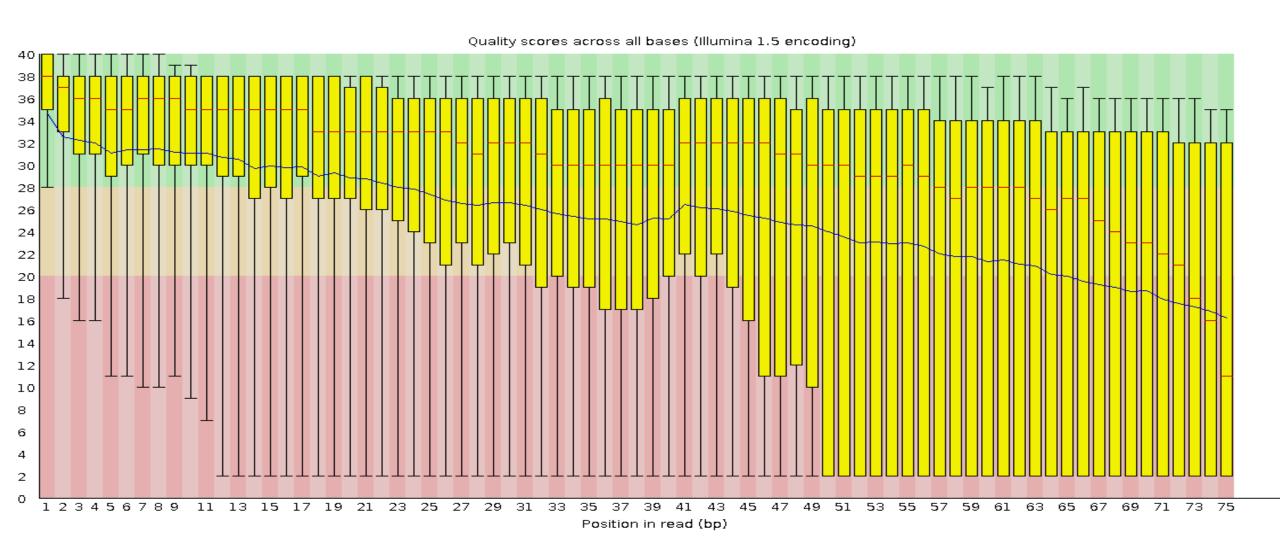
Per position base quality (FastQC)





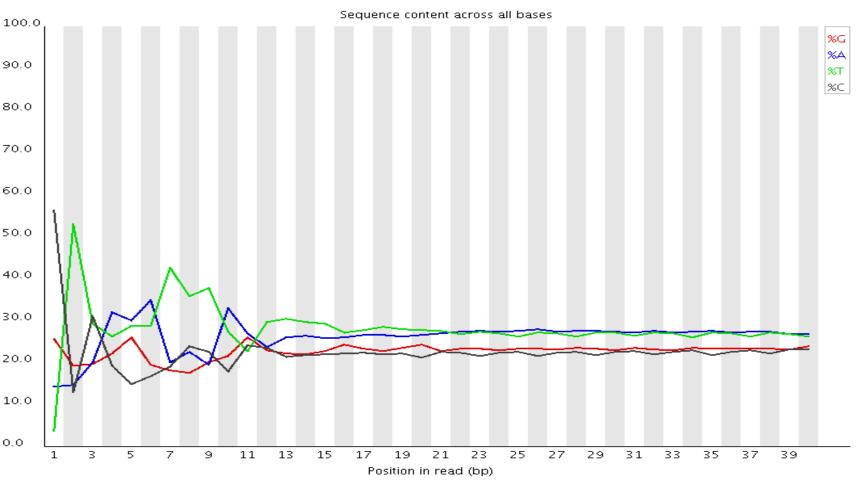
Per position base quality (FastQC)





Per position sequence content (FastQC)





- Enrichment of k-mers at the 5' end due to use of random hexamers or transposases in the library preparation
- Typical for RNA-seq data
- Can't be corrected, doesn't usually effect the analysis
- Now: the barcodes

Check also the length of the reads

- Barcode read: should be 20 bp
- Other length => BAM tagging won't work
- To get rid of shorter barcode reads:
 - o select **BOTH FASTQ files**
 - o tool Preprocessing / Trim reads with Trimmomatic with parameter Minimum length of reads to keep = 20
 - O Continue analysis with the files that are named "--trimmed.fq.gz" (instead of the "--unpaired_trimmed.fq.gz" ones).
- To trim too long barcode reads:
 - o after the previous step: select **read1** file, run **Preprocessing / Trim reads with Trimmomatic** with parameter **Number of bases to keep from the start = 20**

Measure	Value
Filename	reads.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	10157126
Sequences flagged as poor quality	0
Sequence length	20
%GC	51



From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

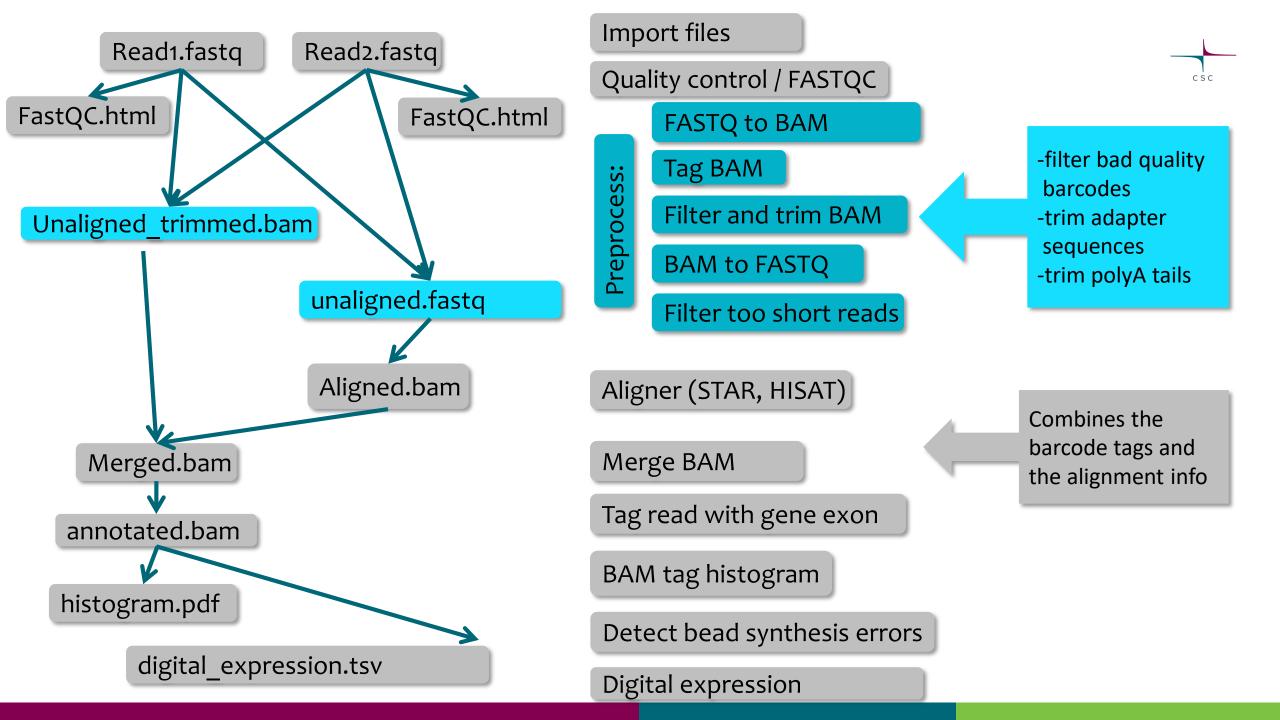


Preprocessing single cell DropSeq FASTQ files

- This tool is a combination of several tools:
 - DropSeq (TagBamWithReadSequenceExtended, FilterBam)
 - o Picard (FASTQ to BAM, BAM to FASTQ) and
 - Trimmomatic (MINLEN)
- The steps are:
 - 1. Convert the FASTO files into an unaligned BAM file
 - 2. Tag the reads in the BAM file with the cellular and molecular barcodes
 - 3. Filter and trim the reads in the BAM file
 - 4. Convert the tagged and trimmed BAM file back into a FASTO file for the alignment
 - 5. Filter out too short reads of the FASTQ file

Why FASTQ-BAM-FASTQ-BAM?

- -FASTQ format cannot hold the information about the cellular and molecular barcodes. BAM format has tag fields which can be used to hold this information
- -In BAM format we can also do some trimming and filtering for the reads.
- -However, the aligners take as input only FASTQ format, which is why we need to transform the trimmed & filtered BAM back to FASTQ format.
- -Later on, after alignment, we merge the two!





Barcoded primer bead

Convert the FASTQ files into an unaligned BAM file

Tag the reads in the BAM file with the cellular and molecular barcodes

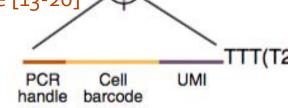
Filter and trim the reads in the BAM file

Convert the tagged and trimmed BAM file back into a FASTQ file for the alignment

• Base range for cell barcode [1-12]

• Base range for molecule barcode [13-20]

Base quality [10]



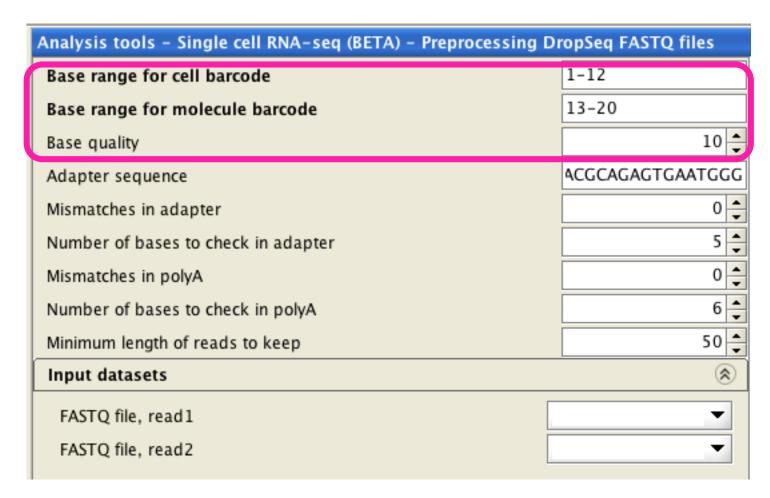
- Adapter sequence [AAGCAGTGGTATCAACGCAGAGTGAATGGG]
- Mismatches in adapter [o]
- Number of bases to check in adapter [5]
- Mismatches in polyA [o]
- Number of bases to check in polyA [6]

Filter out too short reads of the FASTQ file • Minimum length of reads to keep [50]



Preprocessing DropSeq FASTQ files –Tagging reads

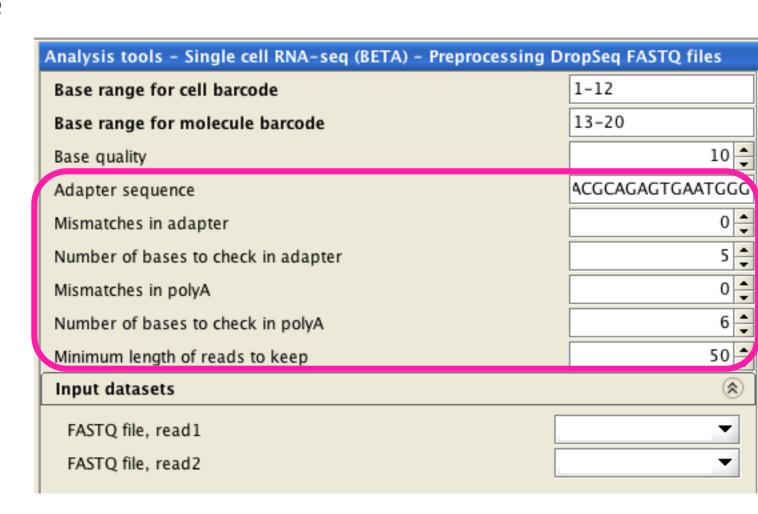
- BAM tags
 - XM = molecular barcode
 - XC = cellular barcode
 - OXQ = number of bases that fall below quality threshold
- After this, we can forget the barcode read
- You need to tell the tool which bases correspond to which barcode



Preprocessing DropSeq FASTQ files –Trimming and filtering



- Quality: filter out reads where more than 1 base have poor quality
- Adapters: Trim away any user determined sequences
 - SMART adapter as default
 - How many mismatches allowed (o)
 - How long stretch of the sequence there has to be at least (5)
- polyA: hard-clip polyA tails
 - How many A's need to be there before clipping happens (6)
 - o mismatches allowed (o)
- Minimum length: filter out short reads from the FASTQ file





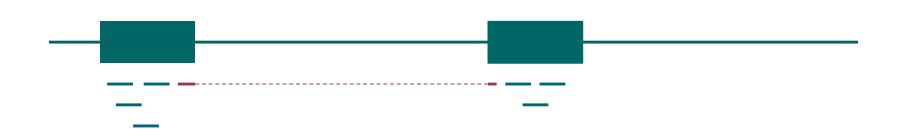
From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

Alignment to reference genome



- Goal is to find out where a read originated from
 - Challenge: variants, sequencing errors, repetitive sequence
- Many organisms have introns, so RNA-seq reads map to genome non-contiguously
 - → spliced alignments needed
 - OBut sequence signals at splice sites are limited and introns can be thousands of bases long
- Splice-aware aligners
 - HISAT2, TopHat
 - STAR



HISAT₂



- HISAT = <u>Hierarchical Indexing for Spliced Alignment of Transcripts</u>
- Fast spliced aligner with low memory requirement
- Reference genome is indexed for fast searching
- Uses two types of indexes
 - One global index: used to anchor each alignment (28 bp is enough)
 - o Thousands of small local indexes, each covering a genomic region of 56 Kbp: used for rapid extension of alignments (good for reads with short anchors over splice sites)
- Uses splice site information found during the alignment of earlier reads in the same run



HISAT2 parameters

Genome	Homo_sapiens.GR ¬
Library type	fr-unstranded •
How many hits to report for a read	5
Base quality encoding used	Sanger - Phred+33 ▼
Minimum intron length	20
Maximum intron length	500000
Disallow soft-clipping	Use soft-clipping
Require long anchor lengths for subsequent assembly	Don't require

8.2.2018

59

STAR



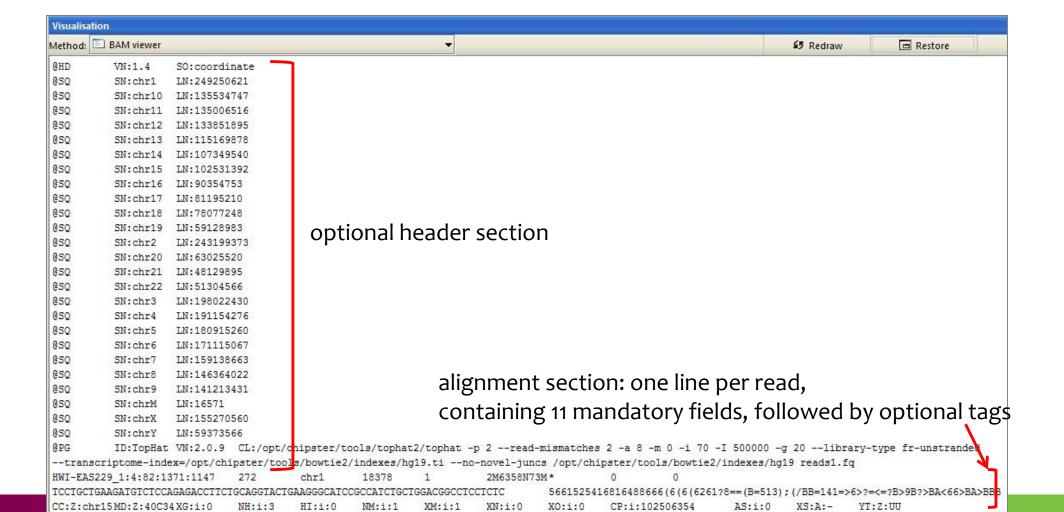
>STAR (Spliced Transcripts Alignment to a Reference) uses a 2-pass mapping process

- splice junctions found during the 1st pass are inserted into the genome index, and all reads are remapped in the 2nd mapping pass
- this doesn't increase the number of detected novel junctions, but it allows more spliced reads mapping to novel junctions.
- Maximum alignments per read -parameter sets the maximum number of loci the read is allowed to map to
 - Alignments (all of them) will be output only if the read maps to no more loci than this. Otherwise no alignments will be output.
- > Chipster offers an Ensembl GTF file to detect annotated splice junctions
 - you can also give your own, for example GENCODE GTFs are recommended
- > Two log files
 - Log_final.txt lists the percentage of uniquely mapped reads etc.
 - Log_progress.txt contains process summary

BAM file format for aligned reads



 SAM (Sequence Alignment/Map) is a tab-delimited text file containing aligned reads. BAM is a binary (and hence more compact) form of SAM.



Fields in BAM file

CSC

- >read name HWI-EAS229_1:2:40:1280:283
- ≻flag 272
- >reference name 1
- >position 18506
- >mapping quality o
- ➤ CIGAR 49M6183N26M
- ≻mate name *
- >mate position o
- ➢insert size o
- **>** sequence

AGGGCCGATCTTGGTGCCATCCAGGGGGCCTCTACAAGGAT AATCTGACCTGCTGAAGATGTCTCCAGAGACCTT

➤ base qualities

ECC@EEF@EB:EECFEECCCBEEEE;>5;2FBB@FBFEEFCF@FFFFCEFFFFEE>FFFFC=@A;@>1@6.+5/5

➤ tags MD:Z:75 NH:i:7 AS:i:-8 XS:A:-

BAM index file (.bai)



- BAM files can be sorted by chromosomal coordinates and indexed for efficient retrieval of reads for a given region.
- The index file must have a matching name. (e.g. reads.bam and reads.bam.bai)
- Genome browser requires both BAM and the index file.
- The alignment tools in Chipster automatically produce sorted and indexed BAMs.
- When you import BAM files, Chipster asks if you would like to preproces them (convert SAM to BAM, sort and index BAM).

Mapping quality

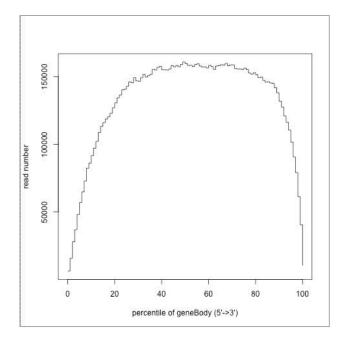


- Confidence in read's point of origin
- Depends on many things, including
 - ouniqueness of the aligned region in the genome
 - olength of alignment
 - o number of mismatches and gaps
- Should be expressed in Phred scores, like base qualities
 - $Q = -10 * log_{10}$ (probability that mapping location is wrong)
- Values differ in different aligners. E. g. unique mapping is
 - 6o in HISAT2
 - 255 in STAR
 - 50 in TopHat
 - https://sequencing.qcfail.com/articles/mapq-values-are-really-useful-but-their-implementation-is-a-mess/

Quality check after alignment with RseQC



- Checks coverage uniformity, read distribution between different genomic regions, novelty of splice junctions, etc.
- Takes a BAM file and a BED file
 - Chipster has BED files available for several organisms
 - You can also use your own BED if you prefer

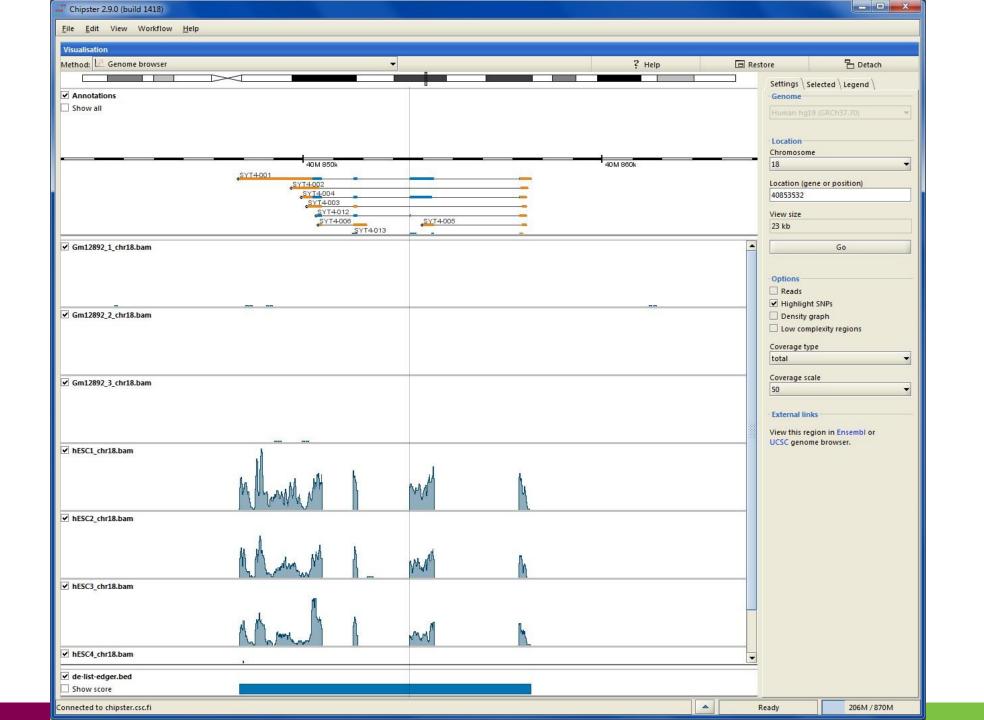


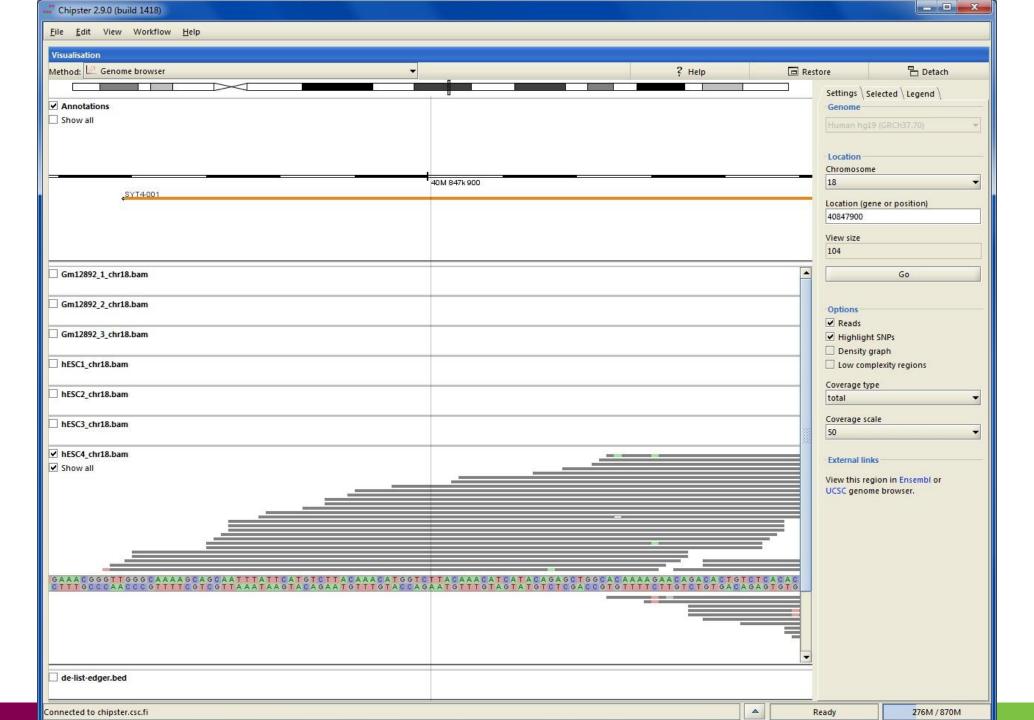
Total Reads	84808 116738			
Total Tags				
Total Assigned Tags	111352			
Group	Total_bases	Tag_count	Tags/Kb	
CDS_Exons	2211343	90961	41.13	
5'UTR_Exons	529860	1662	3.14	
3'UTR_Exons	1415234	12423	8.78	
Introns	25801210	5349	0.21	
TSS_up_1kb	1295771	31	0.02	
TSS_up_5kb	5332522	321	0.06	
TSS_up_10kb	8804879	584	0.07	
TES_down_1kb	1292506	217	0.17	
TES_down_5kb	5108821	344	0.07	
TES down 10kb	8282641	373	0.05	



Visualize alignments in genomic context: Chipster Genome Browser

- Integrated with Chipster analysis environment
- Automatic coverage calculation (total and strand-specific)
- Zoom in to nucleotide level
- Highlight variants
- Jump to locations using BED, GTF, VCF and tsv files
- View details of selected BED, GTF and VCF features
- Several views (reads, coverage profile, density graph)



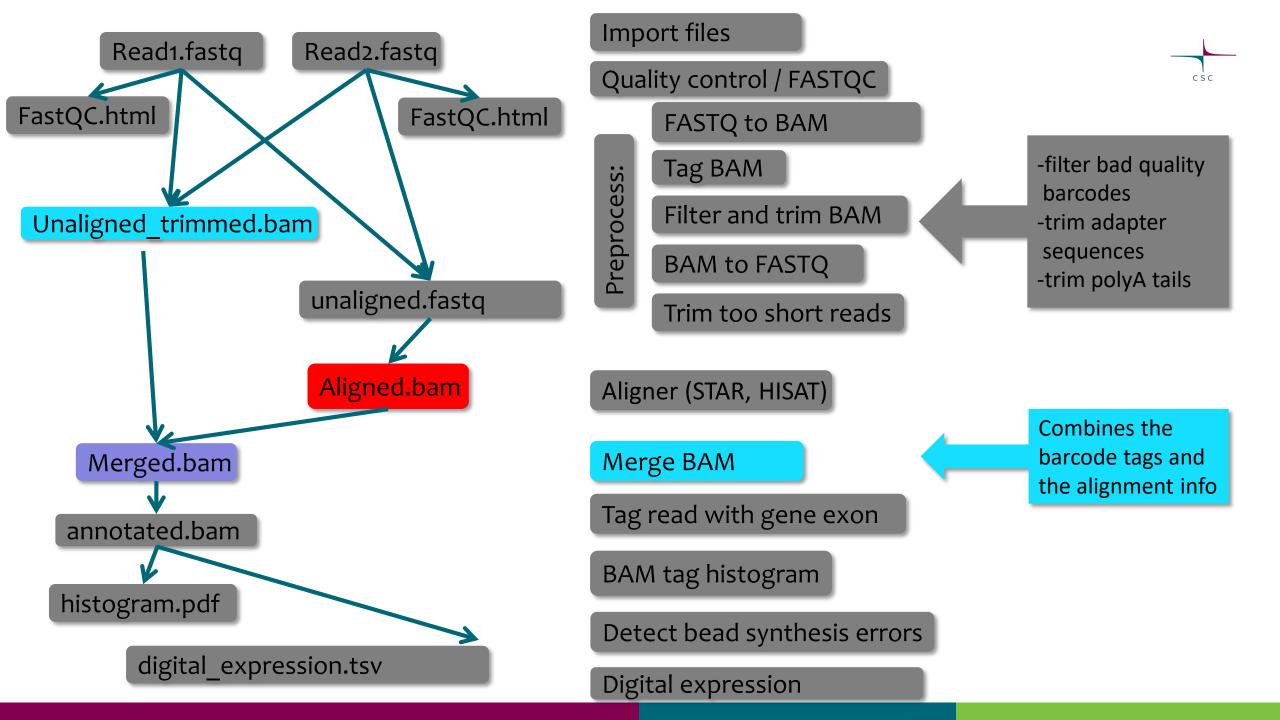


CSC



From FASTQ to expression matrix -preprocessing of DropSeq data

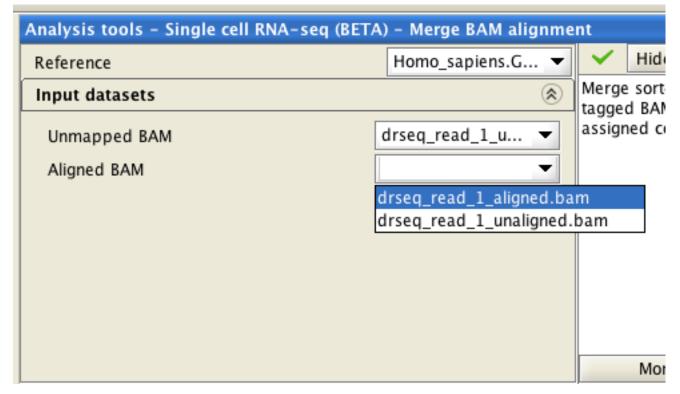
- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix





Merge BAM files

- Why was this again...?
 - o In the alignment, we lost the molecular and cellular tag information (because aligners only eat FASTQ format). We have that info in the unaligned BAM file, so now we combine the two.
- Before merging, the tool sorts the files in queryname order
- Secondary alignments are ignored!
- Remember to...
 - Choose the reference
 - O Make sure the files are correctly assigned!





From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

Annotaate



- GE-tag = name of the gene, if the read overlaps an exon
- XF-tag = location (intron, exon...)
- Choose the annotation file (GTF) OR use your own

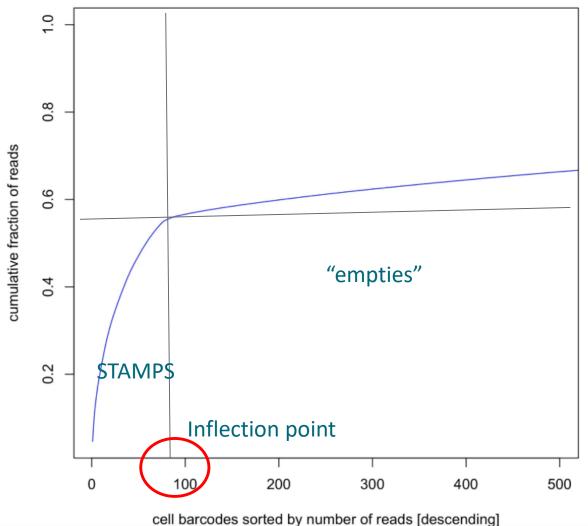


From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

Estimate the number of usable cells

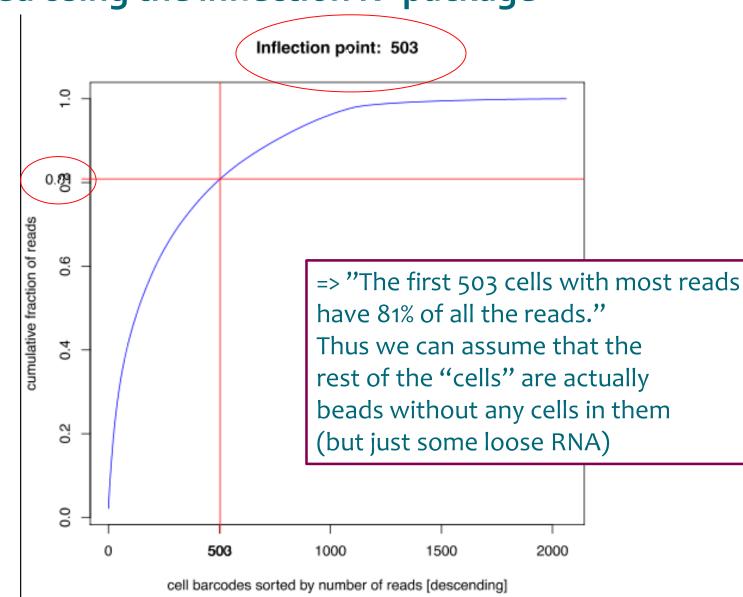
- How many cells do you want to keep?
- To estimate this: the inflection point
 - extract the number of reads per cell (barcode)
 - plot the cumulative distribution of reads
 - select the "knee" of the distribution (knee = inflection point)
 - The number of STAMPs (=beads exposed to a cell in droplets): cell barcodes to the left of the inflection point
 - Empties (=beads only exposed to ambient RNA in droplets): to the right of the inflection point





Inflection point also computed using the inflection R -package

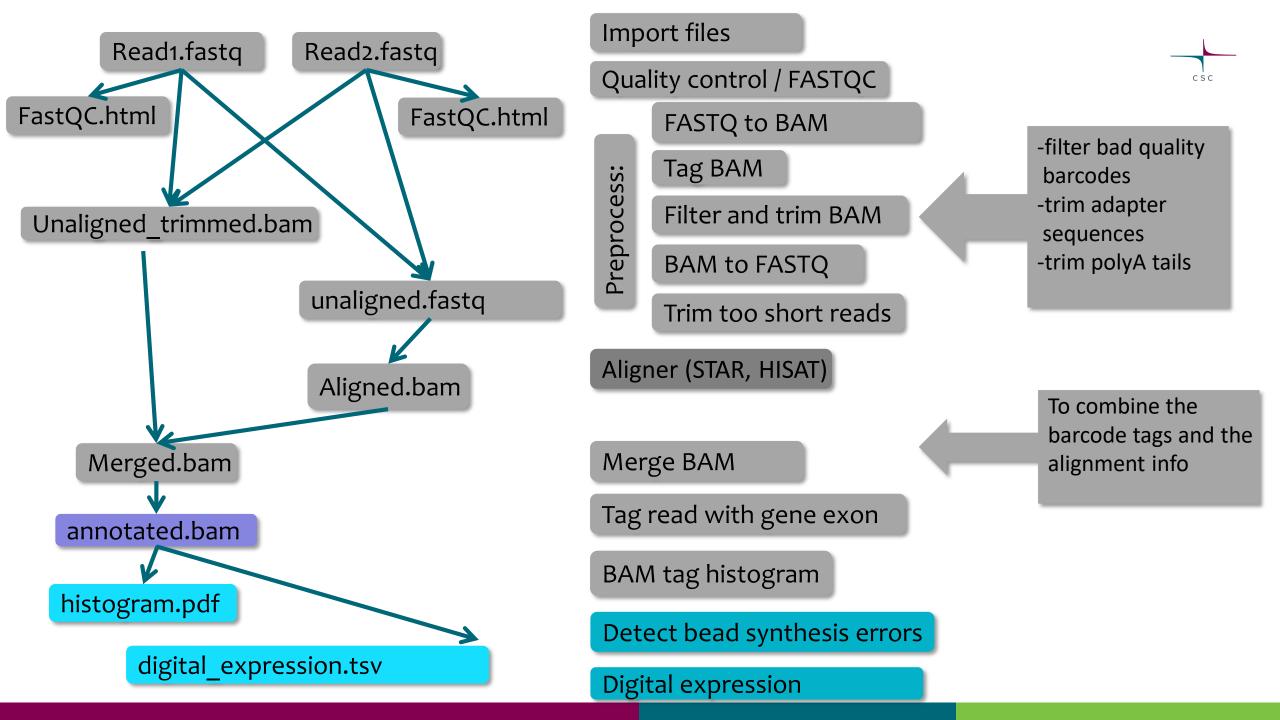
- Sometimes visual estimation can be tricky, so the tool also gives a numerical estimate
- Finds the inflection point using extreme distance estimator (ede) from inflection R-package
- Thank you DawitYohannes!





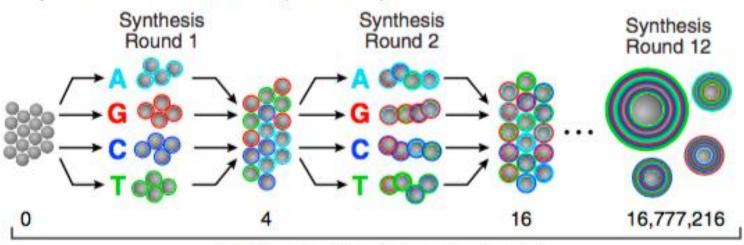
From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix



Detect bead synthesis errors

ATTC GAGT TAT? CAGC GTAAT TTTT Cell (12) UMI (8)



Number of unique barcodes in pool

- Known problem with ChemGenes beads: a percentage of beads did not undergo all twelve split-and-pool bases.
 - o A **mixed base** at base 12 (=actually 1st base of UMI) and a **fixed T** base at base 20 (=actually the 1st base of the polyT segment)
- DropSeq-tool *DetectBeadSynthesisErrors* fixes this:
 - The last base is trimmed off and all cell barcodes with identical sequence at the first 11 bases are merged
 - o If any other UMI base is fixed, the reads with that cell barcode are discarded
- Parameter: number of barcodes on which to perform the correction
 - o roughly 2 times the anticipated number cells (empirically found that this allows to recover nearly every defective cell barcode).

Error types

- SYNTHESIS_MISSING_BASE ~
 - o 1 or more bases missing from cell barcode => T's at the end of UMIs
 - Fix: insert an "N" (reading frame fixed) and merge. If more than 1 missing, discard these reads
- SINGLE_UMI_ERROR
 - At <u>each</u> position of the UMIs, single base appears in >80% of the UMIs for that cell.
 - Fix: cell barcodes with this property are dropped
- PRIMER_MATCH
 - Same as with SINGLE_UMI_ERROR, + the UMI matches one of the PCR primers
 - Fix: these barcodes are dropped
- OTHER
 - OUMIs are extremely skewed towards at least one base, but not at all 8 positions
 - Fix: these barcodes are dropped



• Error: AAAAACGTGGG-CAGCGTAATTT

Fixed: AAAAACGTGGGNCAGCGTAATTT

Synthesis statistics



- synthesis_stats.txt contains a bunch of useful information:
 - 1. CELL_BARCODE the 12 base cell barcode
 - 2. NUM_UMI the number of total UMIs observed
 - 3. FIRST_BIASED_BASE the first base position where any bias is observed (-1 for no detected bias)
 - 4. SYNTH_MISSING_BASE as 3 but specific to runs of T's at the end of the UMI
 - 5. ERROR_TYPE
 - 6. For bases 1-8 of the UMI, the observed base counts across all UMIs. This is a "|" delimited field, with counts of the A,C,G,T,N bases.
- synthesis_stats.summary.txt contains a histogram of the SYNTHESIS_MISSING_BASE errors*, as well as the counts of all other errors, the number of total barcodes evaluated, and the number of barcodes ignored.

*1 or more bases missing from cell barcode => T's at the end of UMIs

CELL_BARCODE TCATTTAGTCGA NUM_UMI 11832

FIRST_BIASED_BASE -1 -1 

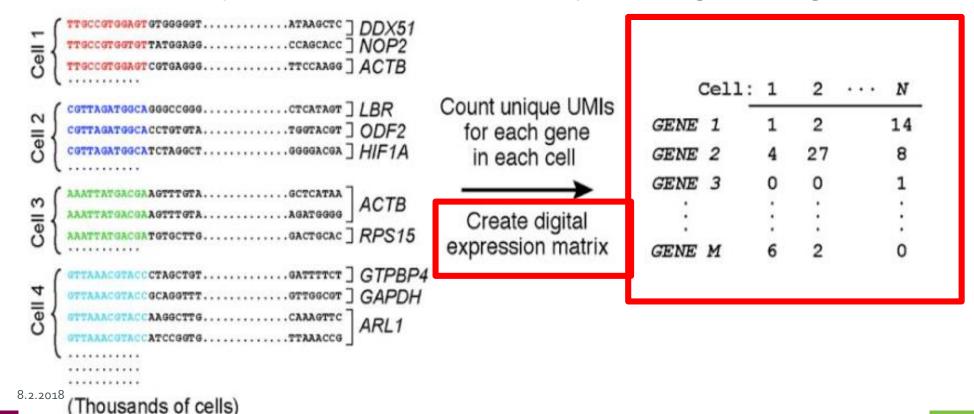
From FASTQ to expression matrix -preprocessing of DropSeq data

- Quality control of raw reads
- Preprocessing: tagging with barcodes & filtering
- Alignment to reference genome
- Merge BAM files
- Annotate with gene names
- Estimate the number of usable cells
- Detect bead synthesis errors
- Generate digital expression matrix

Digital expression matrix



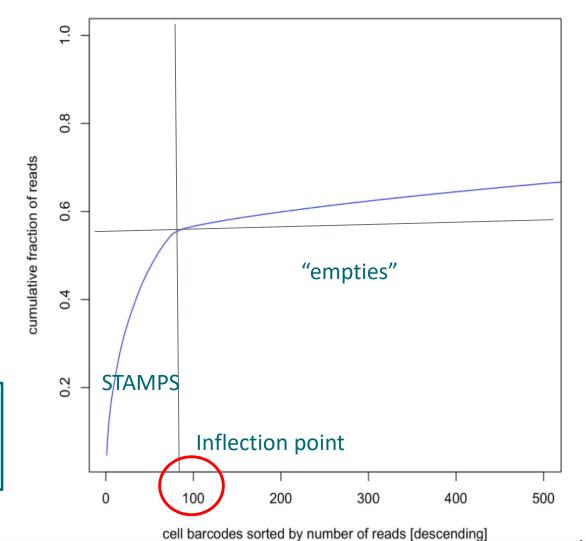
- To digitally count gene transcripts:
 - 1. a list of HQ UMIs in each gene, within each cell, is assembled
 - 2. UMIs within edit distance = 1 are merged together
 - 3. The total number of unique UMI sequences is counted
 - => this number is reported as the number of transcripts of that gene for a given cell



Filtering the data for DGE

- Why don't we just take all the cells?
 - othe aligned BAM can contain hundreds of thousands o cell barcodes
 - Some of them are "empties"
 - Some cell (barcode)s contain just handful of reads
 - o It is painful to deal with a huge matrix.
- Filtering based on:
 - Number of core barcodes (to X) xells with most reads
 - Minimum number of expressed genes per cell

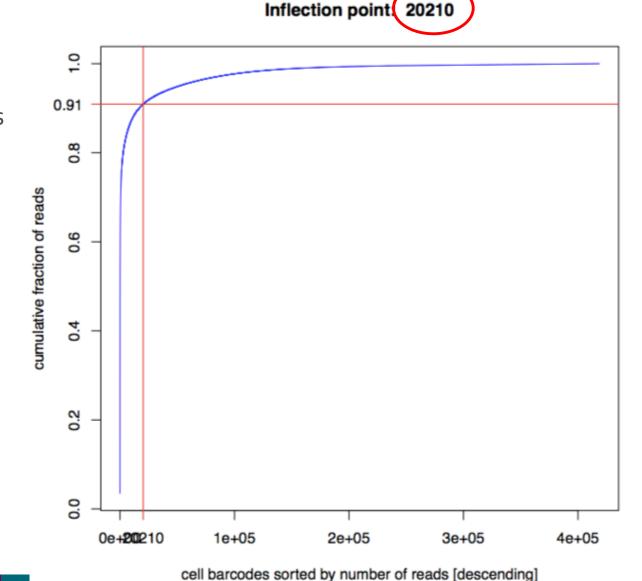
NOTE: You can always also choose to have a bigger number of cells and use that for further analysis (Seurat has it's own filtering tools) NOTE: "read" is a molecule here, which may or may not have the same/almost same UMI as another molecule

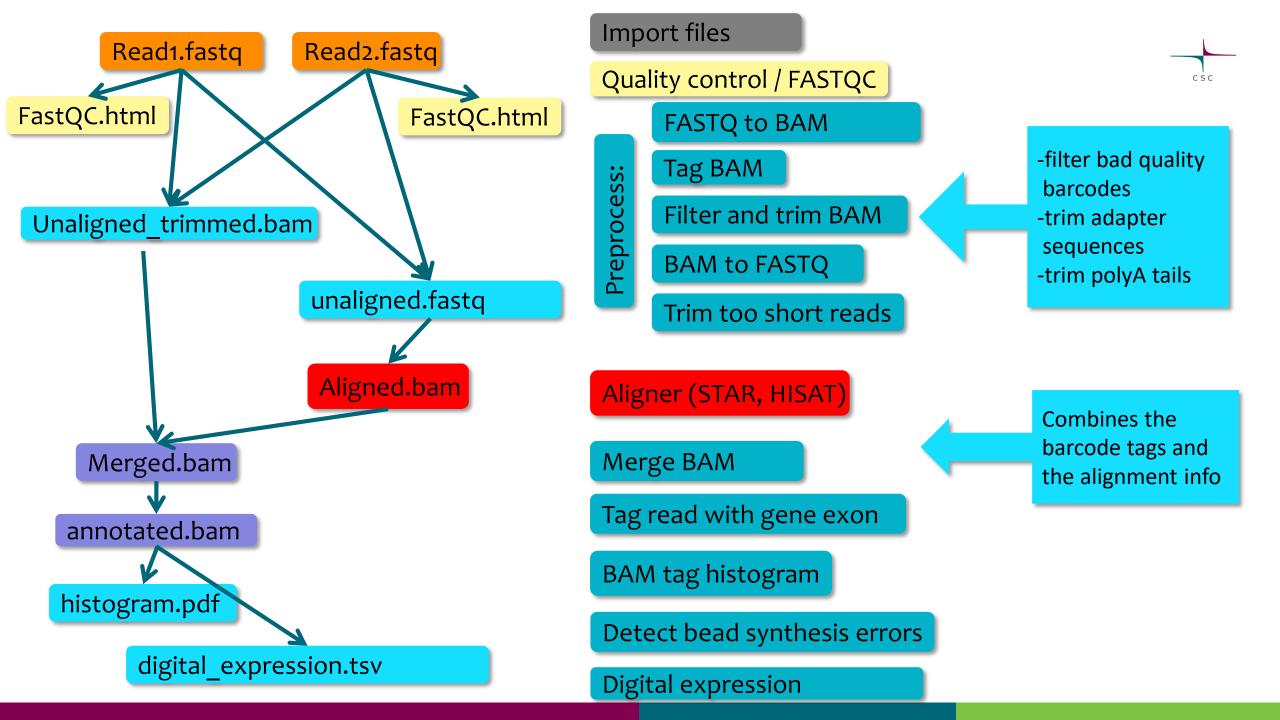




Filtering the data for DGE –huge number of reads

- What do you do if this happens?
 - Use the "minimum number of expressed genes per cell" as a requirement
 - O Do you know how many cells to expect? Use that number!
 - Take a large number of cells and filter in the Seurat tools
- Why this happens?







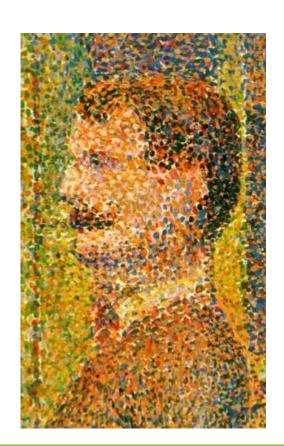
Clustering analysis with Seurat tools



Seurat

http://satijalab.org/seurat

- Seurat combines dimensionality reduction and graph-based partitioning algorithms for unsupervised clustering of single cells.
- The approach can be described briefly:
 - 1. Identification of highly variable genes
 - 2. Linear dimensionality reduction (PCA, principal component analysis) on variable genes
 - 3. Determine significant principal components
 - 4. Graph based clustering to classify distinct groups of cells
 - 5. Non-linear dimensional reduction (t-SNE, t-Distributed Stochastic Neighbor Embedding) for cluster visualization
 - 6. Marker gene discovery, visualization, and downstream analysis





Clustering analysis with Seurat tools



- Setup & preprocessing (filtering, log normalization)
- Quality control
- Filter cells
- Regress unwanted sources of variation
- Detect variable genes
- Linear dimensional reduction (PCA)
- Determine statistically significant principle components (=identify the true dimensions of the data)
- Cluster the cells
- Non-linear dimensional reduction (tSNE)
- Find differentially expressed genes (biomarkers for the clusters)

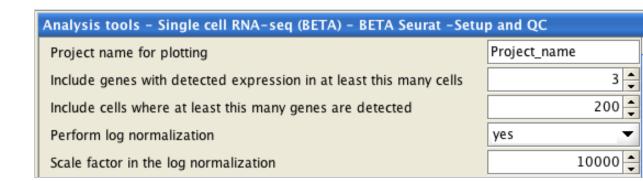
- Setup & QC
- Filtering, regression
 & detection of
 variable genes
- Linear dimensional reduction (PCA)
- Clustering the cells

Visualize biomarkers

Setting up Seurat object



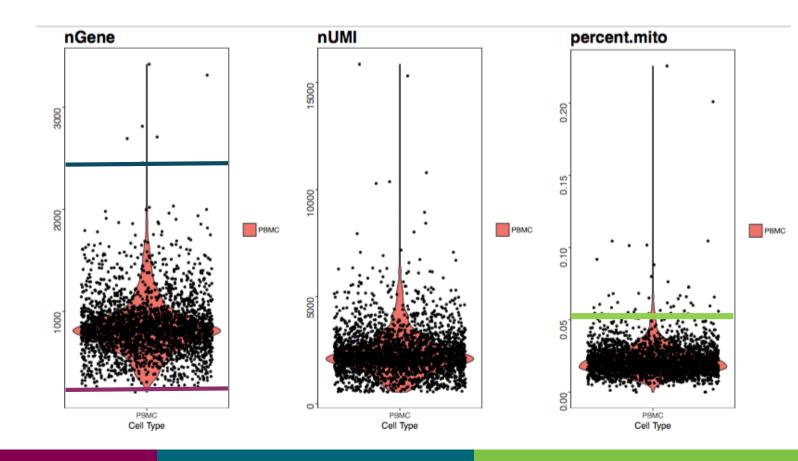
- In these tools, we are working with an R object (.Robj)
 - Can't be opened in Chipster, can be exported & imported to R
- Import one of the following
 - Tar package of three 10X Genomics output files
 - DGE matrix from the DropSeq tools
 - Check that the input file is correctly assigned!
- Log normalization
 - o normalizes the gene expression measurements for each cell by the total expression, multiplies this by a scale factor (10,000 by default), and log transforms the result
- Filtering
 - Keep genes which are detected at least in X cells
 - Keep cells where at least Y genes are detected
- Give a name for the project (used in some plots)





Quality control & filtering for empties, multiplets and broken cells

- Empty (no cell in droplet) →
 low gene count (<200)
- Multiplet (more than one cell in droplet) → large gene count (>2500)
- Broken cell in droplet →
 large percentage of
 mitochondrial transcripts
 (>5%)





Clustering analysis with Seurat tools

- Setup & preprocessing (filtering, log normalization)
- Quality control
- Filter cells
- Regress unwanted sources of variation
- Detect variable genes
- Linear dimensional reduction (PCA)
- Determine statistically significant principle components (=identify the true dimensions of the data)
- Cluster the cells
- Non-linear dimensional reduction (tSNE)
- Find differentially expressed genes (biomarkers for the clusters)

- Setup & QC
- Filtering, regression
 & detection of
 variable genes
- Linear dimensional reduction (PCA)
- Clustering the cells

Visualize biomarkers

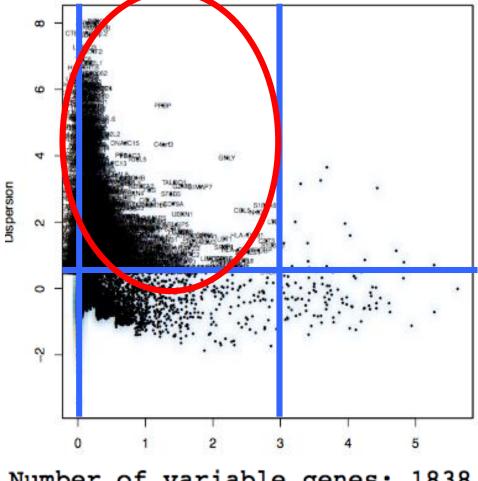


Remove unwanted sources of variation

- Single cell data typically contains 'uninteresting' variation
 - otechnical noise
 - o batch effects
 - o cell cycle stage, etc
- Removing this variation improves downstream analysis
- Seurat constructs linear models to predict gene expression based on user-defined variables
 - o batch, cell alignment rate, no of detected molecules per cell, mitochondrial transcript percentage
 - Seurat regresses the given variables individually against each gene, and the resulting residuals are scaled
 - oscaled z-scored residuals of these models are used for dimensionality reduction and clustering
 - o Chipster currently uses the no of detected molecules per cell and mitochondrial transcript percentage

Detection of variable genes

- Downstream analysis focuses on highly variable genes
- Seurat finds them in the following way
 - calculate the average expression and dispersion for each gene
 - place genes into bins based on expression
 - calculating a z-score for dispersion within each bin
- Set the parameters to mark visual outliers on the dispersion plot
 - o exact parameter settings may vary based on the data type, heterogeneity in the sample, and normalization strategy
 - The default parameters are typical for UMI data that is normalized to a total of 10 000 molecules



Number of variable genes: 1838

Bottom cutoff on x-axis for identifying variable genes	0.0125
Top cutoff on x-axis for identifying variable genes	3.0
Bottom cutoff on y-axis for identifying variable genes	0.5



Clustering analysis with Seurat tools

- Setup & preprocessing (filtering, log normalization)
- Quality control
- Filter cells
- Regress unwanted sources of variation
- Detect variable genes
- Linear dimensional reduction (PCA)
- Determine statistically significant principle components (=identify the true dimensions of the data)
- Cluster the cells
- Non-linear dimensional reduction (tSNE)
- Find differentially expressed genes (biomarkers for the clusters)

- Setup & QC
- Filtering, regression
 & detection of
 variable genes
- Linear dimensional reduction (PCA)
- Clustering the cells

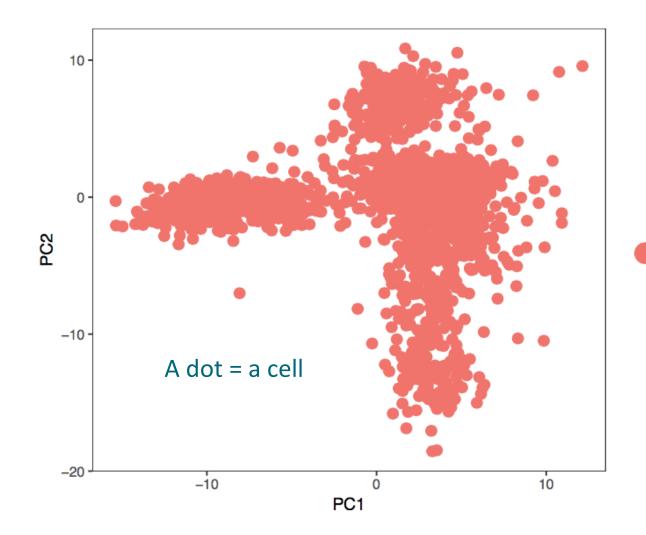
Visualize biomarkers

97



Linear dimensionality reduction (PCA) on variable genes

- Principal Component Analysis = PCA
- Reduce the numerous, possibly correlating variables (=counts for each gene) into a smaller number of linearly uncorrelated dimensions (=principal components)
- Essentially, each PC represents a robust 'metagene' — a linear combination of hundreds to thousands of individual transcripts





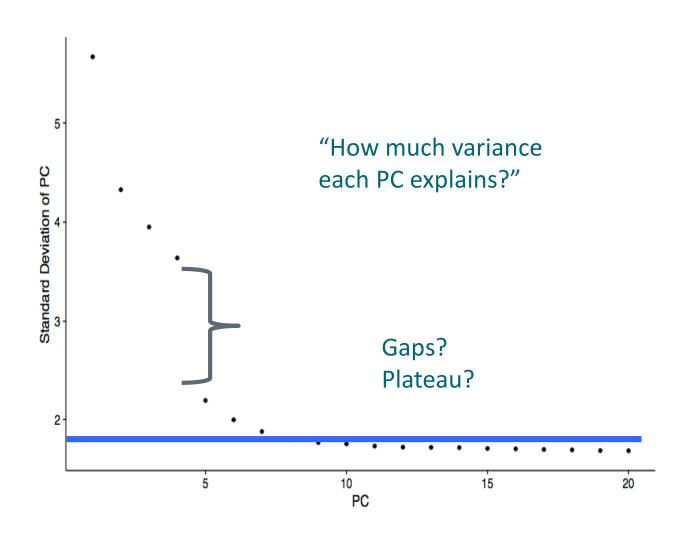
Determine significant principal components

- A key step to this clustering approach involves selecting a set of principal components (PCs) for downstream clustering analysis
- However, estimating the true dimensionality of a dataset is a challenging and common problem in machine learning.
- The tool provides a couple plots to aid in this:
 - Elbow plot
 - PCHeatmap



Determine significant principal components 1: Elbow plot

• The **elbow** in the plot tends to reflect a transition from informative PCs to those that explain comparatively little variance.

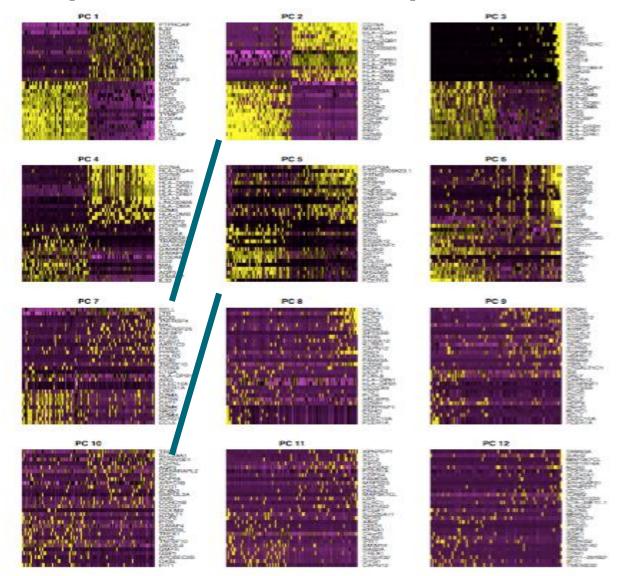




Determine significant principal components 2: PCHeatmap

• Displays the extremes across both genes and cells, and can be useful to help exclude PCs that may be driven primarily by ribosomal/mitochondrial or cell cycle genes.

"Is there still a difference between the extremes?"





Clustering analysis of with Seurat tools

- Setup & preprocessing (filtering, log normalization)
- Quality control
- Filter cells
- Regress unwanted sources of variation
- Detect variable genes
- Linear dimensional reduction (PCA)
- Determine statistically significant principle components (=identify the true dimensions of the data)
- Cluster the cells
- Non-linear dimensional reduction (tSNE)
- Find differentially expressed genes (biomarkers for the clusters)

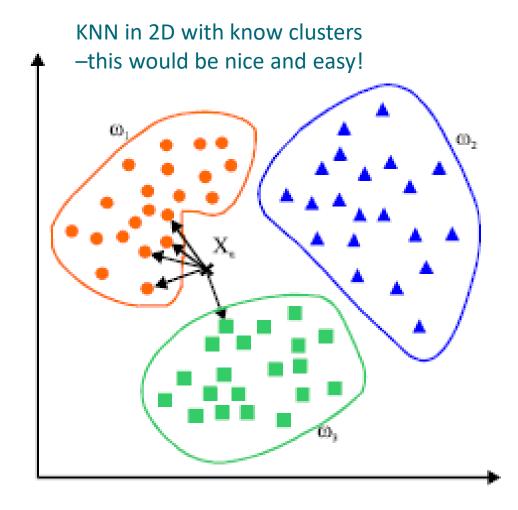
- Setup & QC
- Filtering, regression
 & detection of
 variable genes
- Linear dimensional reduction (PCA)
- Clustering the cells

Visualize biomarkers



Clustering -why is it so tricky?

- Need to use **unsupervised** methods (we don't know beforehand, how many clusters there are)
- Our data is big and complex:
 - Lots of cells
 - Lots of dimensions (=genes)
 - Lots of noise (both technical and biological)
- ...which is why:
 - o the algorithm is a bit tricky, and
 - We reduce and select the dimensions to use in clustering (= we do the PCA and select X components to use, instead of using all thoudans of genes)



Graph based clustering to classify distinct groups of cells

• Seurat clustering = similar to

o SNN-Cliq (C. Xu and Su, Bioinformatics 2015) and

PhenoGraph (Levine et al. Cell 2015)

• ...which are graph based methods:

Identify k-nearest neighbours of each cell

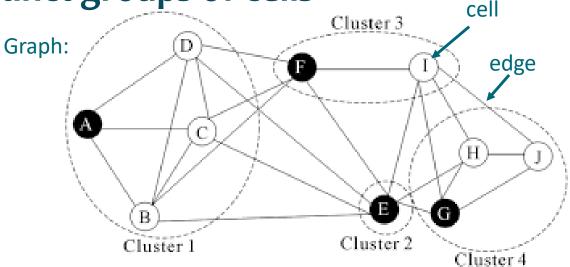
o Distance measure: Euclidean + Jaccard distance

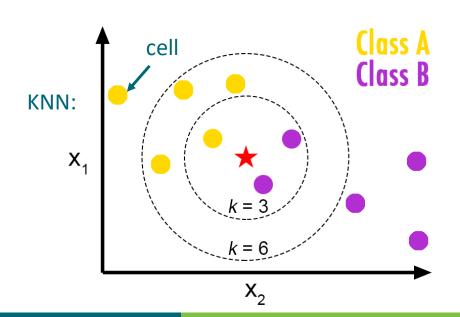
 Calculate the number of Shared Nearest Neighbours (SNN) between each pair of cells

3. Build the graph: add an **edge** between cells, if they have at least one SNN

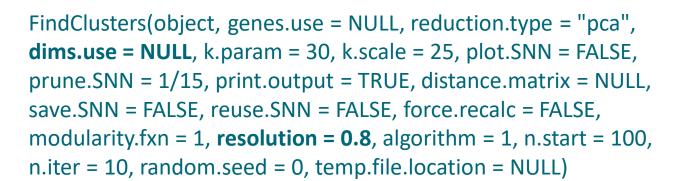
4. Clusters: group of cells with many edges between them

 Smart Local Moving algorithm (SLM, Blondel et al., Journal of Statistical Mechanics)





Clustering parameters



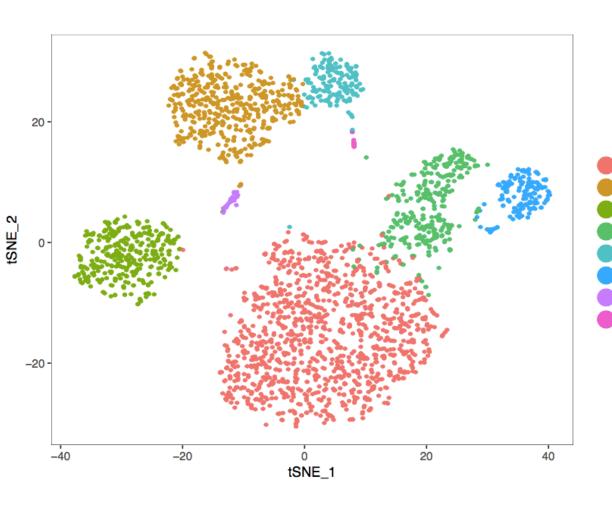


- Lots of parameters!
- **Dims.use** = number of principal components
- **Resolution** = "granularity":
 - o increased values lead to a greater number of clusters.
 - Values o.6-1.2 typically returns good results for single cell datasets of around 3K cells
 - Optimal resolution often increases for larger datasets
 - (If you get very few or very many clusters, try adjusting)

Analysis tools - Single cell RNA-seq (BETA) - BETA Seurat -Clustering				
Number of principal components to use	10			
Resolution for granularity	0.6			
Cluster hiemarker gane has to be expressed in at least this fraction	0.25			

Non-linear dimensional reduction (t-SNE) for cluster visualization





- T-SNE = t-distributed Stochastic Neighbor Embedding
 - Non-linear algorithm, different transformations to different regions
- Output
 Why do we now use t-SNE, why not the good old PCA?
 - We give the selected PCA "pseudo-genes" to t-SNE
 - T-SNE also reduces dimensions, but it is "more faithful to the original data"
 - PCA can find clusters too, but t-SNE does just that –
 it reduced the dimensions so that the clusters
 become visible
- Good text about readingT-SNE's: https://distill.pub/2016/misread-tsne/



Finding differentially expressed genes (biomarkers for the clusters)

• Parameters:

- o min.pct : requires a gene to be <u>detected</u> at least THIS minimum percentage in <u>either</u> of the two groups of cells (default: 0.25)
- o **thresh.test**: requires a gene to be <u>differentially expressed</u> (on average) by THIS amount between the two groups (default: 0.25)
- You can set both of these to zero, but with a dramatic increase in time (this will test a large number of genes that are unlikely to be highly discriminatory)

Options for tests: bimod, roc, Students t-test, Tobit-test, Poisson, negative-binomial distribution

NOTE: We are now comparing a cluster to all other cells. So for example cluster 1 vs all

others.

Analysis tools - Single cell RNA-seq (BETA) - BETA Seurat -Clustering					
Cluster biomarker gene has to be expressed in at least this fractio	0.25				
Differential expression threshold for a cluster biomarker gene	0.25				
Which test to use for finding marker genes	bimod ▼				
Only positive changes	FALSE ▼				



Markers for a particular cluster

• You can filter the resulting list to get only the biomarkers for a certain cluster:

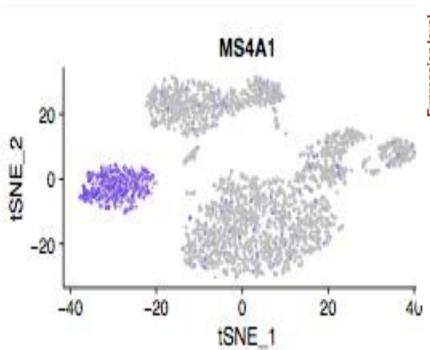
Utilities / Filter table by column value			
Column to filter by	cluster		
Does the first column have a title	no		
Cutoff	2.0		
Filtering criteria	equal-to		

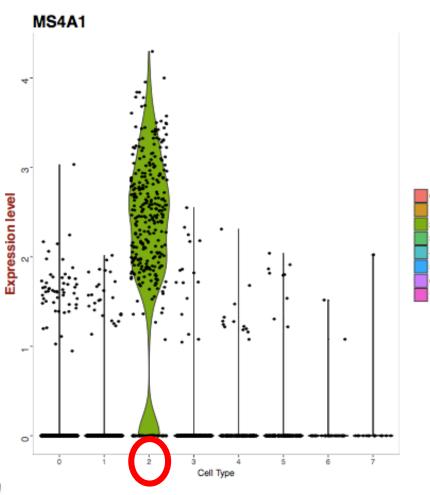
Showing 320 rows	of 320 and all 7 columns				
	p_val	avg_diff	pct.1	pct.2	cluster
TCL1A	1.030004e-191	2.4922363	0.624	0.022	2
HLA-DRB12	1.030181e-176	1.4341113	0.985	0.381	2
LINC00926	2.78283e-175	1.9663931	0.566	0.009	2
HLA-DPA12	6.494113e-166	1.3878677	0.98	0.421	2
VPREB3	1.41582e-150	1.6848535	0.493	0.006	2
S100A41	1.067153e-145	-2.110098	0.367	0.879	2
BANK1	1.547914e-111	1.4315893	0.481	0.037	2
HLA-DRB52	4.812796e-111	1.3080622	0.854	0.276	2
C1004C1	0.353060- 100	1.6531713	0.350	0.046	2

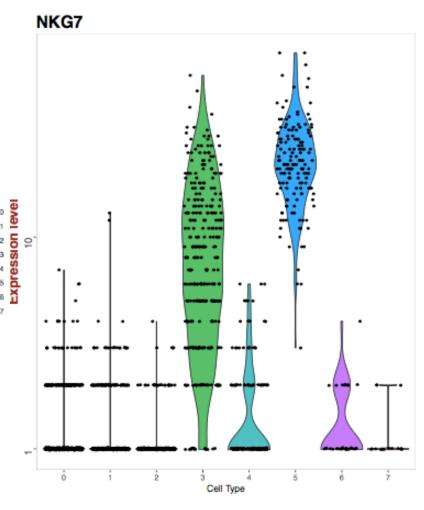


Visualize biomarkers

• Select a marker gene from the lists









Clustering analysis with Seurat tools

- Setup & preprocessing (filtering, log normalization)
- Quality control
- Filter cells
- Regress unwanted sources of variation
- Detect variable genes
- Linear dimensional reduction (PCA)
- Determine statistically significant principle components (=identify the true dimensions of the data)
- Cluster the cells
- Non-linear dimensional reduction (tSNE)
- Find differentially expressed genes (biomarkers for the clusters)

- Setup & QC
- Filtering, regression
 & detection of
 variable genes
- Linear dimensional reduction (PCA)
- Clustering the cells

Visualize biomarkers