

GAME-CIBOG WORKSHOP

RNA-seq
scRNA-seq

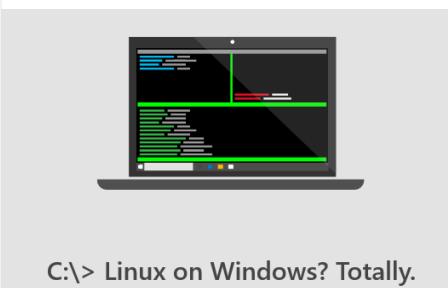
CURRICULUM VITAE

Markus Joppich

- 
- 1989 • Born in Bonn
 - 2012 • B.Sc. Computer Science at **RWTH AACHEN UNIVERSITY**
 - iGEM  Aachen 2014 Team Member
 - Best Measurement Project & Best Supporting Software
 - M.Sc. Computer Science at **RWTH AACHEN UNIVERSITY**
 - 2015 • Started PhD at     
 - 2018 • iGEM  Munich 2018 Supervisor
 - Finalist, Best Manufacturing Project, Best Software
 - 2021 • Defended PhD in November 2021
 - PostDoc + Habilitation on scRNA-seq, causal networks



 [mjoppich](#)
 [@mjoppich](#)



C:\> Linux on Windows? Totally.

Markus
Joppich

Vascular neutrophilic inflammation and immunothrombosis distinguish severe COVID-19 from influenza pneumonia

Leo Nicolai^{1,2,3} | Alexander Leunig^{1,2} | Sophia Brambs¹ | Rainer Kaiser^{1,2,3} | Markus Joppich⁴ | Marie-Louise Hoffknecht¹ | Christoph Gold¹ | Anouk Engel¹ | Vivien Polewka¹ | Maximilian Muenchhoff^{3,5,6} | Johannes C. Hellmuth^{3,7,8} | Adrian Ruhle^{3,5} | Stephan Ledderose⁹ | Tobias Weinberger^{1,2,3} | Heiko Schulz⁹ | Clemens Scherer^{1,2,3} | Martina Rudelius⁹ | Michael Zoller¹⁰ | Oliver T. Keppler^{3,5,6} | Bernhard Zwölfer¹⁰ | Michael von Bergweilt-Baldon^{7,8} | Stefan Kääb^{1,2,3} | Ralf Zimmer⁴ | Roman D. Bülow¹¹ | Saskia von Stillfried¹¹ | Peter Boor^{11,12} | Steffen Massberg^{1,2,3} | Kami Pekayvaz^{1,2,3} | Konstantin Stark^{1,2,3}

RESEARCH ARTICLE

Self-sustaining IL-8 loops drive a prothrombotic neutrophil phenotype in severe COVID-19

Rainer Kaiser,^{1,2} Alexander Leunig,^{1,2} Kami Pekayvaz,^{1,2} Oliver Popp,^{4,5} Markus Joppich,⁴ Vivien Polewka,¹ Raphael Escaig,¹ Afra Anjum,¹ Marie-Louise Hoffknecht,¹ Christoph Gold,¹ Sophia Brambs,¹ Anouk Engel,¹ Sven Stockhausen,¹ Viktoria Knottenberg,¹ Anna Titova,¹ Mohamed Hajj,^{4,5} Clemens Scherer,^{1,2,3} Maximilian Muenchhoff,^{3,8} Johannes C. Hellmuth,^{3,9} Kathrin Saar,^{4,5} Benjamin Schubert,^{10,11,12} Anne Hilgendorff,^{12,13,14} Christian Schulz,^{1,2} Stefan Kääb,^{1,2,3} Ralf Zimmer,⁴ Norbert Hübner,^{4,5,6} Steffen Massberg,^{1,2,3} Philipp Mertins,^{4,5} Leo Nicolai,^{1,2,3} and Konstantin Stark^{1,2,3}

Title Protective immune trajectories in early viral containment of non-pneumonic SARS-CoV-2 infection

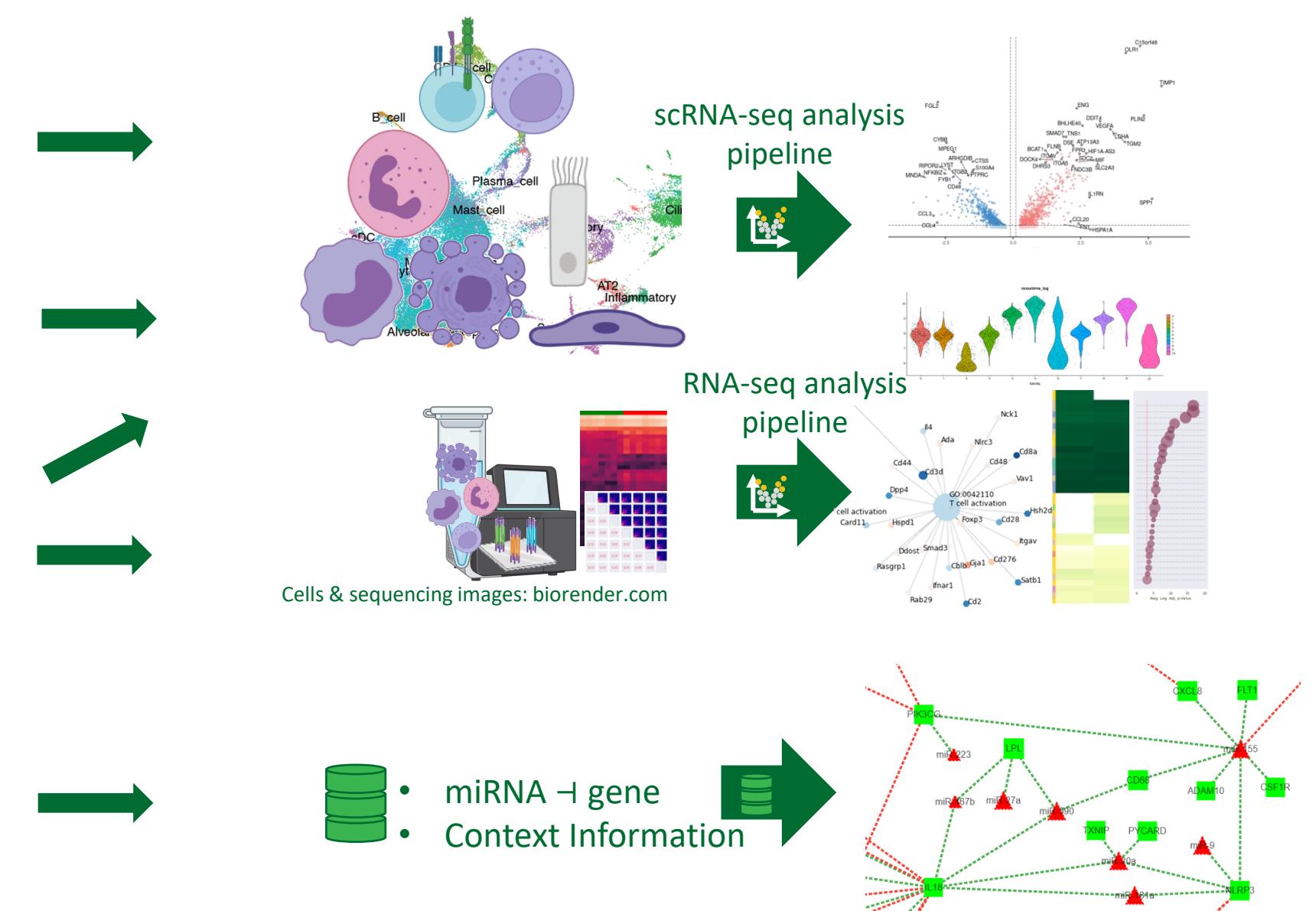
Author list
 Kami Pekayvaz^{1,2,*}, Alexander Leunig^{1,2,#}, Rainer Kaiser^{1,2}, Markus Joppich³, Sophia Brambs¹, Aleksandar Janic⁴, Oliver Popp⁵, Daniel Nixdorf^{6,7}, Valeria Fumagalli^{8,9}, Nora Schmidt¹⁰, Vivien Polewka¹, Afra Anjum¹, Viktoria Knottenberg¹, Luke Eivers¹, Lucas E. Wange⁴, Christoph Gold¹, Marieluise Kirchner⁵, Maximilian Muenchhoff^{11,12,13}, Johannes C. Hellmuth^{7,13}, Clemens Scherer^{1,2,13}, Raquel Rubio-Acerro^{12,14}, Tabita Eser^{12,14}, Flora Deák^{12,14}, Kerstin Puchinger¹⁴, Niklas Kuhl^{15,16}, Andreas Linder^{15,16}, Kathrin Saar⁵, Lukas Tomas^{1,2}, Christian Schulz^{1,2}, Andreas Wieser^{12,14}, Wolfgang Enard⁴, Inge Kroidl^{12,14}, Christof Geldmacher^{12,14}, Michael von Bergweilt-Baldon^{7,13}, Oliver T. Keppler^{11,12,13}, Mathias Munschauer¹⁰, Matteo Iannaccone^{8,9,17}, Ralf Zimmer³, Philipp Mertins⁵, Norbert Hübner⁵, Michael Hoelscher^{12,14}, Steffen Massberg^{1,2}, Konstantin Stark^{1,2,§,*}, Leo Nicolai^{1,2,§,*}

Using Context-Sensitive Text Mining to Identify miRNAs in Different Stages of Atherosclerosis

Markus Joppich¹ Christian Weber² Ralf Zimmer¹¹Department of Informatics, LFE Bioinformatics, Ludwig-Maximilians-Universität München, Munich, Germany²Institute for Cardiovascular Prevention, Ludwig-Maximilians-Universität München, Amalienstr. 17, Munich, Bavaria 80333, Germany (e-mail: joppich@bio.lmu.de).

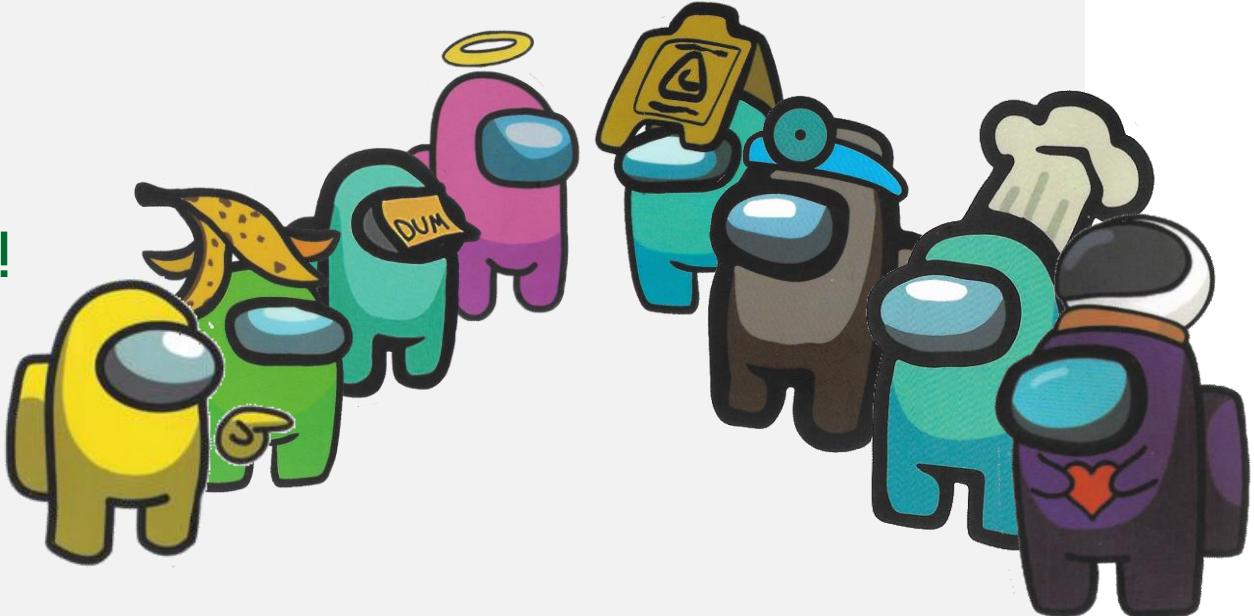
Thromb Haemost 2019;119:1247–1264.

SEQU-INTO: Early detection of impurities, contamination and off-targets (ICOs) in long read/MinION sequencing

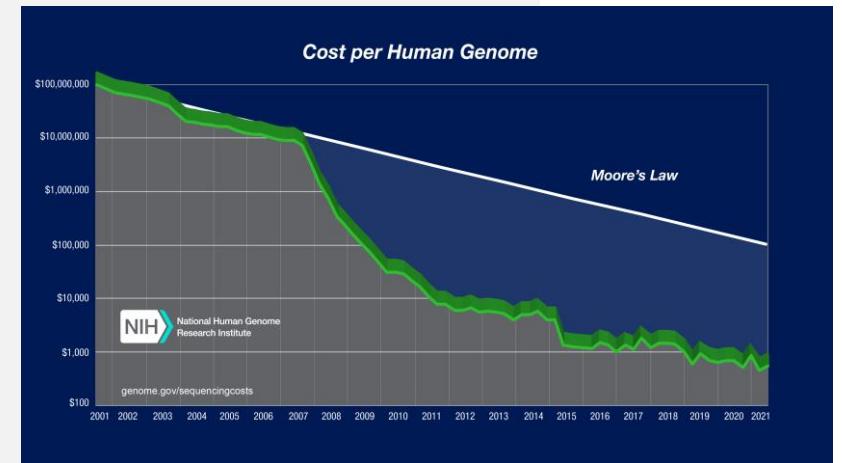
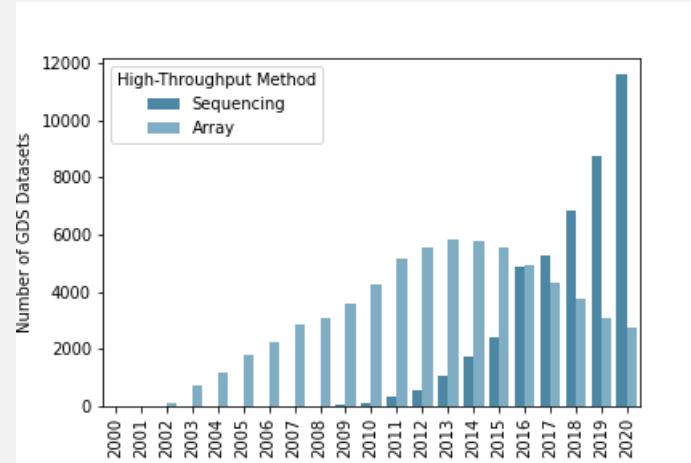
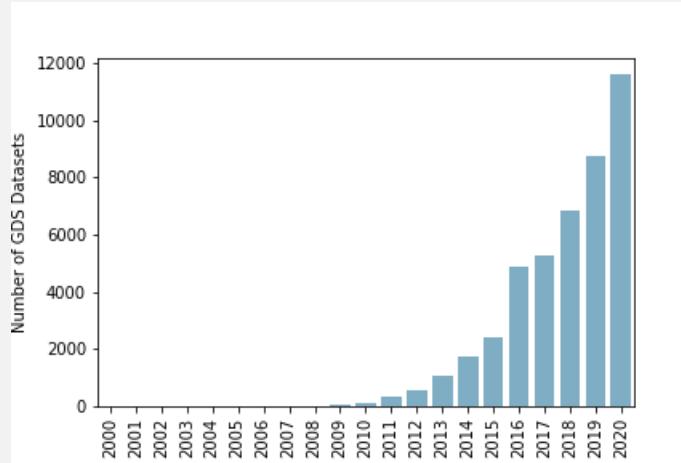
Markus Joppich^{a,1}, Margaryta Olenchuk^{a,1}, Julia M. Mayer^{a,1}, Quirin Emslander^b, Luisa F. Jimenez-Soto^c,Ralf Zimmer^{a,*}^aLFE Bioinformatics, Department of Informatics, Ludwig-Maximilians-Universität München, 80333 München, Germany^bPhysics of Synthetic Biological Systems, Physics Department, Technische Universität München, 85748 Garching, Germany^cWalther Straub Institute for Pharmacology and Toxicology, Ludwig-Maximilians-Universität München, Goethestrasse 33, 80336 München, Germany

ABOUT YOU

- Who are you?
- What are you studying?
- Why are you here?
- What is your experience?
- What are your expectations?
- Collect answers in your team!
 - We'll meet here in ~10min



WORKSHOP ON HIGH-THROUGHPUT SEQUENCING

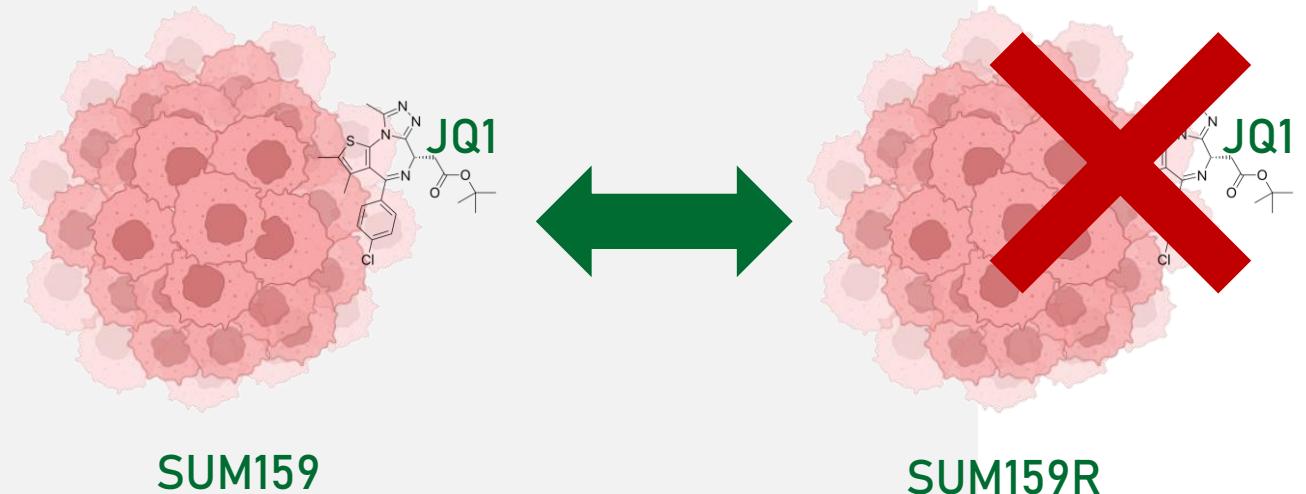


- Sequencing experiments are common and performed on a routinely basis
 - Costs are continuously decreasing
 - Sequencing data is more flexible than array experiments

Resource

Synthetic Lethal and Resistance Interactions with BET Bromodomain Inhibitors in Triple-Negative Breast Cancer

Shaokun Shu,^{1,7,18,19} Hua-Jun Wu,^{2,11,12,18} Jennifer Y. Ge,^{1,2,13} Rhamy Zeid,^{1,20} Isaac S. Harris,^{8,14} Bojana Jovanović,^{1,7,16} Katherine Murphy,¹ Binbin Wang,^{2,21} Xintao Qiu,^{1,5} Jennifer E. Endress,^{8,14} Jaime Reyes,¹ Klothilda Lim,^{1,5} Alba Font-Tello,^{1,5} Sudeepa Syamala,^{1,5} Tengfei Xiao,² Chandra Sekhar Reddy Chilamakuri,¹⁵ Evangelia K. Papachristou,¹⁵ Clive D'Santos,¹⁵ Jayati Anand,¹ Kunihiro Hinohara,^{1,7} Wei Li,^{2,11,22} Thomas O. McDonald,^{2,6,11,12} Adrienne Luoma,^{3,9} Rebecca J. Modiste,¹⁷ Quang-De Nguyen,¹⁷ Brittany Michel,⁴ Paloma Cejas,^{1,5} Cigall Kadoc,^{4,10,16} Jacob D. Jaffe,¹⁶ Kai W. Wucherpfennig,^{3,9} Jun Qi,¹ X. Shirley Liu,^{2,5,11} Henry Long,^{1,5} Myles Brown,^{1,5,7,14} Jason S. Carroll,¹⁵ Joan S. Brugge,^{8,14} James Bradner,^{1,7,23} Franziska Michor,^{2,6,11,12,14,16,*} and Komelia Polyak^{1,5,6,7,14,16,24,*}



- triple-negative breast cancer (TNBC)
 - Does not express estrogen receptor (ESR1), progesterone receptor (PGR), HER2/neu (ERBB2)
 - Cell line SUM159
 - metastatic human breast cancer cell line derived from a patient with anaplastic breast carcinoma
 - Control (SUM159) and BET bromodomain inhibitors (BBDIs) resistant (SUM159R)
 - JQ1 is inhibitor of BET family of bromodomain proteins

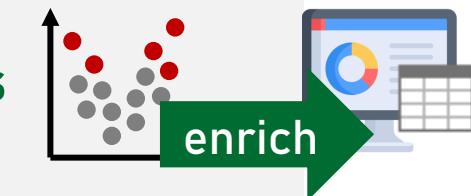
<https://www.sciencedirect.com/science/article/pii/S1097276520302690>

INTEGRATIVE BIOINFORMATICS

- Integrative Bioinformatics
 - Subdiscipline in Bioinformatics
 - Uses tools of computer science and electronic infrastructure applied to Bio*



- Uses existing tools/methods/resources to gain new insights



COMPUTER RESOURCES

- Computer
- Some kind of Linux
 - You use Linux? → you're fine (but need Win/Mac to run Loupe browser)
 - You use Windows? → use Windows Subsystem for Linux → you're fine



- You use Mac OS? → you're mostly fine, but I can't help you
 - My Mac experience dates back to Mac OS X Leopard!

SOFTWARE RESOURCES

- You need some R version

```
mjoppich@spectre3:/mnt/c/Users/mjopp$ R --version
R version 4.0.3 (2020-10-10) -- "Bunny-Wunnies Freak Out"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)
```

- You may want to use jupyter+R kernel

- `install.packages("devtools")`
- `devtools::install_github("IRkernel/IRkernel")`
- `IRkernel::installspec()`

- Some kind of IDE or jupyter-lab in the browser

- Big fan of Visual Code
 - View PDF, PNG, JPEG
 - Supports jupyter notebooks and plots



```
options(jupyter.plot_mimetypes = 'image/png')
options(repr.plot.width = 1, repr.plot.height = 0.75, repr.plot.res = 100)
```

TIMELINE

GAME-CIBoG kickoff workshop

	Day1 (22/01/22)	Day2 (23/01/22)	Between Day2-3	Day3 (29/01/22)	Day4 (30/01/22)
NU staff	Opening remarks Dr. Katusno (10 min) Ice-breaking in each group Dr. Bustos (30min) Introduction/ Presenting the assignment Dr.Hinohara / Dr. Kato (30 min)			Organizing the data and brainstorming between students (3 hour) give advanced advice for the interpretation (Dr.Hinohara/Dr. Kato) → let every team to do interpretation of analyzed data	Presenting the answer from each group (10min per group 1 hour) Presenting the answer Dr.Hinohara / Dr.Kato (1 hour) Awards ceremony (5min) Closing remarks Dr. Katusno (10 min)
LMU staff	RNA-seq lecture (1) Dr. Hinske (1 hour)	RNA-seq lecture (2) Dr. Joppich (3 hours)			
Students	24 students from 5 Unive Divde into 6 groups	Goal: *know how to analyze single-cell RNA-seq *start analysis by themselves	Goal: *finish the analysis by themselves *start interpretation of the analyzed data	*analyze by themselves *interpretation of analyzed data	Goal: Interpret the analyzed data present their result

RNA-seq data analysis Interpretation of analyzed data

GIT REPOSITORY

mjoppich / game_cibog_2022 Private

Code Issues Pull requests Actions Projects Wiki Security Insights Settings

main 1 branch 0 tags Go to file Add file Code

mjoppich preprocessed data a376daa 2 minutes ago 3 commits

alignments/bulk	folder structure	2 minutes ago
bulk	folder structure	2 minutes ago
bulk_own	folder structure	2 minutes ago
code	folder structure	2 minutes ago
images	folder structure	2 minutes ago
reads	folder structure	2 minutes ago
ready_to_use_data	folder structure	2 minutes ago
references	folder structure	2 minutes ago
sc	preprocessed data	2 minutes ago
tools	folder structure	2 minutes ago
README.md	Initial commit	10 minutes ago

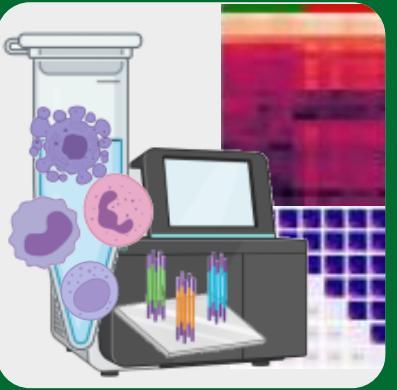
README.md

game_cibog_2022

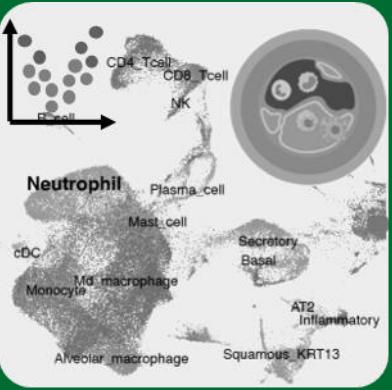
Workshop material for GAME-CIBoG workshop

Data for direct use

- **git clone**
https://github.com/mjoppich/game_cibog_2022.git
- **Folder structure useful (not required though!) for workshop**
- **Send eMail to**
joppich@bio.ifi.lmu.de
with GitHub username to get access

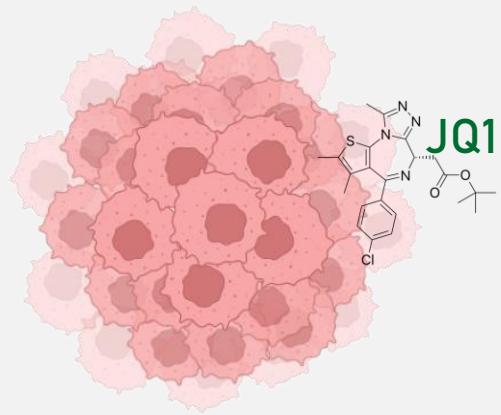


bulk RNA-seq



scRNA-seq

BULK RNA-SEQ



SUM159

Material Isolation

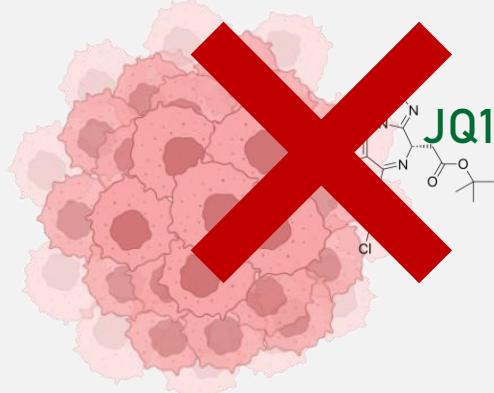
Sequencing

Alignment

Quantification

logFC + pVal

Genes



SUM159R

Quite some time

Matter of hours*

Matter of weeks

GETTING THE DATA

- Gene expression data are collected in large repositories
 - NCBI Gene Expression Omnibus (GEO)
 - EBI European Nucleotide Archive (ENA)
- For publication whole samples are collected
 - Raw sequencing data (shared with SRA, Sequencing Read Archive)
 - *Processed* expression data
- Here: use processed data
 - Saves time
 - Was probably done right anyhow

SEARCHING THE DATA

 CellPress

Molecular Cell

Resource

Synthetic Lethal and Resistance Interactions with BET Bromodomain Inhibitors in Triple-Negative Breast Cancer

Shaokun Shu,^{1,7,18,19} Hua-Jun Wu,^{2,11,12,18} Jennifer Y. Ge,^{1,2,13} Rhamy Zeid,^{1,20} Isaac S. Harris,^{8,14} Bojana Jovanović,^{1,7,16} Katherine Murphy,¹ Binbin Wang,^{2,21} Xintao Qiu,^{1,5} Jennifer E. Endress,^{8,14} Jaime Reyes,¹ Klothilda Lim,^{1,5} Alba Font-Tello,^{1,5} Sudeepa Syamala,^{1,5} Tengfei Xiao,² Chandra Sekhar Reddy Chilamakuri,¹⁵ Evangelia K. Papachristou,¹⁵ Clive D'Santos,¹⁵ Jayati Anand,¹ Kunihiko Hinohara,^{1,7} Wei Li,^{2,11,22} Thomas O. McDonald,^{2,6,11,12} Adrienne Luoma,^{3,9} Rebecca J. Modiste,¹⁷ Quang-De Nguyen,¹⁷ Brittany Michel,⁴ Paloma Cejas,^{1,5} Cigall Kadoch,^{4,10,16} Jacob D. Jaffe,¹⁶ Kai W. Wucherpfennig,^{3,9} Jun Qi,¹ X. Shirley Liu,^{2,5,11} Henry Long,^{1,5} Myles Brown,^{1,5,7,14} Jason S. Carroll,¹⁵ Joan S. Brugge,^{8,14} James Bradner,^{1,7,23} Franziska Michor,^{2,6,11,12,14,16,*} and Komelia Polyak^{1,5,6,7,14,16,24,*}

Series GSE131102

Status Public on Mar 30, 2020

Title Synthetic lethal and resistance interactions with BET bromodomain inhibitors

Organism *Homo sapiens*

Experiment type Expression profiling by high throughput sequencing
Genome binding/occupancy profiling by high throughput sequencing
Other

Summary This SuperSeries is composed of the SubSeries listed below.

Overall design Refer to individual Series

Citation(s) Shu S, Wu HJ, Ge JY, Zeid R et al. Synthetic Lethal and Resistance Interactions with BET Bromodomain Inhibitors in Triple-Negative Breast Cancer. *Mol Cell* 2020 Jun;18(6):1096-1113.e8. PMID: 32416067





<https://www.ncbi.nlm.nih.gov/geo/>

Deposited Data

All raw genomic data	GEO	GSE131102
All raw numeric data and image files	Mendeley	https://doi.org/10.17632/p4ypdxmsk5.1

<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE131102>

GETTING THE RIGHT DATA

- Interested in RNA-seq data
SUM159 cells vs SUM159-JQ1R cells

This SuperSeries is composed of the following SubSeries:

≡ Less...
GSE131026 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [ATAC-Seq]
GSE131030 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [Barcode_seq]
GSE131091 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [CRISPR screen]
GSE131097 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [ChIP-Seq]
GSE131099 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [RNA-Seq]
GSE131135 Synthetic lethal and resistance interactions with BET bromodomain inhibitors [scRNA-Seq]



**SUM159 cells
at 3h and 24h**

**SUM159-JQ1R
cells at 3h
and 24h**

Samples (62)
[≡ Less...](#)

GSM3763447	SUM149_DMSO_12h_Replicate1
GSM3763448	SUM149_DMSO_12h_Replicate2
GSM3763449	SUM149_DMSO_12h_Replicate3
GSM3763450	SUM149_JQ1_10uM_12h_R1
GSM3763451	SUM149_JQ1_10uM_12h_R2
GSM3763452	SUM149_JQ1_10uM_12h_R3
GSM3763453	SUM149_JQ1_2uM_12h_R1
GSM3763454	SUM149_JQ1_2uM_12h_R2
GSM3763455	SUM149_JQ1_2uM_12h_R3
GSM3763456	SUM149R_DMSO_12h_Replicate1
GSM3763457	SUM149R_DMSO_12h_Replicate2
GSM3763458	SUM149R_DMSO_12h_Replicate3
GSM3763459	SUM149R_JQ1_10uM_12h_R1
GSM3763460	SUM149R_JQ1_10uM_12h_R2
GSM3763461	SUM149R_JQ1_10uM_12h_R3
GSM3763462	SUM149R_JQ1_2uM_12h_R1
GSM3763463	SUM149R_JQ1_2uM_12h_R2
GSM3763464	SUM149R_JQ1_2uM_12h_R3
GSM3763465	SUM159_DMSO_12h_R1
GSM3763466	SUM159_DMSO_12h_R2
GSM3763467	SUM159_DMSO_24h_R1
GSM3763468	SUM159_DMSO_24h_R2
GSM3763469	SUM159_DMSO_3h_R1
GSM3763470	SUM159_DMSO_3h_R2
GSM3763471	SUM159_JQ1_12h_R1
GSM3763472	SUM159_JQ1_12h_R2
GSM3763473	SUM159_JQ1_24h_R1
GSM3763474	SUM159_JQ1_24h_R2
GSM3763475	SUM159_JQ1_3h_R1
GSM3763476	SUM159_JQ1_3h_R2
GSM3763477	SUM159_JQ1R_DMSO_24h_R1
GSM3763478	SUM159_JQ1R_DMSO_24h_R2
GSM3763479	SUM159_JQ1R_DMSO_3h_R1
GSM3763480	SUM159_JQ1R_DMSO_3h_R2

GETTING QUANTIFICATION DATA

Submission date May 13, 2019
Last update date Apr 01, 2020
Contact name Kornelia Polyak
E-mail(s) kornelia_polyak@dfci.harvard.edu
Phone 617-632-2106
Organization name Dana-Farber Cancer Institute
Department Medical Oncology
Lab Polyak
Street address 450 Brookline Ave
City Boston
State/province MA
ZIP/Postal code 02215
Country USA

Platform ID Series (2)
GPL11154 GSE131099 Synthetic lethal and resistance bromodomain inhibitors [RNA-seq]
GSE131102 Synthetic lethal and resistance bromodomain inhibitors

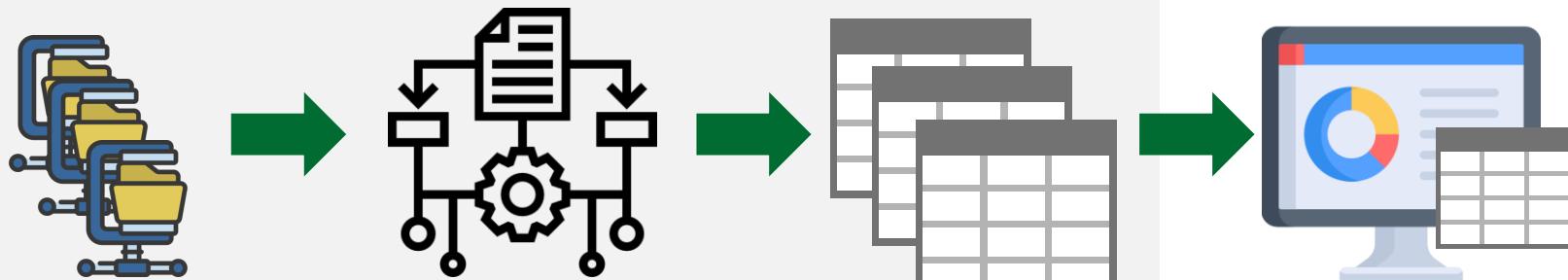
Relations
BioSample SAMN11634886
SRA SRX5824615

Supplementary file
GSM3763467_SUM159_DMSO_2_R1.counts.txt.gz 94.6 Kb
SRA Run Selector

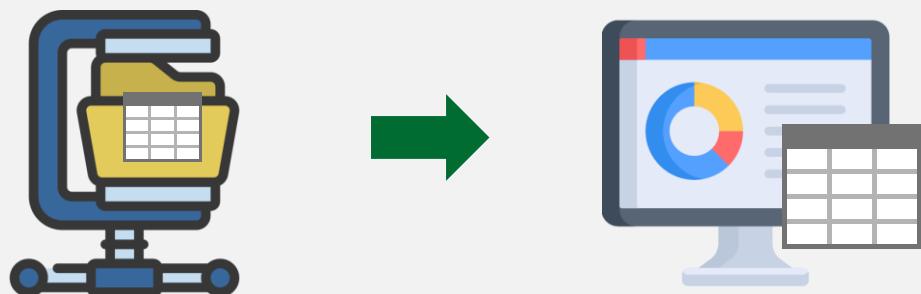
Raw Data (green arrow pointing to the left)

Processed (green arrow pointing to the left)

- Raw Data: you are the boss
 - Much more data, more steps, but most flexible (e.g. new methods)

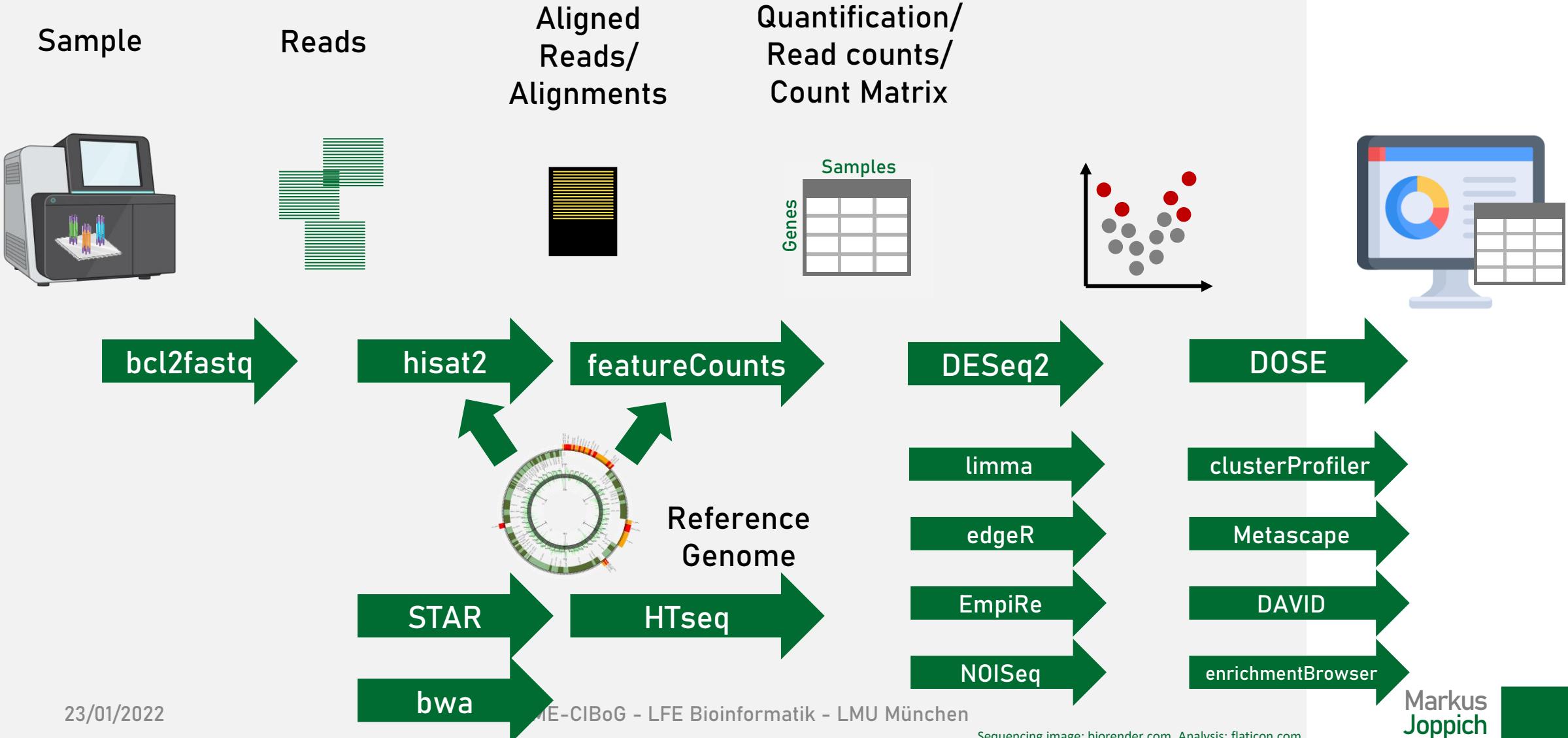


- Processed Data: much easier (download-read-go)



BULK PROCESSING PIPELINE

All the steps needed for bulk data processing



REQUIRED TOOLS FOR DOWNLOADING DATA

Getting the read data

- **sratoolkit** for downloading files
 - wget --output-document tools/sratoolkit.tar.gz <http://ftp-trace.ncbi.nlm.nih.gov/sra/sdk/current/sratoolkit.current-ubuntu64.tar.gz>
 - cd tools && tar -vxzf sratoolkit.tar.gz
 - vdb-config –interactive
 - Make path to sratoolkit visible:
 - Add

```
export PATH=$PATH:/mnt/t/workshop/tools/sratoolkit.2.11.2-ubuntu64/bin
```

to `~/.bashrc` and source `~/.bashrc`
- Then we can call *fastq-dump* as

```
fastq-dump --gzip SRR9048093
```

for downloading specific SRA-file

REQUIRED FILES AND REFERENCES

Doing the alignment

- Hisat2
 - sudo apt install hisat2
 - <http://daehwankimlab.github.io/hisat2/>
 - Index/Reference genome: H. sapiens, GRCh38, https://genome-dbx.s3.amazonaws.com/hisat/grch38_genome.tar.gz
 - tar xfz grch38_genome.tar.gz
- featureCounts / subread
 - sudo apt install subread
 - <https://sourceforge.net/projects/subread/files/subread-2.0.3/>
- Genome Annotation File:
 - http://ftp.ensembl.org/pub/release-105/gtf/homo_sapiens/Homo_sapiens.GRCh38.105.gtf.gz
 - gunzip Homo_sapiens.GRCh38.105.gtf.gz

Binaries

Version: HISAT2 2.2.1

Release Date: 7/24/2020

Source	https://cloud.biohpc.swmed.edu/index.php/s/fE9QCsX3NH4QwBi/download
OSX_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/zMgEtnF6LjnFr/download
Linux_x86_64	https://cloud.biohpc.swmed.edu/index.php/s/oTtGWbWjaxsQ2Ho/download

Gene annotation

What can I find? Protein-coding and non-coding genes, splice variants, cDNA and protein sequences, non-coding RNAs.

 More about this genebuild

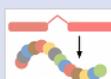
 Download FASTA files for genes, cDNAs, ncRNA, proteins

 Download GTF or GFF3 files for genes, cDNAs, ncRNA, proteins

 Update your old Ensembl IDs

Pax6 INS
FOXP2 BRC2
DMD ssh

Example gene



Example transcript

USING THE TOOLS

How to call the tools?

- General: check websites of tools, `<tool> --help`

```
mjoppich@spectre3:/mnt/w/game_cibog$ hisat2 --help
HISAT2 version 2.1.0 by Daehwan Kim (infphilo@gmail.com, wwwccb.jhu.edu/people/infphilo)
Usage:
  hisat2 [options]* -x <ht2-idx> {-1 <m1> -2 <m2> | -U <r>} [-S <sam>]

  <ht2-idx>  Index filename prefix (minus trailing .X.ht2).
  <m1>        Files with #1 mates, paired with files in <m2>.
              Could be gzip'ed (extension: .gz) or bzip2'ed (extension: .bz2).
  <m2>        Files with #2 mates, paired with files in <m1>.
              Could be gzip'ed (extension: .gz) or bzip2'ed (extension: .bz2).
  <r>         Files with unpaired reads.
              Could be gzip'ed (extension: .gz) or bzip2'ed (extension: .bz2).
  <sam>       File for SAM output (default: stdout)

  <m1>, <m2>, <r> can be comma-separated lists (no whitespace) and can be
  specified many times. E.g. '-U file1.fq,file2.fq -U file3.fq'.
```

DOWNLOADING BULK RAW DATA

```
#24h R1
! cd reads/bulk/ && fastq-dump --gzip SRR9048093
#24h R1
! cd reads/bulk/ && fastq-dump --gzip SRR9048094
#JQ1 24h R1
! cd reads/bulk/ && fastq-dump --gzip SRR9048103
#JQ1 24h R1
! cd reads/bulk/ && fastq-dump --gzip SRR9048104
```

```
Read 17872256 spots for SRR9048093
Written 17872256 spots for SRR9048093
Read 40740514 spots for SRR9048094
Written 40740514 spots for SRR9048094
Read 43554379 spots for SRR9048103
Written 43554379 spots for SRR9048103
Read 28277603 spots for SRR9048104
Written 28277603 spots for SRR9048104
```

```
! mv reads/bulk/SRR9048093.fastq.gz reads/bulk/SUM159_24h_R1.fastq.gz
! mv reads/bulk/SRR9048094.fastq.gz reads/bulk/SUM159_24h_R2.fastq.gz
! mv reads/bulk/SRR9048103.fastq.gz reads/bulk/SUM159R_24h_R1.fastq.gz
! mv reads/bulk/SRR9048104.fastq.gz reads/bulk/SUM159R_24h_R2.fastq.gz
```

- With the tools available just download all files
 - Rename for better understandability!
 - Downloading took ~1h

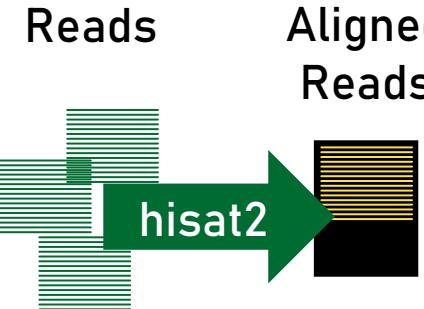
ALIGNING READS

Must be done for all samples!

```
! hisat2 -p 4 -x ./references/grch38/genome -U ./reads/bulk/SUM159_24h_R2.fastq.gz |  
samtools view -@ 4 -bS - |  
samtools sort -@ 4 -o alignments/bulk/SUM159_24h_R2.bam -  
  
40740514 reads; of these:  
40740514 (100.00%) were unpaired; of these:  
1926278 (4.73%) aligned 0 times  
33800194 (82.96%) aligned exactly 1 time  
5014042 (12.31%) aligned >1 times  
95.27% overall alignment rate  
[bam_sort_core] merging from 12 files and 4 in-memory blocks...
```

Alignment performance

- *hisat2* expects reference file and reads, output: sam format
- *samtools view* converts sam format into compressed bam format
 - Primarily done to save disk space!
- *samtools sort* sorts the reads by genome positions



REMINDER: GENE STRUCTURE

Here: CCL2 (http://www.ensembl.org/Homo_sapiens/Gene/Summary?db=core;g=ENSG00000108691;r=17:34255218-34257203)

Gene: CCL2 ENSG00000108691

Ensembl ID

Gene Symbol

Old gene symbols

Description: C-C motif chemokine ligand 2 [Source:HGNC Symbol;Acc:HGNC:10618]
Gene Symbols: GDCF-2, HC11, MCAF, MCP-1, MCP1, MGC9434, SCYA2, SMC-CF
Location: Chromosome 17: 34,255,218-34,257,203 forward strand.
GRCh38:CM000679.2
About this gene
Transcripts

This gene has 4 transcripts (splice variants), 154 orthologues, 26 paralogues and is associated with 4 phenotypes.
Show transcript table

Show/hide columns (1 hidden) Filter

Transcript ID	Name	bp	Protein	Biotype	CCDS	UniProt Match	RefSeq Match	Flags
ENST00000225831.4	CCL2-201	741	99aa	Protein coding	CCDS11277	P13500	NM_002982.4	MANE Select v0.95 Ensembl Canonical GENCODE basic APPRIS P1 TSL:1
ENST00000580907.5	CCL2-202	736	65aa	Protein coding	-	J3KRT7	-	GENCODE basic TSL:2
ENST00000582017.1	CCL2-203	996	No protein	Retained intron	-	-	-	TSL:NA
ENST00000624362.1	CCL2-204	992	No protein	TEC	-	-	-	TSL:NA

gene

transcript

exon

Non-coding

coding

Diagram illustrating the gene structure of CCL2. The top part shows the Ensembl gene summary page for CCL2, highlighting the gene symbol (CCL2), Ensembl ID (ENSG00000108691), and various transcript details. The bottom part shows a genomic track diagram where the gene (CCL2) is represented by a yellow line, and transcripts (CCL2-201, CCL2-202, CCL2-203, CCL2-204) are shown as colored lines (orange, red, blue, green) with their respective biotypes: protein coding, protein coding, retained intron, and TEC. Arrows indicate the flow from gene to transcript, and arrows labeled 'exon', 'non-coding', and 'coding' point to specific regions of the transcripts.

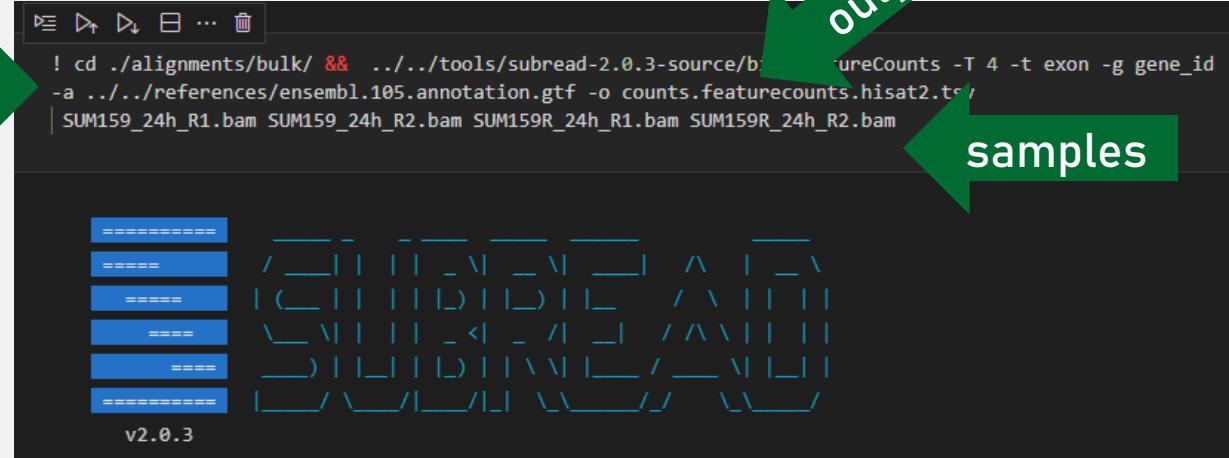
QUANTIFYING GENE COUNTS

From Alignments to count matrices

```
✓ alignments/bulk
  └── SUM159_24h_R1.bam
  └── SUM159_24h_R2.bam
  └── SUM159R_24h_R1.bam
  └── SUM159R_24h_R2.bam
```

reference

```
! cd ./alignments/bulk/ && ../../tools/subread-2.0.3-source/bin/subread featureCounts -T 4 -t exon -g gene_id
-a ../../references/ensembl.105.annotation.gtf -o counts.featurecounts.hisat2.tsv
SUM159_24h_R1.bam SUM159_24h_R2.bam SUM159R_24h_R1.bam SUM159R_24h_R2.bam
```



v2.0.3

output

samples

Aligned
Reads

Counts

Samples



subread



- Run **featureCounts** with **alignments files** and **genome annotation**
 - Faster with more threads: **-T 4**
 - Quantify by exons: **-t exon**
 - Summarize by **gene_id**: **-g gene_id**

	Status	SUM159_24h_R1.bam	SUM159_24h_R2.bam	SUM159R_24h_R1.bam	SUM159R_24h_R2.bam
1	Assigned	12730408	29461647	31964377	28473798
2	Unassigned_Unmapped	1101887	1926278	1944297	1520084
4	Unassigned_Read_Type	0	0	0	0
5	Unassigned_Singleton	0	0	0	0
6	Unassigned_Mappingquality	0	0	0	0
7	Unassigned_Chimera	0	0	0	0
8	Unassigned_FragmentLength	0	0	0	0
9	Unassigned_Duplicate	0	0	0	0
10	Unassigned_MultiMapping	6488761	15144678	15243105	10011071
11	Unassigned_Secondary	0	0	0	0
12	Unassigned_NonSplit	0	0	0	0
13	Unassigned_NoFeatures	765478	1736538	1548873	1132298
14	Unassigned_Overlapping_Length	0	0	0	0
15	Unassigned_Ambiguity	1102693	2602009	2799429	1726337
16					

COUNT MATRIX

Final Output

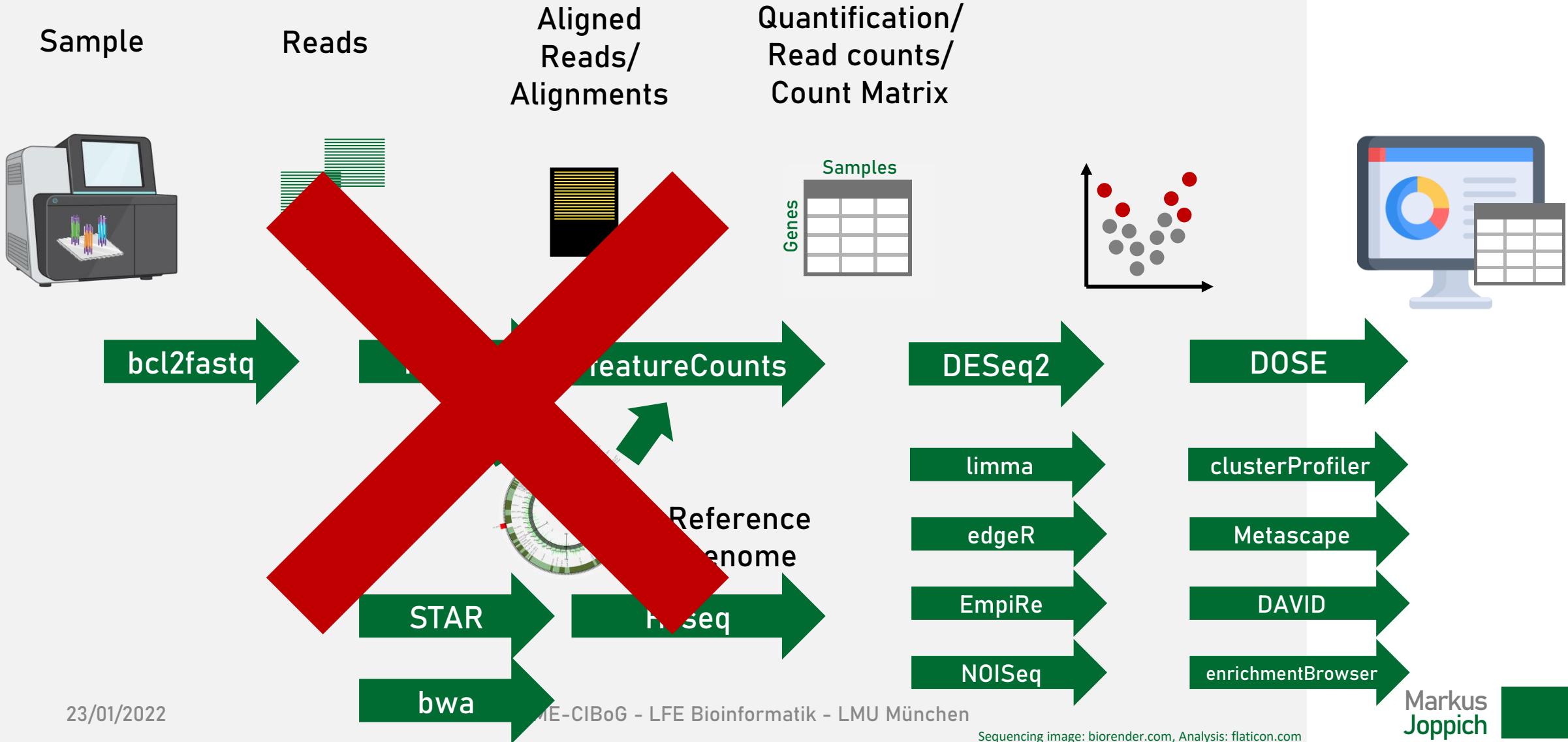
	A	B	C	D	E	F	G	H	I	J	K	L	M
1	Geneid	Chr	Start	End	Strand	Leng	SUM159_24h_R1.bam	SUM159_24h_R2.bam	SUM159R_24h_R1.bam	SUM159R_24h_R2.bam	SUM159_SUM	SUM159R_SUM	ALL_SUM
2	ENSG00000198804	MT	5904	7445	+	1542	141955	327668	468323	317305	469623	785628	1255251
3	ENSG00000198886	MT	10760	12137	+	1378	138689	320222	438344	267733	458911	706077	1164988
4	ENSG00000210082	MT	1671	3229	+	1559	76078	178317	247196	142901	254395	390097	644492
5	ENSG00000198712	MT	7586	8269	+	684	64070	149254	241954	153261	213324	395215	608539
6	ENSG00000198938	MT	9207	9990	+	784	62879	149328	305177	187365	212207	492542	704749
7	ENSG00000198727	MT	14747	15887	+	1141	59983	132070	206999	144340	192053	351339	543392
8	ENSG00000074800	1;1;1;1;1;8861000;88861429;8-;-;-;-;-;-	5216	60100	126070	98074	68928	186170	167002	353172			
9	ENSG00000111640	12;12;12;16534512;66534569;6+;+;+;+;+	2396	53008	117725	94460	64355	170733	158815	329548			
10	ENSG00000156508	6;6;6;6;673489308;73492804;-;-;-;-;-	10063	52679	115930	89297	63060	168609	152357	320966			
11	ENSG00000167658	19;19;19;13976056;33976747;3-;-;-;-;-	4021	44637	114492	119551	61234	159129	180785	339914			
12	ENSG00000163359	2;2;2;2;223732400;237324814-;-;-;-;-	21516	50449	108069	69778	52632	158518	122410	280928			

Total Sums	SUM159_24h_R1.bam	SUM159_24h_R2.bam	SUM159R_24h_R1.bam	SUM159R_24h_R2.bam
	12730442	29461460	31964210	20473685

- Large table (>60000 rows!)
- Sums in counts differ between samples and replicates!
 - Absolute values are not comparable!
- Ensembl Geneids are ... unique, but not interpretable!

AVOIDING THE HASSLE OF DOING ALIGNMENT+COUNTING

Use preprocessed data



DOWNLOADING THE PROCESSED DATA

- Interested in 8 samples
 - Manually get the links, download with command-line
 - `wget -O <filename> <URL>`

The screenshot shows a bioinformatics data download interface. At the top, it displays sample details:

Submission date	May 13, 2019
Last update date	Apr 01, 2020
Contact name	Kornelia Polyak
E-mail(s)	kornelia_polyak@dfci.harvard.edu
Phone	617-632-2106
Organization name	Dana-Farber Cancer Institute
Department	Medical Oncology
Lab	Polyak
Street address	450 Brookline Ave
City	Boston
State/province	MA
ZIP/Postal code	02215
Country	USA

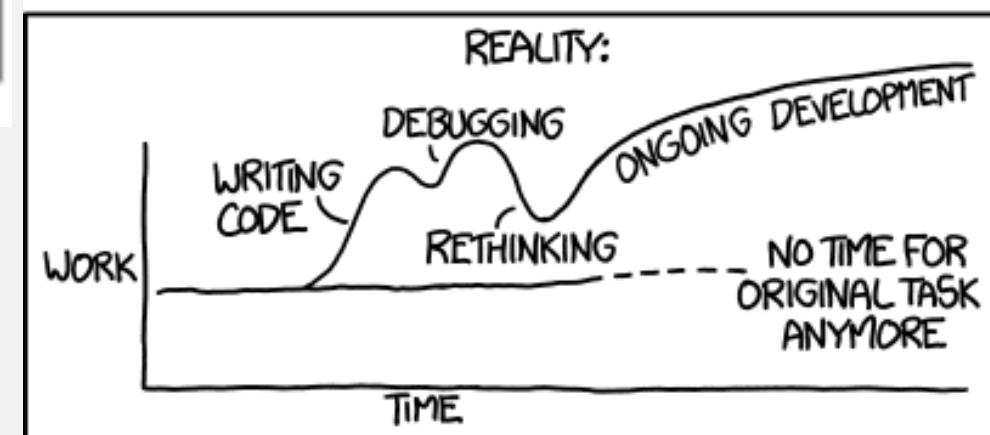
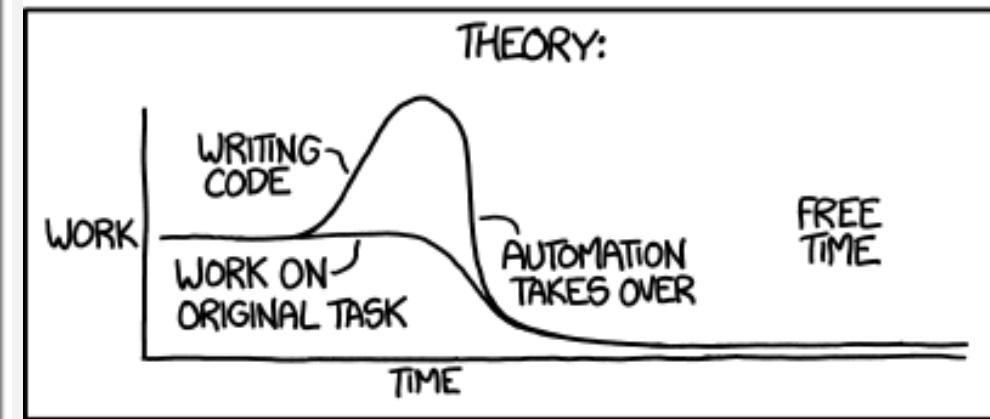
Below this, it lists platform ID (GPL11154) and two series (GSE131099 and GSE131102), both described as "Synthetic lethal and resistance interactions between bromodomain inhibitors [RNA-Seq]".

The "Relations" section shows BioSample (SAMN11634886) and SRA (SRX5824615).

At the bottom, a "Supplementary file" table lists 10 files, all of which are "GSM3763467_SUM159_DMSO_24h_R1.counts.txt.gz" files with varying suffixes (e.g., R1, R2, counts.txt.gz).

A context menu is open over one of the files, with "Copy Link" highlighted.

"I SPEND A LOT OF TIME ON THIS TASK.
I SHOULD WRITE A PROGRAM AUTOMATING IT!"

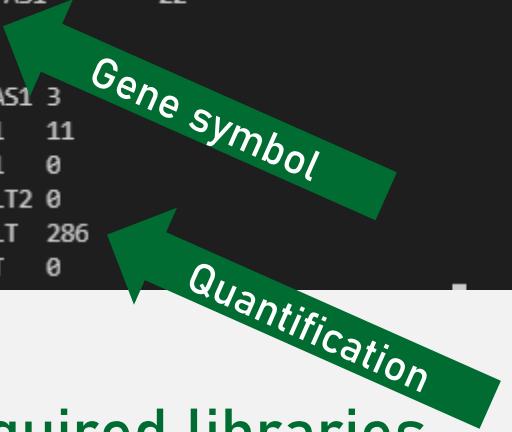


THE FIRST DE ANALYSIS

Reading in the data

- Let's have a peak at the downloaded files: zless

```
mjoppich@spectre3:/mnt/w/game_cibog$ zcat bulk/GSM3763467_SUM159_DMSO_24h_R1.counts.txt.gz | head
A1BG      44
A1BG-AS1    22
A1CF
A2M
A2M-AS1    3
A2ML1     11
A2MP1      0
A3GALT2    0
A4GALT    286
A4GNT      0
```



- Required libraries
 - I/O
 - DE analysis
 - Plotting

```
library(readr)
library('DESeq2')
library(ggplot2)
library(cowplot)
```

THE FIRST DE ANALYSIS

Getting all samples loaded

- Read single samples

```
Let's start with reading in the 24h samples
SUM159_24h_r1 = read_tsv("../bulk/GSM3763468_SUM159_DMSO_24h_R1.counts.txt.gz", col_names=F)
SUM159_24h_r2 = read_tsv("../bulk/GSM3763468_SUM159_DMSO_24h_R2.counts.txt.gz", col_names=F)

SUM159jq1r_24h_r1 = read_tsv("../bulk/GSM3763477_SUM159_JQ1R_DMSO_24h_R1.counts.txt.gz", col_names=F)
SUM159jq1r_24h_r2 = read_tsv("../bulk/GSM3763477_SUM159_JQ1R_DMSO_24h_R2.counts.txt.gz", col_names=F)

Rows: 23841 Columns: 2
```

- Merge all samples into one data frame
 - Name the columns

```
colnames(SUM159_24h_r1) = c("gene", "SUM159_24h_r1")
colnames(SUM159_24h_r2) = c("gene", "SUM159_24h_r2")
colnames(SUM159jq1r_24h_r1) = c("gene", "SUM159jq1r_24h_r1")
colnames(SUM159jq1r_24h_r2) = c("gene", "SUM159jq1r_24h_r2")
```

```
sum159_24h.expr_df <- merge(SUM159_24h_r1, SUM159_24h_r2, by = 'gene')
sum159jq1r_24h.expr_df <- merge(SUM159jq1r_24h_r1, SUM159jq1r_24h_r2, by = 'gene')
```

```
merge_24h.expr_df <- merge(sum159_24h.expr_df, sum159jq1r_24h.expr_df, by = 'gene')
```

THE FIRST DE ANALYSIS

Getting all samples loaded

- Regularly inspect your objects

```
head(merge_24h.expr_df)
```

A data.frame: 6 × 5

	gene	SUM159_24h_r1	SUM159_24h_r2	SUM159jq1r_24h_r1	SUM159jq1r_24h_r2
	<chr>	<dbl>	<dbl>	<dbl>	<dbl>
1	_alignment_not_unique	1946526	4392993	5571795	3575085
2	_ambiguous	385372	884778	815090	514010
3	_no_feature	2015667	4587061	5254348	3574216
4	_not_aligned	0	0	0	0
5	_too_low_aQual	0	0	0	0
6	A1BG	44	214	171	63

- Remove unwanted “genes”

```
expr_df_24h = merge_24h.expr_df[grep("^_", merge_24h.expr_df$gene, invert=T),]  
rownames(expr_df_24h) = expr_df_24h$gene  
expr_df_24h$gene = NULL  
expr_df_24h
```

A data.frame: 23836 × 4

	SUM159_24h_r1	SUM159_24h_r2	SUM159jq1r_24h_r1	SUM159jq1r_24h_r2
	<dbl>	<dbl>	<dbl>	<dbl>
A1BG	44	214	171	63
A1BG-AS1	22	45	69	58
A1CF	2	4	1	0
A2M	4	28	0	0

DESeq2

Brief introduction

Love et al. *Genome Biology* (2014) 15:550
DOI 10.1186/s13059-014-0550-8



METHOD

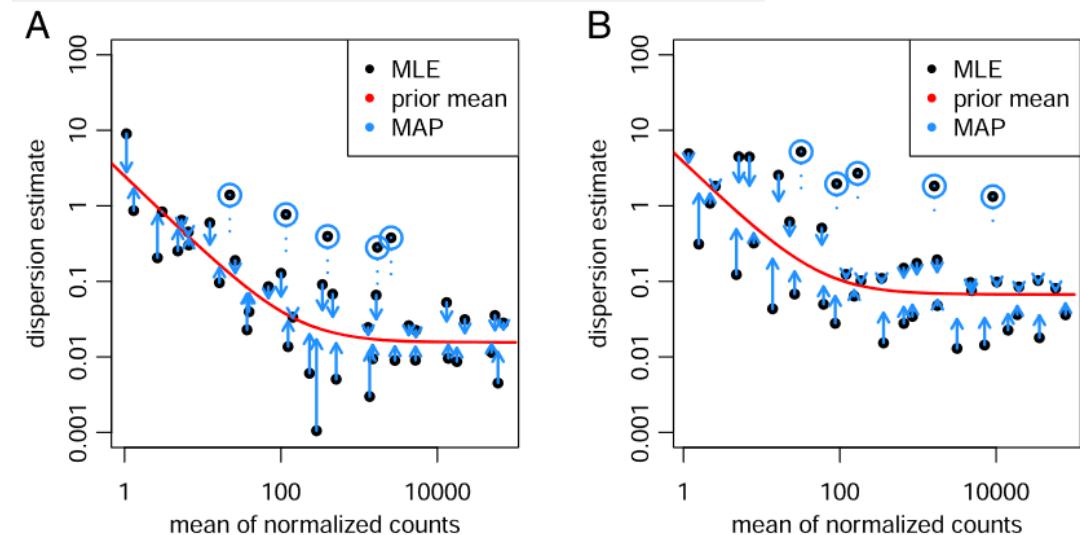
Open Access

Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2

Michael I Love^{1,2,3}, Wolfgang Huber² and Simon Anders^{2*}

Abstract

In comparative high-throughput sequencing assays, a fundamental task is the analysis of count data, such as read counts per gene in RNA-seq, for evidence of systematic changes across experimental conditions. Small replicate numbers, discreteness, large dynamic range and the presence of outliers require a suitable statistical approach. We present *DESeq2*, a method for differential analysis of count data, using shrinkage estimation for dispersions and fold changes to improve stability and interpretability of estimates. This enables a more quantitative analysis focused on the strength rather than the mere presence of differential expression. The *DESeq2* package is available at <http://www.bioconductor.org/packages/release/bioc/html/DESeq2.html>.



If you want to see different examples:

https://hbctraining.github.io/DGE_workshop_salmon_online/schedule/links-to-lessons.html

DESeq2

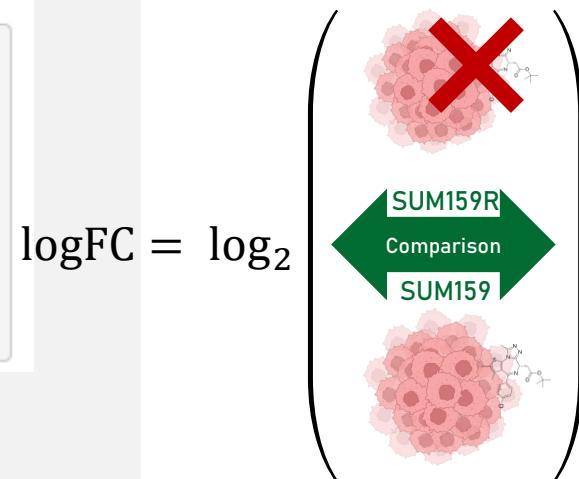
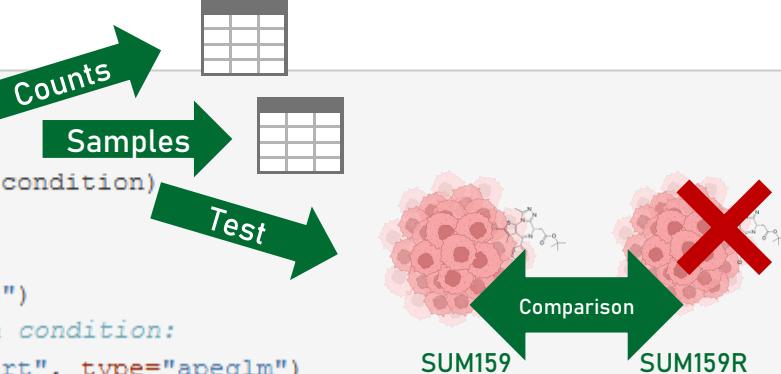
Quick Start

Quick start

Here we show the most basic steps for a differential expression analysis. There are a variety of steps upstream of DESeq2 that result in the generation of counts or estimated counts for each sample, which we will discuss in the sections below. This code chunk assumes that you have a count matrix called `cts` and a table of sample information called `coldata`. The `design` indicates how to model the samples, here, that we want to measure the effect of the condition, controlling for batch differences. The two factor variables `batch` and `condition` should be columns of `coldata`.

```
dds <- DESeqDataSetFromMatrix(countData = cts,
                               colData = coldata,
                               design= ~ batch + condition)

dds <- DESeq(dds)
resultsNames(dds) # lists the coefficients
res <- results(dds, name="condition_trt_vs_untrt")
# or to shrink log fold changes association with condition:
res <- lfcShrink(dds, coef="condition_trt_vs_untrt", type="apeglm")
```



<https://www.bioconductor.org/packages/devel/bioc/vignettes/DESeq2/inst/doc/DESeq2.html>

DESEQ2 PREPARATIONS

- Preparing sample sheet

```
pdat_24h = data.frame( id=c("SUM159_24h_R1", "SUM159_24h_R2", "SUM159R_24h_R1", "SUM159R_24h_R2"),  
                      group=c("SUM159_24h", "SUM159_24h", "SUM159R_24h", "SUM159R_24h") )  
  
A data.frame: 4 × 2  
  id      group  
  <chr>    <chr>  
1 SUM159_24h_R1  SUM159_24h  
2 SUM159_24h_R2  SUM159_24h  
3 SUM159R_24h_R1  SUM159R_24h  
4 SUM159R_24h_R2  SUM159R_24h
```

- Assigns each sample/replicate a group
 - We will then compare two groups!

DESEQ2

DESeq2 analysis for 24h

```
dds_24h <- DESeqDataSetFromMatrix(countData=expr_df_24h,  
                                    colData=pdat_24h,  
                                    design=~group, tidy = FALSE)
```

converting counts to integer mode

```
Warning message in DESeqDataSet(se, design = design, ignoreRank):  
“some variables in design formula are characters, converting to factors”
```

- Create DESeqDataSet (DDS)
 - Count data/quantifications
 - Sample sheet
 - Design pattern (groups from sample sheet)
- Call DESeq
 - Perform shrinkage*
 - Order result by p-value (BH corrected)



```
dds_24h <- DESeq(dds_24h)  
res_lfc_24h <- lfcShrink(dds_24h, coef=resultsNames(dds_24h)[[2]], type="apeglm")  
res_lfc_24h <- res_lfc_24h[order(res_lfc_24h$padj, decreasing=FALSE),]
```

```
res_24h <- results(dds_24h)  
res_24h <- res_24h[order(res_24h$padj, decreasing=FALSE),]
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

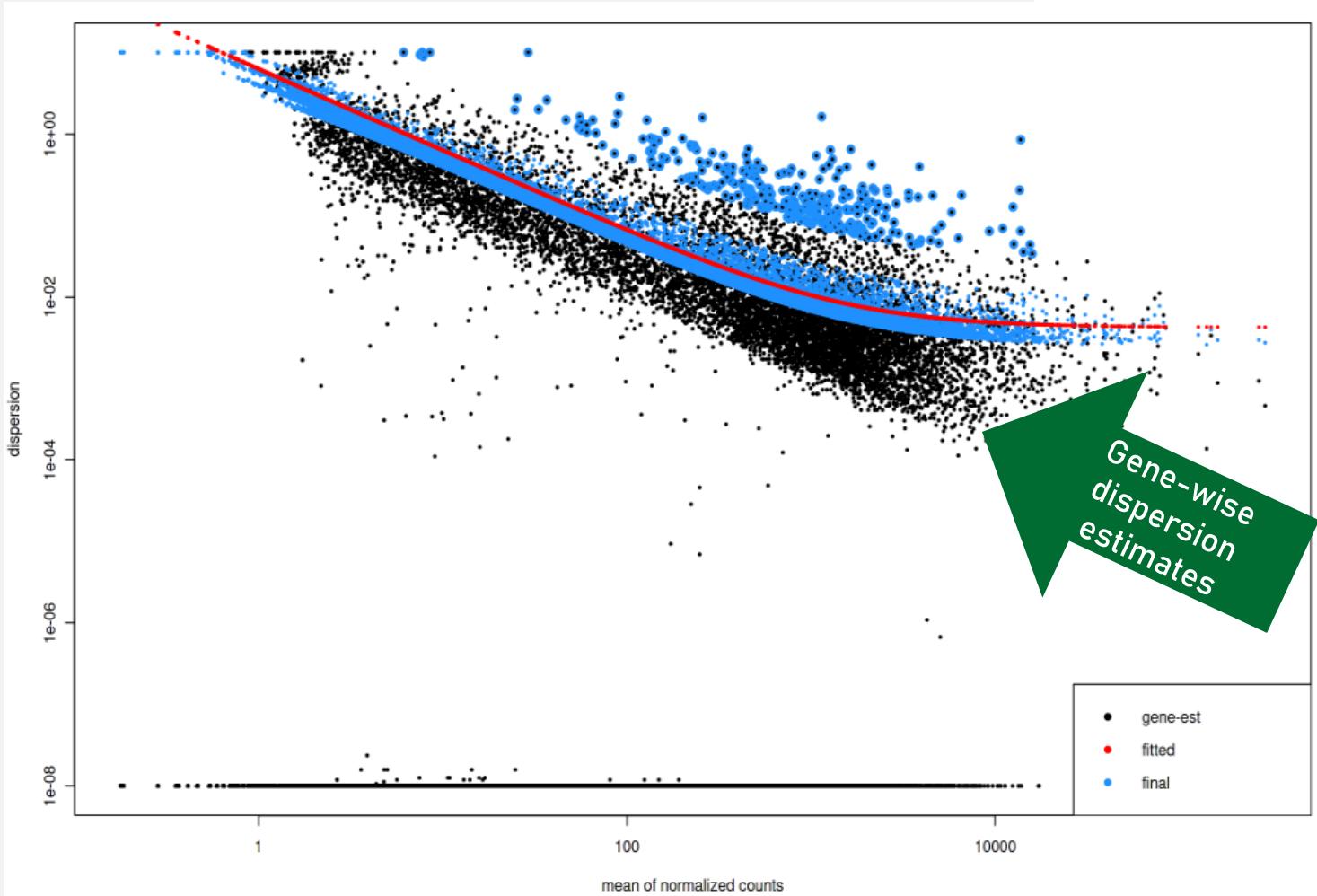
fitting model and testing

using 'apeglm' for LFC shrinkage. If used in published research, please cite:
Zhu, A., Ibrahim, J.G., Love, M.I. (2018) Heavy-tailed prior distributions for
sequence count data: removing the noise and preserving large differences.
Bioinformatics. <https://doi.org/10.1093/bioinformatics/bty895>

DESEQ2

Interpreting the run

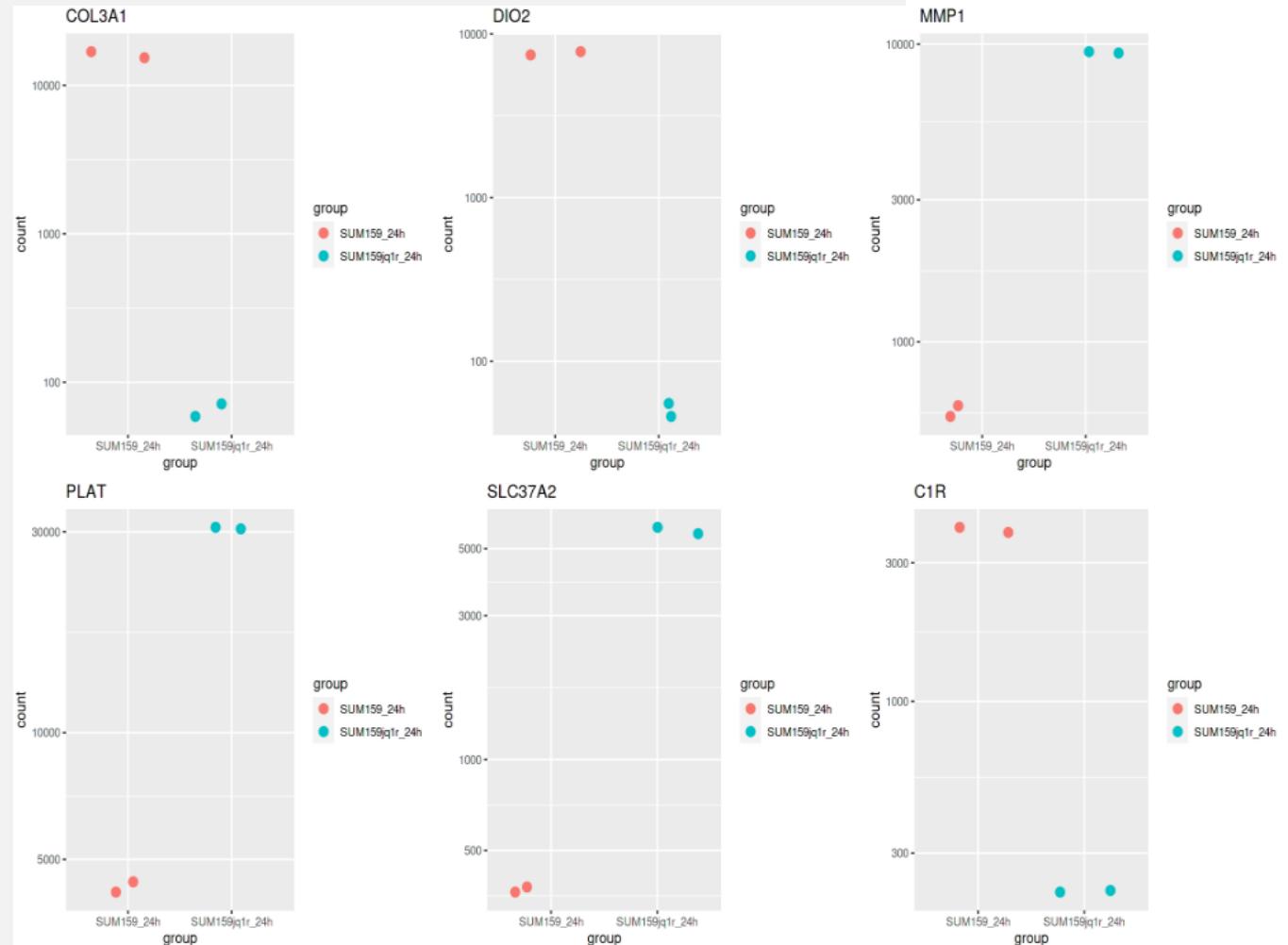
- Dispersion is a measure of spread or variability in the data
 - Here: related to the mean and variance of the data
- Plots expected dispersion value for genes of a given expression strength
- Expectation: generally scattering around the curve
 - dispersion decreasing with increasing mean expression levels



DESEQ2

Interpreting the run

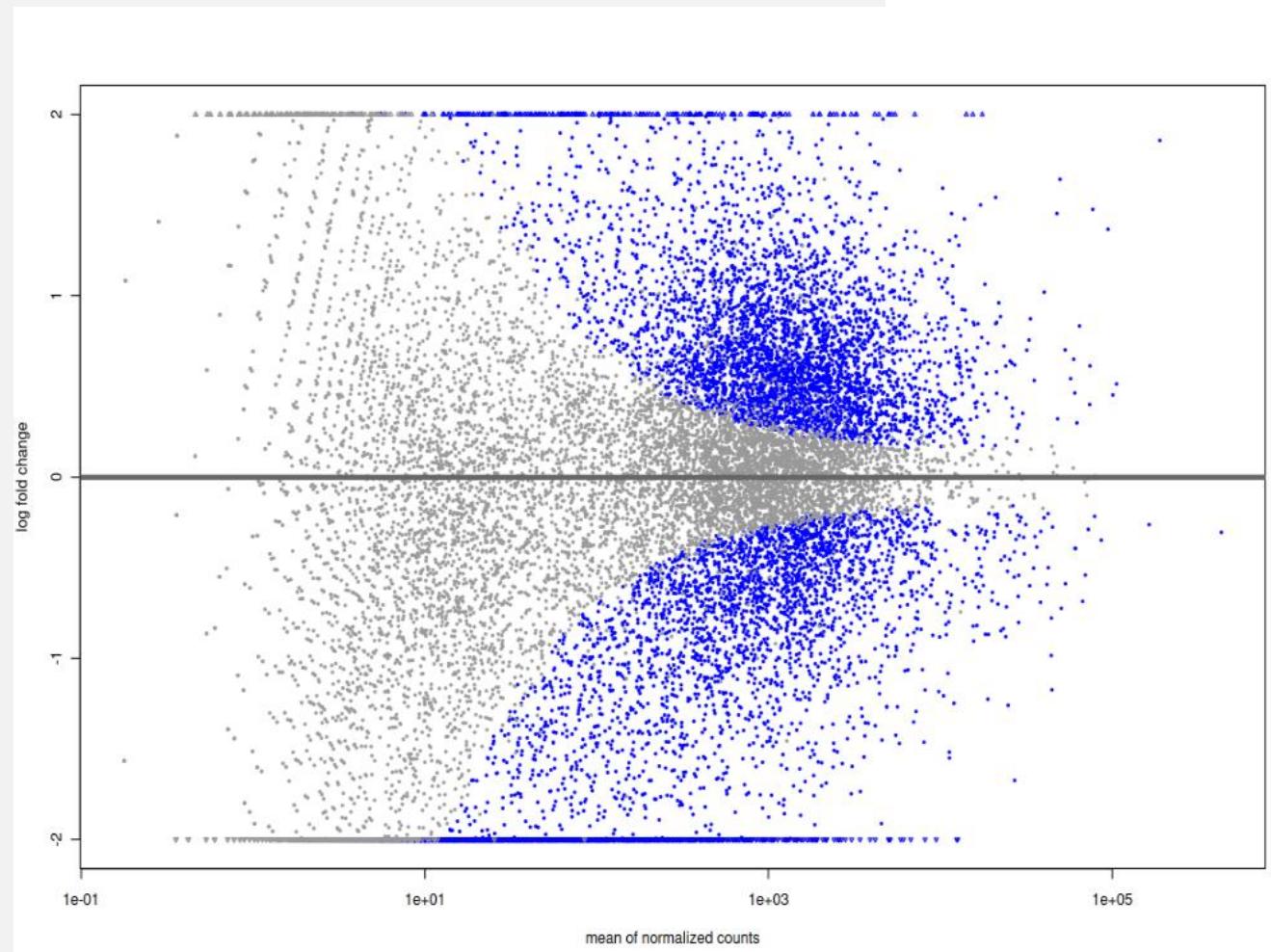
- Take a look at the most differentially expressed genes
 - Normalized! Counts per replicate
 - COL3A1 not highly expressed in SUM159R
 - But in SUM159



DESEQ2

Interpreting the run

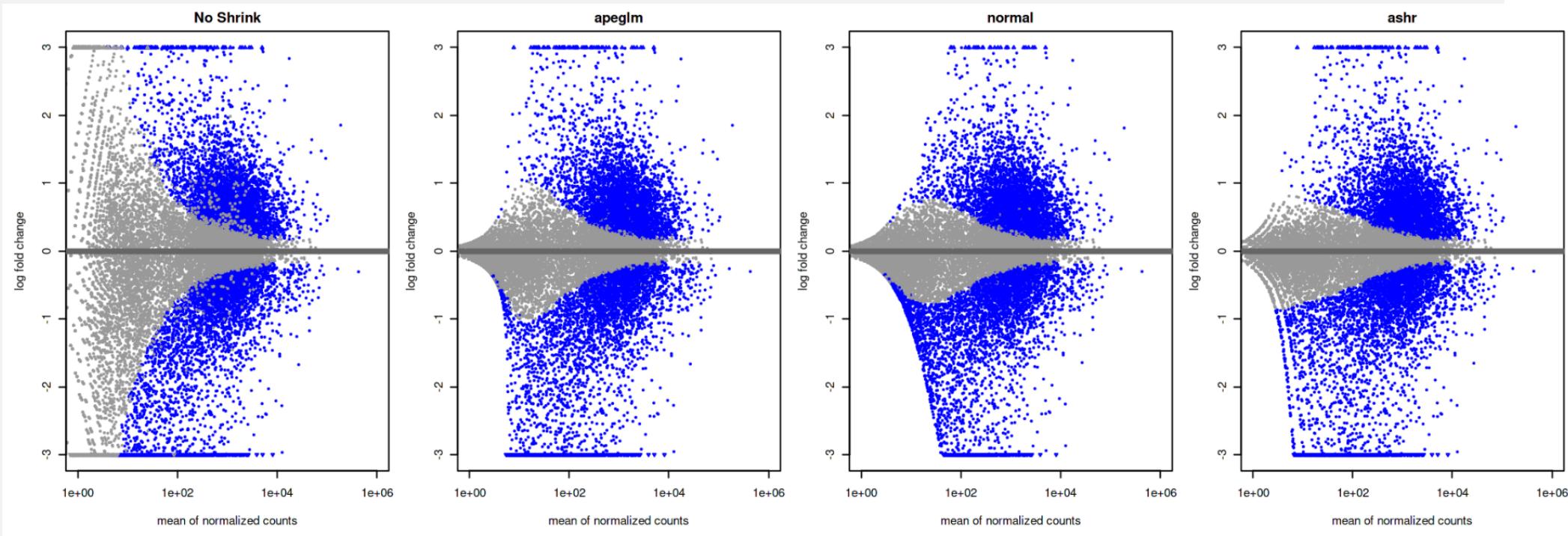
- MA plot
 - Log Fold Changes by mean of expression (in normalized counts)
 - Notice: high fold changes particularly for low counts
 - $12 / 3 \Rightarrow \text{logFC} = 2$
 - ± 2 read accuracy
 - $14/1 \Rightarrow \text{logFC} = 3.81$



DESEQ2

Interpreting the run

- Different shrinkage methods
 - No influence on statistical evaluation!
 - Shrinkage removes the noise associated with log2FCs from low count genes without requiring arbitrary filtering thresholds



ANNOTATING RESULTS

Make it more readable!

```
. log2 fold change (MAP): group SUM159R 24h vs SUM159 24h  
Wald test p-value: group SUM159R 24h vs SUM159 24h  
DataFrame with 6 rows and 5 columns
```

	baseMean	log2FoldChange	lfcSE	pvalue
ENSG00000168542	7623.80	-8.01821	0.1590666	0.00000e+00
ENSG00000254166	7798.00	-5.74974	0.1165984	0.00000e+00
ENSG00000196611	4798.88	4.00955	0.1011819	0.00000e+00
ENSG00000211448	3935.12	-7.20920	0.1654849	0.00000e+00
ENSG00000104368	17205.64	2.84687	0.0806986	1.03324e-274
ENSG00000134955	2884.91	4.37507	0.1265195	9.71242e-263
				4.61547e-271

```
gencodeAnnot = as.data.frame(read_tsv("../references/gencode.grch38.human.35.gtfout.list"))  
colnames(gencodeAnnot) = c("gene_stable_id", "gene_name", "biotype")  
gencodeAnnot$gene_stable_id = gsub("\\.[0-9]+$", "", gencodeAnnot$gene_stable_id)  
head(gencodeAnnot)
```

gene_stable_id	gene_name	biotype	
ENSG000000000003	TSPAN6	protein_coding	
2	ENSG000000000005	TNMD	protein_coding
3	ENSG000000000419	DPM1	protein_coding
4	ENSG000000000457	SCYL3	protein_coding
5	ENSG000000000460	C1orf112	protein_coding
6	ENSG000000000938	FGR	protein_coding

Row.names	baseMean	log2FoldChange	lfcSE	pvalue	padj	gene_name	biotype	
<1<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<chr>	<chr>	
10002	ENSG00000168542	7623.801	-8.018207	0.1590657	0.000000e+00	0.000000e+00	COL3A1	protein_coding
12941	ENSG00000196611	4798.877	4.009554	0.10118186	0.000000e+00	0.000000e+00	MMP1	protein_coding
13923	ENSG00000211448	3935.122	-7.209199	0.16548490	0.000000e+00	0.000000e+00	DIO2	protein_coding
18508	ENSG00000254166	7797.995	-5.749736	0.11659844	0.000000e+00	0.000000e+00	CASC19	lncRNA
2473	ENSG00000104368	17205.636	2.846869	0.08069864	1.033238e-274	4.615475e-271	PLAT	protein_coding
5773	ENSG00000134955	2884.914	4.375074	0.12651954	9.712416e-263	3.615447e-259	SLC37A2	protein_coding

DESEQ2

Creating outputs

- Filter for only entries with valid adjusted p-values
- Significance only decided on adjusted p-value (BH)
- Often also filtered by logFC
 - Clear signal (p-value)
 - Strong signal (logFC)

```
outres_24h = as.data.frame(res_lfc_24h[! is.na(res_lfc_24h$padj),])
outres_24h = merge(outres_24h, gencodeAnnot, by.x=0, by.y="gene_stable_id")

outres_24h = outres_24h[order(outres_24h$padj), ]

normExpr_24h = counts(dds_24h, normalized=T)
cn_24h = colnames(normExpr_24h)
normExpr2_24h = cbind(rownames(normExpr_24h), normExpr_24h)
colnames(normExpr2_24h) = c("Geneid", cn_24h)

write.table(outres_24h, file="../bulk_own/24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")
write.table(normExpr2_24h, file="../bulk_own/normed_expr.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")

printres_24h = outres_24h[outres_24h$padj < 0.001, ]
write.table(printres_24h, file="../bulk_own/sig.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")

printres_24h.up = printres_24h[printres_24h$log2FoldChange > 1, ]
write.table(printres_24h.up , file="../bulk_own/sig.up.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")

printres_24h.down = printres_24h[printres_24h$log2FoldChange < -1, ]
write.table(printres_24h.down , file="../bulk_own/sig.down.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")
```

```
top24h.de = outres_24h[order(abs(outres_24h$log2FoldChange),decreasing=TRUE),]

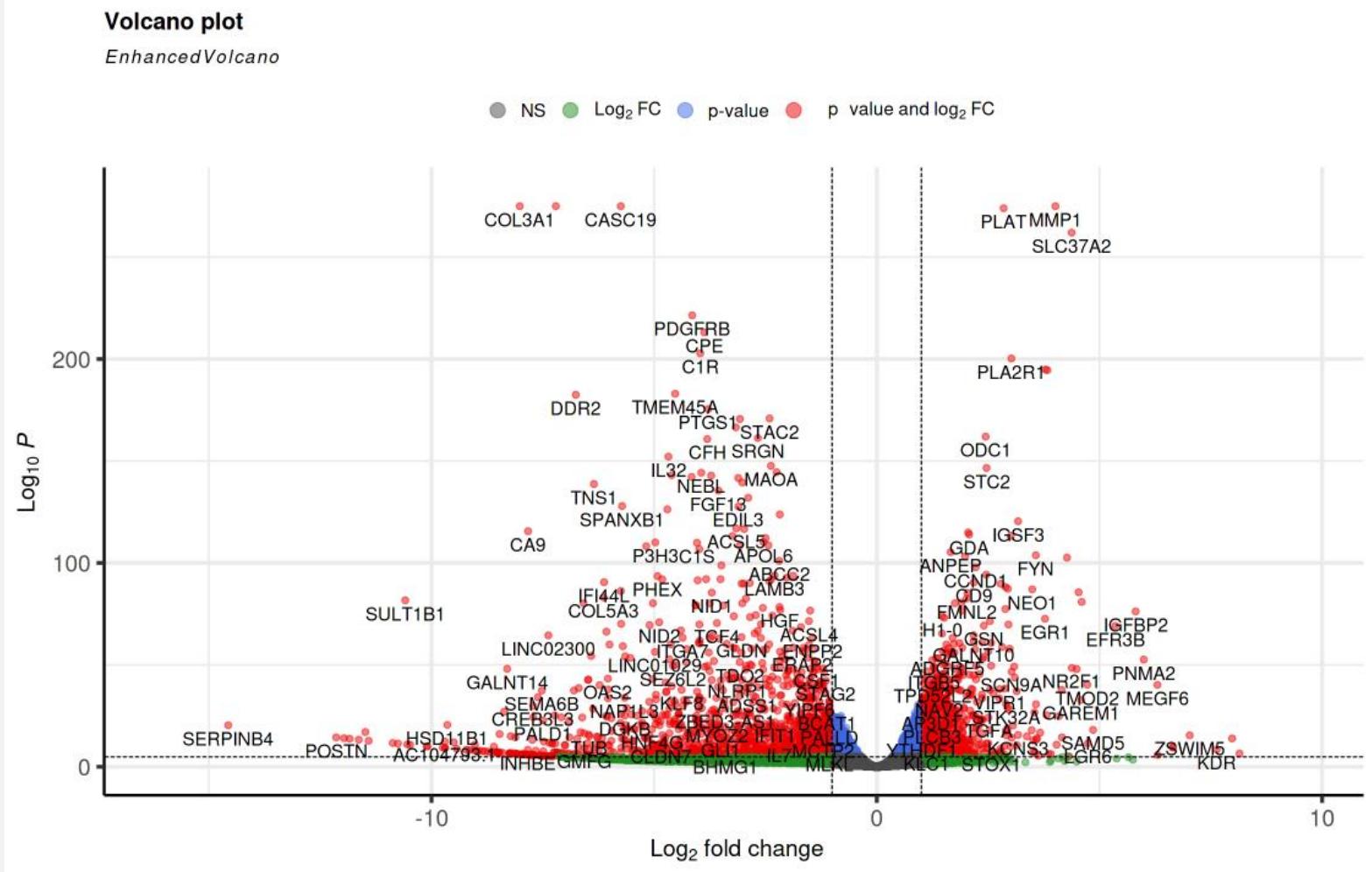
top24h.de.500.down = head(top24h.de[top24h.de$log2FoldChange < -1,], n=500)
top24h.de.500.up = head(top24h.de[top24h.de$log2FoldChange > 1,], n=500)

write.table(top24h.de.500.down , file="../bulk_own/top.down.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")
write.table(top24h.de.500.up , file="../bulk_own/top.up.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")

write.table(top24h.de , file="../bulk_own/top.all.24h_sum159_vs_sum159_jq1r.tsv", row.names=F, quote=F, sep="\t")
```

DESEQ2

Results



- Enhanced Visualization, e.g. Volcano Plot
 - library(ggrepel)
 - library(EnhancedVolcano)

```
EnhancedVolcano(outres_24h,  
  lab = outres_24h$gene_name,  
  x = 'log2FoldChange',  
  y = 'pvalue')
```

SET ENRICHMENT

Trying to identify common structures

estrogen receptor-negative breast cancer, C4733092		
Source: ALL		
Gene	Uniprot	Gene Full Name
TNF	P01375	tumor necrosis factor
TP53	P04637	tumor protein p53
IL6	P05231	interleukin 6
VEGFA	P15692	vascular endothelial growth factor A
PIK3CA	P42336	phosphatidylinositol-4,5-bisphosphate 3-kinase catalytic subunit alpha
BCL2	P10415	BCL2 apoptosis regulator
EGFR	P00533	epidermal growth factor receptor
CTNNB1	P35222	catenin beta 1
PTEN	P60484	phosphatase and tensin homolog
CDKN2A	P42771 Q8N726	cyclin dependent kinase inhibitor 2A
CXCL8	P10145	C-X-C motif chemokine ligand 8
AKT1	P31749	AKT serine/threonine kinase 1
PTGS2	P35354	prostaglandin-endoperoxide synthase 2
IGF1	P05019	insulin like growth factor 1
STAT3	P40763	signal transducer and activator of transcription 3
TLR4	O00206	toll like receptor 4
CCL2	P13500	C-C motif chemokine ligand 2
PIK3CD	O00329	phosphatidylinositol-4,5-bisphosphate 3-kinase catalytic subunit delta

- Gene sets are just lists of genes
 - Sometimes called gene signatures
- Genes were found to be associated with gene set description
 - Here: estrogen receptor-negative breast cancer
 - Association crucial!
 - Automatic retrieval? Evidence? Context?
 - Gene sets/terms not independent!
 - Overlapping genes!
- Command-line tools or web service?
 - disgenet2r+DOSE+clusterProfiler
 - Metascape
- Method?
 - ORA/hypergeometric test
 - GSEA



DISGENET2R: DIRECTLY FINDING

- Follow manual:

<https://disgenet.org/static/disgenet2r/disgenet2r.html>

- Create account
- Install disgenet2r

Selected gene symbols

```
input_genes = head(top24h.de.500.up$gene_name, n=500)
print(input_genes)

disgenet_input <- gene2disease( gene = input_genes, vocabulary = "HGNC", database = "ALL")

common_genes = intersect(disgenet_input@qresult$gene_symbol, input_genes)

print(paste("Remaining Gene Count: ", length(common_genes), "of", length(input_genes)))

res_enrich <- disease_enrichment( entities = common_genes, vocabulary = "HGNC", database = "ALL")
```

Perform enrichment

- Performs ORA on DisGeNET gene sets

A data.frame: 19 × 7								
ID <chr>	Description <chr>	source <chr>	Ratio <chr>	BgRatio <chr>	pvalue <dbl>	FDR <dbl>		
							<chr>	<dbl>
373 C0009402	Colorectal Carcinoma	ALL	165/424	5473/21666	5.346205e-10	2.049735e-06		
2247 C0178874	Tumor Progression	ALL	128/424	3865/21666	3.904144e-10	2.049735e-06		
4343 C0686619	Secondary malignant neoplasm of lymph node	ALL	101/424	2825/21666	9.763622e-10	2.495582e-06		
6593 C2939419	Secondary Neoplasm	ALL	90/424	2492/21666	6.356075e-09	1.218459e-05		
1445 C0033578	Prostatic Neoplasms	ALL	69/424	1722/21666	9.714469e-09	1.440937e-05		
3706 C0376358	Malignant neoplasm of prostate	ALL	138/424	4502/21666	1.127494e-08	1.440937e-05		
5421 C1458155	Mammary Neoplasms	ALL	96/424	2780/21666	1.682902e-08	1.843499e-05		
	Schizophrenia	ALL	98/424	2872/21666	2.084294e-08	1.997796e-05		
	Prostate carcinoma	ALL	134/424	4388/21666	2.623150e-08	2.141940e-05		
	Liver carcinoma	ALL	164/424	5725/21666	3.072684e-08	2.141940e-05		
	Oestrogen receptor positive breast cancer	ALL	31/424	510/21666	2.989114e-08	2.141940e-05		
	Malignant neoplasm of lung	ALL	215/424	21566/21666	4.860294e-08	3.105728e-05		
	Carcinoma	ALL	156/424	21566/21666	5.0906e-08	4.046879e-05		
	Squamous cell carcinoma of esophagus	ALL	75/424	2053/21666	3.7174e-08	5.278418e-05		
	Primary malignant neoplasm of lung	ALL	120/424	3894/21666	1.197581e-07	6.122033e-05		
	ovarian neoplasm	ALL	87/424	2542/21666	1.407022e-07	6.743155e-05		
	Adenocarcinoma	ALL	79/424	2235/21666	1.628121e-07	7.309775e-05		
	Malignant neoplasm of endometrium	ALL	52/424	1235/21666	1.754802e-07	7.309775e-05		
	Lung Neoplasms	ALL	59/424	1486/21666	1.811238e-07	7.309775e-05		

DOSE/CLUSTERPROFILER

```
sel_df = outliers_24h[outliers_24h$log2FoldChange > 0,]

nsamples=1500
geneList = head(sel_df$log2FoldChange, n=nsamples)
names(geneList) = head(sel_df$gene_name, n=nsamples)
names(geneList) = mapIds(org.Hs.eg.db, head(sel_df$gene_name, n=nsamples), 'ENTREZID', 'SYMBOL')
geneList = geneList[!is.na(names(geneList))]

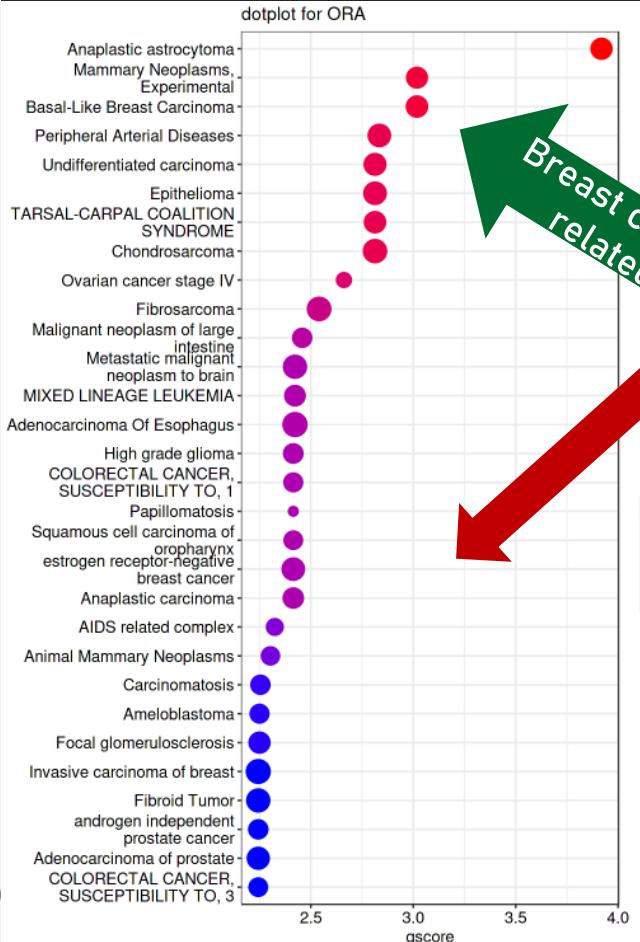
geneList = geneList[order(geneList, decreasing = TRUE)]
head(geneList)
```

- Requires list-objects with log-Foldchanges as value, and ENTREZID as name!
 - mapIds transfers gene symbols to entrez ID
 - Conversion is not exact: 1:m mapping, might not find EntrezID for all gene symbols
- Extremely well done tutorial available: <https://yulab-smu.top/biomedical-knowledge-mining-book/dose-enrichment.html>

DOSE/CLUSTERPROFILER

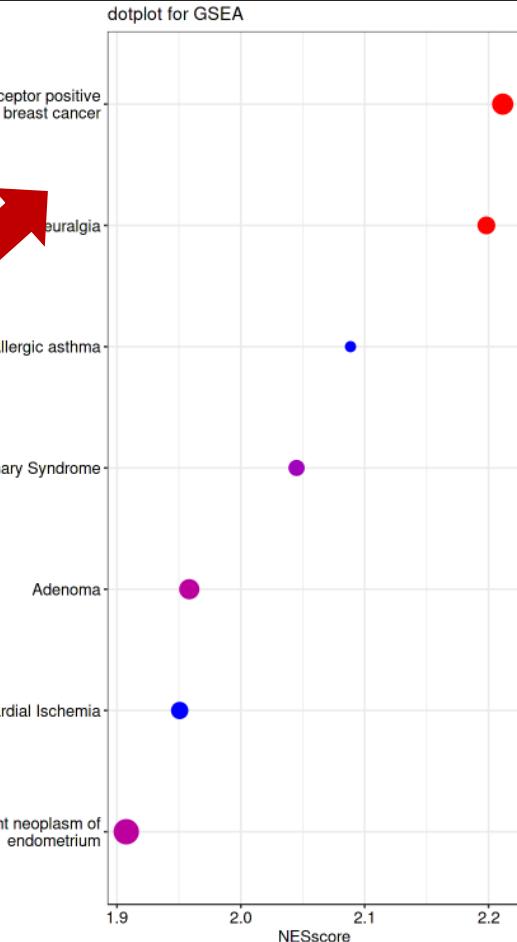
Overrepresentation analysis (ORA)
→ hypergeometric test

```
dgnEnrich <- enrichDGN(names(geneList),  
|   pvalueCutoff = 0.5,  
|   pAdjustMethod = "BH")  
dgnEnrich <- setReadable(dgnEnrich, 'org.Hs.eg.db')
```



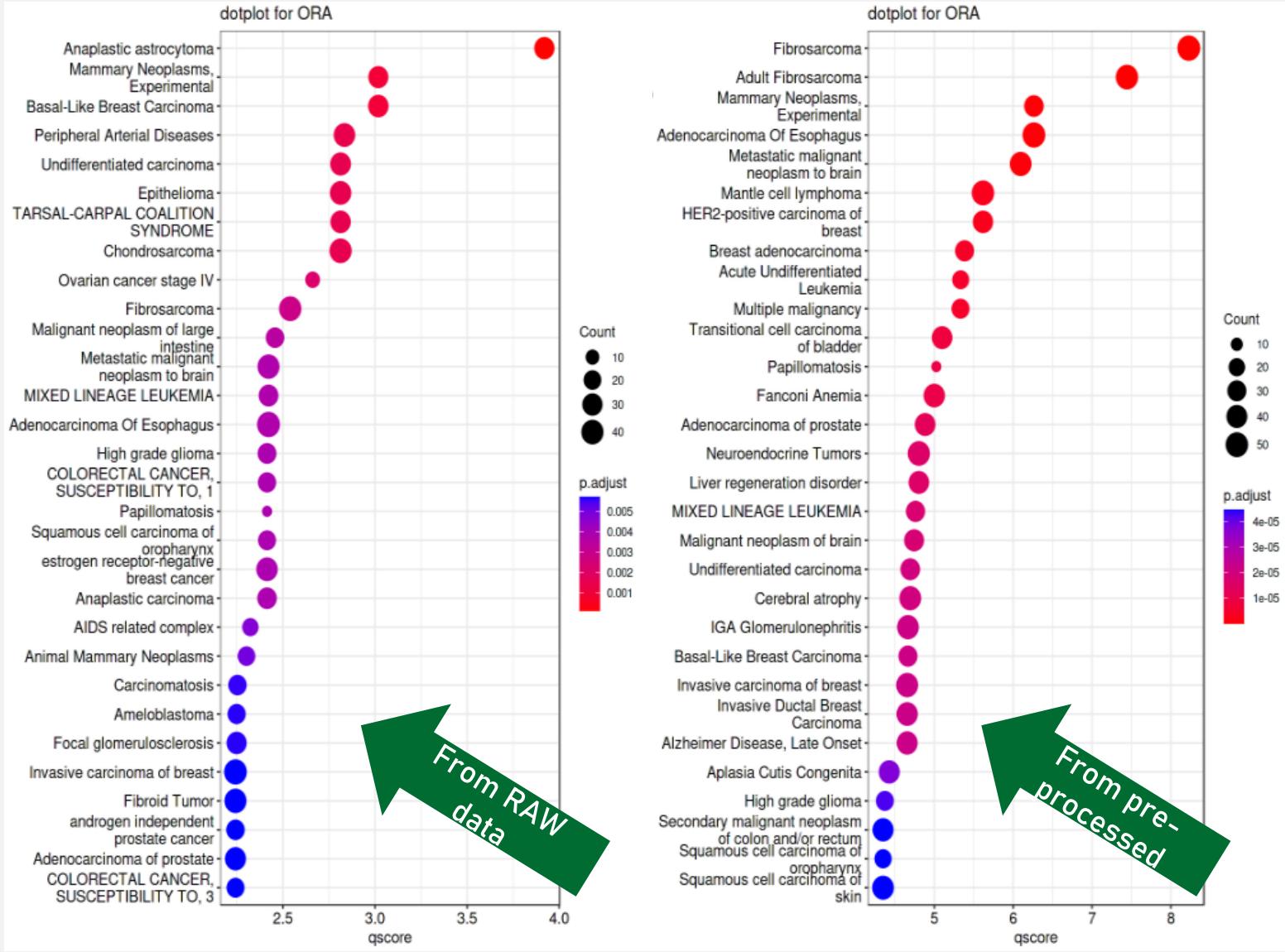
Gene Set Enrichment Analysis
→ GSEA

```
dgnGSE <- gseDGN(geneList,  
|   pvalueCutoff = 0.5,  
|   pAdjustMethod = "BH")  
dgnGSE <- setReadable(dgnGSE, 'org.Hs.eg.db')
```



Breast cancer related

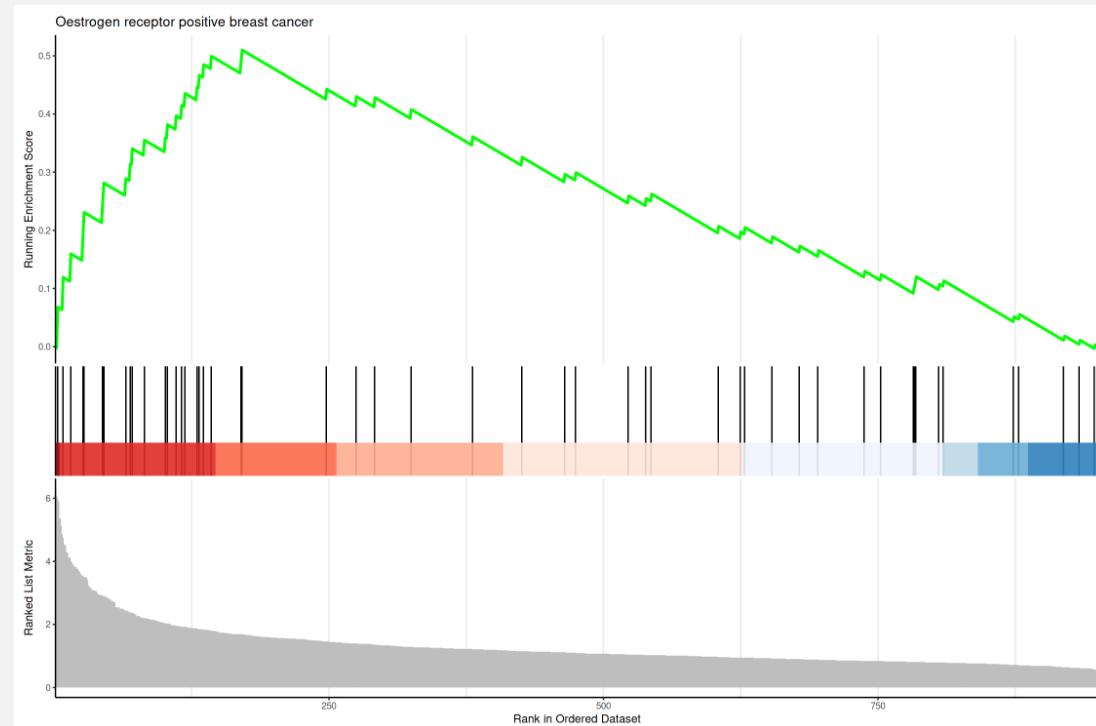
THE INPUT MATTERS



- Same data-set, different results!
- Why?!
 - Different gene annotations?
 - Different alignment?
 - Different quantification strategy?
- Exclude non-protein-coding genes?

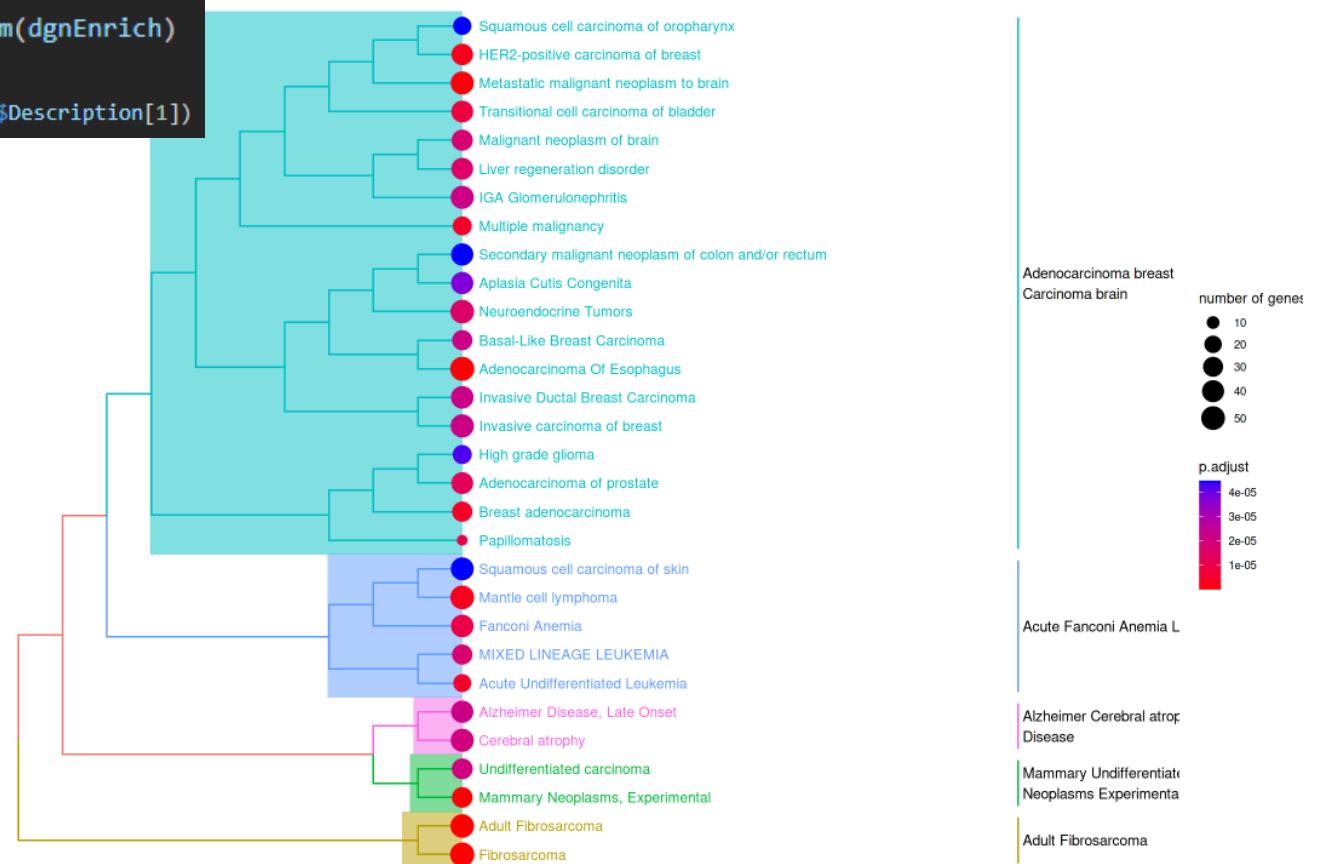
ADVANCED PLOTS

GSEA / Tree clustering



```
dgnEnrich2 <- pairwise_termsim(dgnEnrich)
treeplot(dgnEnrich2)
```

```
gseaplot2(dgnGSE, geneSetID = 1, title = dgnGSE$Description[1])
```



- The `treeplot()` function performs hierarchical clustering of enriched terms/gene sets
 - Jaccard's similarity index across enriched terms

GENE SET ENRICHMENT

Using Metascape

- <http://metascape.org/gp/index.html>
- Use one of the results produced before
 - Here: Top500 regulated genes (logFC), up-regulated
- Select gene identifier
 - PROBEID is a gene synonym
 - H. sapiens is correctly recognized
- Start Express Analysis

The screenshot shows the Metascape web application interface. It consists of three main sections: Step 1, Step 2, and Step 3.

Step 1: A "Multiple Gene Lists" section with a "Select File..." button and a "File Upload" window. The upload window shows a file tree from "This PC > Workshop (W:) > game_cibog > bulk_c" and a list of files including "metascape", "24h_sum159_vs_sum159_jq1r.tsv", and several other tsv files. A "Cancel" button is visible.

Step 2: A "Select Columns in your Excel file." dropdown menu. The options listed are: Row_names (Type: Ensembl Gene), baseMean (Type: Gene ID), log2FoldChange (Type: Gene Synonym), lfcSE (Type: Gene Synonym), pvalue, padj, gene_name (Type: Gene Synonym), and biotype (Type: Gene Synonym). The "gene_name" option is currently selected.

Step 3: Buttons for "Express Analysis" and "Custom Analysis". Below these buttons, there are two sets of input fields:

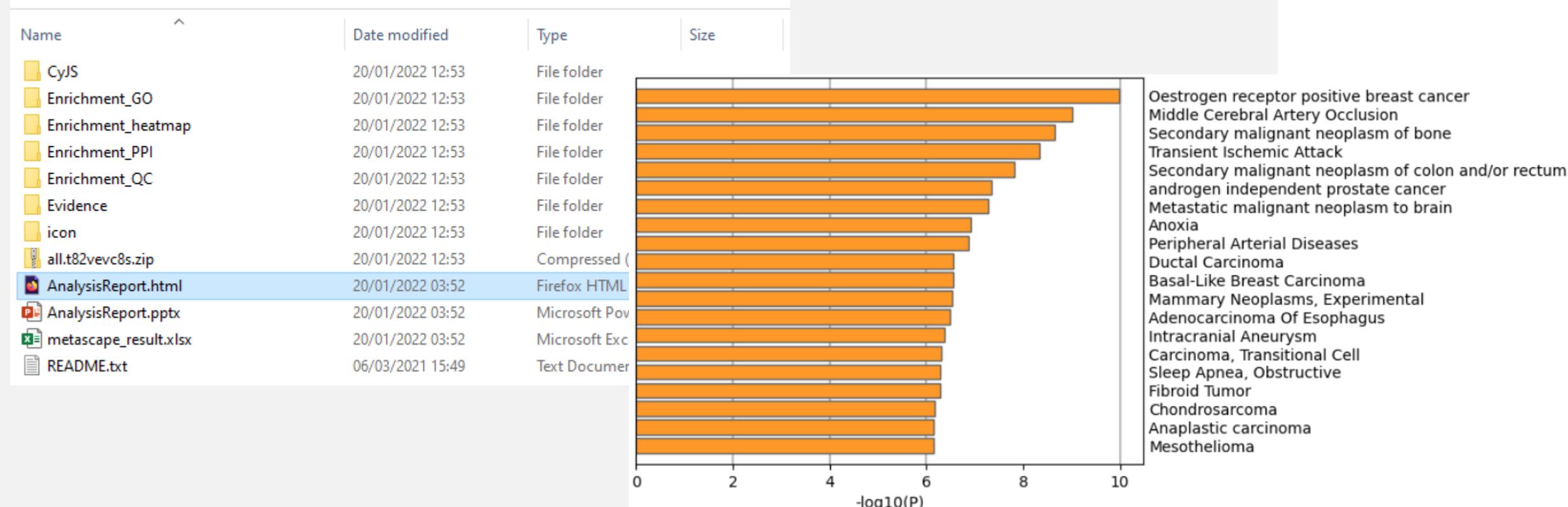
- Step 2: "Optional if you only consider human species in your study." Input as species: Any Species, Analysis as species: H. sapiens (471).
- Step 3: Express Analysis and Custom Analysis buttons.

METASCAPE RESULTS

Save them or they are gone

Gene list enrichments are identified in the following ontology categories: COVID, Cell_Type_Signatures, DisGeNET, PaGenBase, TRRUST, Transcription_Factor_Targets. All genes in the genome have been used as the enrichment background. Terms with a p-value < 0.01, a minimum count of 3, and an enrichment factor > 1.5 (the enrichment factor is the ratio between the observed counts and the counts expected by chance) are collected and grouped into clusters based on their membership similarities. The top few enriched clusters (one term per cluster) are shown in the Figure 4-9. The algorithm used here is the same as that is used for pathway and process enrichment analysis.

PC > Workshop (W:) > game_cibog > bulk_own > metascape_top24h_up



GENE SET ENRICHMENT RESULT COMPARISON

- Within zip-Download all tables available as excel file

Description	PARENT_C	LogP	Enrichmer	Z-score	#TotalGeneInLibrary	#GeneInGO	#GeneInHitList	#GeneInGOAndHitList	%InGO	STDV	%InGenesID	Hits
Oestrogen receptor positive breast cancer		-10	4	8.4	30211	510	463		31	6.7	1.2	374 595 5AREG CCND1 BCL2 KLF5
Middle Cerebral Artery Occlusion		-9	3.4	7.7	30211	626	463		33	7.1	1.2	134 290 3ADORA1 ANPEP AREG
Secondary malignant neoplasm of bone		-8.7	3.3	7.5	30211	647	463		33	7.1	1.2	374 596 1AREG BCL2 CCN2 DPYSL2

- Does not occur in disgenet2r/DOSE analysis!

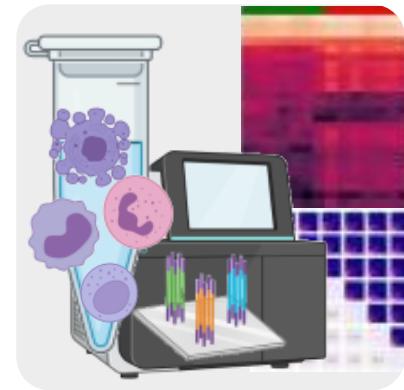
- Different inputs!
- Adj. p-Value sorted gene list in R
- FC-sorted gene list for Metascape

- Not only statistical test highly relevant, also selected genes!

```
metascapeHits = c("AREG", "CCND1", "BCL2", "KLF5", "CHKA", "EPHB2", "EREG", "FYN",  
"GLI2", "GSN", "FOXA1", "CXCL8", "KDR", "MAPT", "MMP1", "RAB27B", "SLC20A1", "TRPM2",  
"FOSL1", "STC2", "MAP4K4", "ADAMTS2", "SEMA4D", "TRIM29", "SGK3", "PAG1", "LGR6",  
"UBASH3B", "DNER", "E2F7", "CD24")  
doseInputGenes_pval = head(outres_24h$gene_name, n=nsamples)  
pvalSortedOverlap = length(intersect(metascapeHits, doseInputGenes_pval))  
doseInputGenes_fc = head(top24h.de.500.up$gene_name, n=nsamples)  
fcSortedOverlap = length(intersect(metascapeHits, doseInputGenes_fc))  
  
print(paste("Metascape set size:", length(metascapeHits)))  
print(paste("Pval Sorted Overlap:", pvalSortedOverlap))  
print(paste("logFC Sorted Overlap:", fcSortedOverlap))  
✓ 0.5s  
[1] "Metascape set size: 31"  
[1] "Pval Sorted Overlap: 22"  
[1] "logFC Sorted Overlap: 31"
```

SUMMARY

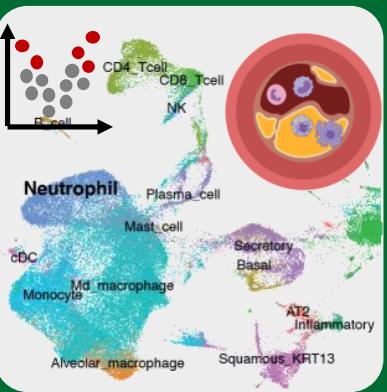
- Analyse public data from raw data (reads)
 - Or from pre-processed data
- Perform differential expression analysis using DESeq2
- Perform Gene Set Enrichment Analyses of significant genes
 - R-based: clusterProfiler, DOSE, disgenet2r, ...
 - Web-based: Metascape
- At each “pipeline” step the choice of tools, methods and parameters affects the final result
 - Sorting of lists, number of genes for set enrichment, etc.



bulk RNA-seq

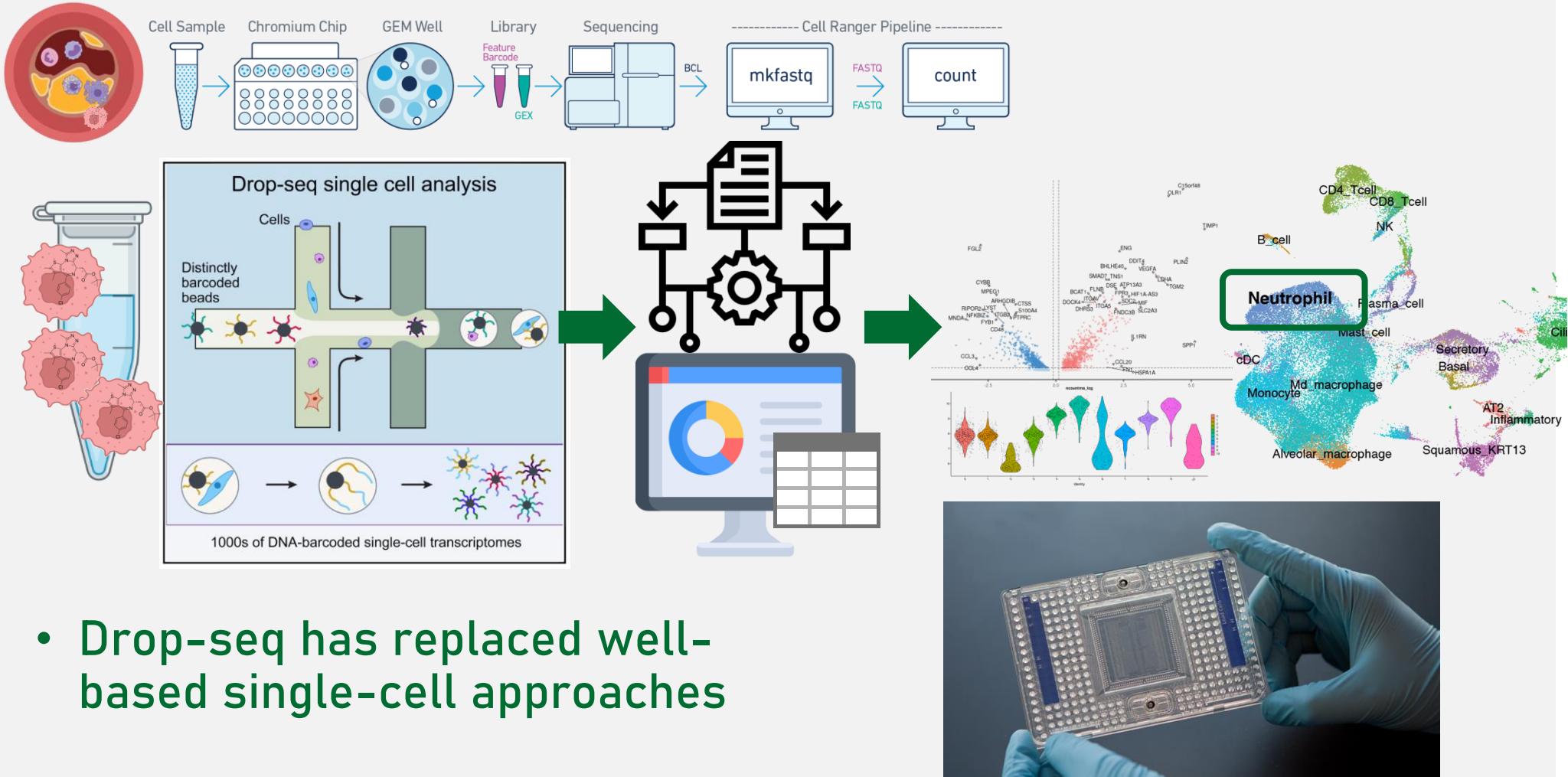


bulk RNA-seq



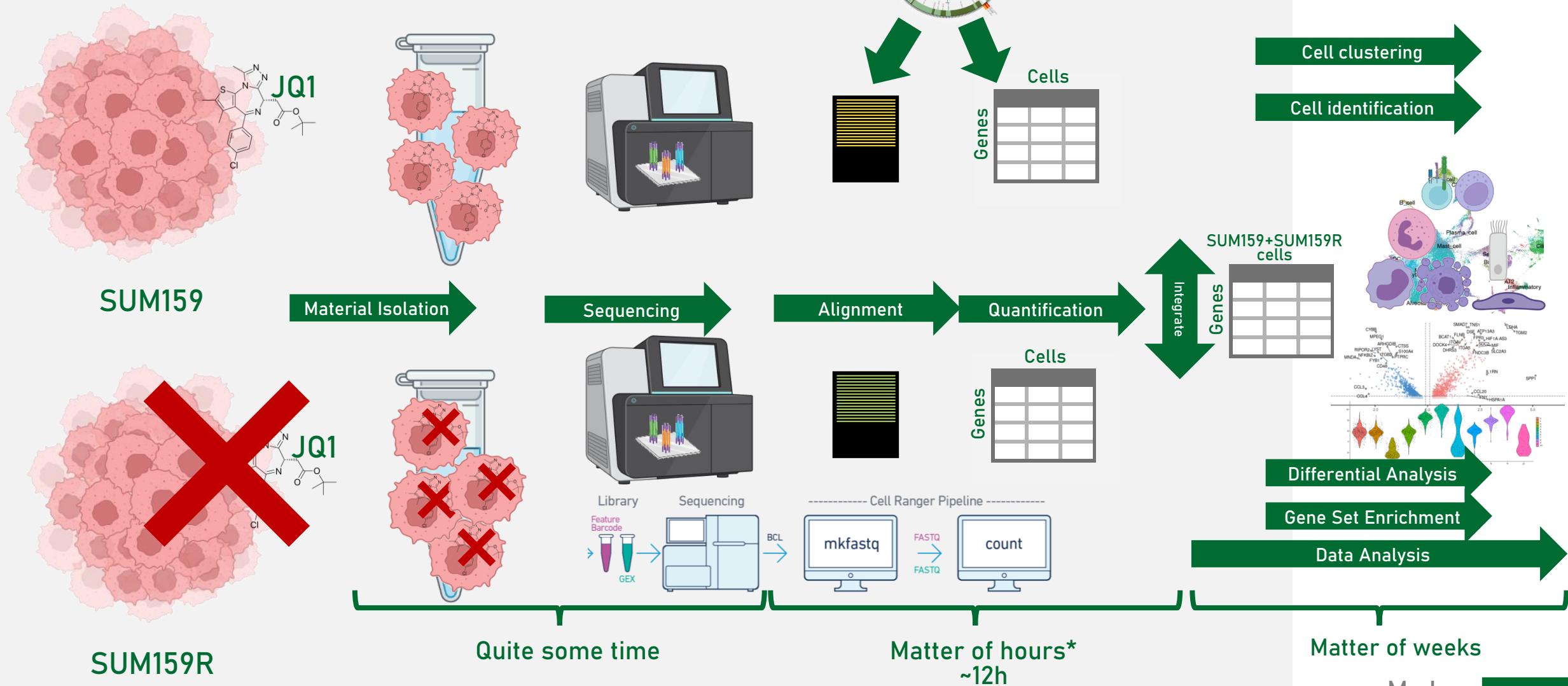
scRNA-seq

SINGLE CELL TRANSCRIPTOMICS



- Drop-seq has replaced well-based single-cell approaches

SINGLE CELL RNA-SEQ



DOWNLOADING CELLRANGER+LOUPE

Cell Ranger - 6.1.2 (October 25, 2021)

- Self-contained, relocatable tar file. Does not require centralized installation.
- Contains binaries pre-compiled for CentOS/RedHat 6.0+ and Ubuntu 12.04+.
- [Download - Linux 64-bit - 768 MB](#) - md5sum: 310d4453acacf0eec52e76aded14024c

curl

wget

```
curl https://cf.10xgenomics.com/releases/cell-exp/cellranger-6.1.2.tar.gz?Expires=1642730466&Policy=eyJTdGF0ZW1lbQi0ltI1l1c291cmN1ijoiaHR0cHM6Ly9jZi4xMhnnZW5vbWljcy5jb20vcmVsZWfzZXMyY2VsbC1leHAvY2VsbHJhbmdlci02LjEuMi50YXIUzoiLCJDb25kaXRpb24iOnsiRGF0ZUxlc3NUjaGFuIjp7IkFXUzpFcG9jaFRpbkUiOjE2NDI3MzA0NjZ9fx1dfQ_&Signature=USFWTdaZIygJwvTt1CSo7xZdWZmDq2CH-w~0t1Wg1CuDTVNLKbGcjYXLNBfIAWeaCckQu4AJrw- bbNRGwF0MS9oB5WheoRAopZyKXuNaCt97tqKs9NN1UpwAbRdgD0JFz1QdKMCT3cg3~AXRLOxFnxlnNVvfjcpcjrnN1dXv7qcS862t VezERTC0OKrwCR1syjFnCKxmHhi04xPRAdoQs8UhoFl0XLqB-WsLM7f~CSNj4~dV-RfRGS- 0ePVuU-s2mOGdDtbOF66je612kZo433Z51eX11GWK6ktwpJEBSqkzNoJFYCDjTw0HCXU~F4X8nTh6rpGNEv5a8QEHa__&Key-Pair-Id=APKA17S6A5RYOXBWRPDA"
```

- Human reference (GRCh38) dataset required for Cell Ranger.
- [Download - 11 GB](#) - md5sum: dfd654de39bff23917471e7fcc7a00cd
- [Build steps](#)

curl

wget

```
curl https://cf.10xgenomics.com/supp/cell-exp/refdata-gex-GRCh38-2020-A.tar.gz
```

- <https://support.10xgenomics.com/single-cell-gene-expression/software/downloads/latest>
- Might require brief registration



Cell Ranger 6.1.2

Single Cell Analysis Pipelines



Loupe Browser 6.0.0

Interactive Analysis

Loupe Browser 6.0.0 (October 7, 2021)

Please follow the [install instructions](#) after downloading one of the installers below.

[Read about what's new in Loupe Browser 6.0](#)

Download for Windows

filesize: 694 MiB

md5sum: 8fce42ab47c3d040d7a1538be8d33900

Download for macOS

filesize: 784 MiB

md5sum: aee779daa98bc101092169666fd1c2aa

PREPARING SINGLE CELL ANALYSES

Good folder structure is crucial!

```
└─ SC
    └─ calls
        ├─ SUM159DMSO
        ├─ SUM159RDMSO
        └─ calls.sh
            └─ SUM159DMSO.err
            └─ SUM159DMSO.out
            └─ SUM159RDMSO.err
            └─ SUM159RDMSO.out
    └─ samples
        └─ SUM159DMSO.csv
        └─ SUM159RDMSO.csv
            └─ SUM159DMSO.err
            └─ SUM159DMSO.out
            └─ SUM159RDMSO.err
            └─ SUM159RDMSO.out
    └─ slides
    └─ tools
        └─ cellranger-6.1.2
        └─ sratoolkit.2.11.2-ubuntu64
            └─ cellranger-6.1.2.tar.gz
            └─ sratoolkit.tar.gz
```

- **sc-dir**

- **Calls subfolder:** script to run cellranger+output
 - **Samples:** sample description files

- **Tools**

- **Cellranger executable**

THE SAMPLE SPECIFICATION

```
> samples > SUM159DMSO.csv
1 fastqs, sample, library_type
2 /mnt/w/game_cibog/reads/SUM159DMSO/, SRR9050597, Gene Expression
3 /mnt/w/game_cibog/reads/SUM159DMSO/, SRR9050598, Gene Expression
4 /mnt/w/game_cibog/reads/SUM159DMSO/, SRR9050599, Gene Expression
5 /mnt/w/game_cibog/reads/SUM159DMSO/, SRR9050600, Gene Expression
```

```
> samples > SUM159RDMSO.csv
1 fastqs, sample, library_type
2 /mnt/w/game_cibog/reads/SUM159RDMSO/, SRR9050605, Gene Expression
3 /mnt/w/game_cibog/reads/SUM159RDMSO/, SRR9050606, Gene Expression
4 /mnt/w/game_cibog/reads/SUM159RDMSO/, SRR9050607, Gene Expression
5 /mnt/w/game_cibog/reads/SUM159RDMSO/, SRR9050608, Gene Expression
```

- Contains paths to all reads for a sample
- Specifies whether reads are from gene expression or antibody capture (e.g. CITE-seq)

THE CALLS SCRIPT

```
3 CURDIR=`pwd`  
4  
5 #path to cellranger binary/executable  
6 CELLRANGER=$CURDIR/tools/cellranger-6.1.2/bin/cellranger  
7 #path to the reference folder  
8 REFFOLDER_HUMAN=$CURDIR/references/refdata-gex-GRCh38-2020-A  
9  
10 #this is for parallelizing cellranger. choose according to your infrastructure!  
11 OPTS="--localcores 6 --localmem 30"  
12  
13 ls $CELLRANGER  
14  
15 BASE=../sc/  
16  
17 #for each sample run cellranger  
18 for SAMPLE in SUM159DMSO SUM159RDMSO;  
19 do  
20  
21 #path to sample file for current sample  
22 SAMPLEFILE=$CURDIR/$BASE/samples/$SAMPLE.csv  
23  
24 #verify samplefile exists  
25 ls $SAMPLEFILE  
26  
27 #if we rerun this sample, delete existing results  
28 rm -rf $BASE/$SAMPLE  
29  
30 # run cellranger for the selected sample, with given reference and given read-files and optional parameters.  
31 # save output in $SAMPLE.out and redirect error-messages to $SAMPLE.err  
32 # wrapping $CELLRANGER into /usr/bin/time --verbose will print how long the run took and how much RAM was used lateron  
33 (cd $BASE/calls/ && /usr/bin/time --verbose $CELLRANGER count --id=$SAMPLE --transcriptome=$REFFOLDER_HUMAN --libraries=$SAMPLEFILE $OPTS > $SAMPLE.out 2> $SAMPLE.err)  
34  
35 done
```

RUNNING CELLRANGER

Organisation matters

Running cellranger

```
! bash sc/calls/calls.sh
```

CELLRANGER OUTPUT

During the run

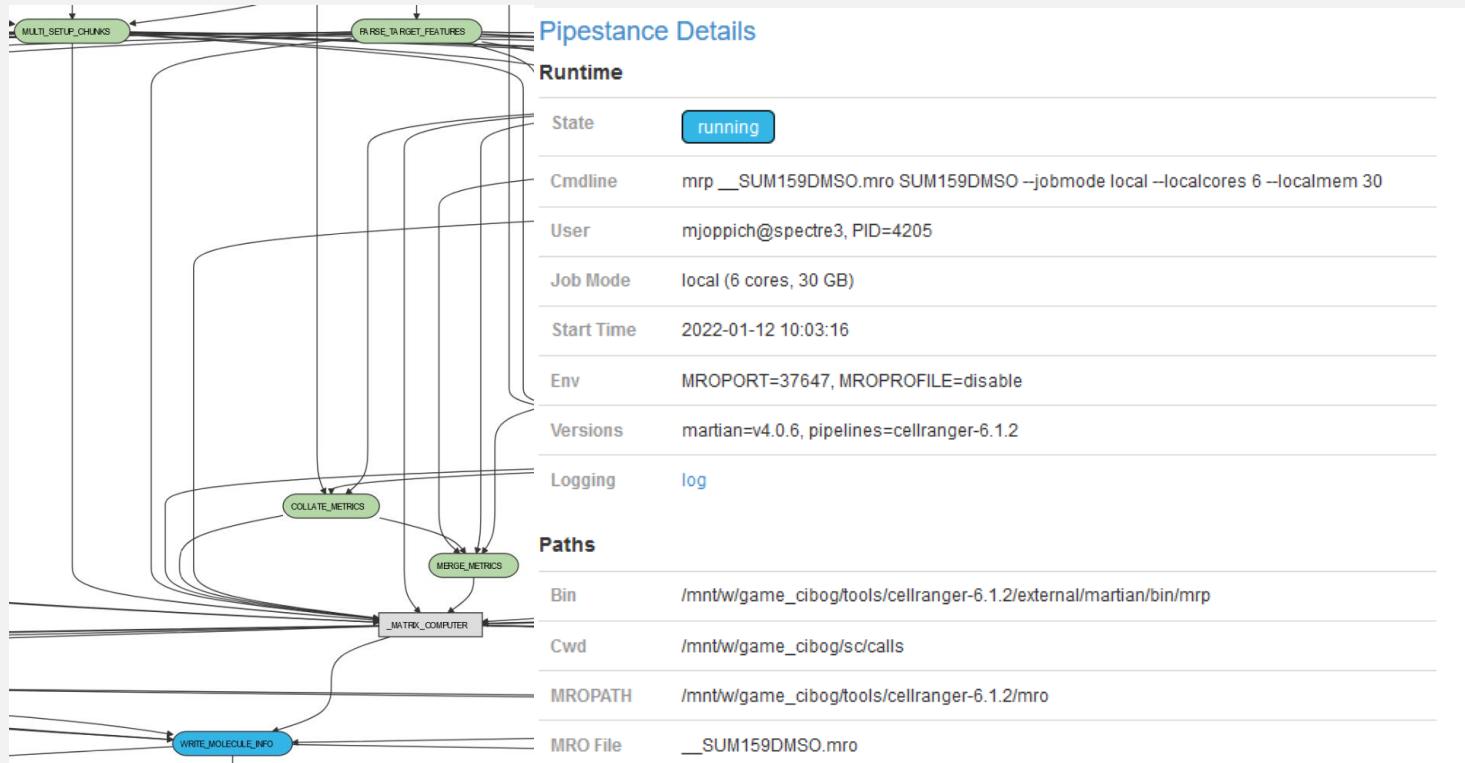
```
1 Martian Runtime - v4.0.6
2 Serving UI at http://spectre3:37647?auth=2WXiTSo22fyUhBGH5BoHt-G9wr50RVX096wuG0F9yCE
3
4 Running preflight checks (please wait)...
5 Checking sample info...
6 Checking FASTQ folder...
7 Checking reference...
8 Checking reference_path (/mnt/w/game_cibog/references/refdata-gex-GRCh38-2020-A) on spectre3...
9 Checking optional arguments...
10 mrc: v4.0.6
11
12 mrp: v4.0.6
13
14 Anaconda: Python 3.8.2
15
16 numpy: 1.19.2
17
18 scipy: 1.6.2
19
20 pysam: 0.16.0.1
21
22 h5py: 3.2.1
23
24 pandas: 1.2.4
25
26 STAR: 2.7.2a
27
28 samtools: samtools 1.10
29 Using htseq 1.10.2
30 Copyright (C) 2019 Genome Research Ltd.
31
32 2022-01-12 10:03:28 [runtime] (ready)
33 2022-01-12 10:03:28 [runtime] (run:local)
34 2022-01-12 10:03:28 [runtime] (ready)
```

Elapsed (wall clock) time (h:mm:ss or m:ss): 5:10:07
Average shared text size (kbytes): 0
Average unshared data size (kbytes): 0
Average stack size (kbytes): 0
Average total size (kbytes): 0
Maximum resident set size (kbytes): 15888584

Web Interface

Logging output

CELLRANGER PROGRESSION



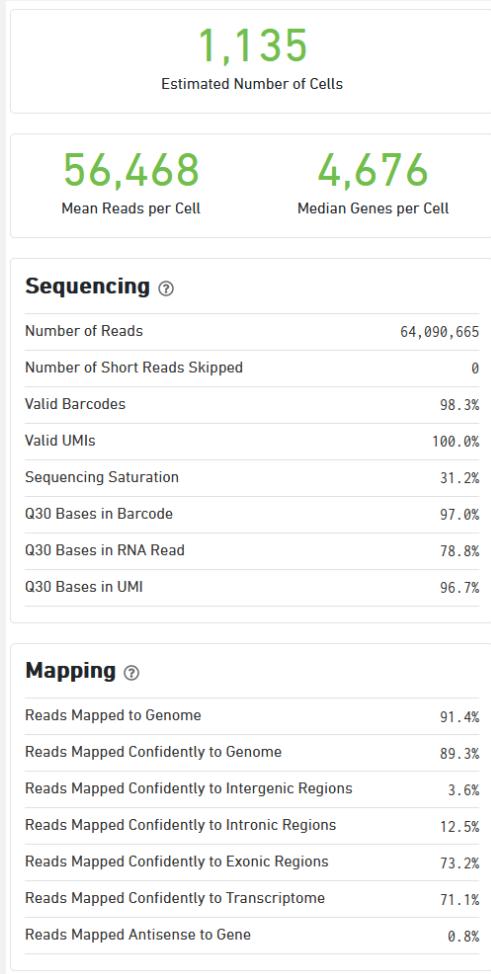
- You can also keep track on progression with web interface
 - But you gotta wait long enough anyway ;)

CELLRANGER OUTPUT

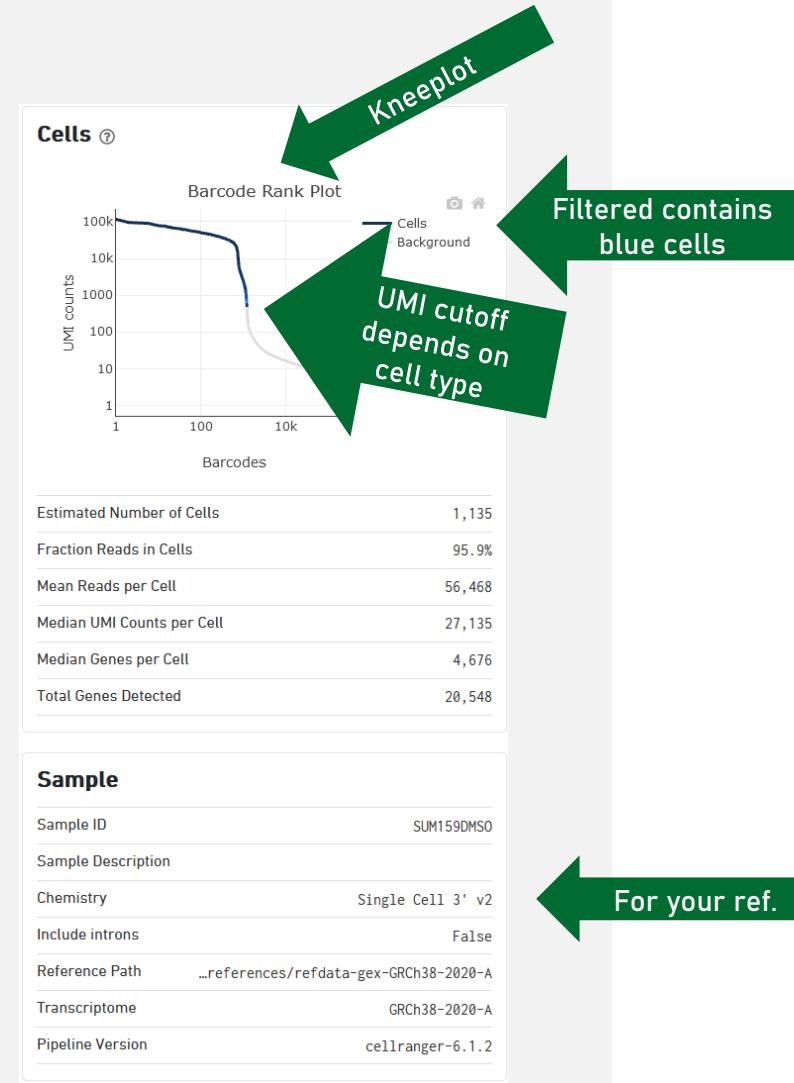
PC > Workshop (W:) > game_cibog > sc > calls > SUM159DMSO > outs >			
Name	Date modified	Type	Size
analysis	12/01/2022 15:11	File folder	
filtered_feature_bc_matrix	12/01/2022 12:18	File folder	
raw_feature_bc_matrix	12/01/2022 12:09	File folder	
cloupe.cloupe	12/01/2022 15:12	CLOUPE File	28.213 KB
filtered_feature_bc_matrix.h5	12/01/2022 12:18	H5 File	6.017 KB
metrics_summary.csv	12/01/2022 15:12	Microsoft Excel C...	1 KB
molecule_info.h5	12/01/2022 15:06	H5 File	94.008 KB
possorted_genome_bam.bam	12/01/2022 12:14	BAM File	4.820.195 KB
possorted_genome_bam.bam.bai	12/01/2022 12:16	BAI File	3.650 KB
raw_feature_bc_matrix.h5	12/01/2022 12:09	H5 File	10.976 KB
web_summary.html	12/01/2022 15:12	Firefox HTML Doc...	3.246 KB

- **web_summary.html** for quick look at the data
 - Alignment rate, etc.
- **cloupe.cloupe** file can be used for analysis
- **raw_feature_bc_matrix**
 - File and folder
 - All barcodes
- **Filtered_feature_bc_matrix**
 - Cellranger filtered barcodes

WEB SUMMARY



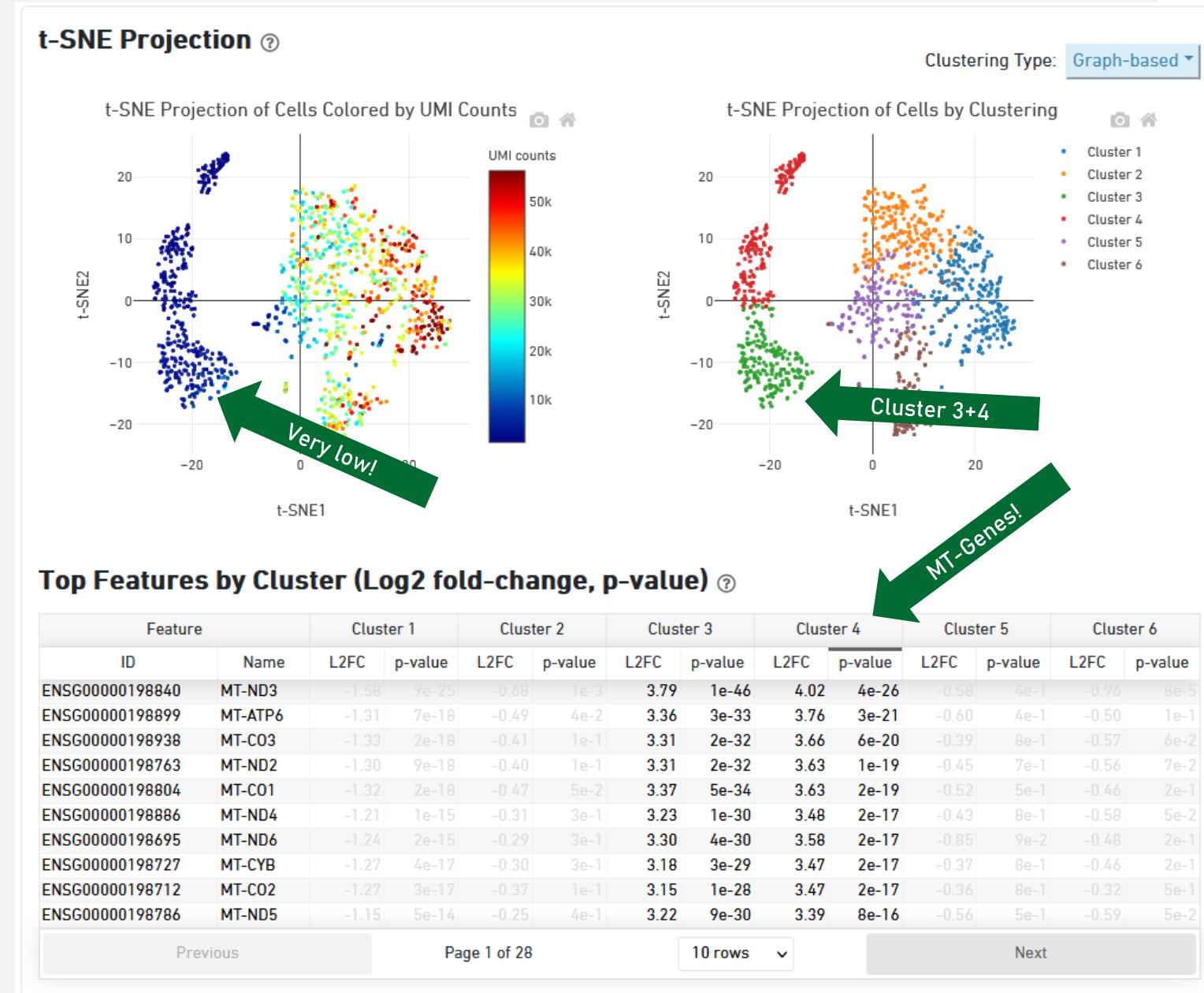
- Estimated number of cells
- Mean reads/genes per cell
- Sequencing information
 - About 25.000-50.000 reads per cell should be planned
- Mapping statistics
 - Most reads should map to exon and transcriptome!



WEB SUMMARY

Very brief analysis

- Initial clustering of all filtered cells
- Marker genes per cluster
- Spot strange clusters which might need further attention
 - High MT-content = damaged cells



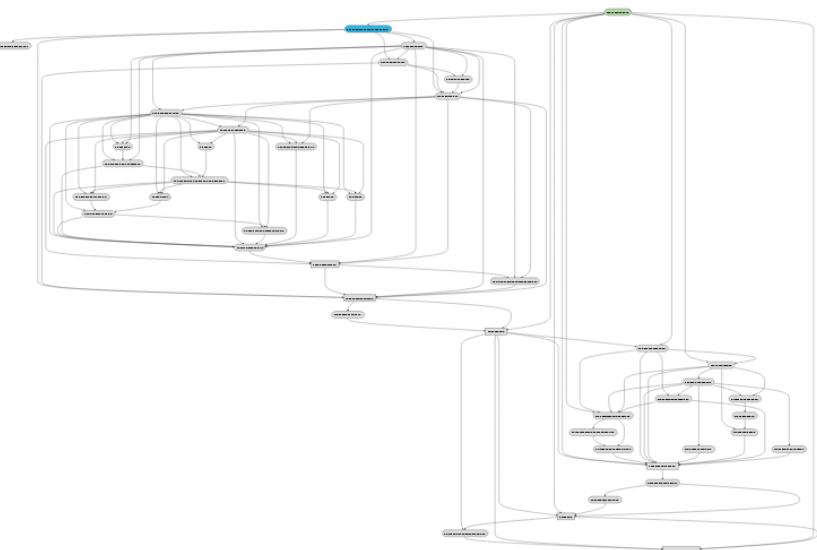
CELLRANGER AGGREGATE

Combines two separate libraries

After having executed cellranger, we must merge both datasets:

```
! cd sc/aggr_calls/ && ../../tools/cellranger-6.1.2/bin/cellranger aggr --localmem 30 --localcores 8 --csv aggr_samples.csv --id SUM159DMSO_RDMSO
```

10X Martian Pipeline Runner / mjoppich / SUM159DMSO_RDMSO / SC_RNA_AGGREGATOR_CS



sc > aggr_calls > aggr_samples.csv

1 sample_id,molecule_h5

2 SUM159DMSO,/mnt/w/game_cibog/sc/calls/SUM159DMSO/outs/molecule_info.h5

3 SUM159RDMSO,/mnt/w/game_cibog/sc/calls/SUM159RDMSO/outs/molecule_info.h5

Absolute paths

Pipestance Details

Runtime

State running

Cmdline mrp __SUM159DMSO_RDMSO.mro SUM159DMSO_RDMSO --jobmode local --localcores 8 --localmem 30

User mjoppich@spectre3, PID=20549

Job Mode local (8 cores, 30 GB)

Start Time 2022-01-12 21:18:27

Env MROPORT=33385, MROPROFILE=disable

Versions martian=v4.0.6, pipelines=cellranger-6.1.2

Logging log

Paths

Bin /mnt/w/game_cibog/tools/cellranger-6.1.2/external/martian/bin/mrp

Cwd /mnt/w/game_cibog/sc/aggr_calls

MROPATH /mnt/w/game_cibog/tools/cellranger-6.1.2/mro

MRO File __SUM159DMSO_RDMSO.mro

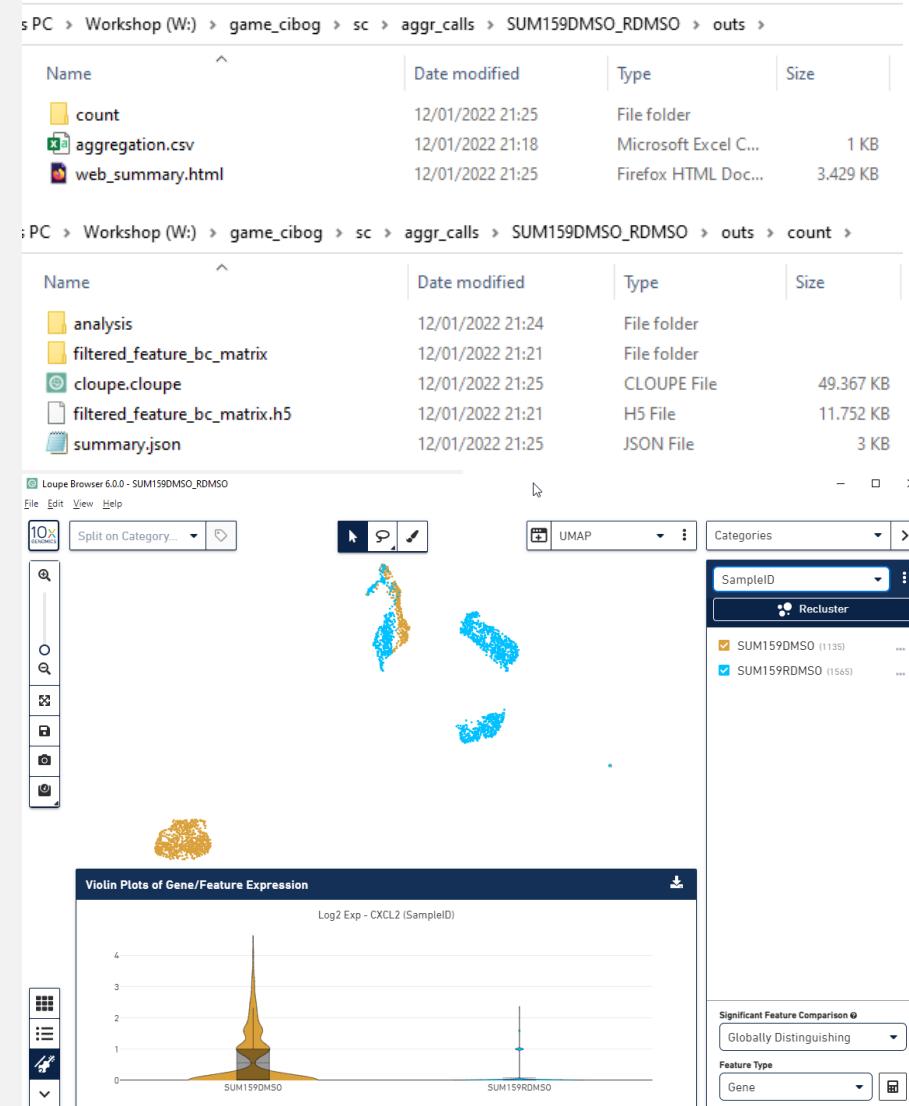
```
@include "rna/sc_rna_aggregator_cs.mro"
```

```
call SC_RNA_AGGREGATOR_CS(
    sample_id      = "SUM159DMSO_RDMSO",
    sample_desc    = "",
    pipestance_root = "/mnt/w/game_cibog/sc/aggr_calls",
    aggregation_csv = "/mnt/w/game_cibog/sc/aggr_calls/aggr_samples.csv",
    normalization_mode = "mapped",
    no_secondary_analysis = false,
)
```

CELLRANGER AGGREGATE OUTPUT

Similar to other outputs

- All “filtered” cells from two cellranger count runs
 - With Seurat or scanpy cell selection possible
 - Maybe the low UMI-count cells are actual cells?
- Integration procedure? Unknown.
 - With Seurat or scanpy multiple ways possible
 - SCTransform, bbknn, rpca, etc.
- cloupe file allows further analysis without using command line tools
 - Usage of scanpy/Seurat strongly encouraged!
 - But worth another workshop!



SCRNA-SEQ ANALYSIS WITH LOUPE

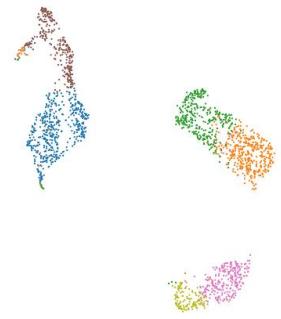
Just the most important things

- Complete tutorial:

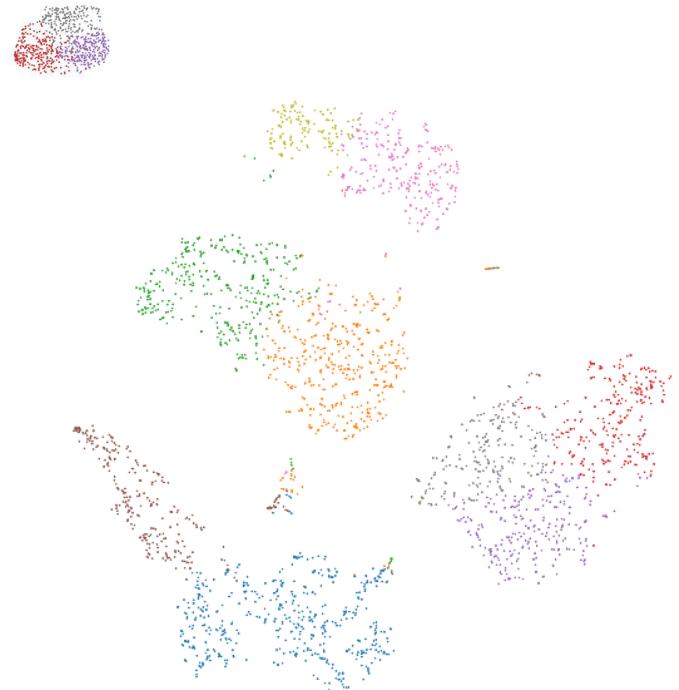
<https://support.10xgenomics.com/single-cell-gene-expression/software/visualization/latest/tutorial>

- Loupe allows

- Visualizations
 - UMAP, t-SNE (whereas UMAP preserves local and global distances better)
 - Violin plots, Heatmaps, Feature plots
- Re-analysis

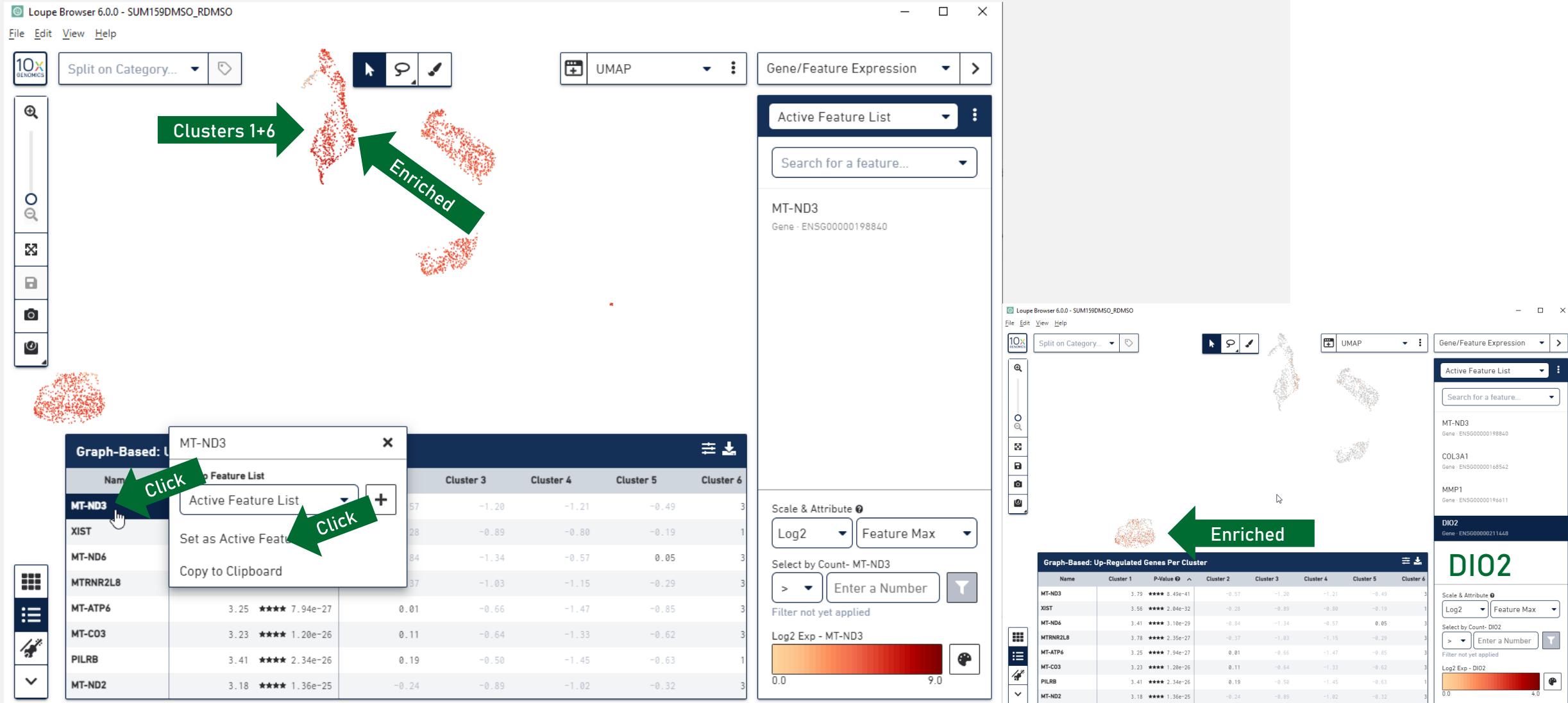


UMAP

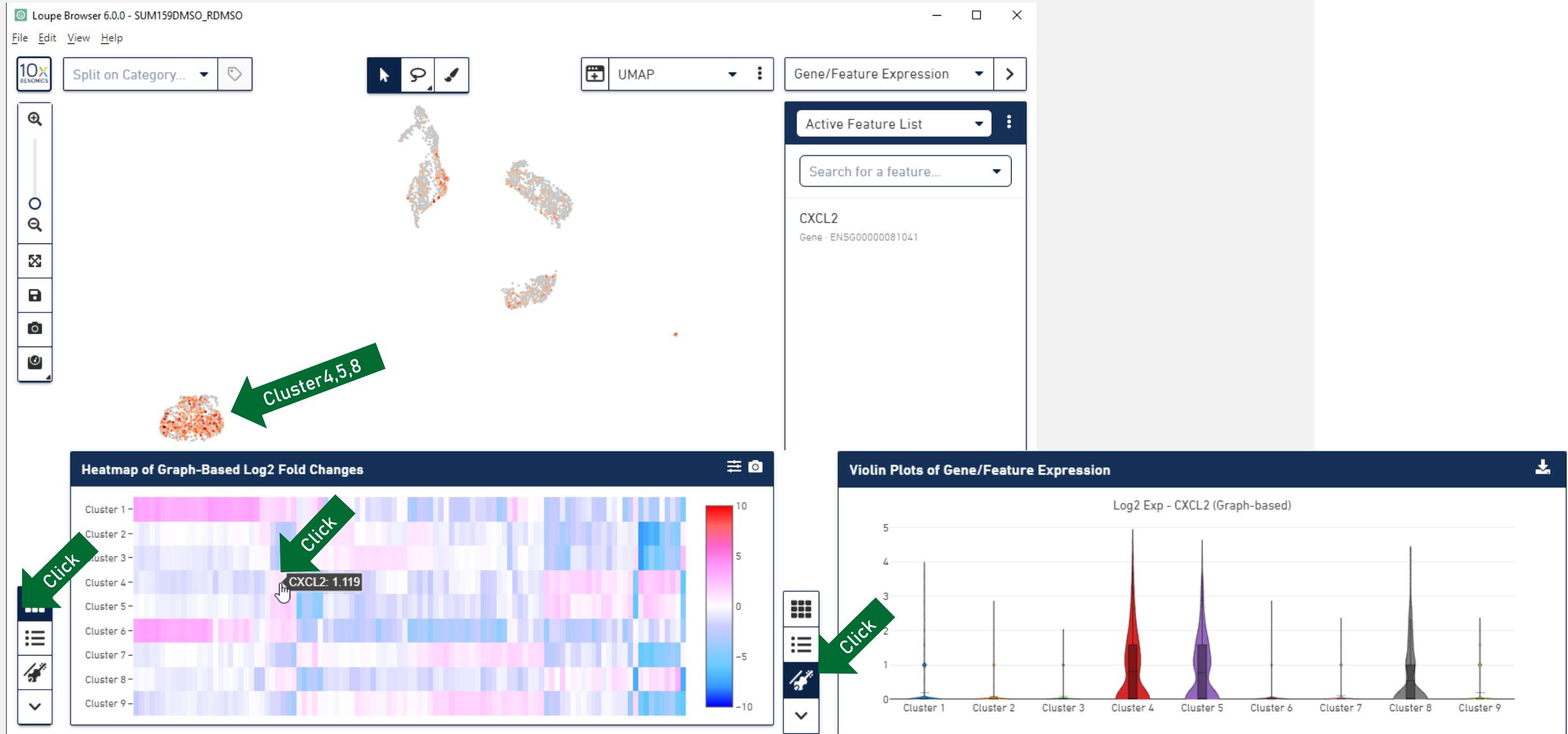


t-SNE

LOUPE: INSPECT FEATURES



LOUPE: HEATMAP



LOUPE: RECLUSTER

- Performs a new analysis
- Let's you remove unwanted cells
 - Such as cells with too high mt-Content

Loupe Browser 6.0.0

Categories
Using the clusters selected below, launch a wizard
to generate a reclustered t-SNE projection

Recluster

Each of the reclustering steps below are optional:

- Cluster 1 (496)
- Cluster 2 (426)
- Cluster 3 (360)
- Cluster 4 (272)
- Cluster 5 (271)
- Cluster 6 (269)
- Cluster 7 (260)
- Cluster 8 (218)
- Cluster 9 (128)

Review Barcodes

Review the barcodes to filter or upload a .csv of barcodes. After this step, you'll have the option to set thresholds to remove poor-quality cells and/or adjust parameters like the number of PCAs. You will then have the option to generate a new t-SNE and/or UMAP projection.

At this time, you may recluster up to 100,000 barcodes.

To adjust which clusters are included in the recluster, select or de-select clusters on the primary Loupe window.

Current selected barcodes

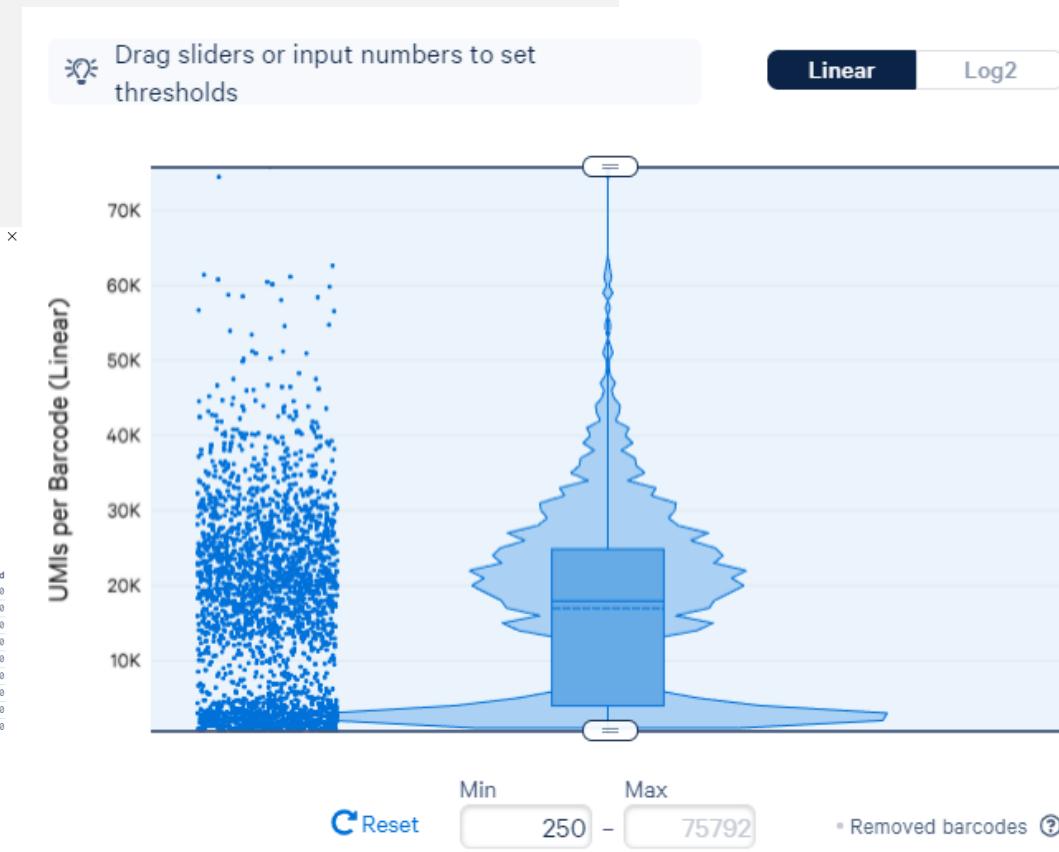
Removed barcodes [?](#)

You may also add to or remove barcodes from the current barcode selection [?](#)

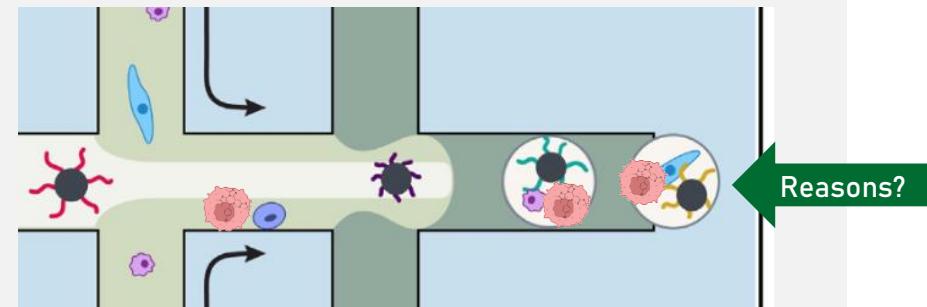
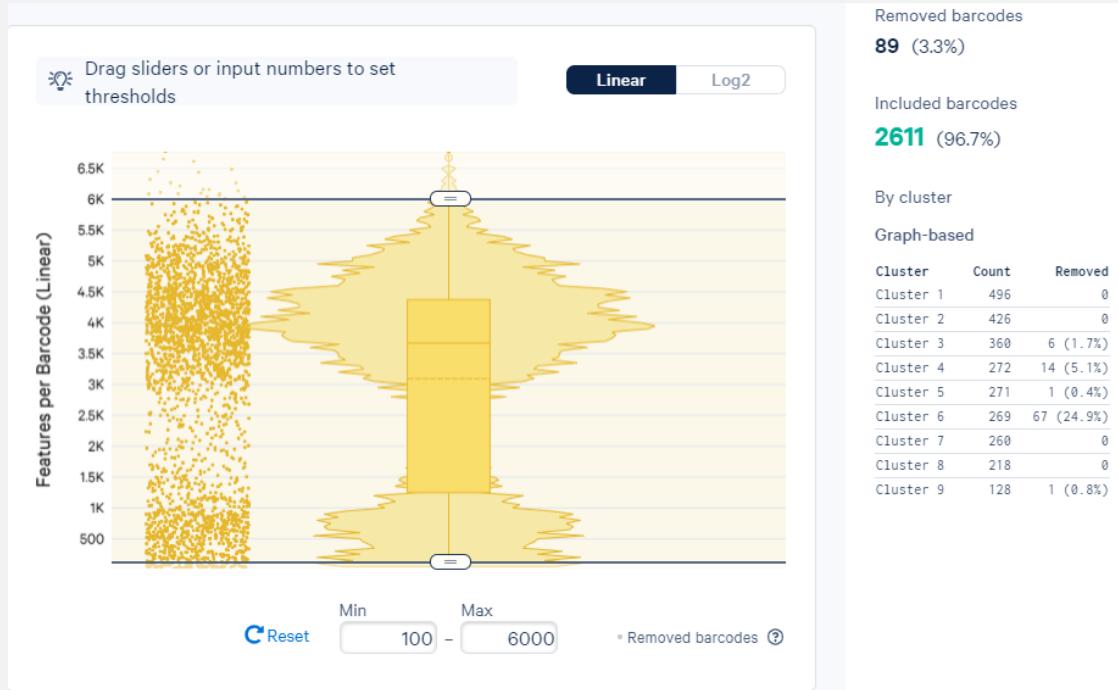
Add barcodes in CSV Remove barcodes in CSV [Upload CSV](#)

Skip to final step [Next](#)

[View our Recluster FAQs](#)

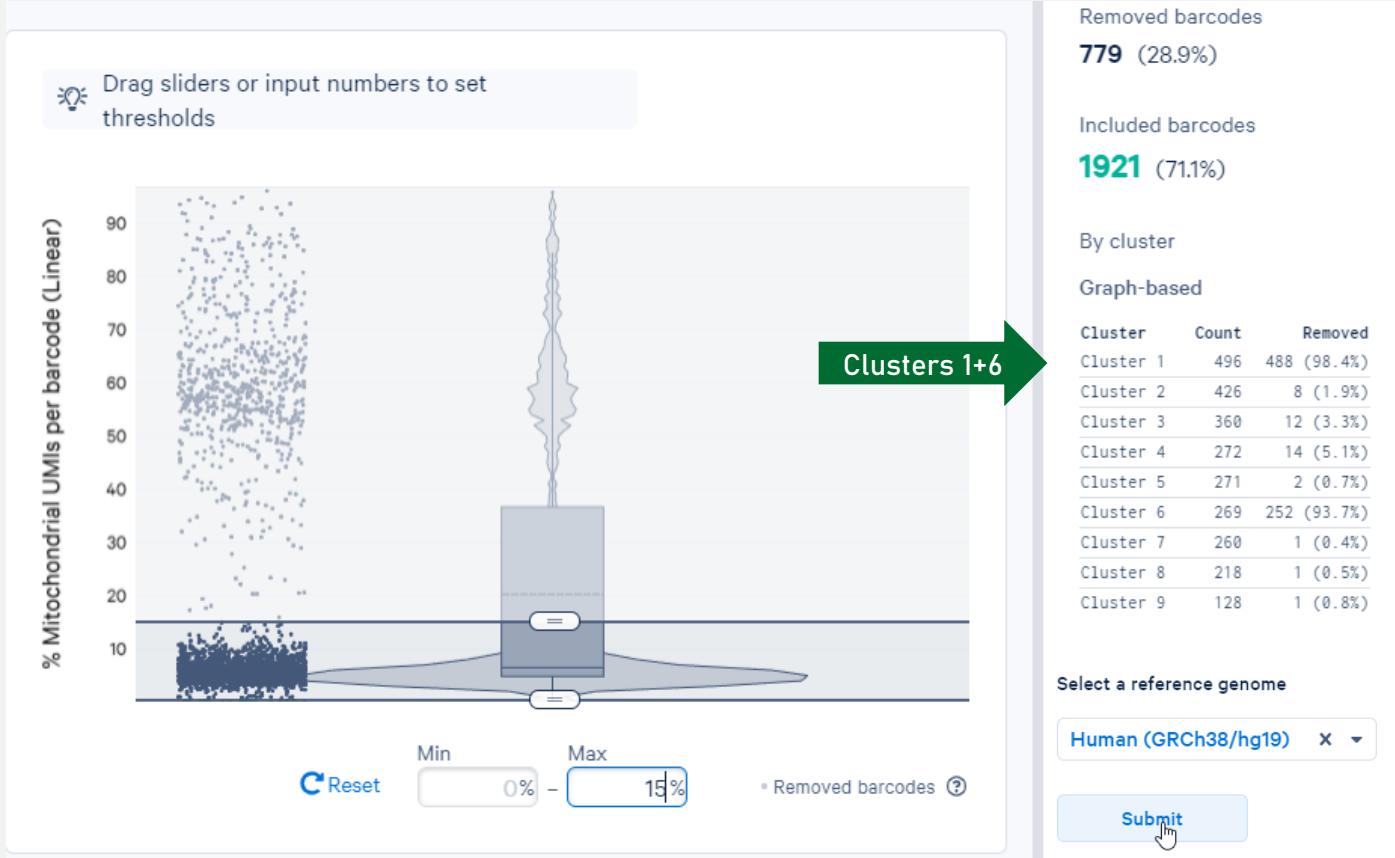


LOUPE: RECLUSTER



- Less than 100 genes might just be fragmented cells
- More than 6000 cells might be doublets
 - 2 cells per droplet

LOUPE: RECLUSTER



- 15% mt-Content is maximum
 - But depends on cell type
- Often also just 5%

LOUPE: RECLUSTER

Recluster

Select the type of plot to generate *

t-SNE UMAP

Select at least one type of plot to generate.

Adjust reanalyze parameters (for advanced users) ▾

Name the recluster *

Give it a good name

250UMI_100_6000FEATURE_MT15

This will generate a new view on the same Loupe file.
You'll still be able to refer back to the original view.

You may click on any of the steps in the left panel to go back and make changes.

Total Results

Starting point

2700

Removed barcodes

779 (28.9%)

Included barcodes

1921 (71.1%)

By cluster

Graph-based

Cluster	Count	Removed
Cluster 1	496	488 (98.4%)
Cluster 2	426	8 (1.9%)
Cluster 3	360	12 (3.3%)
Cluster 4	272	14 (5.1%)
Cluster 5	271	2 (0.7%)
Cluster 6	269	252 (93.7%)
Cluster 7	260	1 (0.4%)
Cluster 8	218	1 (0.5%)
Cluster 9	128	1 (0.8%)

 Processing, do not close this window

This might take several minutes.

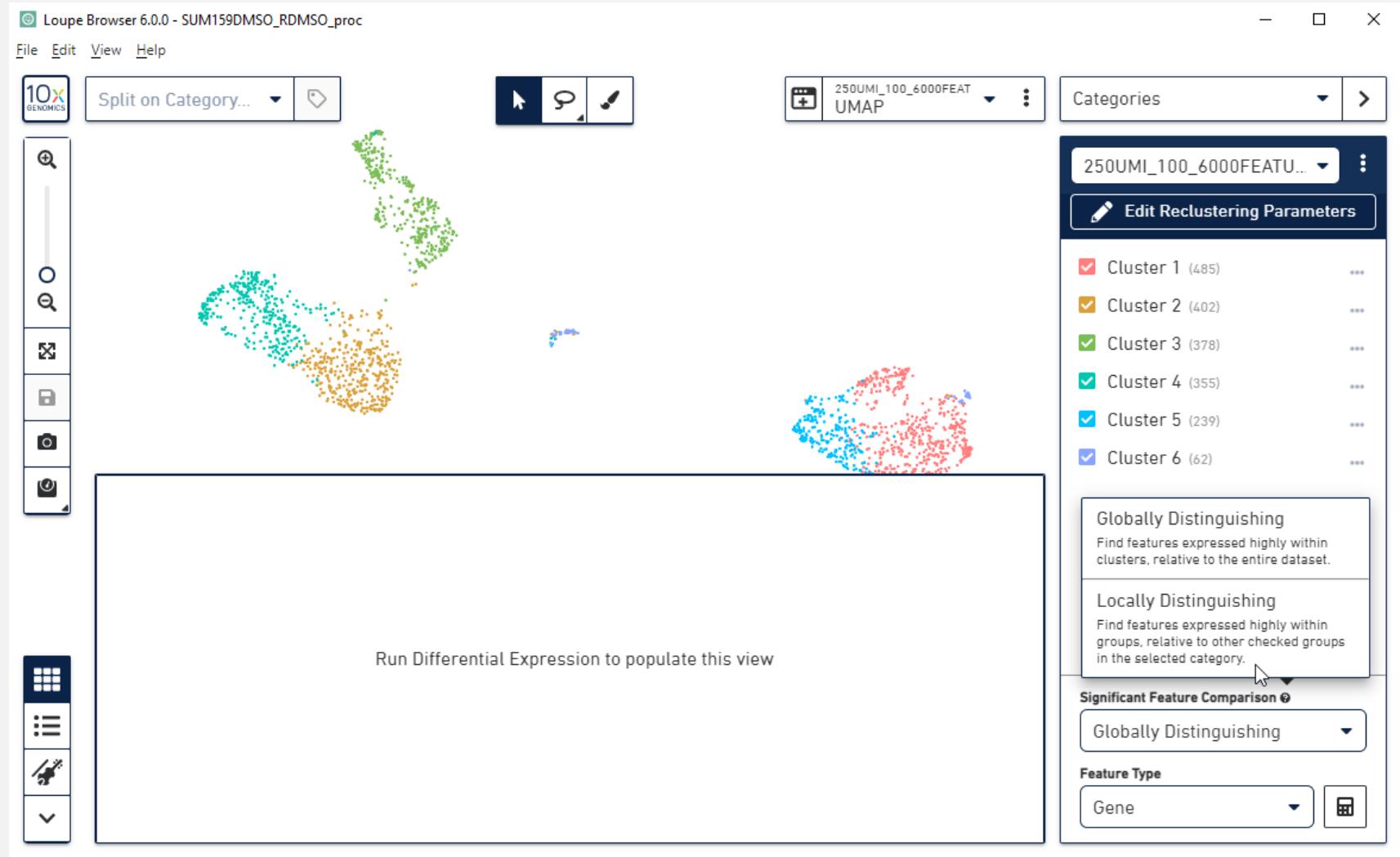


Waiting for your reclustering?

[View our reclustering FAQs](#)

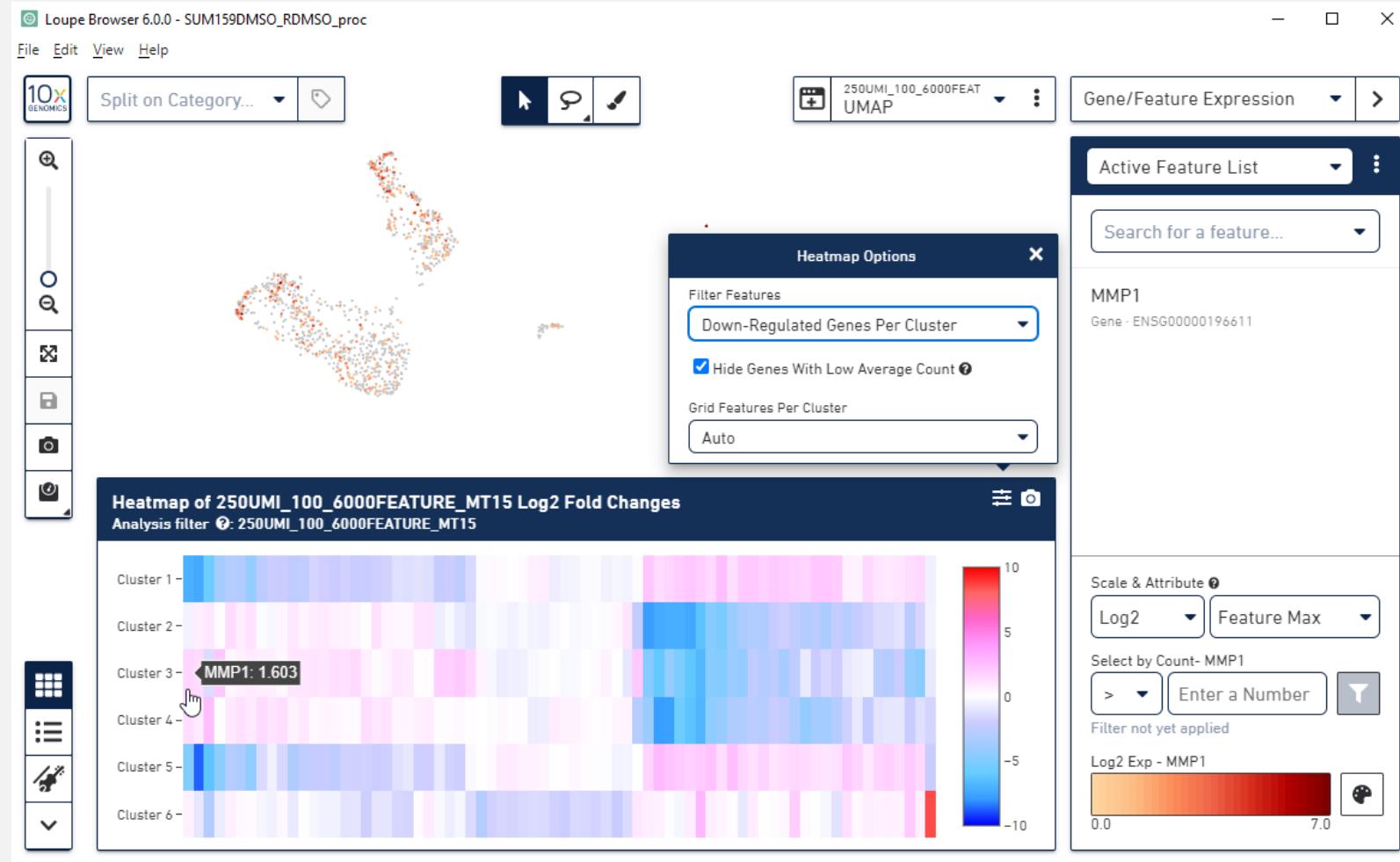
[Cancel Recluster](#)

LOUPE: FILTERED DATASET

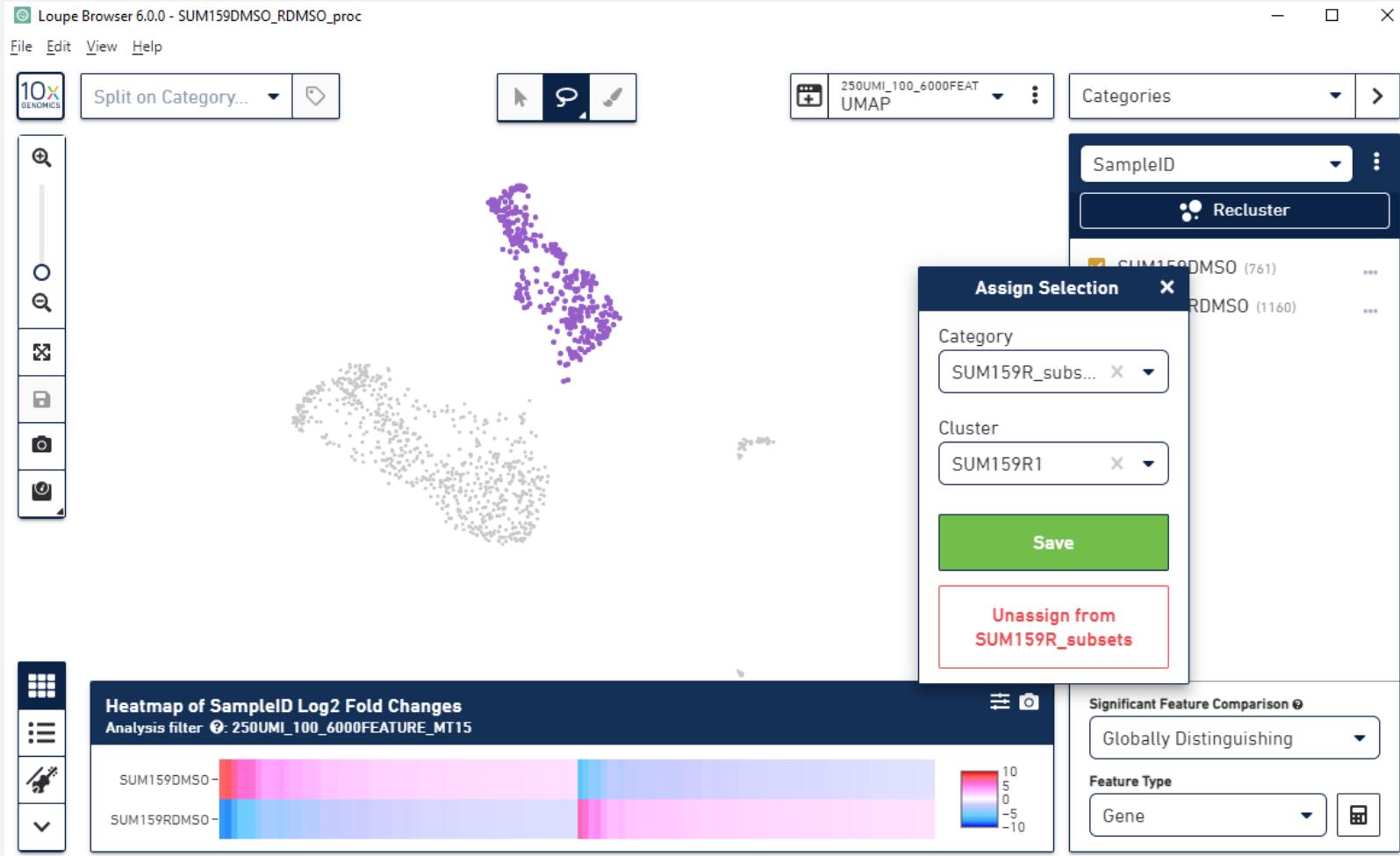


Find Markers

LOUPE: HEATMAP



LOUPE: CREATE SUB-CLUSTERING



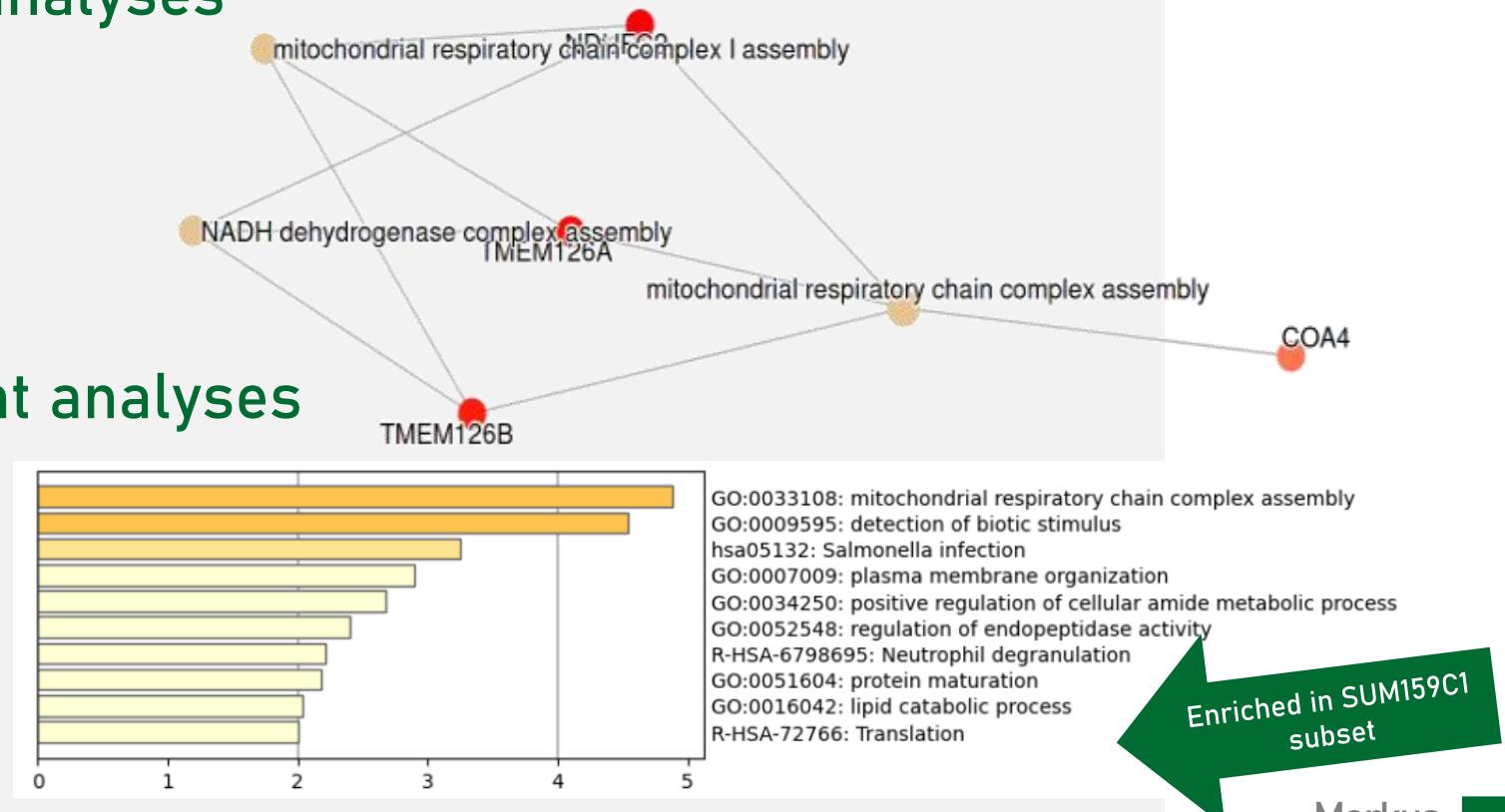
LOUPE: DATA EXPORT

The image shows the Loupe Browser 6.0.0 interface. At the top, there is a navigation bar with tabs for 'Loupe Browser 6.0.0 - SUM159DMSO_RDMSO_proc', 'File', 'Edit', 'View', and 'Help'. Below the navigation bar are several tool icons: 'Split on Category...', a search icon, a zoom icon, a 'Categories' dropdown, and a 'Recluster' button. The main area features two UMAP plots: one on the left with blue points and one on the right with orange points. To the left of the plots is a vertical toolbar with various icons for filtering, selecting, and analyzing data. A 'Table Options' dialog box is open in the center, showing settings for 'Displayed Numeric Value' (Log2 Fold Change), 'Filter Features' (Up-Regulated Genes Per Cluster (Default)), and 'CSV Export Count' (Top 100 Genes). An 'Export Table to CSV' button is at the bottom of this dialog. To the right of the UMAP plots is a 'SUM159R_subsets' panel titled 'Recluster' which lists 'SUM159R1 (379)' and 'SUM159R2 (730)'. Below this is a table titled 'SUM159R_subsets: Up-Regulated Genes Per Cluster' with an analysis filter of '250UMI_100_6000FEATURE_MT15'. The table lists genes with their Log2 Fold Change values and P-values. At the bottom right is a large table of gene expression data with columns for FeatureID, FeatureName, SUM159R1 Average, SUM159R1 Log2 Fold Change, and SUM159R1 P-Value.

FeatureID	FeatureName	SUM159R1 Average	SUM159R1 Log2 Fold Change	SUM159R1 P-Value
ENSG00000253377	AC068672	1.349142878	3.530195793	2.47E-22
ENSG00000137691	CFAP300	1.402430388	2.150597426	3.39E-11
ENSG00000137500	CCDC90B	3.887566102	1.912991489	2.97E-10
ENSG00000152558	TMEM123	11.0402033	1.852589647	7.46E-10
ENSG00000246067	RAB30-DT	4.275111632	1.880286258	1.04E-09
ENSG00000225383	SFTA1P	1.005196219	2.098340267	1.31E-09
ENSG00000137509	PRCP	1.649490664	1.682944272	1.94E-07
ENSG00000223784	LINP1	1.361253676	1.759323228	2.36E-07
ENSG00000078124	ACER3	1.22319058	1.629331154	1.06E-06
ENSG00000137692	DCUN1D5	13.12810484	1.499185498	2.08E-06
ENSG00000151366	NDUFC2	22.81432094	1.482994504	2.87E-06
ENSG00000109321	AREG	1.717311132	1.705556692	3.07E-06
ENSG00000171202	TMEM126	6.45	ENSG00000 PGM2L1	1.634958 1.004291 0.01039
ENSG00000087884	AAMDC	3.46	ENSG00000 VMP1	2.315585 0.9686 0.013977
ENSG00000166435	XRRA1	1.47	ENSG00000 BEX1	1.2789 1.036113 0.017087
ENSG00000171204	TMEM126	5.48	ENSG00000 WDR72	2.797594 0.950319 0.017254
ENSG00000137693	YAP1	1.49	ENSG00000 JMJD8	1.215924 0.963404 0.018445
ENSG00000149273	RPS3	1.50	ENSG00000 TMEM134	1.569559 0.906085 0.029937
ENSG00000137492	THAP12	1.51	ENSG00000 PAFAH1B3	1.254679 0.879472 0.04437
ENSG00000159063	ALG8	1.52	ENSG00000 ANAPC15	3.875455 0.840493 0.045953
ENSG00000245694	CRNDE	2.53	ENSG00000 DKK1	1.840841 0.861719 0.063495
FNSG00000149196	HIKFSHI	5		

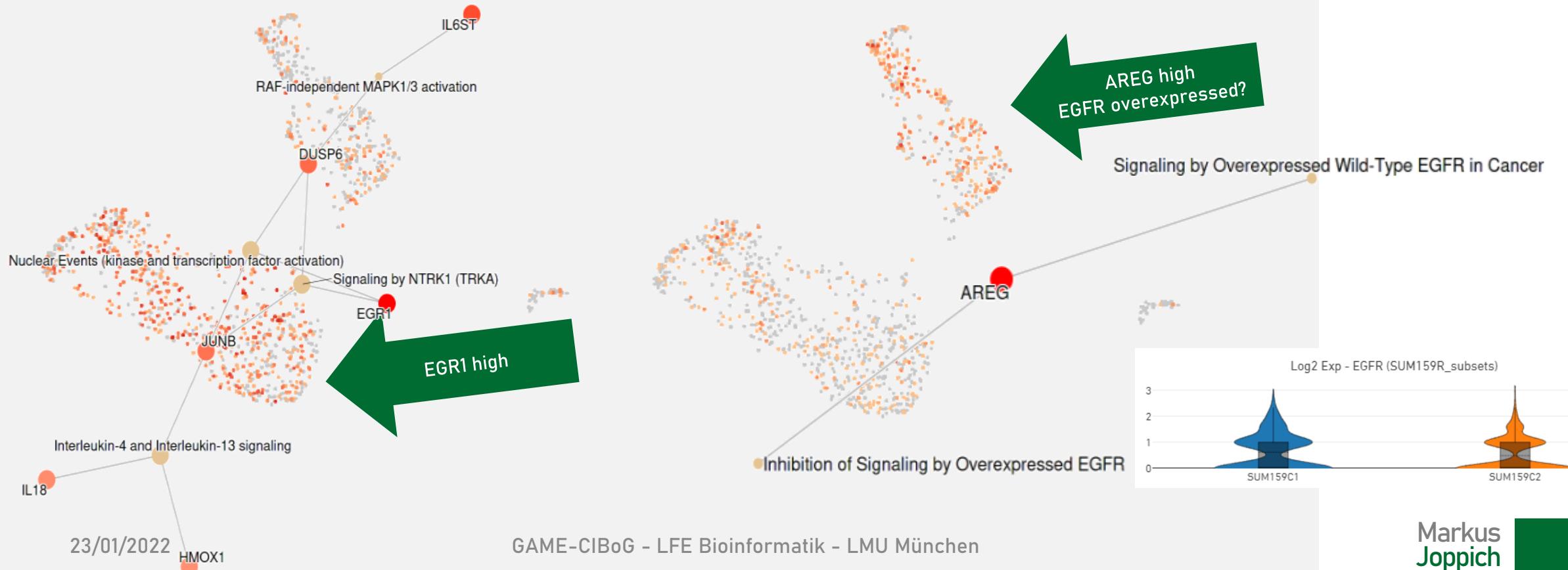
LOUPE: SET ENRICHMENT

- Export data from Loupe
 - Warning! Contains also insignificantly regulated genes!
- Run R-based enrichment analyses
 - Gene Ontology
 - Reactome
 - Etc.
- Run web-based enrichment analyses
 - Filter genes for significance before uploading!

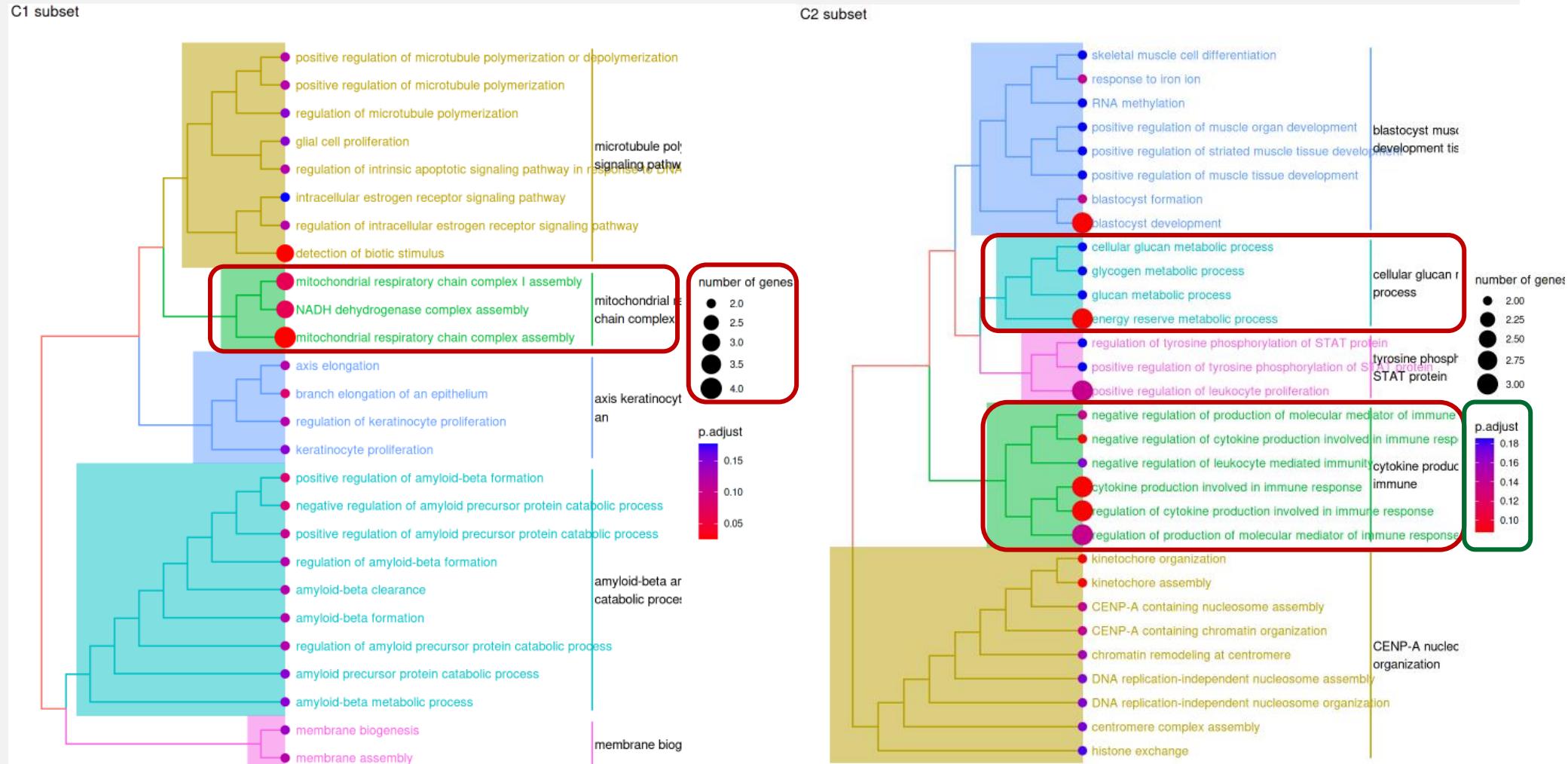


SET ENRICHMENT RESULTS: COMPARING C1+C2

- In C2: EGR1 transcription factor increased
- In C1: AREG increased which is relevant to EGFR/HER1



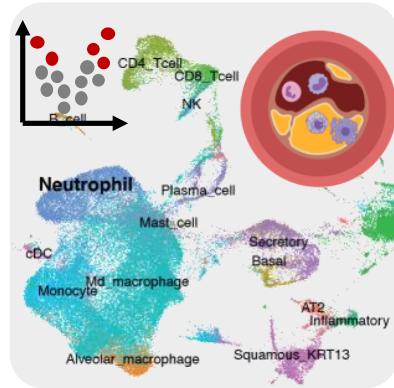
SET ENRICHMENT RESULTS: COMPARING C1+C2



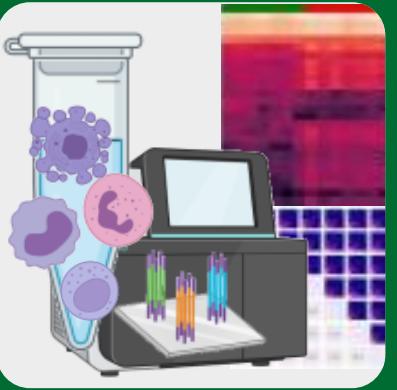
SUMMARY

scRNA-seq analysis

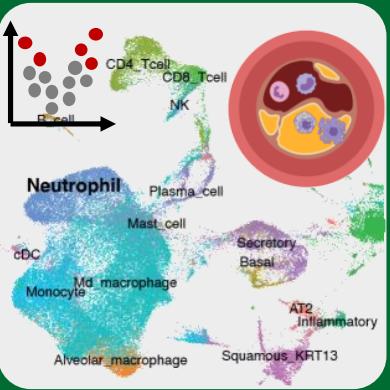
- Perform filtering and clustering of scRNA-seq data from within Loupe
- Analyse scRNA-seq data with Loupe
 - Create new clusters/categories
 - Watch for specific genes
 - Perform Differential Expression analysis
- Perform Gene Set Enrichment Analyses of significant genes
 - R-based: clusterProfiler, DOSE, disgenet2r, ...
 - Web-based: Metascape



scRNA-seq



bulk RNA-seq



scRNA-seq

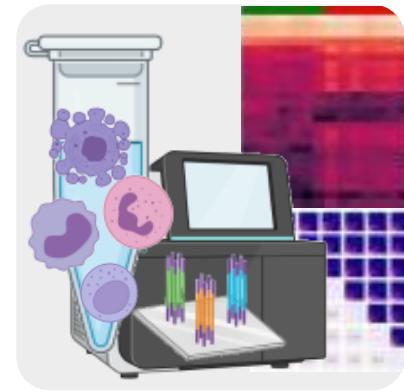
TIMELINE

GAME-CIBoG kickoff workshop

	Day1 (22/01/22)	Day2 (23/01/22)	Between Day2-3	Day3 (29/01/22)	Day4 (30/01/22)
NU staff	Opening remarks Dr. Katusno (10 min) Ice-breaking in each group Dr. Bustos (30min) Introduction/ Presenting the assignment Dr.Hinohara / Dr. Kato (30 min)			Organizing the data and brainstorming between students (3 hour) give advanced advice for the interpretation (Dr.Hinohara/Dr. Kato) → let every team to do interpretation of analyzed data	Presenting the answer from each group (10min per group 1 hour) Presenting the answer Dr.Hinohara / Dr.Kato (1 hour) Awards ceremony (5min) Closing remarks Dr. Katusno (10 min)
LMU staff	RNA-seq lecture (1) Dr. Hinske (1 hour)	RNA-seq lecture (2) Dr. Joppich (3 hours)			
Students	24 students from 5 Unive Divde into 6 groups	Goal: *know how to analyze single-cell RNA-seq *start analysis by themselves	Goal: *finish the analysis by themselves *start interpretation of the analyzed data	*analyze the analysis by themselves *interpretation of the analyzed data	Goal: Interpret the analyzed data present their result
	RNA-seq data analysis Interpretation of analyzed data				

SUMMARY

- Analyse public data from raw data (reads)
 - Or from pre-processed data
- Perform differential expression analysis using DESeq2
- Perform Gene Set Enrichment Analyses of significant genes
 - R-based: clusterProfiler, DOSE, disgenet2r, ...
 - Web-based: Metascape
- At each “pipeline” step the choice of tools, methods and parameters affects the final result
 - Sorting of lists, number of genes for set enrichment, etc.

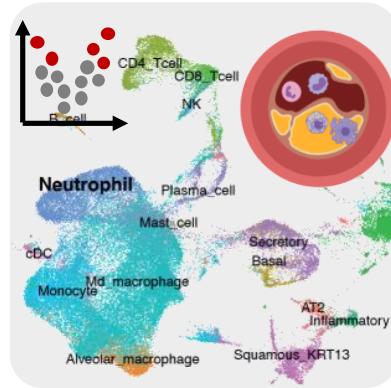


bulk RNA-seq

SUMMARY

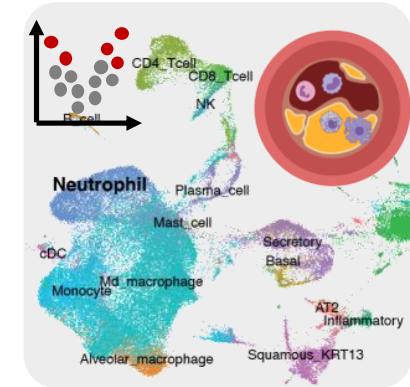
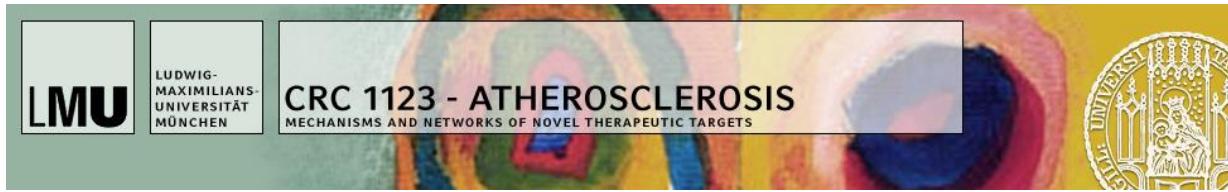
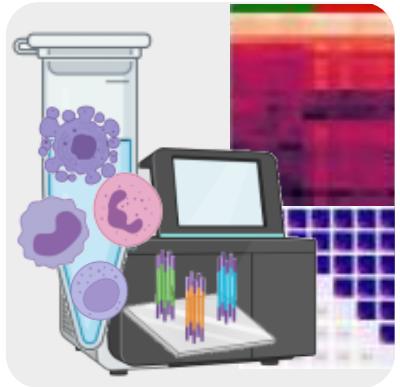
scRNA-seq analysis

- Analyse public data from raw data (reads)
- Run cellranger count and aggregate pipelines
- Perform sub-clustering of scRNA-seq data from within Loupe
 - Cell filtering and
 - Differential Expression analysis
- Perform Gene Set Enrichment Analyses of significant genes
 - R-based: clusterProfiler, DOSE, disgenet2r, ...
 - Web-based: Metascape



scRNA-seq

THANK YOU



FOR YOUR ATTENTION