Sergej Ruff, AG Jung, 01.06.2022
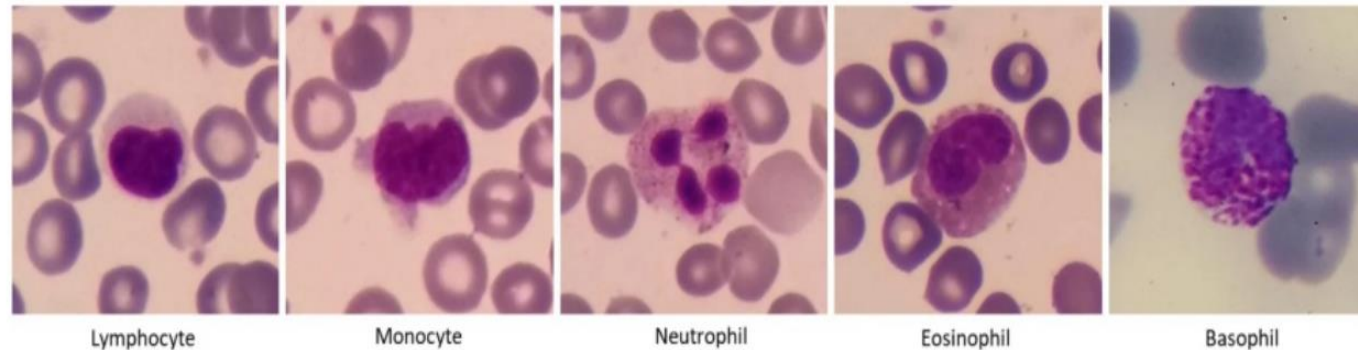
# Forschungskonzeption: single-cell RNA-seq analysis

# Inhalt

-Einleitung: Droplet und Information zum Paper

- Methoden und Daten: Datenanalyse in Seurat

-Forschungshypothesen für die Bachelorarbeit

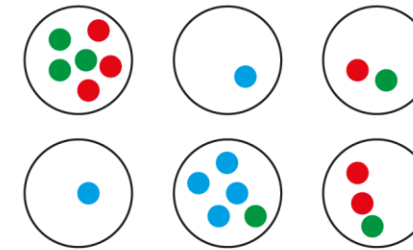# Wozu Single-Cell RNA-seq?

- Heterogene Zellpopulationen
  - Charakterisierung
  - Identifizierung
  - Krankheiten
- Untersuchung von Co-Expression-Mustern
- Expressionsunterschiede
- Untersuchung seltener Zellpopulationen



**B** Single cell transcriptome analysis
**C** Bulk analysis
**D** Coexpression Matrix (single cell)
**E** Coexpression Matrix (bulk analysis)

Macaulay IC, Voet T (2014) Single Cell Genomics: Advances and Future Perspectives. PLoS Genet 10(1): e1004126. https://doi.org/10.1371/journal.pgen.1004126



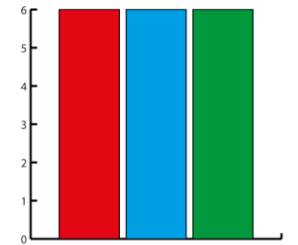Lymphocyte    Monocyte    Neutrophil    Eosinophil    Basophil

Five types of white blood cells in the normal peripheral blood.

Kouzehkanan, Z.M., Saghari, S., Tavakoli, S. *et al.* A large dataset of white blood cells containing cell locations and types, along with segmented nuclei and cytoplasm. *Sci Rep* **12,** 1123 (2022). https://doi.org/10.1038/s41598-021-04426-x

# Droplet basierte Methoden

1. Vorbereitung der Zellsuspension

2. Zellsortierung

3. Einkapseln einzelner Zellen in Droplets

4. CDNA-Synthese und Amplifikation

5. Bibliothek

6. Sequenzierung, data alignment und Interpretation

Salomon, Robert, et al. "Droplet-based single cell RNAseq tools: a practical guide." *Lab on a Chip* 19.10 (2019): 1706-1727.

# Beads

**Oligo-d(T) Primer** bindet den Poly-A-Schwanz der RNA.

**Unique Molecular Identifiers (UMI):** Zufällige, einzigartige Sequenzen auf den Beads, die als Tags die einzelnen Moleküle markieren.
- Transcript counting, normalisation of amplification artifacts

**Cell Barcode:** identische, mehrmals auftauchende Sequenz auf den Beads. Für Identifikation einzelner Zellen.

**Primer Region** für Amplifikationen - **Template Switching Oligonucleotides (TSO)**

Salomon, Robert, et al. "Droplet-based single cell RNAseq tools: a practical guide." *Lab on a Chip* 19.10 (2019): 1706-1727.

# *Entwicklung eines Forschungskonzepts*

## Maturation trajectories and transcriptional landscape of plasmablasts and autoreactive B cells in COVID-19

Christoph Schultheiß,[1] Lisa Paschold,[1] Edith Willscher,[1] Donjete Simnica,[1] Anna Wöstemeier,[2] Franziska Muscate,[2] Maxi Wass,[1] Stephan Eisenmann,[3] Jochen Dutzmann,[4] Gernot Keyßer,[5] Nicola Gagliani,[2,6,7] and Mascha Binder[1,8,*]

- doi: 10.1016/j.isci.2021.103325
- Veröffentlichung: 19.11.2021
- Online Veröffentlicht: 23.10.2021

-ArrayExpress: E-MTAB-11011

Schultheiß, Christoph, et al. "Maturation trajectories and transcriptional landscape of plasmablasts and autoreactive B cells in COVID-19." *Iscience* 24.11 (2021): 103325.

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

7

# Wozu?

## Infection-induced plasmablasts are a nutrient sink that impairs humoral immunity to malaria

Rahul Vijay,[1,10] Jenna J. Guthmiller,[2,8,10] Alexandria J. Sturtz,[1] Fionna A. Surette,[1,3] Kai J. Rogers,[1] Ramakrishna R. Sompallae,[4] Fengyin Li,[1,9] Rosemary L. Pope,[2] Jo-Anne Chan,[5] Fabian de Labastida Rivera,[6] Dean Andrew,[6] Lachlan Webb,[6] Wendy J. Maury,[1,3] Hai-Hui Xue,[1,3,7] Christian R. Engwerda,[6] James S. McCarthy,[6] Michelle J. Boyle,[5,6] and Noah S. Butler[1,2,3]

Vijay, Rahul et al. "Infection-induced plasmablasts are a nutrient sink that impairs humoral immunity to malaria." *Nature immunology* vol. 21,7 (2020): 790-801. doi:10.1038/s41590-020-0678-5

## Rapid and Massive Virus-Specific Plasmablast Responses during Acute Dengue Virus Infection in Humans

Jens Wrammert,[✉a,b] Nattawat Onlamoon,[c,f] Rama S. Akondy,[a,b] Guey C. Perng,[a] Korakot Polsrila,[c] Anmol Chandele,[a] Marcin Kwissa,[a] Bali Pulendran,[a] Patrick C. Wilson,[d] Orasri Wittawatmongkol,[e] Sutee Yoksan,[f,g] Nasikarn Angkasekwinai,[h] Kovit Pattanapanyasat,[c,f] Kulkanya Chokephaibulkit,[e] and Rafi Ahmed[a,b]

‣ Author information  ‣ Article notes  ‣ Copyright and License information    Disclaimer

Wrammert, Jens et al. "Rapid and massive virus-specific plasmablast responses during acute dengue virus infection in humans." *Journal of virology* vol. 86,6 (2012): 2911-8. doi:10.1128/JVI.06075-11

COVID-19 als <u>Modell</u>:

- Kein Exposure= kein Memory.

- Frühere Verweise auf...

   *... Vermeidung von GC-Reaktionen.*

   *... Hohe Konzentration an PB.*

Modell, um B-Zell-Antworten und ihre Konsequenzen auf das Immunologische Gedächtnis und Immunpathologie zu untersuchen.

- Untersuchung der Zellpopulationen!

Schultheiß, Christoph, et al. "Maturation trajectories and transcriptional landscape of plasmablasts and autoreactive B cells in COVID-19." *Iscience* 24.11 (2021): 103325.

# Daten

## Beschreibung

Sc-Seq Daten von B-Lymphozyten aus PBMC
für gesunde Patienten („HD"), "active" COVID-19 Patienten und "recovered" Patienten.

## Datenset

E-MTAB-11011

## Ordner

E-MTAB-11011.processed.1

## Counts (Pre-Processed)

pbmc.HD_gex_and_vdj.rds
pbmc.active.2.5.3.8_gex_and_vdj.rds
pbmc.recovered.14.16.26_gex_and_vdj.rds



*Patienten für scRNA-seq.*

"recovered" — 3 Patienten

"Active" — 4 Patienten

"HD" — 1 Patient

Schultheiß, Christoph, et al. "Maturation trajectories and transcriptional landscape of plasmablasts and autoreactive B cells in COVID-19." *Iscience* 24.11 (2021): 103325.

# SingleCellExperiment Class (Bioconductor)



**Slots**

Assays: Primäre Daten (Matrix mit Seq.-Counts).

rowData: Information zu den Genen (Reihen des Assays)

ColData: Information zu den Zellen (Spalten des Assays).

Besonderheiten:

reducedDims: dimensionale Reduktionen.

Alternative Experimente

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor." *Nature methods* 17.2 (2020): 137-145.

# Seurat

The `Seurat` object ...

## Slots

| Slot | Function |
|------|----------|
| `assays` | A list of assays within this object |
| `meta.data` | Cell-level meta data |
| `active.assay` | Name of active, or default, assay |
| `active.ident` | Identity classes for the current object |
| `graphs` | A list of nearest neighbor graphs |
| `reductions` | A list of DimReduc objects |
| `project.name` | User-defined project name (optional) |
| `tools` | Empty list. Tool developers can store any internal data from their methods here |
| `misc` | Empty slot. User can store additional information here |
| `version` | Seurat version used when creating the object |

## Slots  Assays

| Slot | Function |
|------|----------|
| `counts` | Stores unnormalized data such as raw counts or TPMs |
| `data` | Normalized data matrix |
| `scale.data` | Scaled data matrix |
| `key` | A character string to facilitate looking up features from a specific `Assay` |
| `var.features` | A vector of features identified as variable |
| `meta.features` | Feature-level meta data |

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor."
*Nature methods* 17.2 (2020): 137-145.

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

14

# Packages

```r
library("Seurat")
library("celldex")
library("SingleR")
library("biomaRt")
library("clusterProfiler")
library("org.Hs.eg.db")

library("SingleCellExperiment")
```

## Daten importieren

readSparseCounts(): *scuttle* package

read10xCounts(): *DropletUtils* package

readRDS()

readH5AD(): *zellkonverter* package

## Erschaffen eines sce-Objektes oder Seurat-Objektes

SingleCellExperiment():*SingleCellExperiment* package

as.SingleCellExperiment(): *Seurat* package

As.Seurat: *Seurat* package

## Datenzugriff

Counts(sce),assay(sce,"counts"),assays(sce)

```
> active
An object of class Seurat
17786 features across 10050 samples within 2 assays
Active assay: integrated (2000 features, 2000 variable features)
 1 other assay present: RNA
 2 dimensional reductions calculated: pca, umap
> SingleCellExperiment(active)
class: SingleCellExperiment
dim: 2000 10050
metadata(0):
assays(1): ''
rownames(2000): IGKV3-15 IGKV3-11 ... AP000345.2 MKLN1-AS
rowData names(0):
colnames(10050): sc2_AAAGATGCACTTACGA-1 sc2_AAAGCAACAGTAGAGC-1 ... sc8_TTTGTCAGTTAAGAAC-1
  sc8_TTTGTCAGTTGTGGAG-1
colData names(0):
reducedDimNames(0):
mainExpName: NULL
altExpNames(0):
```

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor."
Nature methods 17.2 (2020): 137-145

# Quality Control

## Motivation
Finden und Entfernen von low-quality libraries

## Gründe für Low-quality Libraries
Zellschäden (Mitochondriale RNA "↑",endogene RNA↓)
Fehler in der Reverse Transkriptase oder PCR-Amplifikation

## Wie beeinflussen Sie die Analyse?
Bilden eigene Cluster
Clustern verschiedene Zelltypen zusammen (induzierte Expressionsprofile)
Einfluss auf PCA

## QC-Metrics

-Anzahl einzigartiger Gene pro Zelle
-Spike-Ins
-Anreicherung an mt-RNA

Bsp. Aus Tutorial          pbmc= Seurat-Objekt

```
pbmc[["percent.mt"]] <-PercentageFeatureSet(pbmc, pattern = "^MT-")
```

"The calculation here is simply the column sum of the matrix present in the counts slot for features belonging to the set divided by the column sum for all features times 100."

```
##                 orig.ident nCount_RNA nFeature_RNA percent.mt
## AAACATACAACCAC-1    pbmc3k       2419          779  3.0177759
## AAACATTGAGCTAC-1    pbmc3k       4903         1352  3.7935958
## AAACATTGATCAGC-1    pbmc3k       3147         1129  0.8897363
## AAACCGTGCTTCCG-1    pbmc3k       2639          960  1.7430845
## AAACCGTGTATGCG-1    pbmc3k        980          521  1.2244898
```

Visualisierung der QC-Metrics als Violin-Plot und Scatter

```
VlnPlot(pbmc, features = c("nFeature_RNA", "nCount_RNA", "percent.mt"), ncol = 3)

plot2 <- FeatureScatter(pbmc, feature1 = "nCount_RNA", feature2 = "nFeature_RNA")

pbmc <- subset(pbmc, subset = nFeature_RNA > 200 & nFeature_RNA < 2500 & percent.mt < 5)
```

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor." Nature methods 17.2 (2020): 137-145

**Violinplot**

Featurescatter
Median absolute Deviation (MAD) für QC.
3*MAD=Ausreißer

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor."
Nature methods 17.2 (2020): 137-145

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

https://satijalab.org/seurat/articles/pbmc3k_tutorial.html Letzter
Zugriff: 23.05.22 17:28

18

**Verunreinigung durch andere Zelltypen?**

Lösung: Zelltyp-Annotation mit *SingleR* package!

Automatisches Annotationsmethode

Abgleich mit Referenz-Datenset mit bekannten Labeln
hpca.se=celldex::HumanPrimaryCellAtlasData()#reference data

Braucht ein SummarizedExperiment-Object
sce = as.SingleCellExperiment(SeuratObj) #convert Seurat to SingleCellExperiment
se = as(sce, "SummarizedExperiment") #convert SingleCellExperiment to SummarizedExperiment

SingleR
pred = SingleR(test = se, ref = hpca.se, assay.type.test=1,
      labels = hpca.se$label.main)

Entfernen der falschen Zellen

```
> as.vector(active$CellAnnotation)
  [1] "B_cell"            "B_cell"            "B_cell"   "Neutrophils"       "Pre-B_cell_CD34-"
  [6] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [11] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [16] "Neutrophils"       "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [21] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [26] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [31] "B_cell"            "Pro-B_cell_CD34+"  "B_cell"   "B_cell"            "B_cell"
 [36] "B_cell"            "Pre-B_cell_CD34-"  "B_cell"   "B_cell"            "B_cell"
 [41] "B_cell"            "B_cell"            "B_cell"   "Astrocyte"         "Neutrophils"
 [46] "Neutrophils"       "T_cells"           "B_cell"   "B_cell"            "B_cell"
 [51] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [56] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [61] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [66] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [71] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [76] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [81] "Pro-B_cell_CD34+"  "B_cell"            "B_cell"   "Neutrophils"       "B_cell"
 [86] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [91] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
 [96] "B_cell"            "B_cell"            "B_cell"   "Pre-B_cell_CD34-"  "B_cell"
[101] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
[106] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "B_cell"
[111] "B_cell"            "B_cell"            "B_cell"   "B_cell"            "Neutrophils"
[116] "B_cell"            "B_cell"            "B_cell"   "Neutrophils"       "B_cell"
```

```
> active = remWC(active)
[1] "Percentage B-cells: 92.22 %"
> recovered = remWC(recovered)
[1] "Percentage B-cells: 99.34 %"
> hd = remWC(hd)
[1] "Percentage B-cells: 100 %"
```

# Daten Normalisierung

Motivation

Entfernen von systematischen, nicht-biologischen Variationen

Fehlerquellen

Fehler in cDNA-capture und PCR-amp.

Methoden

Library Size Normalisierung

Spike-In Normalisierung

Deconvolution

Aus Tutorial

```
pbmc <- NormalizeData(pbmc, normalization.method = "LogNormalize", scale.factor = 10000)
```

- "LogNormalize: Feature counts for each cell are divided by the total counts for that cell and multiplied by the scale.factor. This is then natural-log transformed using log1p."

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor." Nature methods 17.2 (2020): 137-145

Motivation

Auswahl an Genen, die nützliche Information über die Biologie des Systems enthalten für Downstreamanalysen

Hoch-Variable Gene suchen

SeuratObj = FindVariableFeatures(SeuratObj, selection.method = "vst", nfeatures = 2000)

Skalierung der Daten

all.genes_SeuratObj = rownames(SeuratObj)
SeuratObj = ScaleData(SeuratObj, features = all.genes_SeuratObj)

- "Shifts the expression of each gene, so that the mean expression across cells is 0"
- "Scales the expression of each gene, so that the variance across cells is 1"
    - ➔ "This step gives equal weight in downstream analyses, so that highly-expressed genes do not dominate"
    
    https://satijalab.org/seurat/articles/pbmc3k_tutorial.html Letzter Zugriff : 23.05.22 17:28

Amezquita, Robert A., et al. "Orchestrating single-cell analysis with Bioconductor." Nature methods 17.2 (2020): 137-145

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

https://satijalab.org/seurat/articles/pbmc3k_tutorial .html Letzter Zugriff: 23.05.22 17:28

21

# PCA, Clustering und UMAP

Nutzen skalierte, selektierte variable Features für PCA.

SeuratObj = RunUMAP(SeuratObj, dims = 1:10)

Clustering
SeuratObj = FindNeighbors(SeuratObj, dims = 1:10)
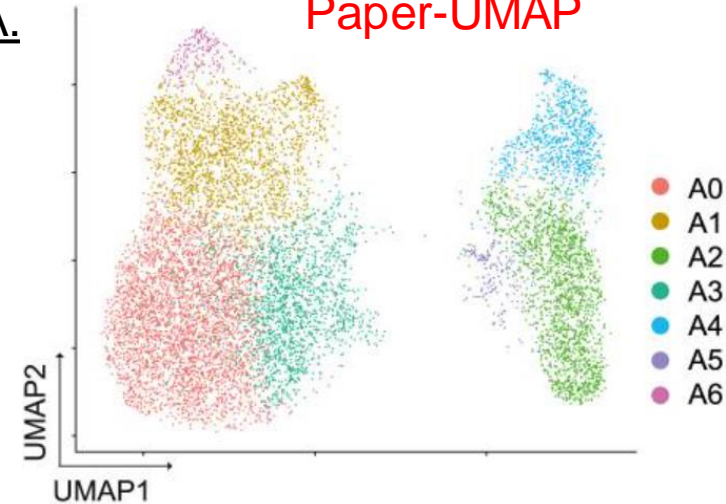SeuratObj = FindClusters(SeuratObj, resolution = 0.5)

UMAP
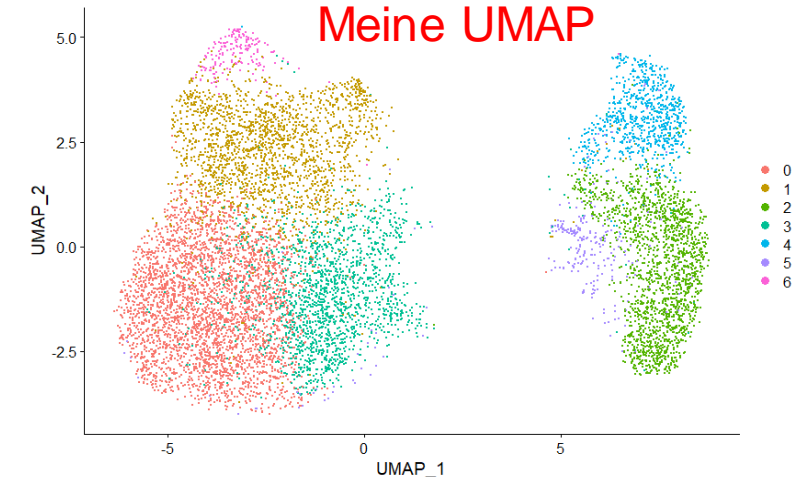SeuratObj = RunUMAP(SeuratObj, dims = 1:10)

Visualisierung
dp=DimPlot(SeuratObj, reduction = "umap")



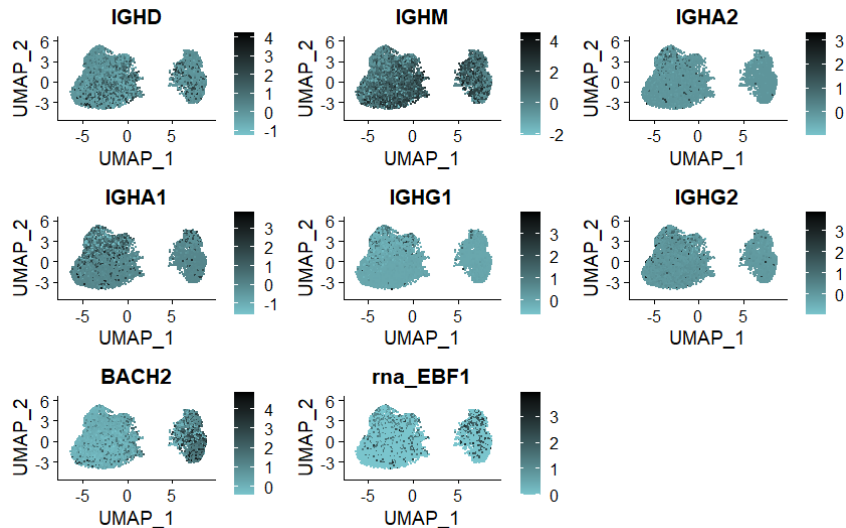scRNA-seq - CD19⁺ B cells - active

Paper-UMAP

Meine UMAP

Stiftung Tierärztliche Hochschule Hannover
University of Veterinary Medicine Hannover, Foundation

https://satijalab.org/seurat/articles/pbmc3k_tutorial
.html Letzter Zugriff: 23.05.22 17:28

## Motivation

Suche nach Markern, die Cluster durch eine differenzielle Expression definieren

Cluster0.active <- FindMarkers(active, ident.1 = 0, min.pct = 0.25)

#findet alle Marker des Clusters 0 der active Covid Patienten

## Visualisierung FeaturePlot()



| | p_val | avg_log2FC | pct.1 | pct.2 | p_val_adj |
|---|---|---|---|---|---|
| IGHV2-26 | 6.316417e-307 | -0.2865865 | 0.819 | 0.314 | 1.263283e-303 |
| CALHM6 | 2.800654e-273 | -0.3398360 | 0.836 | 0.378 | 5.601308e-270 |
| IGHV4-39 | 4.252269e-236 | -0.3198256 | 0.840 | 0.417 | 8.504537e-233 |
| AC243960.1 | 9.955206e-236 | -0.2516593 | 0.804 | 0.319 | 1.991041e-232 |
| IGKV1-9 | 1.004380e-217 | 0.3951521 | 0.803 | 0.350 | 2.008759e-214 |
| DBNDD1 | 1.303543e-208 | -0.2584789 | 0.807 | 0.332 | 2.607087e-205 |
| ENTPD1 | 7.980514e-206 | -0.4229162 | 0.235 | 0.294 | 1.596103e-202 |
| EIF2AK1 | 1.039076e-198 | -0.3036104 | 0.834 | 0.382 | 2.078151e-195 |
| GIHCG | 6.523331e-195 | -0.2553940 | 0.222 | 0.252 | 1.304666e-191 |
| IGKV2-30 | 1.551499e-194 | 0.2681073 | 0.770 | 0.329 | 3.102998e-191 |
| IGHV1-2 | 2.228663e-191 | -0.7594778 | 0.882 | 0.438 | 4.457326e-188 |

Marker von Active_Cluster_2

# Gene Enrichment (KEGG)

## Packages

library("biomaRt")
library("clusterProfiler")
library("org.Hs.eg.db")

Entrez-IDs

Log-fold Changes

| | ID | Description | GeneRatio | BgRatio | pvalue | p.adjust | qvalue | geneID |
|---|---|---|---|---|---|---|---|---|
| hsa03010 | hsa03010 | Ribosome | 21/88 | 158/8146 | 9.414375e-18 | 1.704002e-15 | 1.704002e-15 | 6222/6176/6206/6208/6134/6202/6158/6191/6 |
| hsa05171 | hsa05171 | Coronavirus disease - COVID-19 | 22/88 | 232/8146 | 2.129562e-15 | 1.927253e-13 | 1.927253e-13 | 6222/6176/6206/6208/6134/6202/6158/6191/6 |

Active_kegg_1

```
kegg = gseKEGG(geneList= ngl,organism= "hsa",minGSSize= 5,
pvalueCutoff= pval,verbose=TRUE)
```

| | ID | Description | GeneRatio | BgRatio | pvalue | p.adjust | qvalue | geneID |
|---|---|---|---|---|---|---|---|---|
| hsa04640 | hsa04640 | Hematopoietic cell lineage | 5/74 | 99/8146 | 0.001988967 | 0.1689545 | 0.1467236 | 3123/931/100133941/3566/3115 |
| hsa05416 | hsa05416 | Viral myocarditis | 4/74 | 60/8146 | 0.002083728 | 0.1689545 | 0.1467236 | 3123/71/5880/3115 |
| hsa04726 | hsa04726 | Serotonergic synapse | 5/74 | 115/8146 | 0.003815597 | 0.1689545 | 0.1467236 | 240/59345/2787/3708/1843 |
| hsa04115 | hsa04115 | p53 signaling pathway | 4/74 | 73/8146 | 0.004256724 | 0.1689545 | 0.1467236 | 5728/143686/900/4193 |

→ pt2: Cardiomyopathy

Active_kegg_2

Nachdem ich die Daten auf ihre Reproduzierbarkeit überprüft habe, teste ich die Stabilität der Daten.

Aufgabe der Bachelorarbeit wäre es jetzt die Stabilität der Resultate zu überprüfen.

Dafür würden Bootstrapanalysen und Parameteranpassung in Frage kommen.

# Literaturverzeichnis

-Kouzehkanan, Z.M., Saghari, S., Tavakoli, S. *et al.* A
large dataset of white blood cells containing cell locations and types, along with segmented nuclei and cytoplasm. *Sci Rep* **12,** 1123 (2022).
https://doi.org/10.1038/s41598-021-04426-x

-Macaulay IC, Voet T (2014) Single Cell Genomics: Advances and Future Perspectives. PLoS Genet 10(1): e1004126. https://doi.org/10.1371/journal.pgen.1004126

-Salomon, Robert, et al. "Droplet-based single cell RNAseq tools: a practical guide." *Lab on a Chip* 19.10 (2019): 1706-1727.

-Schultheiß, Christoph, et al. "Maturation trajectories and transcriptional landscape of
plasmablasts and autoreactive B cells in COVID-19." *Iscience* 24.11 (2021): 103325.

-Portugal, Silvia et al. "Atypical memory B cells in human chronic infectious diseases: An interim report." *Cellular immunology* vol. 321 (2017): 18-25.
doi:10.1016/j.cellimm.2017.07.003

-Vijay, Rahul et al. "Infection-induced plasmablasts are a nutrient sink that impairs humoral immunity to malaria." *Nature immunology* vol. 21,7 (2020): 790-801.
doi:10.1038/s41590-020-0678-5

-Wrammert, Jens et al. "Rapid and massive virus-specific plasmablast responses during acute dengue virus infection in humans." *Journal of virology* vol. 86,6 (2012):
2911-8. doi:10.1128/JVI.06075-11

-Abbildung Seite 14:https://github.com/satijalab/seurat/wiki/Assay (Letzter Zugriff: 23.05.2022, 13:50)

-https://satijalab.org/seurat/articles/pbmc3k_tutorial.html Letzter Zugriff: 23.05.22 17:28

-https://bioconductor.org/packages/release/bioc/html/SingleR.html Letzter Zugriff: 23.05.2022 19:31