



**Vilniaus
universitetas**

Doktoranto Karolio Šablausko ataskaita už 2022/2023 mokslo metų antrąjį pusmetį

Darbo vadovas: prof. Audronė Jakaitienė

Disertacijos pavadinimas: Characterization of genetic changes using deep neural networks

Doktorantūros pradžios ir pabaigos metai: 2022 – 2026

Studijų metai: 1.

1 lentelė: Doktorantūros studijų planas

Studijų metai	Egzaminai	
	Planas	Ivykdyta
I (2022/2023)	1	1
II (2023/2024)	2	0
III (2024/2025)	1	0
IV (2025/2026)	0	0
Iš viso:	4	0

Studijų metai	Dalyvavimas konferencijose				Publikacijos					
	Tarptautinėse		Nacionalinėse		Su citavimo rodikliu			Be citavimo rodiklio		
	Planas	Įvykdyta	Planas	Įvykdyta	Planas	Įvykdyta	Būklė	Planas	Įvykdyta	Būklė
I (2022/ 2023)	0	0	0	0	0	0		0	0	
II (2023/ 2024)	1	0	0	0	0	0		0	0	
III (2024/ 2025)	1	0	0	0	1	0		0	0	
IV (2025/ 2026)	1	0	0	0	1	0		0	0	
Iš viso:	3	0	0	0	2	0		0	0	

Publikacijos

Priimta spaudai (2023-09-16):
British Journal of Haematology
IF: 8.615

TITLE: Immunogenicity and clinical effectiveness of BNT162b2 booster against SARS-CoV-2 Omicron in patients with haematological malignancies: a national prospective cohort study

SHORT TITLE: Booster dose against Omicron in haematological patients

KEYWORDS: BNT162b2, booster, Omicron, haematologic patients, protection

Lukas Kevličius^{1, 2*}, Karolis Šablauskas^{1, 2*} Kazimieras Maneikis^{1,2}, Dovilė Juozapaitė¹, Ugnė Ringelevičiūtė¹, Vilmantė Vaitekėnaitė^{1, 2}, Birutė Davainienė^{1, 2}, Guoda Daukėlaitė^{1, 2}, Dominika Vasilevska², Mindaugas Stoškus¹, Ieva Narkevičiūtė¹, Violeta Sivickienė¹, Kęstutis Rudaitis¹, Birutė Masaitytė – Vitlipienė¹, Mantas Minkauskas¹, Daniel Naumovas¹, Tumas Beinortas^{3, 4#}, Laimonas Griškevičius^{1, 2#}

Publikacijos

Pateikta spaudai (2023-08-31)

Nature Communications

IF: 16.6

Title

Comprehensive reanalysis of CNVs in unsolved rare disease patients results in new diagnoses

Authors

German Demidov*, Burcu Yaldiz*, José Garcia-Pelaez*, Elke de Boer*, Nika Schuermans*, Liedewei Van de Vonde*, Ida Paramonov, Lennart Johansson, Francesco Musacchia, Elisa Benetti, Gemma Bullich, **Karolis Sablauskas**, Sergi Beltran, Christian Gilissen, Alex Hoischen, Stephan Ossowski, Richarda M. de Voer, Katja Lohmann, Carla Oliveira, Ana Topf, Lisenka E.L.M. Vissers, Steven Laurie

1 Thesis design

1.1 Thesis aim

- To contribute to the advancement of deep learning techniques for the analysis of next generation sequencing data, with a focus on single cell RNA sequencing (scRNA-seq) data.

1.2 Thesis objectives

- **Data preprocessing and feature engineering:** create efficient data preprocessing pipeline suited for scRNA-seq data, including normalization and batch effect correction to prepare the data for further analysis.
- **Deep differential expression analysis:** develop and implement deep learning-based approach for identifying differentially expressed genes and pathways.
- **Evaluation and benchmarking:** conduct extensive benchmarking and cross-validation experiments to assess the performance and generalizability of deep learning models, comparing them to traditional methods.
- **Biological case study:** apply said techniques in the interpretation of biological data gathered during a biomedical study.

Data structure

	GENE_ 1	GENE_ 2	GENE_ 3	...	GENE_ N
CELL_ 1	10	0	15		0
CELL_ 2	0	12	8		9
CELL_ 3	9	0	0		0
...					
CELL_ M	8	7	1		0

	barcode	library	total_counts	pct_counts_mito	patient	part	timepoint	response	n
55	AACAGCGAAGCACTGT	AC3_2	1705	2.05279	AC	2	3	RD	
66	AACAGCGAAGCTCGTA	AC3_2	618	2.42718	AC	2	3	RD	
97	AACAGCGAATGACCTG	AC3_2	1655	0.060423	AC	2	3	RD	
134	AACAGCGACATGAGTC	AC3_2	606	3.13531	AC	2	3	RD	
135	AACAGCGACATGCACA	AC3_2	509	2.75049	AC	2	3	RD	

Associated information (e.g. timepoint, patient, response about each barcode / cell is saved in a separate dataframe.

Single cell data is represented in M x N matrix

Cell counts by patient by timepoint

Vilniaus
universitetas

		Timepoint			
	patient	1	2	3	4
0	AC	3707	4932	5360	7486
1	AG	3171	6170	3575	3643
2	AA	4700	4453	4681	1457
3	AL	4649	3429	2884	5124
4	AB	1673	3264	2249	3287
5	AM	3885	2320	961	4845
6	V	777	681	547	605
7	AD	5014	938	3510	5336
8	AI	4691	3210	6033	10919
9	N	3288	3262	1700	1881
10	AH	7985	4449	0	4543
11	AE	8190	1372	2234	3126
12	AK	7069	16162	2933	0
13	S	1320	754	265	209
14	Z	3864	1853	3687	4581
15	T	484	441	550	336
16	P	1226	799	1053	0

Analysis steps:

1. Calculate normalized Z scores:
 1. Normalize values per cell
 2. Log-scale values
 3. Scale values
2. Select highly informative genes (37k -> 3k):
 1. Remove known biological artefacts
 2. Select variable genes
3. Principal component analysis (PCA)
4. Integrate using Harmony (reduces batch effect)
5. Compute neighborhood graph
6. Dimensionality reduction:
 1. UMAP
 2. Force directed layout
7. Find doublets and remove them
8. Repeat 1-6 without doublets
9. Find clusters / cell types

Before filtering doublets: (223 782, 37 733)

After filtering doublets: (222 932, 37 733)

37,733

222,782

	GENE_ 1	GENE_ 2	GENE_ 3	...	GENE_ N
CELL_ 1	10	0	15		0
CELL_ 2	0	12	8		9
CELL_ 3	9	0	0		0
...					
CELL_ M	8	7	1		0

Analysis steps:

- 1. Calculate normalized Z scores:**
 - 1. Normalize values per cell**
 - 2. Log-scale values**
 - 3. Scale values**

	GENE_ 1	GENE_ 2	GENE_ 3	...	GENE_ N	
CELL_ 1	10	0	15		0	SUM=10,000

- Linear -> logarithmic effect
- Scale data to unit variance and zero mean.

Vilnius

Analysis steps:

Select highly informative genes (37k -> 3k):

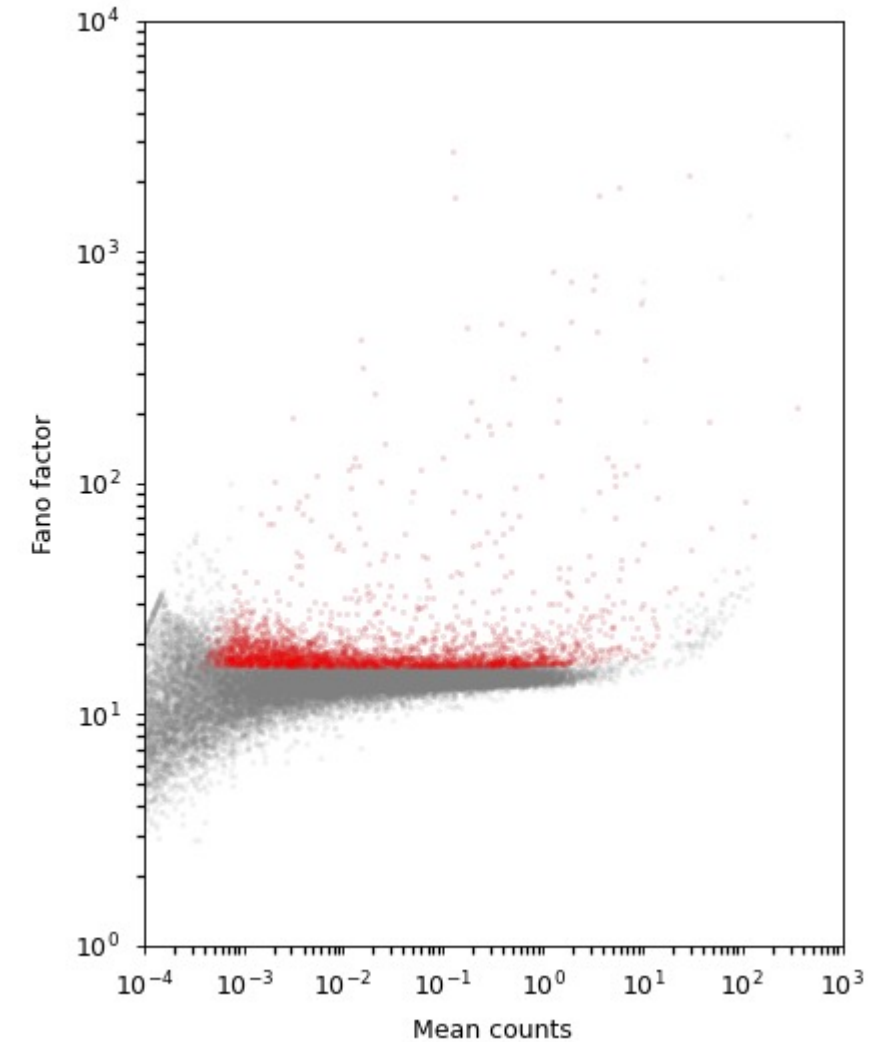
Remove known biological artefacts

Select variable genes

Remove known biological artefacts:

- Exclude mitochondrial, ribosomal and hemoglobin genes

Find and select the most variable genes

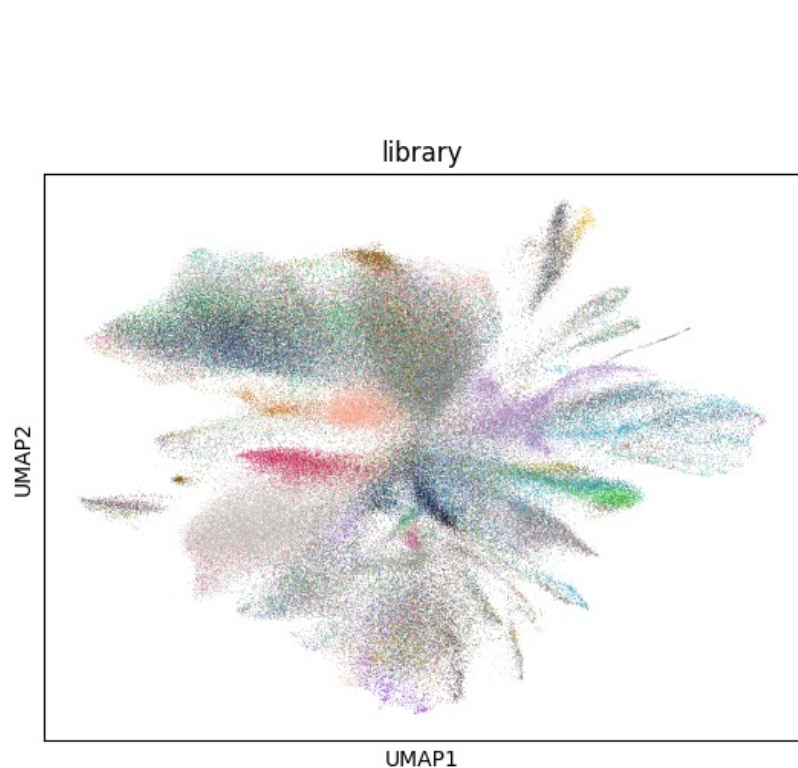


Principal component analysis (PCA) Integrate using Harmony (reduces batch)

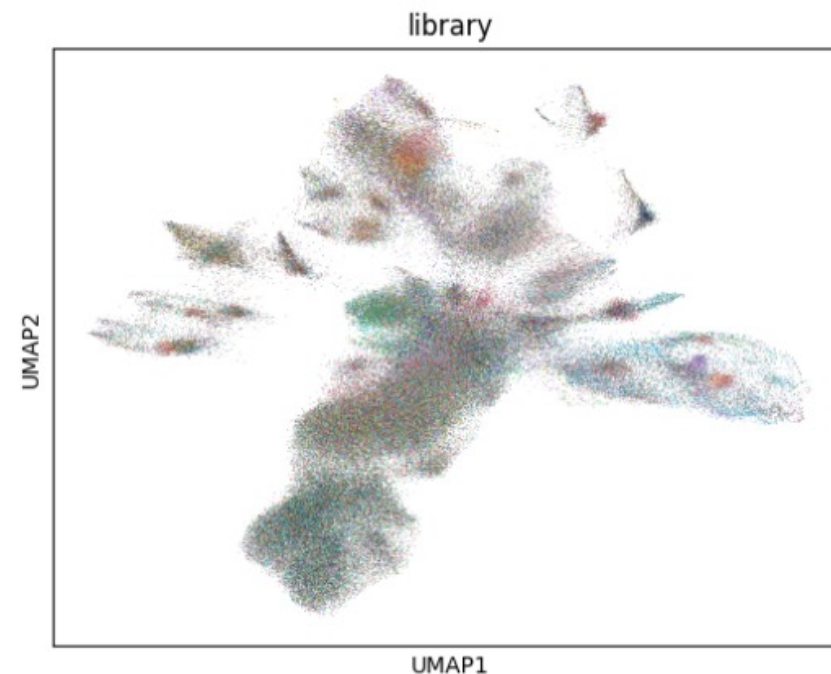
Vilniaus
universitetas

Before Harmony
integrate

After Harmony
integrate

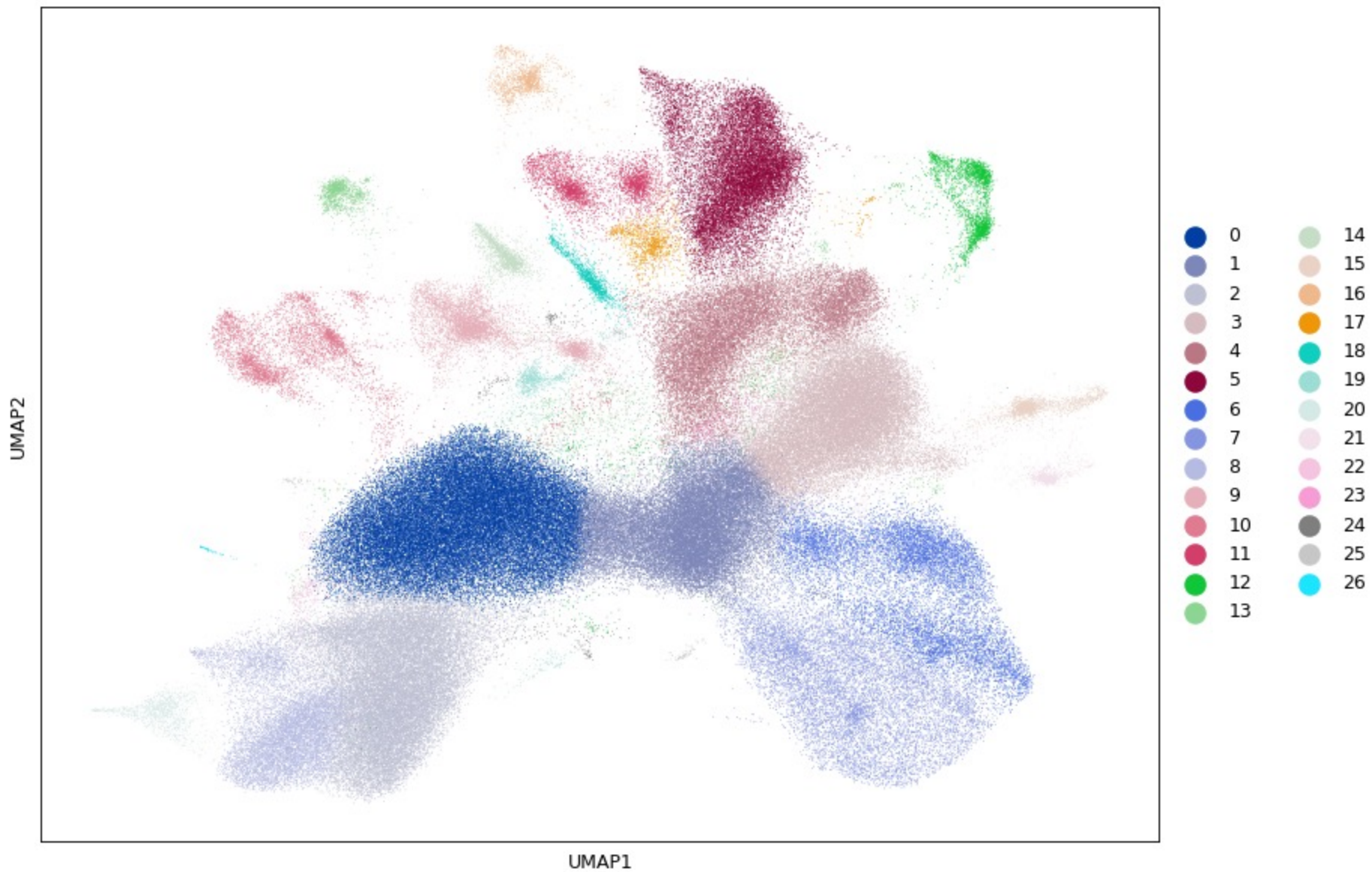


- | | | |
|---------|---------|--------|
| ● AA1_1 | ● AE2_1 | ● N3_1 |
| ● AA1_2 | ● AE3_1 | ● N4_1 |
| ● AA2_1 | ● AE3_2 | ● P1_1 |
| ● AA2_2 | ● AE4_1 | ● P2_1 |
| ● AA3_1 | ● AG1_1 | ● P3_1 |
| ● AA3_2 | ● AG2_1 | ● S1_1 |
| ● AA4_1 | ● AG3_1 | ● S1_2 |
| ● AB1_1 | ● AG4_1 | ● S2_1 |
| ● AB2_1 | ● AH1_1 | ● S3_1 |
| ● AB3_1 | ● AH2_1 | ● S4_1 |
| ● AB4_1 | ● AH4_1 | ● T1_1 |
| ● AC1_1 | ● AI1_1 | ● T2_1 |
| ● AC1_2 | ● AI2_1 | ● T3_1 |
| ● AC2_1 | ● AI3_1 | ● T4_1 |
| ● AC2_2 | ● AI4_1 | ● V1_1 |
| ● AC3_1 | ● AK1_1 | ● V1_2 |
| ● AC3_2 | ● AK2_1 | ● V2_1 |
| ● AC4_1 | ● AK3_1 | ● V3_1 |
| ● AC4_2 | ● AL1_1 | ● V4_1 |
| ● AD1_1 | ● AL2_1 | ● Z1_1 |
| ● AD1_2 | ● AL3_1 | ● Z1_2 |
| ● AD2_1 | ● AL4_1 | ● Z2_1 |
| ● AD3_1 | ● AM1_1 | ● Z2_2 |
| ● AD3_2 | ● AM2_1 | ● Z3_1 |
| ● AD4_1 | ● AM3_1 | ● Z3_2 |
| ● AD4_2 | ● AM4_1 | ● Z4_1 |
| ● AE1_1 | ● N1_1 | ● Z4_2 |
| ● AE1_2 | ● N2_1 | |



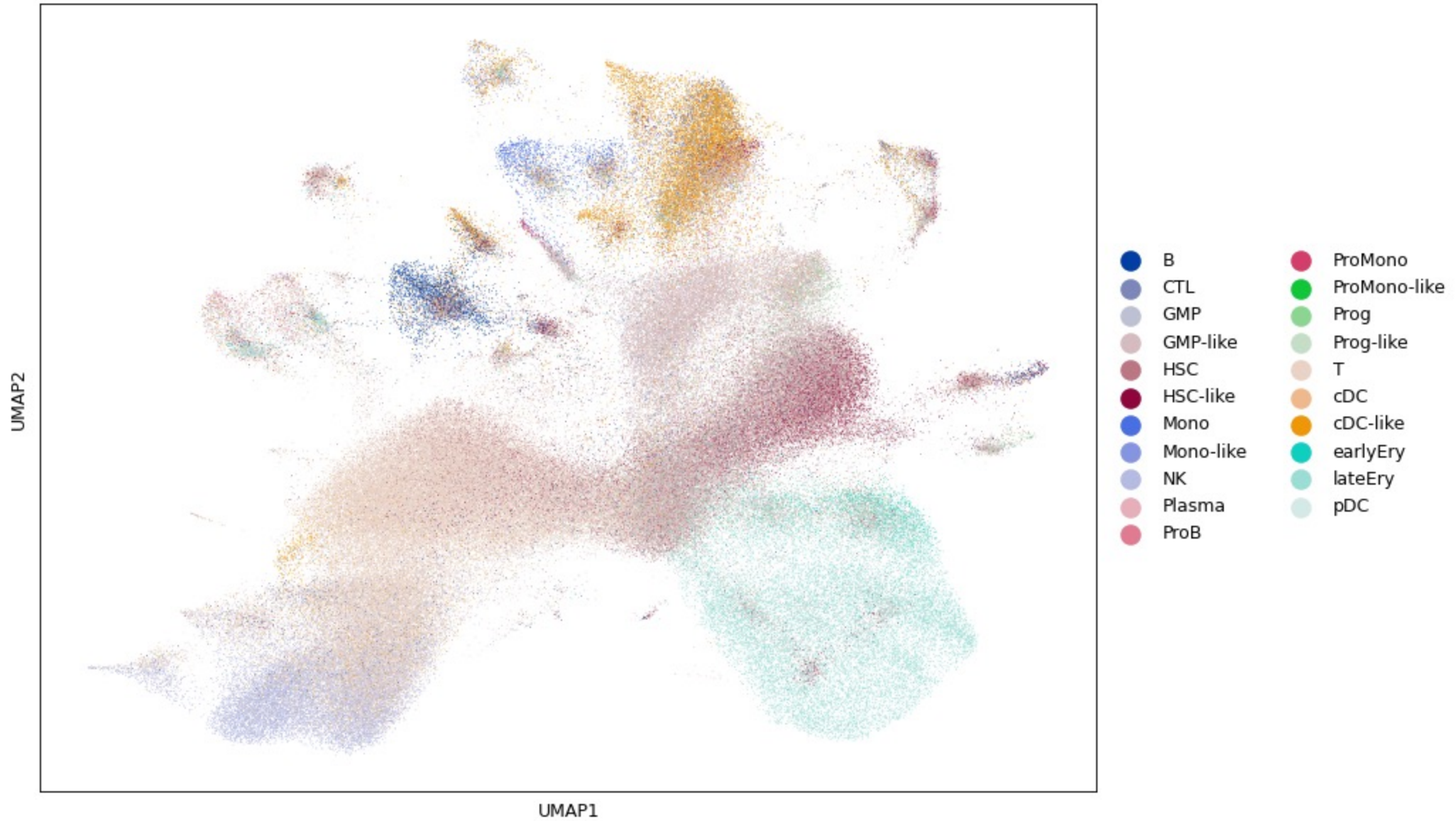
leiden

Universitas



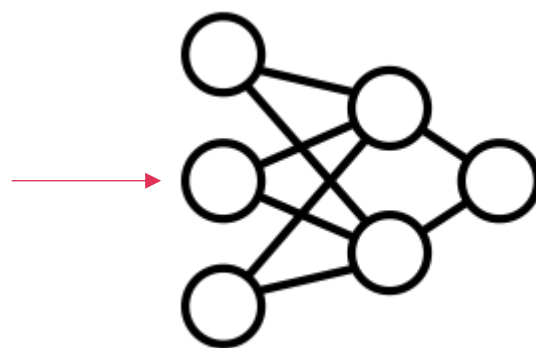
closest_Galen_2019_long

tetas



	barcode	library	total_counts	pct_counts_mito	patient	part	timepoint	response	n
55	AACAGCGAAGCACTGT	AC3_2	1705	2.05279	AC	2	3	RD	
66	AACAGCGAAGCTCGTA	AC3_2	618	2.42718	AC	2	3	RD	
97	AACAGCGAATGACCTG	AC3_2	1655	0.060423	AC	2	3	RD	
134	AACAGCGACATGAGTC	AC3_2	606	3.13531	AC	2	3	RD	
135	AACAGCGACATGCACA	AC3_2	509	2.75049	AC	2	3	RD	

	GENE_ 1	GENE_ 2	GENE_ 3	...	GENE_ N
CELL_ 1	10	0	15		0



REFRACTORY

SENSITIVE

CR vs RD timepoint 1 only

Rapidly dividing cells only

nias
iversitetas

