



**Vilnius  
University**

---

Doktoranto Karolio Šablausko ataskaita už  
2025/2026 mokslo metų pirmąjį pusmetį

Darbo vadovė: prof. Audronė Jakaitienė

**Disertacijos pavadinimas:** Genetinių pokyčių charakterizavimas naudojant giliojo mokymo neuroninius tinklus (angl. *Characterization of genetic changes using deep neural networks*)

**Darbo vadovas:** prof. Audronė Jakaitienė

**Doktorantūros pradžios ir pabaigos metai:** 2022 – 2027  
(akademinės atostogos 2024-10-01 – 2025-01-31)

**Studijų metai:** 4.

1 lentelė: Doktorantūros studijų planas

Studijų metai	Egzaminai	
	Planas	Įvykdyta
I (2022/2023)	1	1
II (2023/2024)	2	2
Akademinės atostogos 2024-10-01 – 2025-01-31		
<b>III (2025/2026)</b>	<b>1</b>	<b>1</b>
IV (2026/2027)	0	0
<b>Iš viso:</b>	<b>4</b>	<b>4</b>

# Visų doktorantūros studijų ir mokslinių tyrim planas bei jo vykdymo suvestinė

Studijų metai	Dalyvavimas konferencijose				Publikacijos			Be citavimo rodiklio
	Tarptautinėse		Nacionalinėse		Su citavimo rodikliu			
	Planas	Įvykdyta	Planas	Įvykdyta	Planas	Įvykdyta	Būklė	
I (2022/2023)	0	0	0	0	0	0	0	Publikacijų be citavimo rodiklio teikti nenumatoma
II (2023/2024)	1	0	0	0	0	0	0	
Akademinės atostogos 2024-10-01 – 2025-01-31								
<b>III (2025/2026)</b>	<b>1</b>	<b>2 pristatymai + 1 santrauka pateikta, laukiama atsakymo</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>1 publikacija + 1 rankraštis (pateikta Nature Communications – atmetė 2026-03-17, 2026-03-18 pateikta Science Translational Medicine)</b>	
IV (2026/2027)	1	0	0	0	1	0	0	
<b>Iš viso:</b>	<b>3</b>	<b>2</b>	<b>0</b>	<b>0</b>	<b>2</b>	<b>1</b>	<b>1</b>	

# Informacija apie tarptautinius renginius ir publikacijas (1)

Pristatytas stendinis pranešimas tarptautinėje konferencijoje „American Society of Hematology Annual Meeting 2025“

Pranešimo tema: „Machine learning framework for molecular classification of hematologic malignancies using transcriptome data“

803A. EMERGING TOOLS, TECHNIQUES, AND ARTIFICIAL INTELLIGENCE IN HEMATOLOGY | NOVEMBER 3, 2025

## Machine learning framework for molecular classification of hematologic malignancies using transcriptome data

Karolis Sablauskas, Egija Berga-Svitina, Tatjana Kiselova, Livija Bardina, Audrone Jakaitiene, Dmitrijs Rots

[Crossmark: Check for Updates](#)

*Blood* (2025) 146 (Supplement 1): 2580.

<https://doi.org/10.1182/blood-2025-2580>

 Split-Screen  Share  Tools  PDF

### Abstract

**BACKGROUND** Accurate characterization of hematologic malignancy is a first step in correct and tailored treatment. Tumor whole transcriptome sequencing (WTS) has recently emerged as a universal technique allowing for the accurate identification of not only all possible fusion transcripts, but also point mutations, copy number alterations, and gene overexpression, being especially useful for diagnosing cases of B-cell acute lymphoblastic leukemias (ALL). Since WTS provides gene expression landscape of the tumor sample, we decided to investigate whether the expression signatures could be used to also accurately classify a full spectrum of pediatric hematologic malignancies.

# Informacija apie tarptautinius renginius ir publikacijas (2)

Žurnale „Science“ (IF 49,7) publikuota „*High-throughput single cell -omics using semi-permeable capsules*“ kartu su autoriais iš VU Gyvybės mokslo centro.

The screenshot shows the Science journal website. At the top, the Science logo is on the left, and navigation links for 'Current Issue', 'First release papers', 'Archive', and 'About' are on the right, along with a 'Submit manuscript' button. Below the navigation, a breadcrumb trail reads 'HOME > SCIENCE > VOL. 391, NO. 6790 > HIGH-THROUGHPUT SINGLE-CELL OMICS USING SEMIPERMEABLE CAPSULES'. The article is categorized as a 'RESEARCH ARTICLE' and 'RESEARCH METHODS'. Social media sharing icons for Facebook, X, Twitter, LinkedIn, YouTube, WeChat, and Email are visible. The article title is 'High-throughput single-cell omics using semipermeable capsules'. The authors listed are DENIS BARONAS, SIMONAS NORVAISIS, JUSTINA ZVIRBLYTE, GRETA LEONAVICIENE, VINCENTA MIKULENAITE, KAROLIS GODA, VYTAUTAS KASETA, KAROLIS SABLAUSKAS, LAIMONAS GRISKEVICIUS, [...], and LINAS MAZUTIS. There is a '+1 authors' button and a link to 'Authors Info & Affiliations'. Below the authors, the journal information is: 'SCIENCE • 18 Dec 2025 • Vol 391, Issue 6790 • pp. 1138-1145 • DOI: 10.1126/science.ady7227'. At the bottom, there are icons for download (9,660), quote (1), and a 'CHECK ACCESS' button.

# Informacija apie tarptautinius renginius ir publikacijas (3)

Paruoštas rankraštis „FRANK: a pan-cancer RNA-seq classifier for childhood tumors with data-efficient learning” kartu su autoriais iš Latvijos.

2026-03-02 Pateikta Nature Communications – atmetė  
2026-03-17. Redaktorius nurodė, jog galima perkelti į  
Communications Medicine (Nature grupės žurnalas, IF  
6.3) ir straipsnis bus siunčiamas recenzijai.

2026-03-18 rankraštis pateiktas Science Translational  
Medicine (IF 14.6), laukiama atsakymo. Jeigu bus  
gautas neigiamas atsakymas, bus sprendžiama dėl  
perkėlimo Communications Medicine.

1 **Title: FRANK: a pan-cancer RNA-seq classifier for childhood tumors with**  
2 **data-efficient learning**

3

4 **Authors:** Karolis Sablauskas<sup>1,2,\*†</sup>, Tatjana Kiselova<sup>3,4</sup>, Livija Bardina<sup>3</sup>, Inga Nartisa<sup>3</sup>, Andrius Zucenka<sup>2,5</sup>,

5 Agnese Viluma<sup>6</sup>, Linda Gailite<sup>6</sup>, Audrone Jakaitiene<sup>1</sup>, Egija Berga-Svitina<sup>3</sup>, Dmitrijs Rots<sup>3,6,7,\*†</sup>

6 <sup>1</sup>Institute of Data Science and Digital Technologies, Vilnius University; Vilnius, Lithuania

7 <sup>2</sup>Hematology, Oncology and Transfusion Medicine Center, National Cancer Center, Vilnius  
8 University Hospital Santaros Klinikos; Vilnius, Lithuania

9 <sup>3</sup>Genetics Laboratory, Children's Clinical University Hospital; Riga, Latvia

10 <sup>4</sup>Bioinformatics Laboratory, Riga Stradins University; Riga, Latvia

11 <sup>5</sup>Faculty of Medicine, Hematology and Oncology Department, Institute of Clinical Medicine,  
12 Vilnius University; Vilnius, Lithuania

13 <sup>6</sup>Institute of Oncology and Molecular Genetics, Riga Stradins University; Riga, Latvia

14 <sup>7</sup>Department of Clinical Genetics, Erasmus Medical Center; Rotterdam, The Netherlands

15 \*Corresponding authors: karolis.sablauskas@santa.lt; d.rots@erasmusmc.nl

16 †These authors contributed equally to this work.

# Informacija apie tarptautinius renginius ir publikacijas (3)

Paruošto rankraščio duomenimis pateikta santrauka „European Society of Human Genetics 2026“ konferencijai.

**Konferencijos potėmė:** *“Bioinformatics, statistical methods and AI”*

**Control/Tracking Number:** 2026-A-2583-ESHG

**Activity:** ESHG Abstract

**Current Date/Time:** 3/20/2026 3:54:45 AM

**FRANK: a pan-cancer RNA-seq classifier for childhood tumors with data-efficient learning**

# Visų doktorantūros mokslinių tyrimų ir disertacijos rengimo etapai

3.	<p>Atskirų daktaro disertacijos dalių (tyrimo metodikos, rezultatų, ginamų teiginių, išvadų, ir kt.) parengimas:</p> <ol style="list-style-type: none"><li>1. Tyrimų apžvalga ir analizė.</li><li>2. Tyrimo metodikos sudarymas.</li><li>3. Teorinis tyrimas.</li><li>4. Empirinis tyrimas.</li><li>5. Gautų duomenų analizė, apibendrinimas.</li><li>6. Išvados, įvadas, literatūros sąrašas.</li></ol>	2025 gruodžio mėn. – 2026 rugsėjo mėn.	<p>Atliekama modelio hiperparametrų paieška taikant optuna, atliekama modelio validacija naudojant išsamius išorinius duomenis (viso 16 398 testavimo mėginiai), atliekama modelio “tvirtumo” (angl. robustness) analizė, testavimas naudojant klinikinius duomenis. Paruošiamas ir publikacijai pateikiamas rankraštis. Remiantis rankraščio duomenimis pateikiama santrauka konferencijai.</p>
----	--	---	--



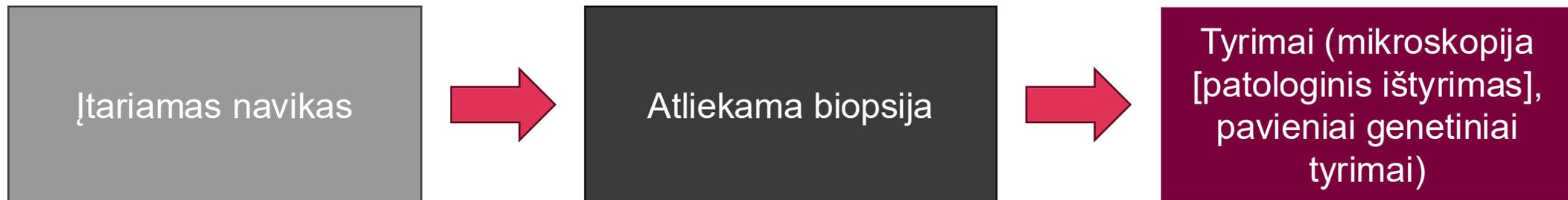
## 1.1 Thesis aim (Tyrimo tikslas)

To contribute to the advancement of deep learning techniques for the analysis of next generation sequencing data, with a focus on RNA sequencing data.

## 1.2 Thesis objectives (Tyrimo uždaviniai)

- **Data preprocessing and feature engineering:** create efficient data preprocessing pipeline suited for RNA sequencing 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.

# Mokslinių rezultatų pristatymas: tyrimo problematika



Tikslas: konkrečiai  
įvardinti naviko tipą

Trukmė: 1-4 sav.

# Tyrimo problematika

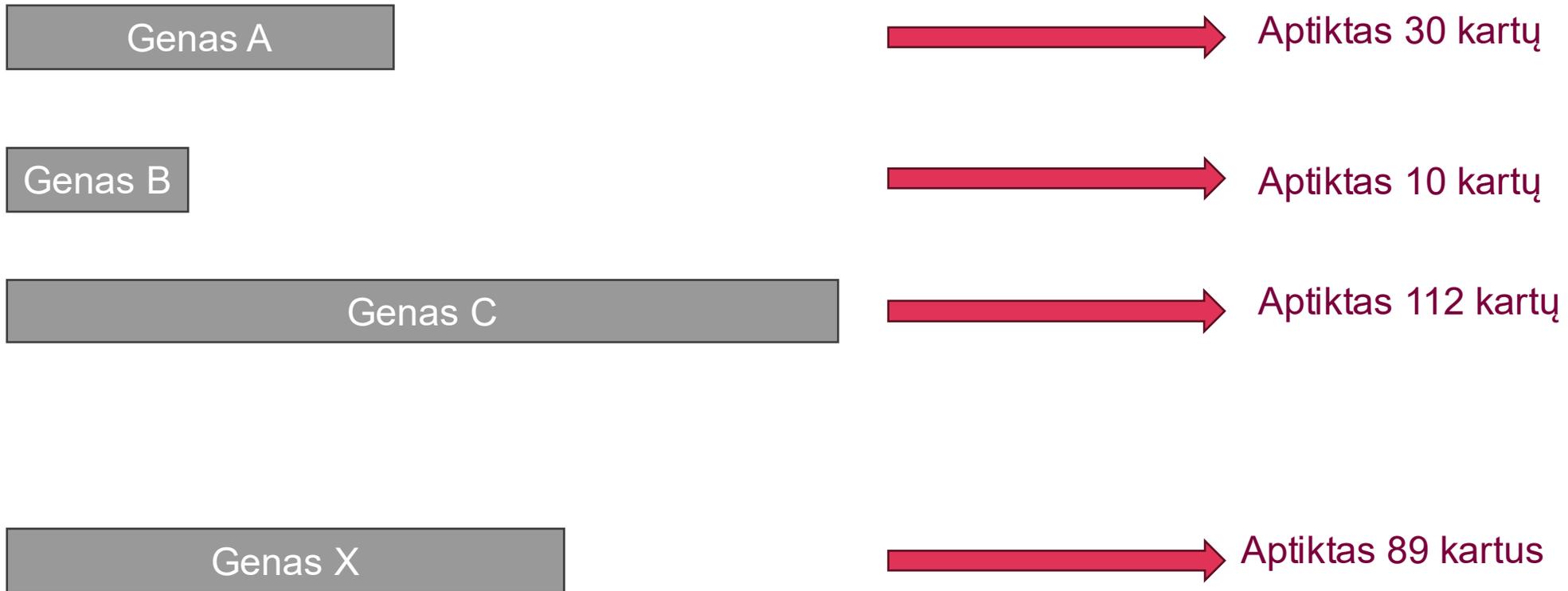
**Pavienių diagnostinių tyrimų trūkumai:**

- Sudėtinga iš infrastruktūros pusės (kiekvienam navikui reikia atskiro lab. metodo)
- Laiko trukmė gauti rezultata.

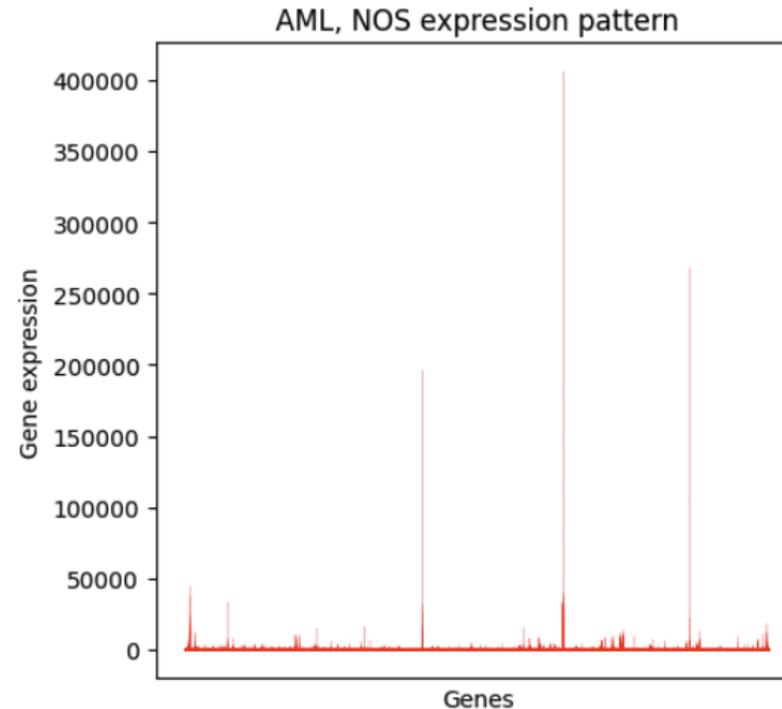
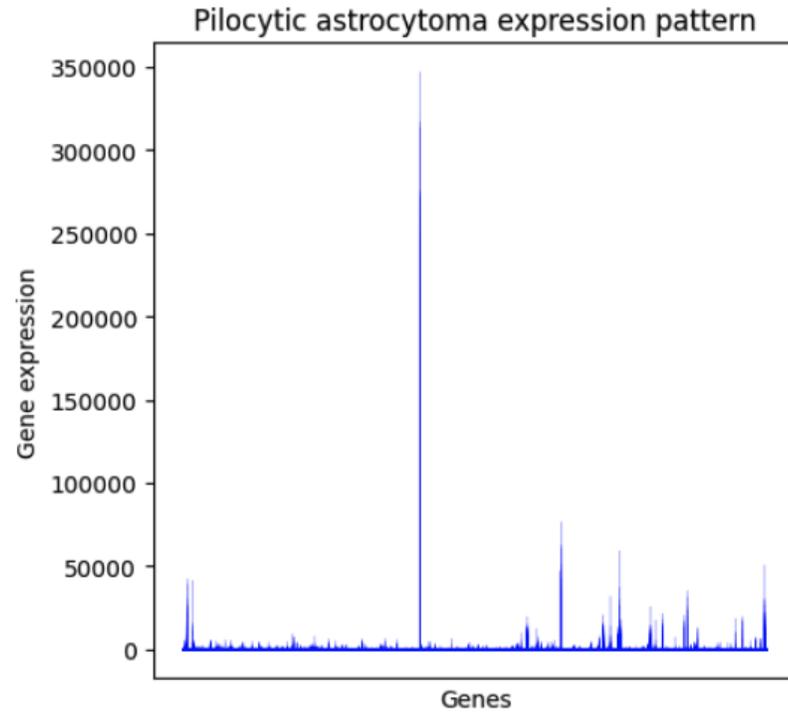
**Galimas sprendimas:**

- “omikos” metodai – didelio kiekio genų ištyrimas

# Viso transkriptomu sekoskaita (angl. Whole transcriptome sequencing - WTS)



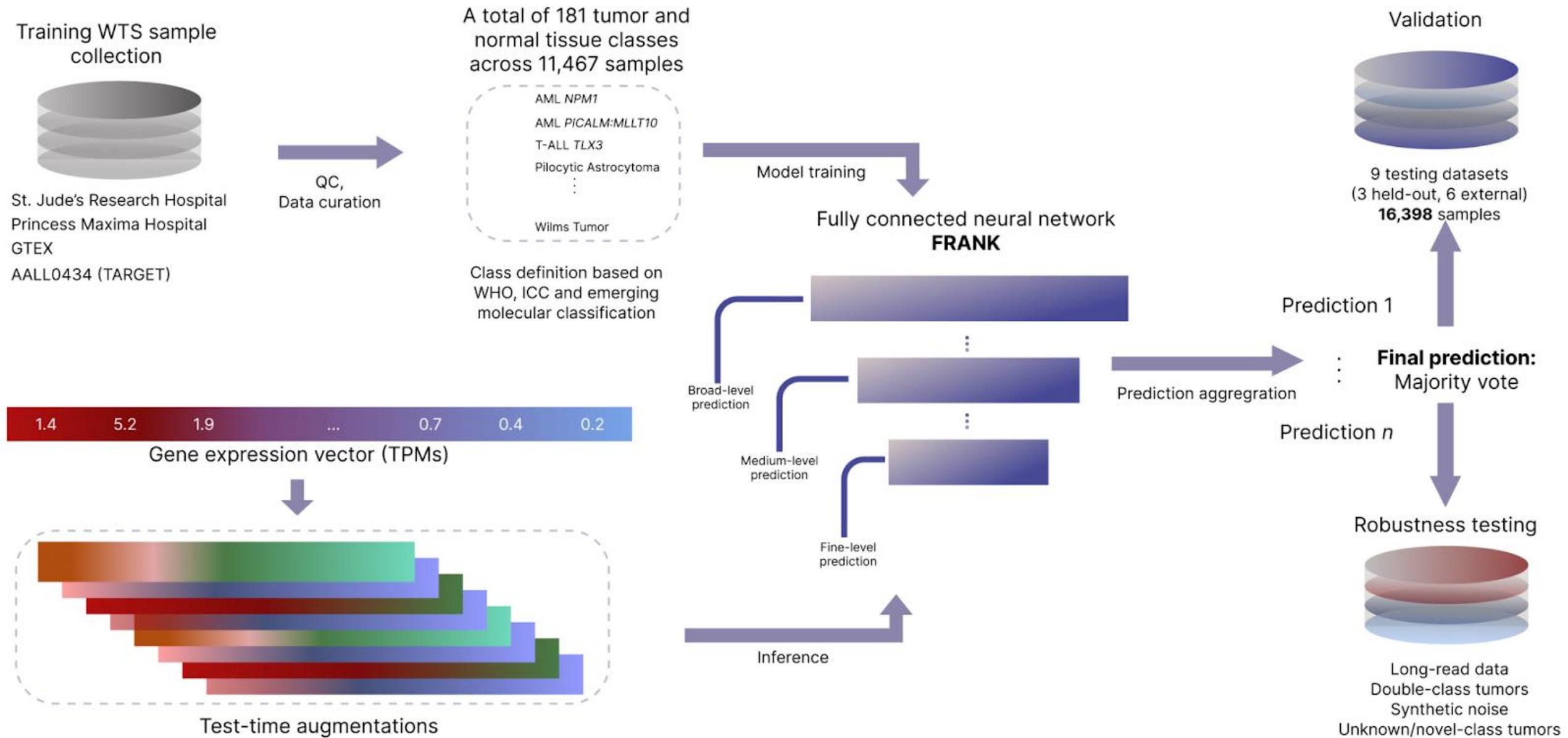
# Genų raiškos palyginimas



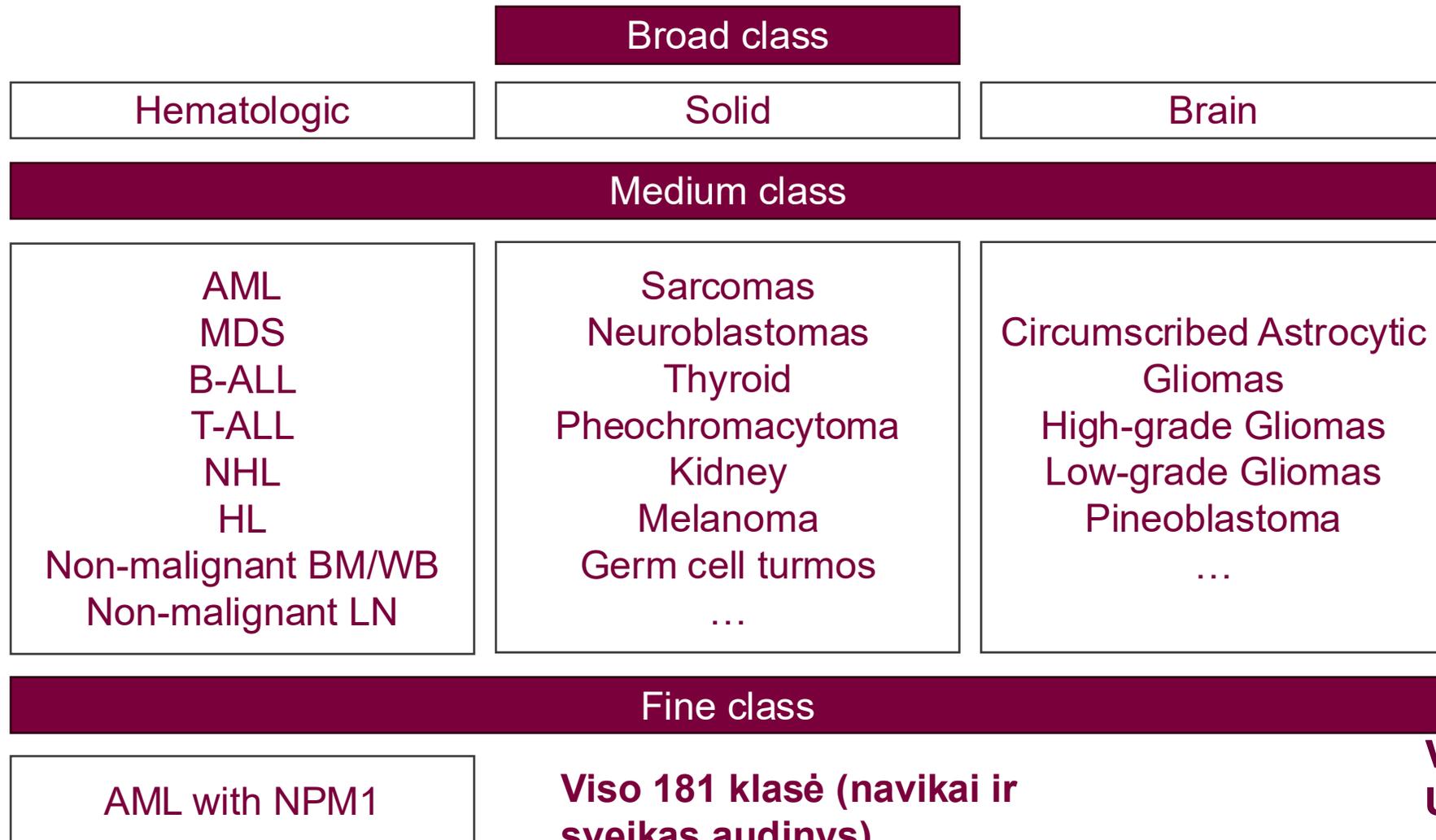
[ 0. , 20.16973919, 0. , ..., 0.3646838 ,

Using WTS gene expression to predict a specific malignancy subtype

No direct detection of fusion, mutations, etc.



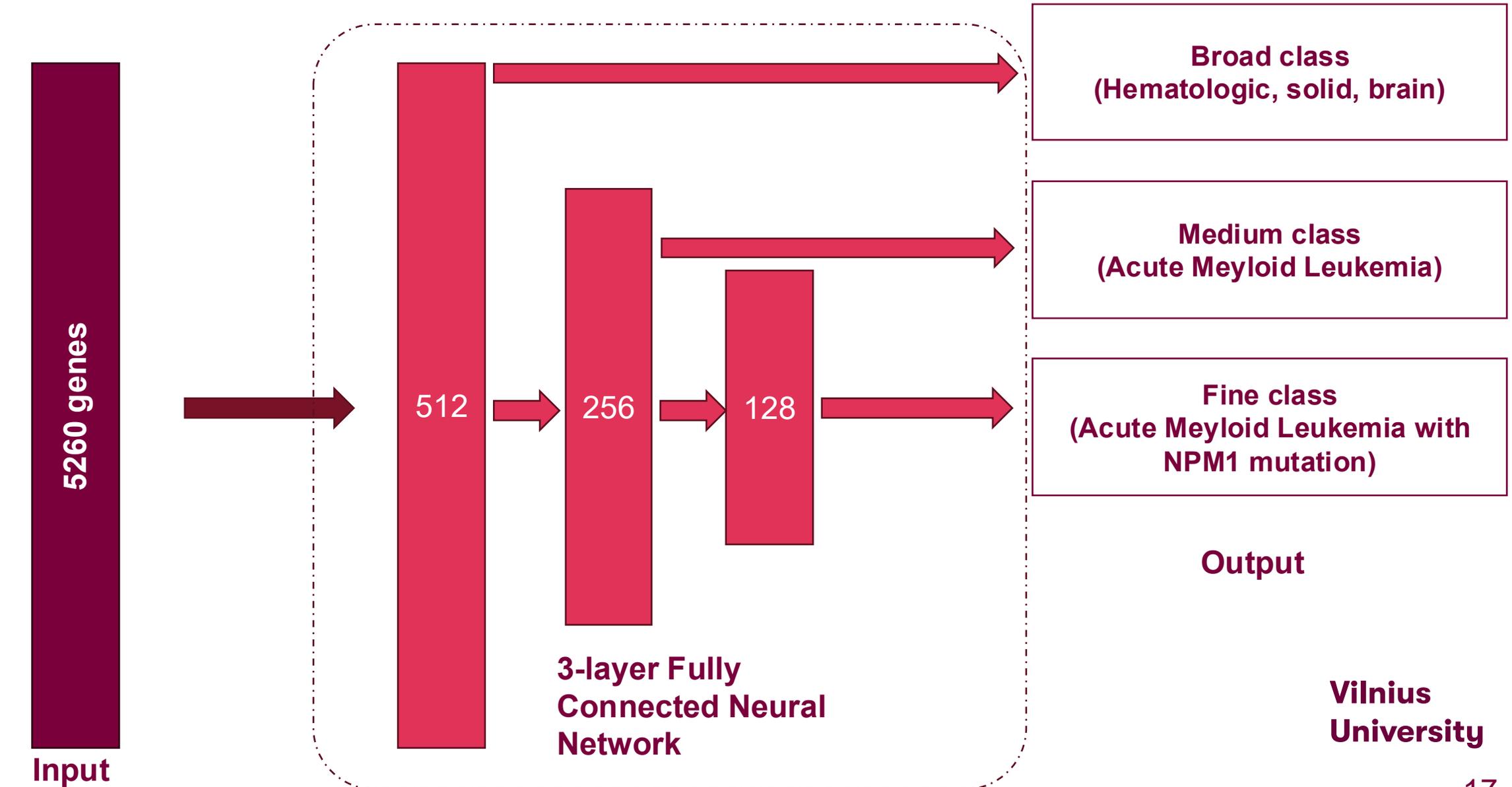
# Hierarchinė navikų klasifikacija

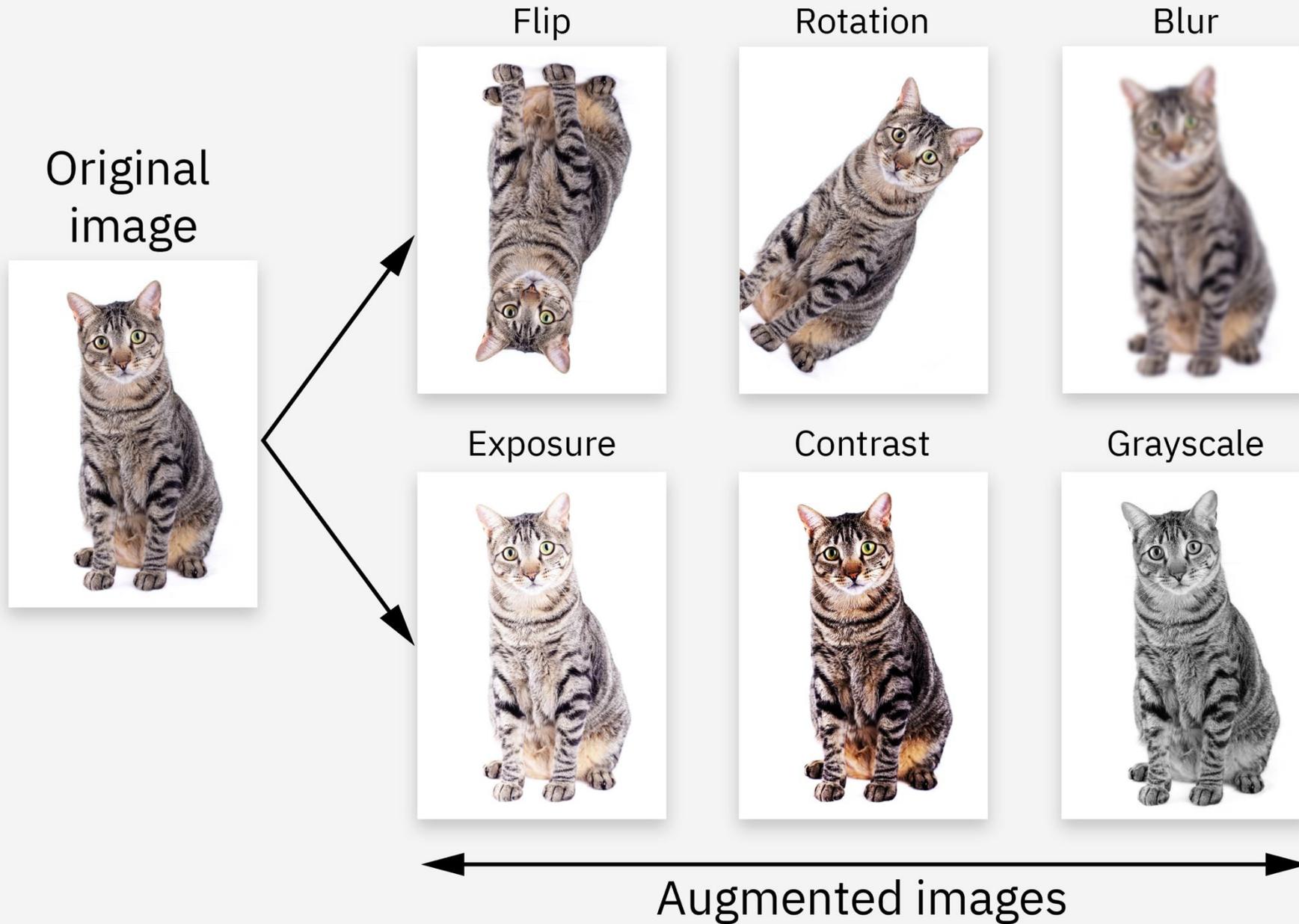


# Mokymo duomenų apžvalga

Cohort	N individuals	N samples	N classes (fine level)	Tissue Preservative	Enrichment Method	GENCODE Version
St. Jude	3,911	4,319 (Tumor $n=4319$ ) (Normal $n=0$ )	144 (Pan-tumor)	FFPE ( $n=805$ ) Fresh/Stabilized ( $n=2,973$ ) Not Available ( $n=541$ )	Not Available ( $n=5$ ) PolyA ( $n=889$ ) Total RNA ( $n=3,425$ )	v31
Princess Maxima	1,136	1,149 (Tumor $n=1110$ ) (Normal $n=39$ )	80 (Pan-tumor)	Fresh/Stabilized ( $n=1149$ )	Total RNA ( $n=1,149$ )	v29
GTEX v10 release	283	5,714 (Tumor $n=0$ ) (Normal $n=5,714$ )	31 (All normal tissue)	Fresh/Stabilized ( $n=5714$ )	Total RNA ( $n=5,714$ )	v39
AALL0434	285	285 (Tumor $n=285$ ) (Normal $n=0$ )	7 (T-ALL)	Fresh/Stabilized ( $n=285$ )	Total RNA ( $n=285$ )	N.A.
Total	5,615	11,467 (Tumor $n=5,714$ ) (Normal $n=5,753$ )	181			

# Klasifikatorius





# — Duomenų papildymas

Triukšmas viso mėginio lygiu

120	455	...	13
119	457	...	15
Gene A	Gene B	...	Gene X

Trūkstami genai

120	455	...	13
120	0	...	13
Gene A	Gene B	...	Gene X

Triukšmas vieno geno lygiu

120	455	...	13
120	677	...	13
Gene A	Gene B	...	Gene X

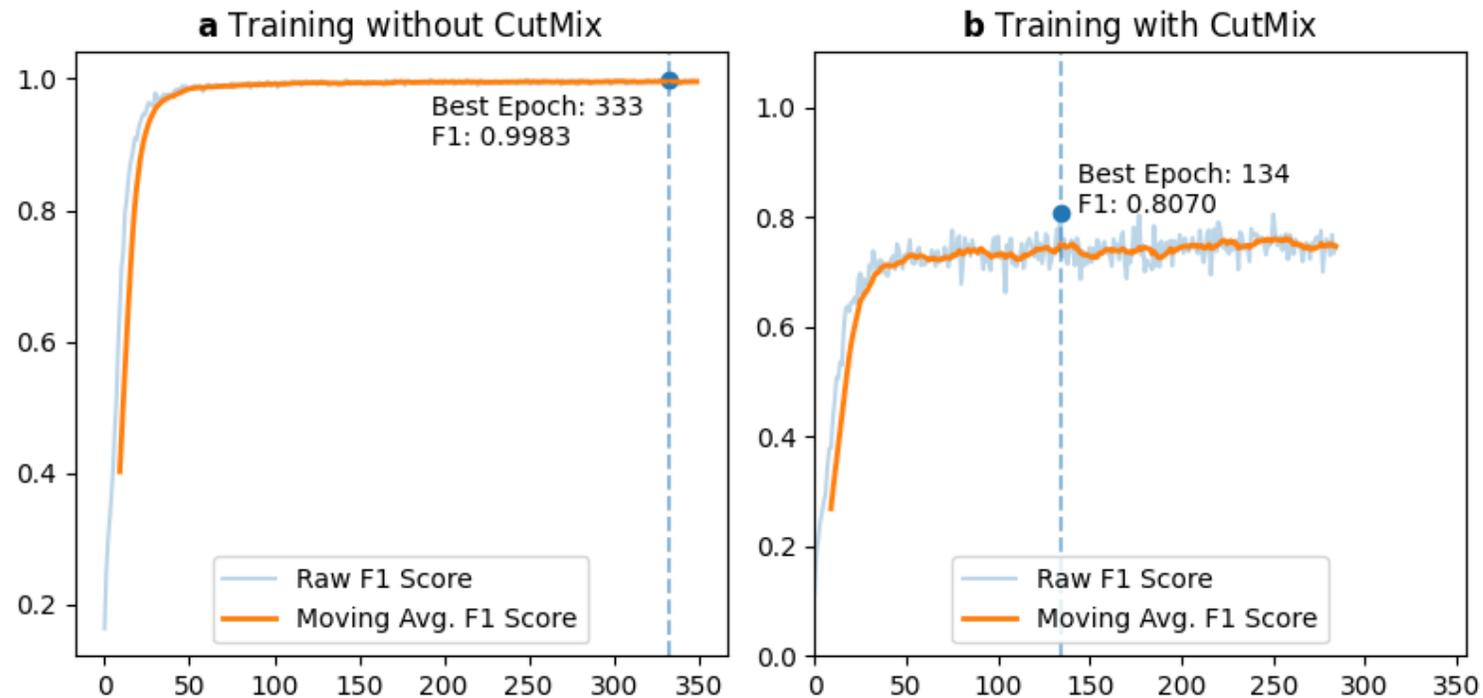
*CutMix*

120	455	...	13
256	567	...	8

120	455	...	8
-----	-----	-----	---

# Mokymo pasirinkimas

- Daugiaklasės klasifikacijos uždavinys - 181 klasė
- Klasių disbalansas (3 – 303 mokymų pavyzdžių klasei)
- Mokymas naudojant visus duomenis – *CutMix* padeda išvengti persimokymo.

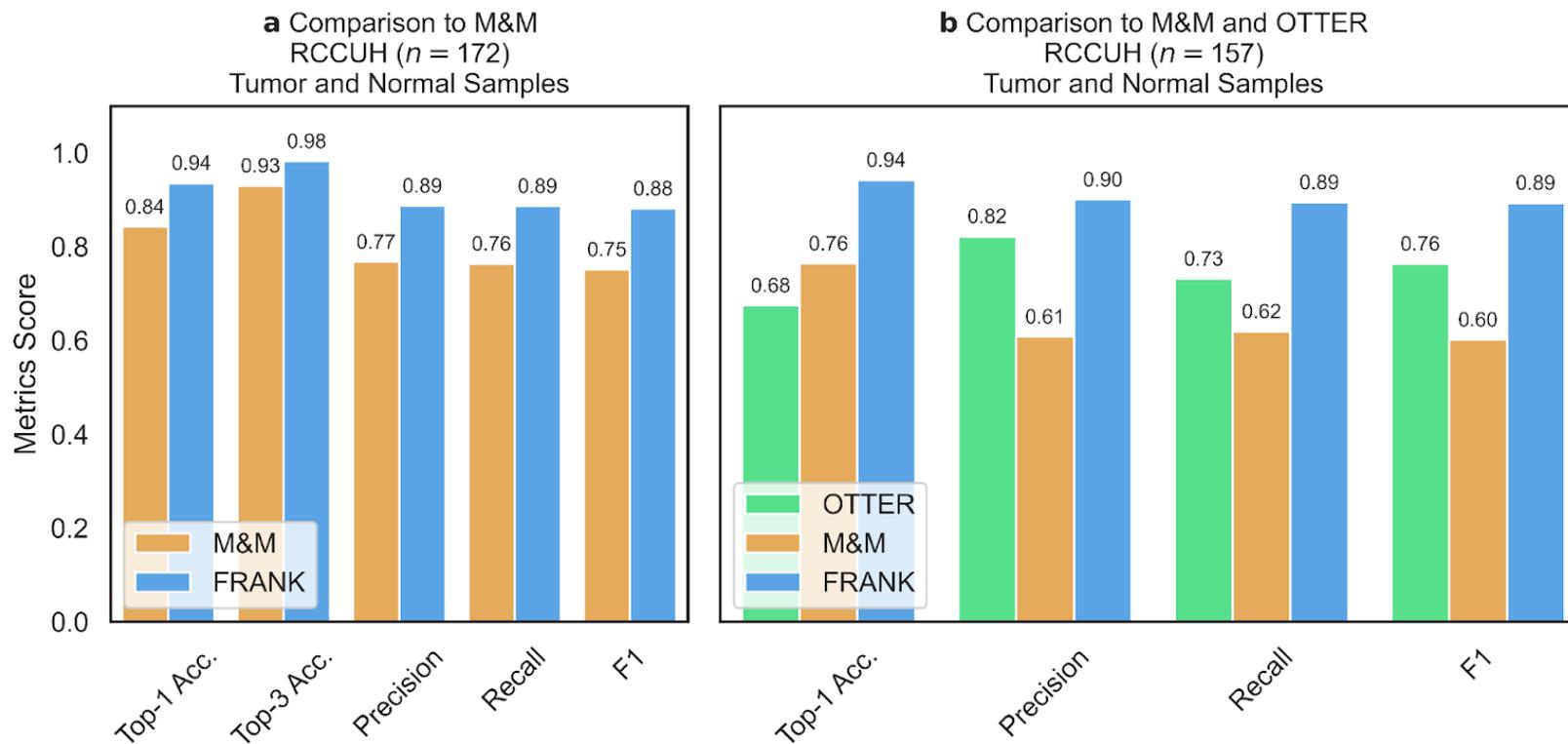


# Rezultatai testavimo kohortose

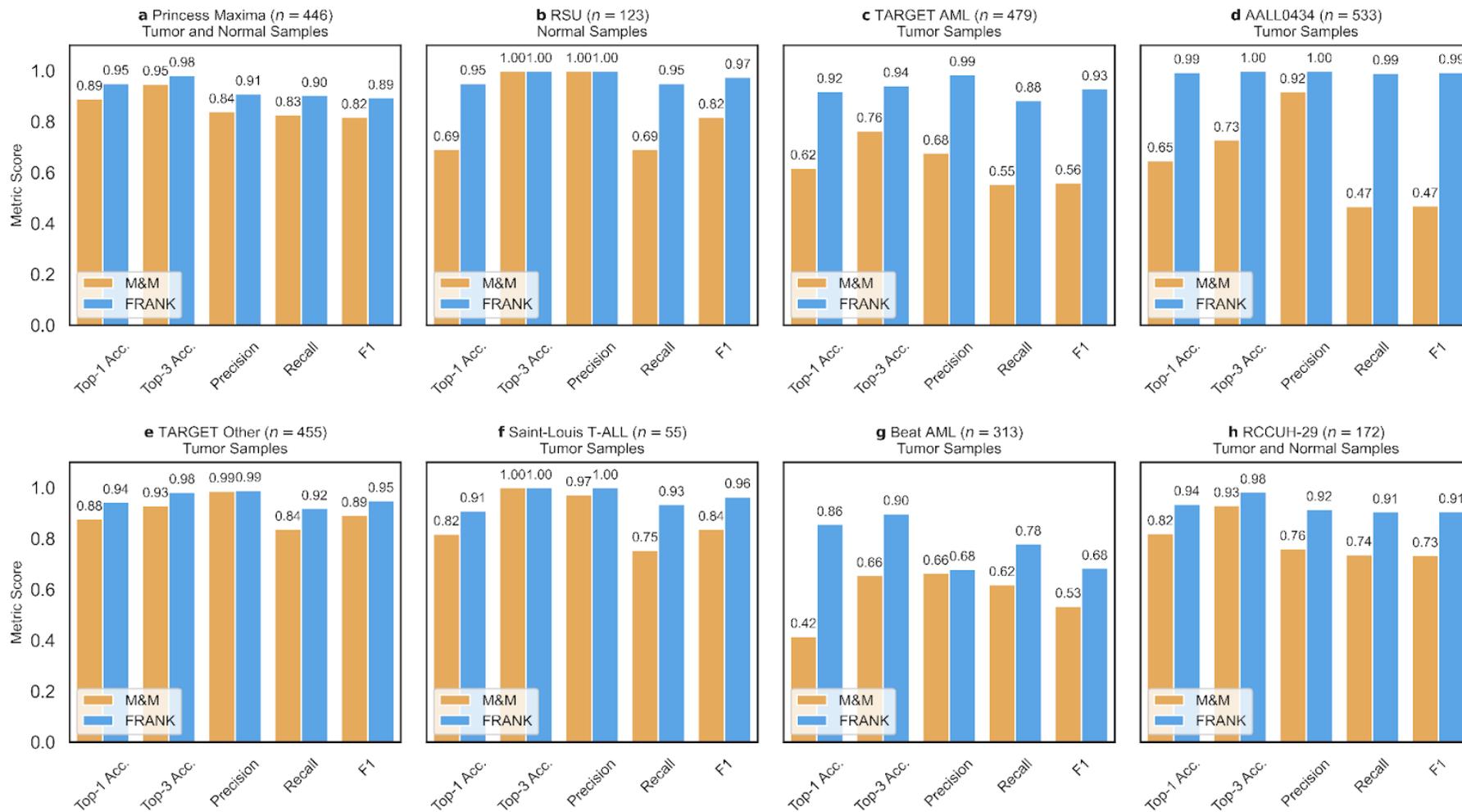
Dataset	External dataset	N samples	N errors	N tumor/normal mismatches	Accuracy (%)	F1 score	Precision	Recall
RCCUH	Yes	178	13	1	92.7	0.87	0.88	0.88
Princess Maxima	No	430	21	4	95.1	0.90	0.92	0.90
TARGET AML	Yes	489	42	7	91.4	0.80	0.87	0.77
TARGET Other	Yes	468	32	1	93.2	0.90	0.99	0.84
AALL0434	No	707	37	0	94.8	0.85	0.90	0.82
Saint-Louis T-ALL	Yes	72	7	0	90.3	0.87	0.86	0.91
Beat AML	Yes	329	49	5	85.1	0.72	0.72	0.79

# Palyginimas su kitais literatūroje aprašytais modeliais

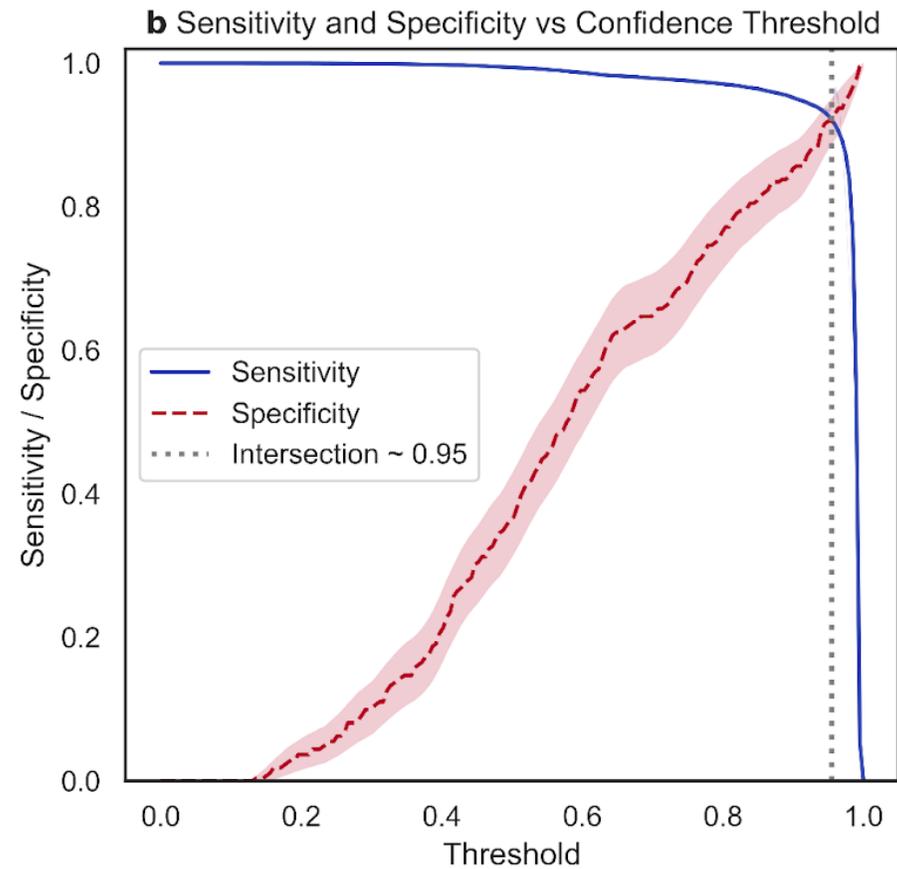
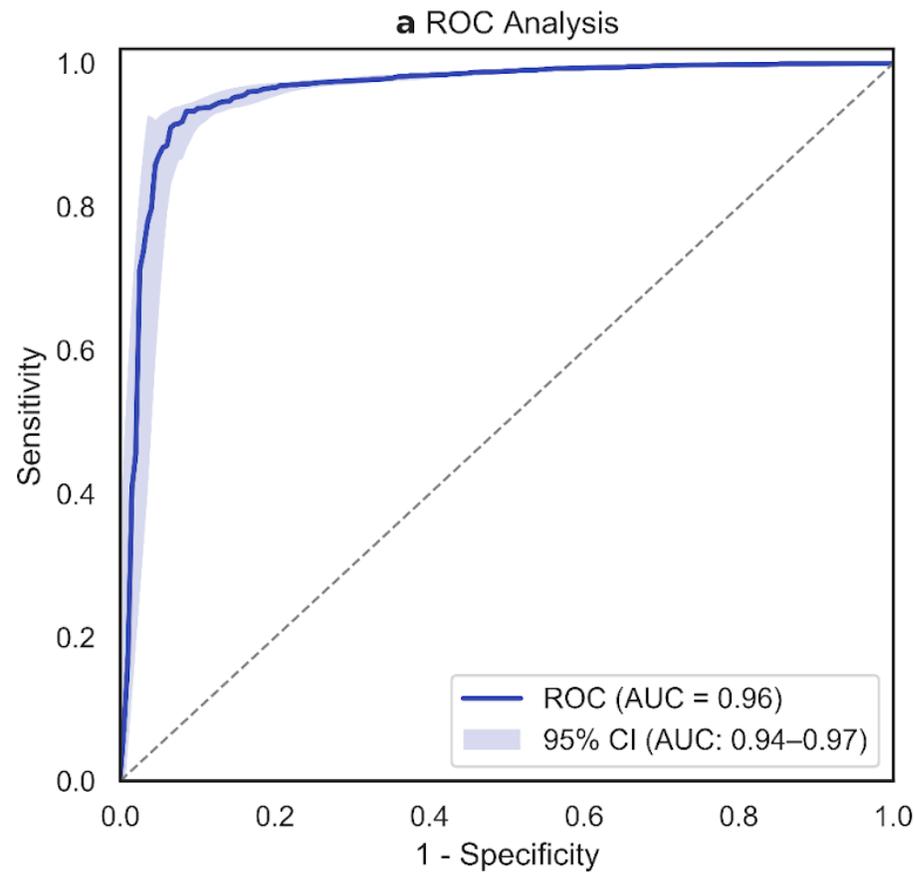
- M&M ir OTTER
- Rygos ligoninės (RCCUH) duomenų rinkinys (nenaudotas publikacijose, patikimi ir “išvalyti” duomenys)
- M&M 172/189 mėginiai tinkami testavimui, OTTER 157/189 mėginiai tinkami testavimui (visi tinkami OTTER, tinkami ir M&M)



# FRANK vs M&M



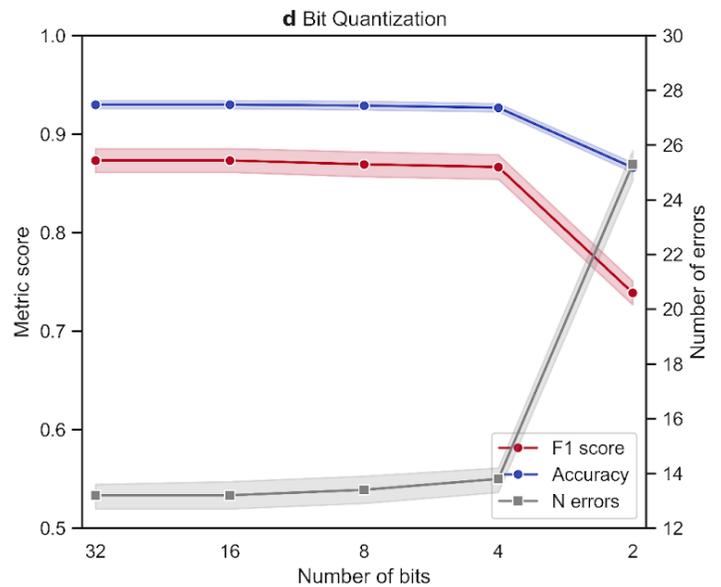
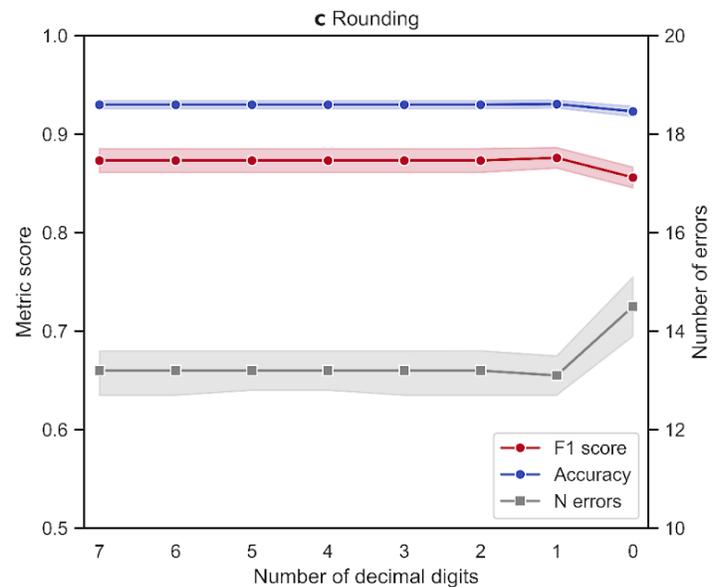
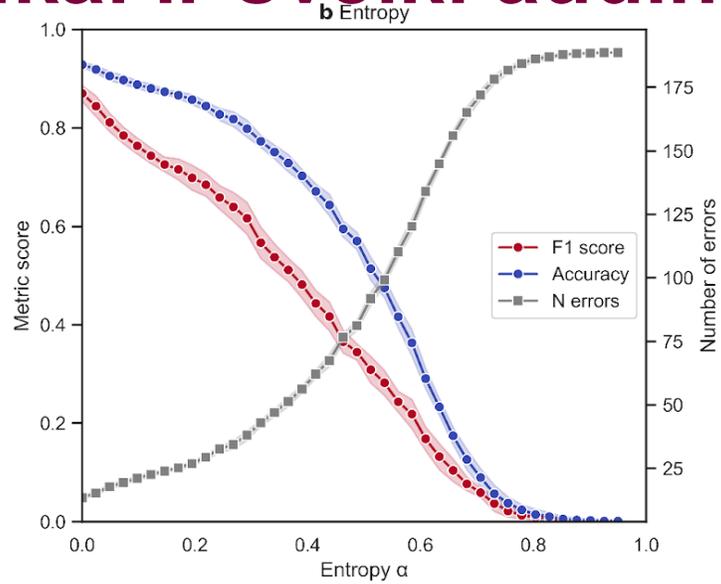
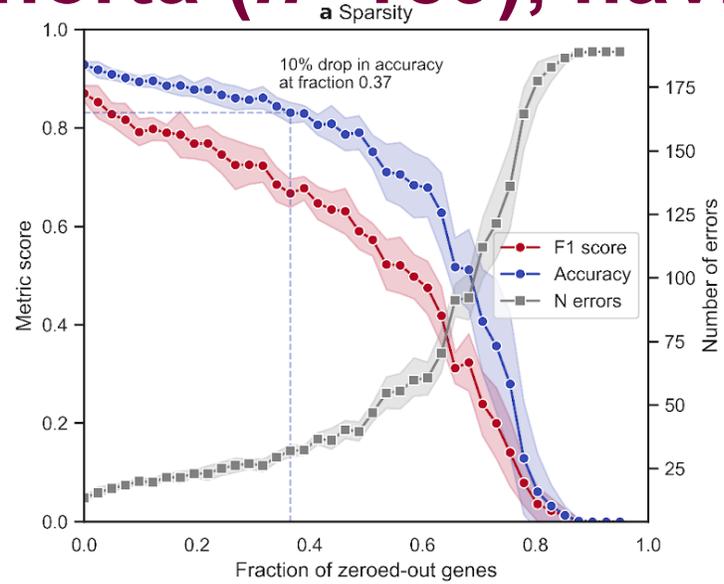
# ROC analizė ( $n=16,398$ , tumor and normal)



# Aggregated metrics on scorable samples

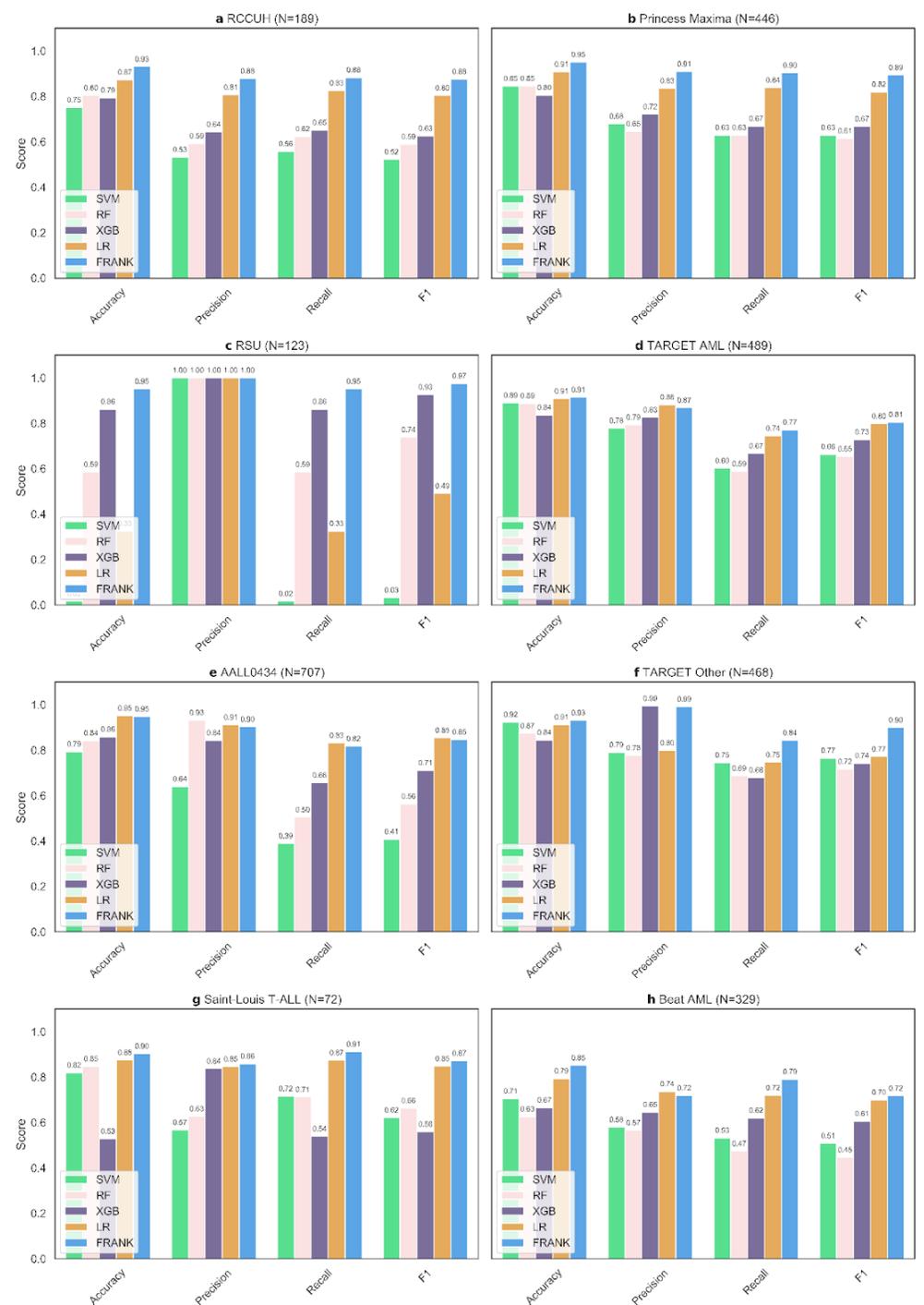
Sample selection	N	N errors	N tumor/normal mismatch	Accuracy (%)	F1 score	Precision	Recall
Tumor and normal tissue - no confidence threshold	16398	272	26	98,3	0,86	0,87	0,88
Tumor and normal tissue - high confidence	14878	22	0	99,9	0,99	0,99	0,99
Tumor only - no confidence threshold	2673	201	18	92,5	0,82	0,84	0,85
Tumor only - high confidence	1785	11	0	99,4	0,98	0,98	0,98
Normal tissue only - no confidence threshold	13725	71	8	99,5	0,98	0,99	0,97
Normal tissue only - high confidence	13093	11	0	99,9	0,99	1	0,99

# Modelio tvirtumo (robustness) analizė, RCCUH kohorta ( $n=189$ ), navikai ir sveiki audiniai



# Palyginimas su kitais klasifikavimo metodais

- Logistinė regresija - panašiausi rezultatai

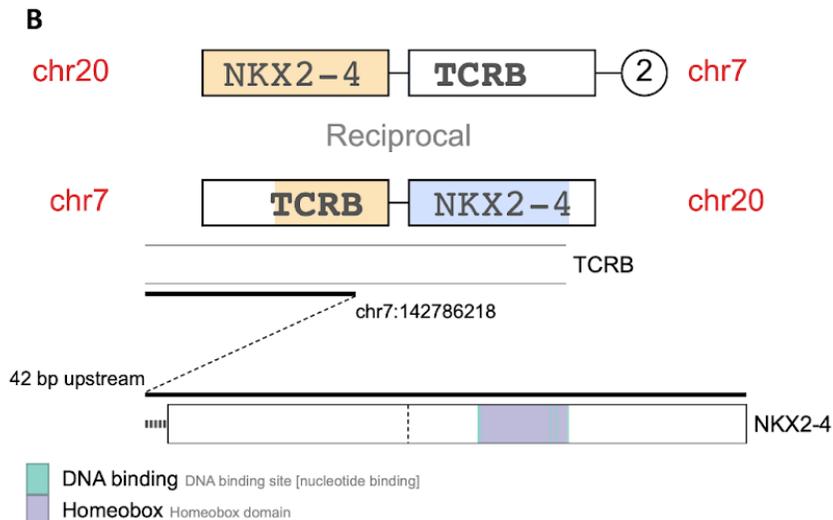
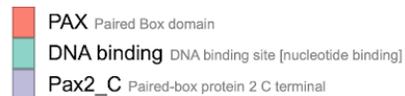
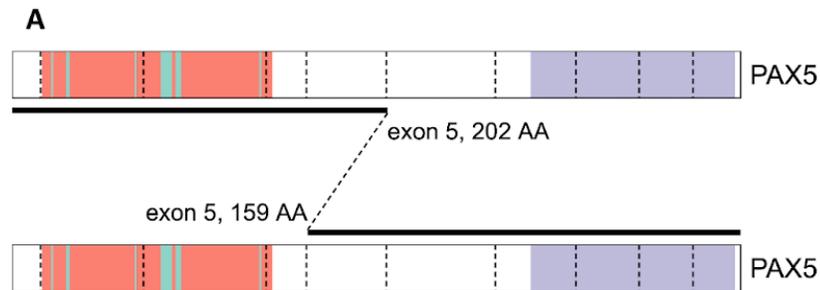


dataset	sample_id	ground_truth	fine_top1_preds	fine_top1_confs
BKUS	220135 45RN	T-cell Acute Lymphoblastic Leukemia, NKX2 Activation	T-cell Acute Lymphoblastic Leukemia, NKX2 Activation	0,856143594

## Initially labeled as “T-ALL, NOS”

dataset	sample_id	ground_truth	fine_top1_preds	fine_top1_confs
BKUS	230137 68DN	B-cell Acute Lymphoblastic Leukemia, PAX5 Alteration	B-cell Acute Lymphoblastic Leukemia, PAX5 Alteration	0,98133 7

## Initially labeled as “B-ALL, NOS”, confirmed to have PAX5 duplication



# Naujumas

- Aptinkamų navikų skaičius (181 vs 96 prieš tai aprašyti literatūroje).
- Duomenų augmentavimo strategija, CutMix panaudojimas, hierarchinės klaidos nuostolio funkcija.
- Išsamus testavimas (M&M – pagrindinė testavimo kohorta iš tos pačios institucijos kaip ir mokymo duomenys).

# Kito semestro planas

- Publikuotas pirmas rankraštis (Science Translational Medicine arba Communications Medicine žurnalai)
- Pradėtas antras rankraštis
- Pranešimo skaitymas tarptautinėje konferencijoje (European Society of Human Genetics)