



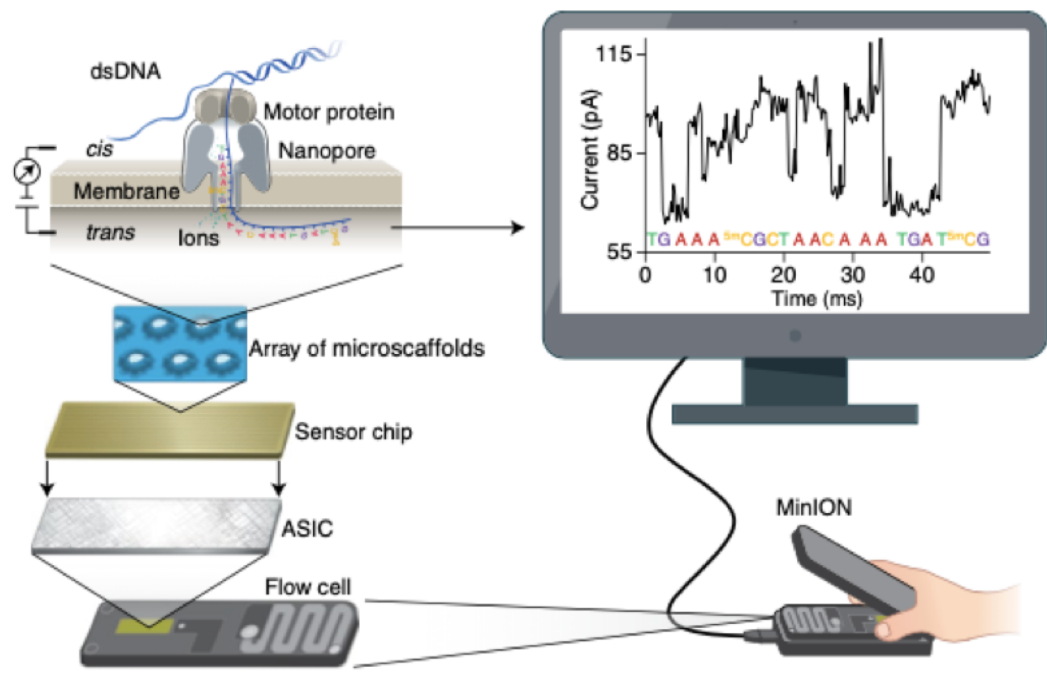
Challenges in the Next Generation (Oxford Nanopore) direct RNA Sequencing (dRNA-seq) Data Processing: Coping With Normalisation of Gene Expression Levels by Principal Component Analysis

A.Kriščiukaitis^{1,2}, R.Petrolis^{1,2}, R.Dragunaitė², R.Stakaitis², D.Skiriutė²

¹ Dept. Physics, Mathematics and Biophysics, Lithuanian University of Health Sciences

² Neuroscience Institute, Lithuanian University of Health Sciences

Nanopore sequencing technology



Modified from Wang, Y., Zhao, Y., Bollas, A. et al. Nanopore sequencing technology, bioinformatics and applications. Nat Biotechnol 39, 1348–1365 (2021)

Introduction: Next generation sequencing (NGS) is getting widely applicable for the detection of molecular markers in modern medical diagnostics. However, processing and evaluation of RNA-seq data is challenging due to some inescapable physical factors causing certain bias in the analyzed data.

Problem: Technical differences (“batch effects”) caused by differences in sample processing (up to RNA extraction, RNA-seq library preparation or the number of live pores) may significantly affect the ability to draw generalizable conclusions from such studies.

The Aim: Elaboration of the method for gene expression level analysis avoiding bias caused by inescapable physical factors in RNA-seq data.

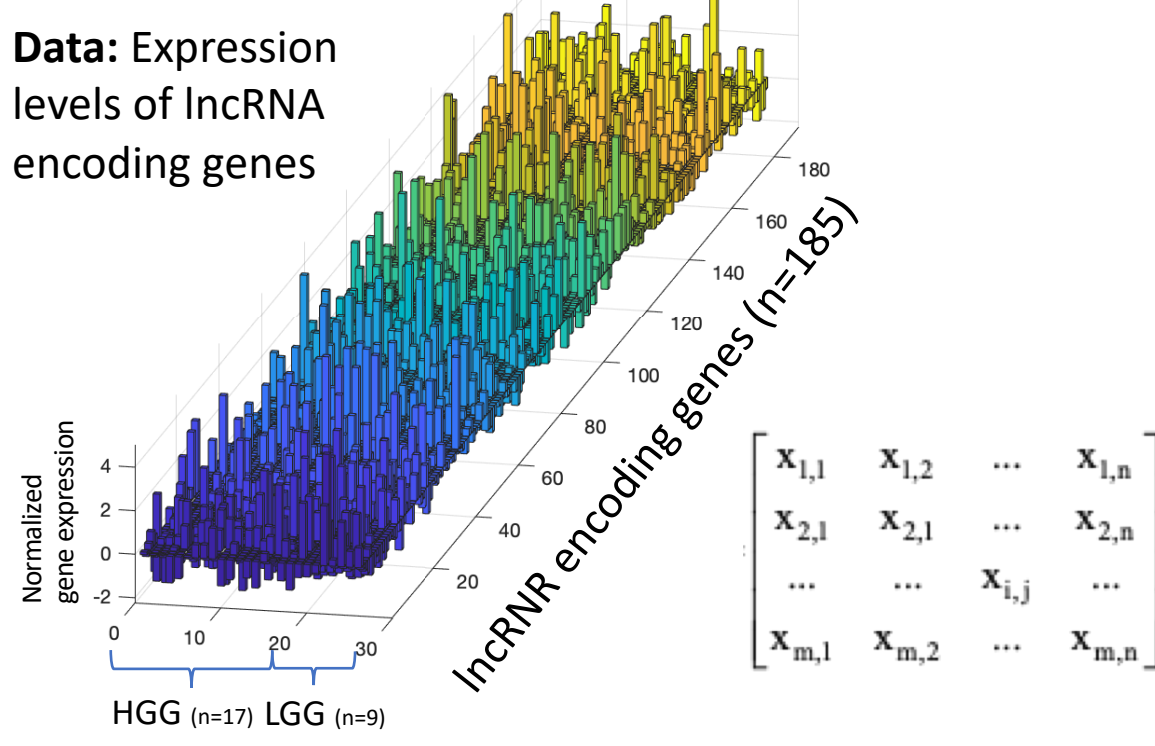
Solution: The key solution to normalization of gene expression levels estimated across different samples lays in so called **housekeeping genes (HKG)**, which are required for the maintenance of the basal cell functions. Thus, they are expected to be equally expressed in all cells.

The list of HKG: [Trends in Genetics 29 (2013), 569–574]

Gene Name	Gene description
C1orf43	chromosome 1 open reading frame 43
CHMP2A	charged multivesicular body protein 2A
EMC7	ER membrane protein complex subunit 7
GPI	glucose-6-phosphate isomerase
PSMB2	proteasome subunit, beta type, 2
PSMB4	proteasome subunit, beta type, 4
RAB7A	member RAS oncogene family
REEP5	receptor accessory protein 5
SNRPD3	small nuclear ribonucleoprotein D3
VCP	valosin containing protein
VPS29	vacuolar protein sorting 29 homolog

Example of method application – questions to answer:

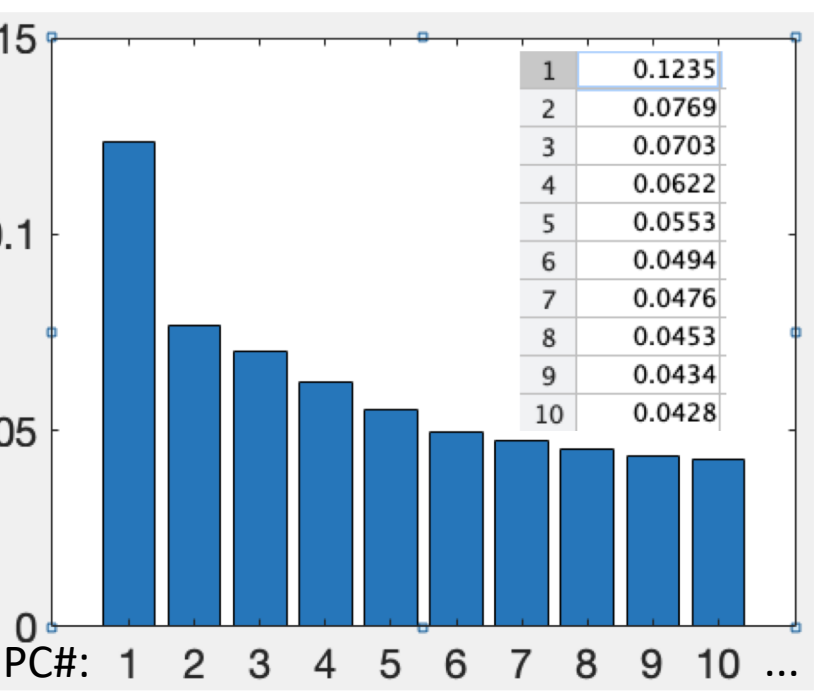
- Does expression of non-coding-RNA (ncRNA) encoding genes is different in Glioblastoma vs Low-grade-glioma (GBM vs LGG) cases?
- If yes, of which genes?



Data: Expression levels of ncRNA encoding genes

Covariation matrix of normalized gene expression levels data set: $R_x = E[X \cdot X^T]$

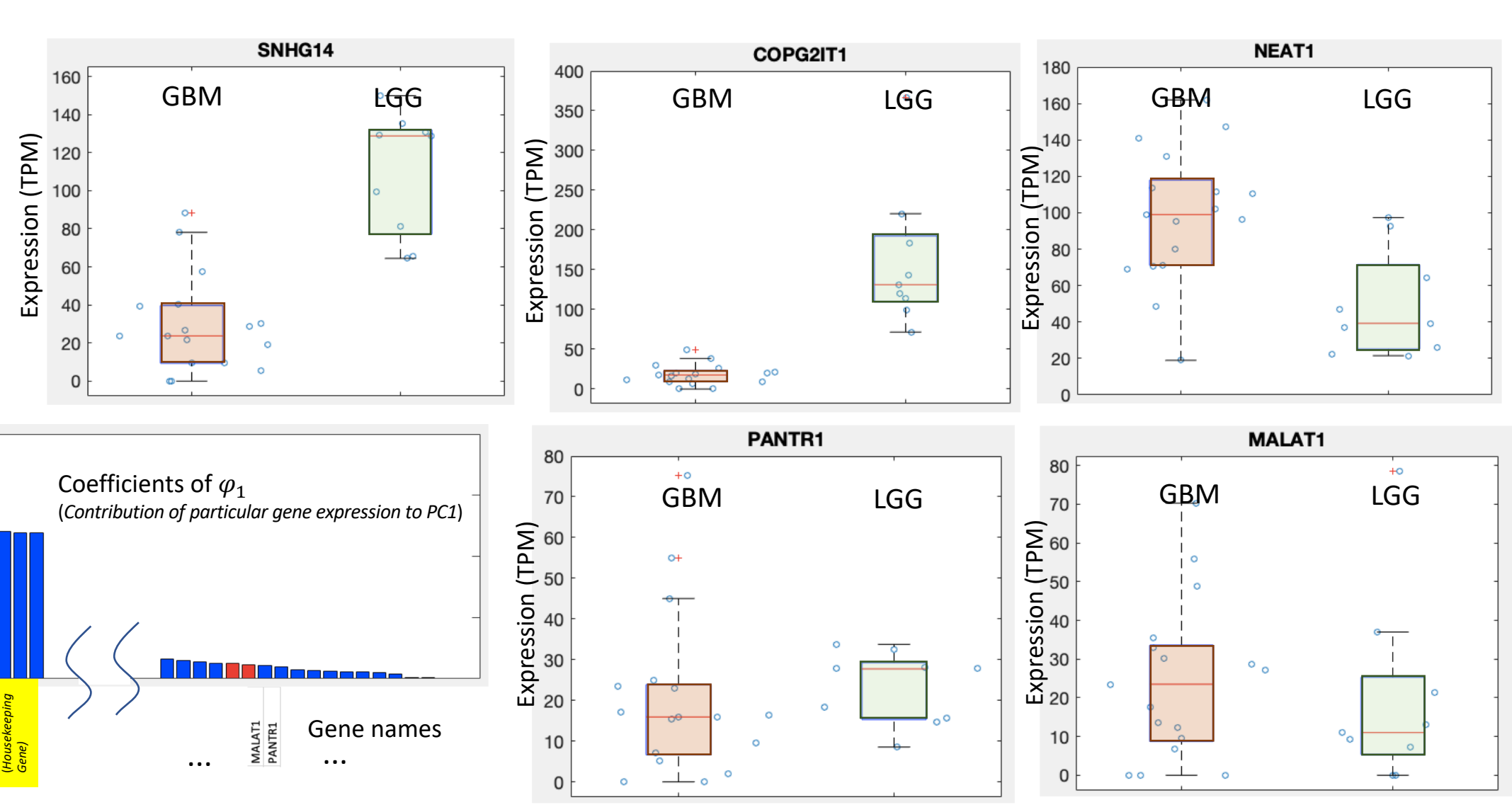
Contribution of first eigenvectors of covariation matrix in representation of total variance in analyzed data



Results: ncRNAs encoding genes which expression differs between GBM and LGG cases.

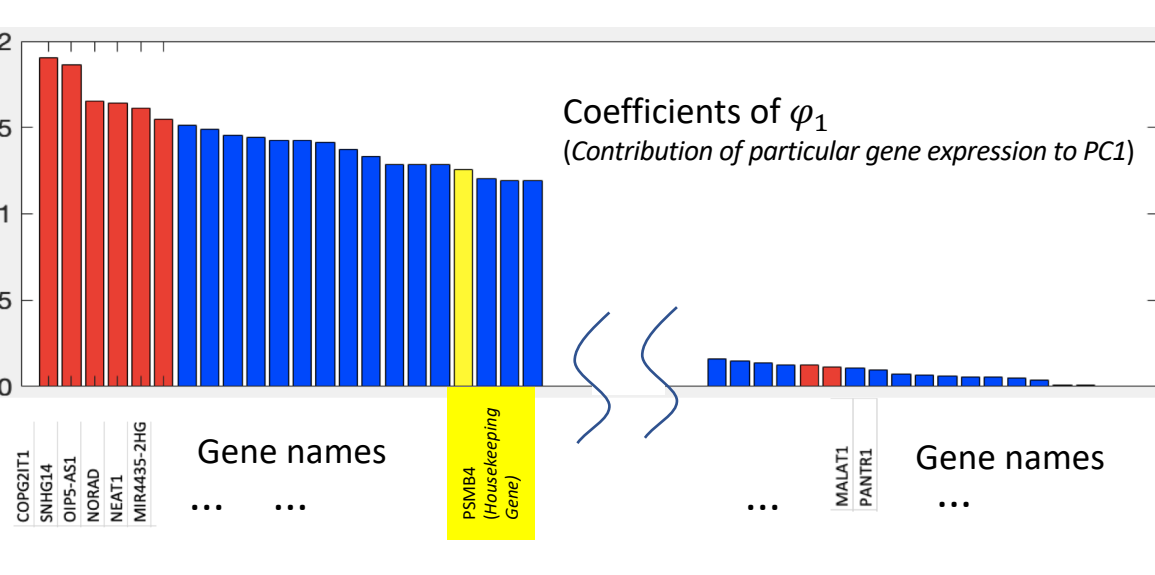
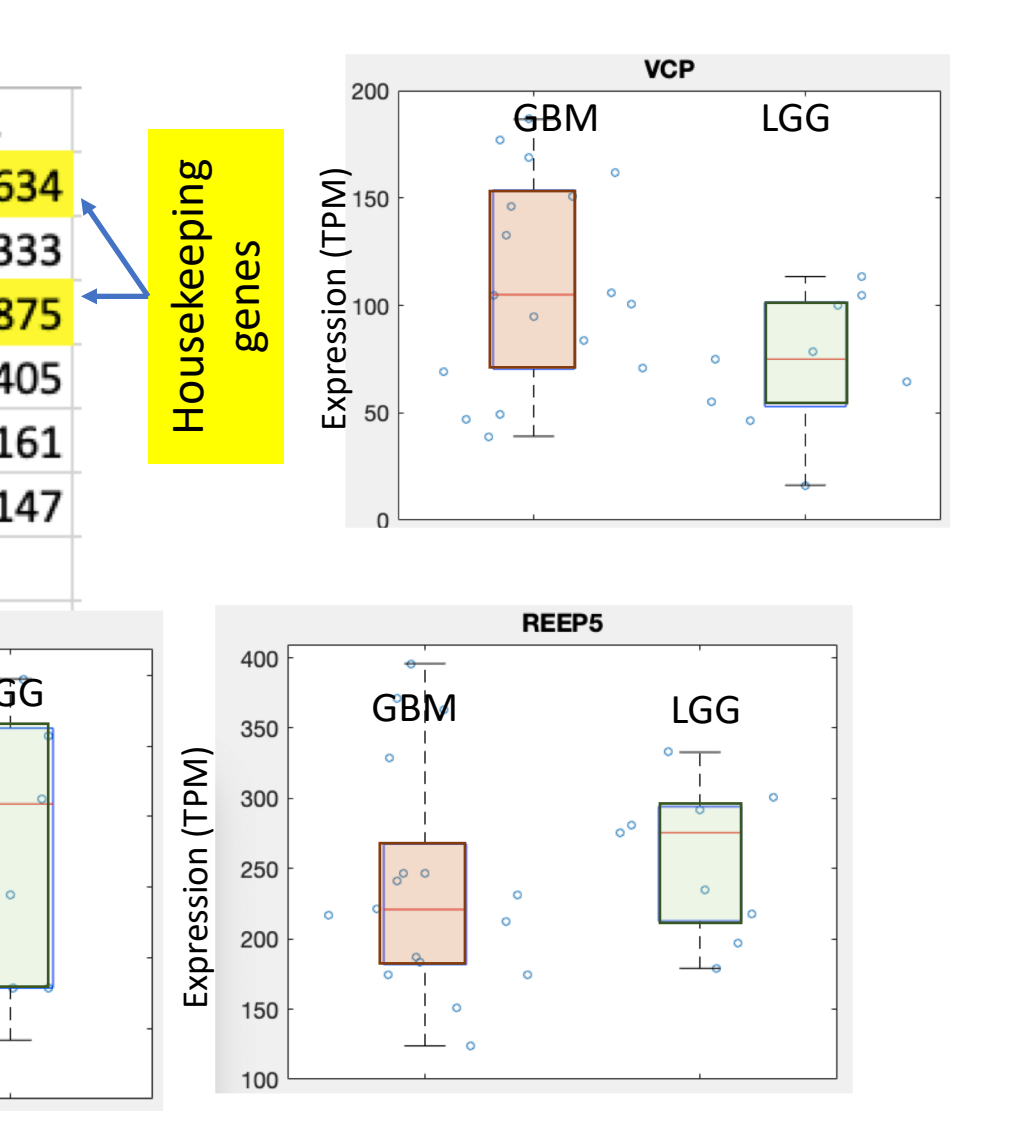
Content of φ_1

Gene Name	PC1 Koef.
COPG2IT1	-0,190234
SNHG14	-0,1861258
OIP5-AS1	-0,1650589
NORAD	-0,1639362
NEAT1	0,1611103
MIR4435-2HG	0,15489642
...	...
MALAT1	0,0124793
PANTR1	-0,011318



Content of φ_2

Gene Name	PC2 Koef.
REEP5	0,21475634
SNHG32	-0,190333
VCP	-0,1828875
CHASERR	-0,175405
SNHG1	-0,1743161
MIR22HG	-0,174147
...	...



Conclusions: Application of Principal Component Analysis to RNA-seq data revealed particular ncRNA encoding genes most differently expressed in the samples of Glioblastoma and Low-grade-glioma. The ncRNAs encoded by COPG2IT1 and SNHG14, together with OIP5-AS1, NORAD, NEAT1, and MIR4435-2HG could be the candidates for biomarkers differentiating Glioblastoma and Low-grade-glioma cases.

Aknowledgements: This study was funded from Research Council of Lithuania (LMT) Grants: S-SEN-20-7 and MIP-20-51.

