

MSc thesis project: A model for statistical inference of gene expression changes in *in vivo* RNA labeling experiments

Background

Analysis of gene expression with methods of Next Generation Sequencing (<https://en.wikipedia.org/wiki/RNA-Seq>) is a key research tool in biomedical research and depends heavily on data science. Typical transcriptome sequencing data produces negative binomial-distributed gene counts [1]. However, more optimal models are required for newly developed methods, for example to account for data (multi)dimensionality or variables that cannot be empirically obtained from a biological system.

Required student background

MSc thesis project for a person keenly interested in applying statistical inference to big data in biomedical research. Basic knowledge of biology and familiarity with transcriptome sequencing data (Next Generation Sequencing) will be an advantage, but is not required.

Research objectives

- Develop a statistical model for inferring gene expression changes from count data obtained using *in vivo* RNA labeling experiments such as SLAM-seq (thiol (**SH**)-linked **alkylation** for the **metabolic sequencing** of RNA) [2,3].
- To improve inference in *in vivo* experiments, such as those performed using Tagger transgenic mouse line [4], different labeling kinetics models will be considered to account for changing label concentration (concentration of any drug in a biological system will change with time, except in cases when a steady state is reached). Other relevant literature, describing the current statistical frameworks used in similar kind of experiments, includes [5–7].

Short data description

- **Data type:** RNA (transcriptome) sequencing data and nucleotide (<https://en.wikipedia.org/wiki/RNA>) conversion rates (count data, see: https://kasperdanielhansen.github.io/genbioconductor/html/Count_Based_RNAseq.html)
- **Tentative data volume:** Approx. total volume of currently available data is 10^8 sequencing reads x 300 bases x 30 samples. For this project a subset of data may be used to reduce computation time. Also, if need be, it is possible to access an HPC cluster

Contact person

Lech Kaczmarczyk (lecka48@liu.se), Krzysztof Bartoszek (krzysztof.bartoszek@liu.se) Project will be done in collaboration with the Department of Biomedical and Clinical Sciences of LiU.

References

- [1] Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biology* 2010;11. <https://doi.org/10.1186/gb-2010-11-10-r106>.
- [2] Herzog VA, Reichholf B, Neumann T, Rescheneder P, Bhat P, Burkard TR, et al. Thiol-linked alkylation of RNA to assess expression dynamics. *Nature Methods* 2017;14:1198–204. <https://doi.org/10.1038/nmeth.4435>.
- [3] Muhar M, Ebert A, Neumann T, Umkehrer C, Jude J, Wieshofer C, et al. SLAM-seq defines direct gene-regulatory functions of the BRD4-MYC axis. *Science* 2018;360:800–5. <https://doi.org/10.1126/science.aa02793>.
- [4] Kaczmarczyk L, Bansal V, Rajput A, Rahman R, Krzyżak W, Degen J, et al. Tagger a swiss army knife for multiomics to dissect cell typespecific mechanisms of gene expression in mice. *PLOS Biology* 2019;17:e3000374. <https://doi.org/10.1371/journal.pbio.3000374>.

- [5] Jürges C, Dölken L, Erhard F. Dissecting newly transcribed and old RNA using GRAND-SLAM. *Bioinformatics* 2018;34:i218–26. <https://doi.org/10.1093/bioinformatics/bty256>.
- [6] Uvarovskii A, Vries ISN, Dieterich C. On the optimal design of metabolic RNA labeling experiments. *PLOS Computational Biology* 2019;15:e1007252. <https://doi.org/10.1371/journal.pcbi.1007252>.
- [7] Neumann T, Herzog VA, Muhar M, Haeseler A von, Zuber J, Ameres SL, et al. Quantification of experimentally induced nucleotide conversions in high-throughput sequencing datasets. *BMC Bioinformatics* 2019;20. <https://doi.org/10.1186/s12859-019-2849-7>.