

Integrative proteo-transcriptomics workflows within the Galaxy framework to explore the correlation between the expression of RNA and proteins

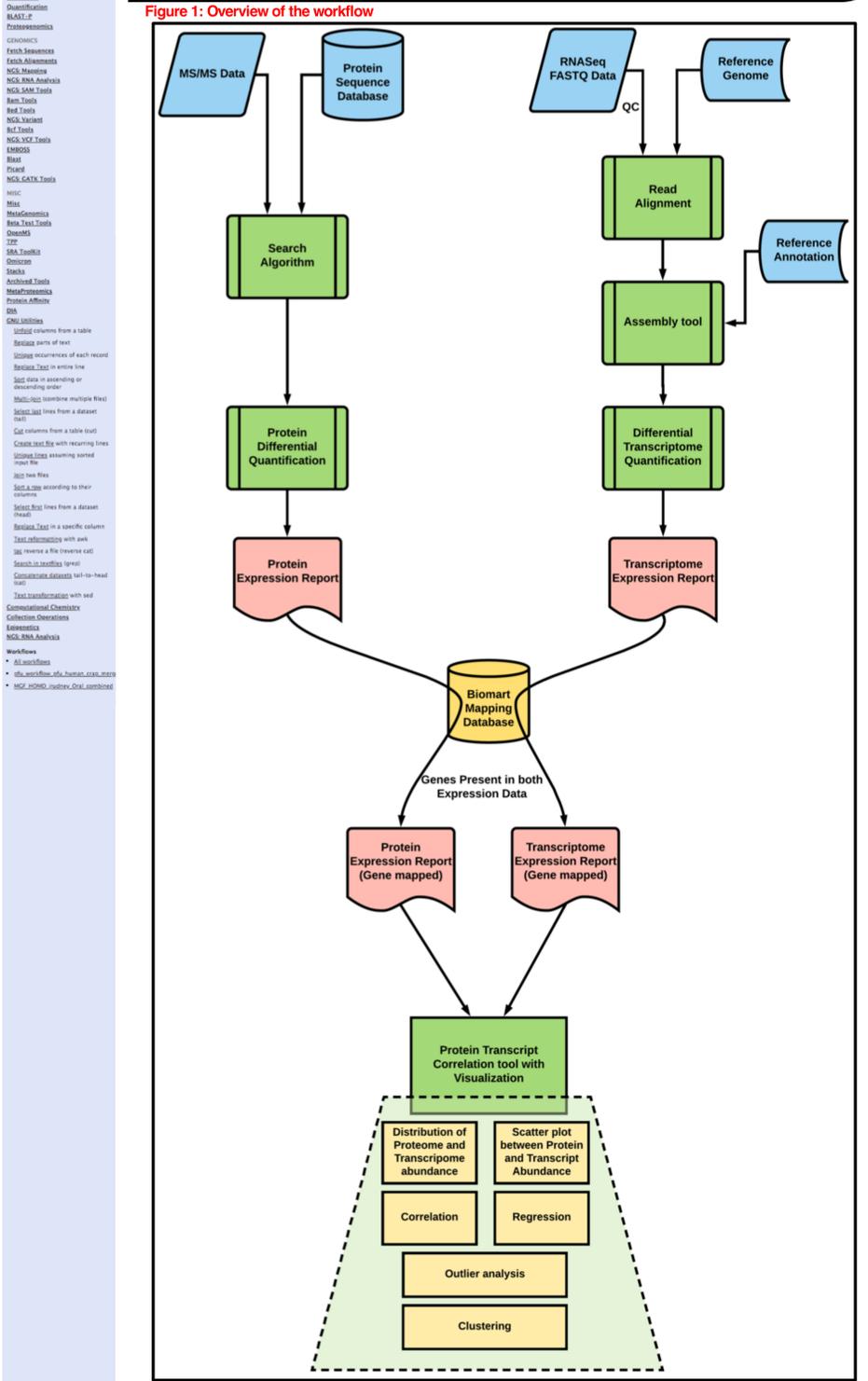


Praveen Kumar^{1,2*}, Priyabrata Panigrahi^{3*}, James Johnson⁴, Caleb Easterly², Subina Mehta², Andrew Rajczewski², Ray Sajulga², Mohammad Heydarian⁵, Krishanpal Anamika³, Timothy Griffin², Pratik Jagtap²

¹Bioinformatics and Computational Biology, University of Minnesota, Minneapolis, USA; ²Biochemistry, Molecular Biology, and Biophysics, University of Minnesota, Minneapolis, USA; ³LABS, Persistent Systems, Pingala-Aryabhata, Pune, India; ⁴Minnesota Supercomputing Institute, University of Minnesota, Minneapolis, USA; ⁵Department of Biology, Johns Hopkins University, Baltimore, USA

Introduction

- Technological advancement in the area of RNA/DNA sequencing techniques, protein/peptide identification techniques along with advances in bioinformatics have led researchers to explore multi-omics approaches.
- Several studies have used differential transcriptomic analysis to catalog gene expression in perturbed conditions, but there are other post-transcriptional regulatory mechanisms may lead to discordant mRNA and protein expression levels.
- Given that proteins are the cell's functional molecules, there has been a considerable interest in comparing protein expression with the cognate mRNA expression.
- In order to facilitate systems-biology analyses, we have developed accessible and user-friendly Galaxy tools and workflows.
- Using the mass-spectrometry-based proteomic data and RNASeq data, the workflow calculates the differential expression of proteins and transcripts respectively.
- The workflow also enables correlation study between expression of proteins and transcripts along with visualizations that will help users in interpretation of the data.



Dataset Used for Testing

- For testing, we have used a published dataset from mouse developmental B-cell samples (DOI:10.4172/jpb.1000302). RNASeq data and protein data were derived from the pre-pro-B cells and pro-B cells (two developmental stages of B-cell).
- We have used protein ITRAQ labeled quantitation ratios from the ProteinPilot results and FPKM ratios from Cuffdiff (after using HISAT2, Stringtie) to correlate the protein and mRNA expression.

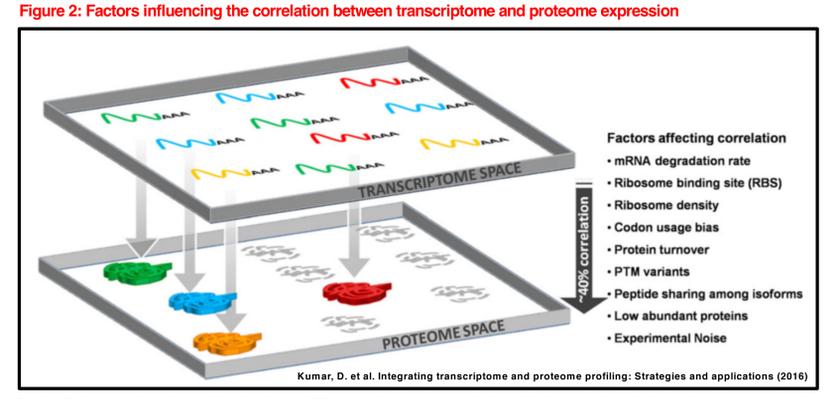


Figure 2: Even though proteins are translated from the mRNA, the correlation between transcriptome and proteome expressions are often low. Apart from some experimental limitations, there are other biological factors that can lead to observing discordant behavior between protein and transcript expression.

Discussion and Plans for Future Versions

- We generated a comprehensive workflow that generates a HTML output with the visualizations (Figure 3) using quantitative RNASeq data (FASTQ files) and proteomic data (MGF files) as inputs.
- The correlation tool explores the association of protein and transcript abundances at multiple levels. After correlating between both and exploring regression analysis, it also looks for outliers which can be vital in investigating the mechanistic cause behind the disparity in protein and mRNA expression.
- Clustering techniques can reveal the set of genes that show a similar pattern of mRNA versus protein expression. Performing a further functional enrichment analysis on genes that either show high mRNA abundance with low protein abundance or vice versa could reveal a class of genes that are being regulated differently.
- Currently, our tool enables performing correlation studies only on a single sample. We are working on including options that can enable users to perform a similar analysis on multiple samples with replicates, including multiple timepoints data.
- We are also exploring the compatibility of interactive visualization in the Galaxy which will provide a competent means to visualize and infer the disproportion of protein expression with transcript abundance.

Figure 3: Glimpse of outputs from the protein transcript correlation tool

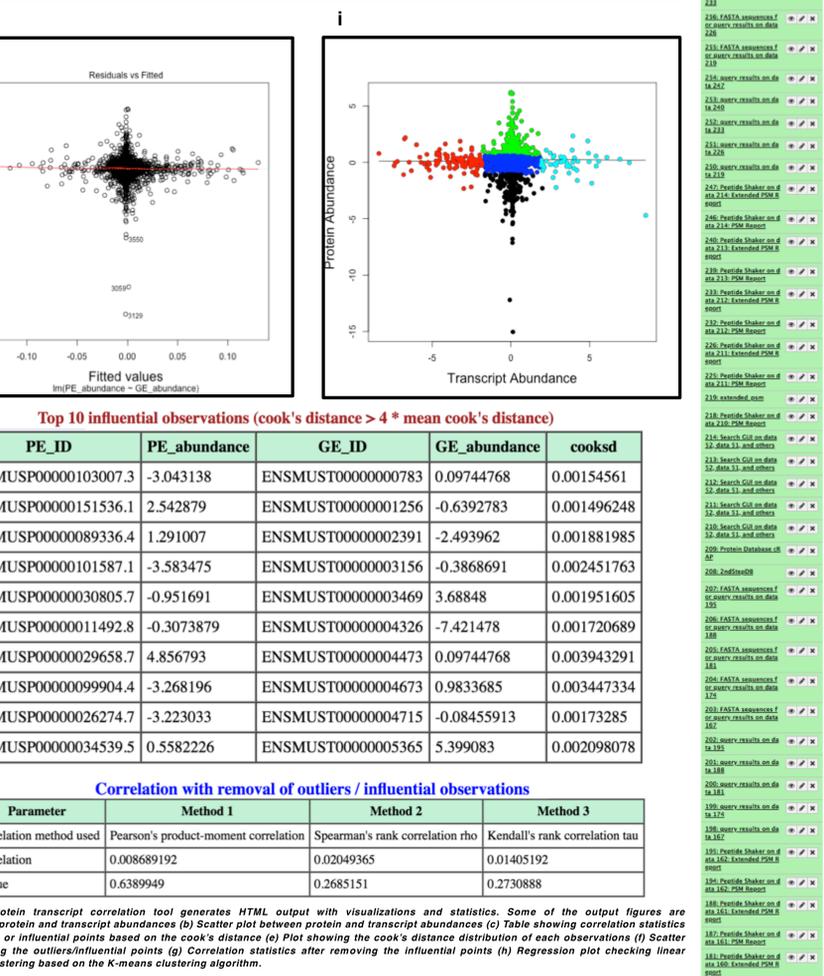
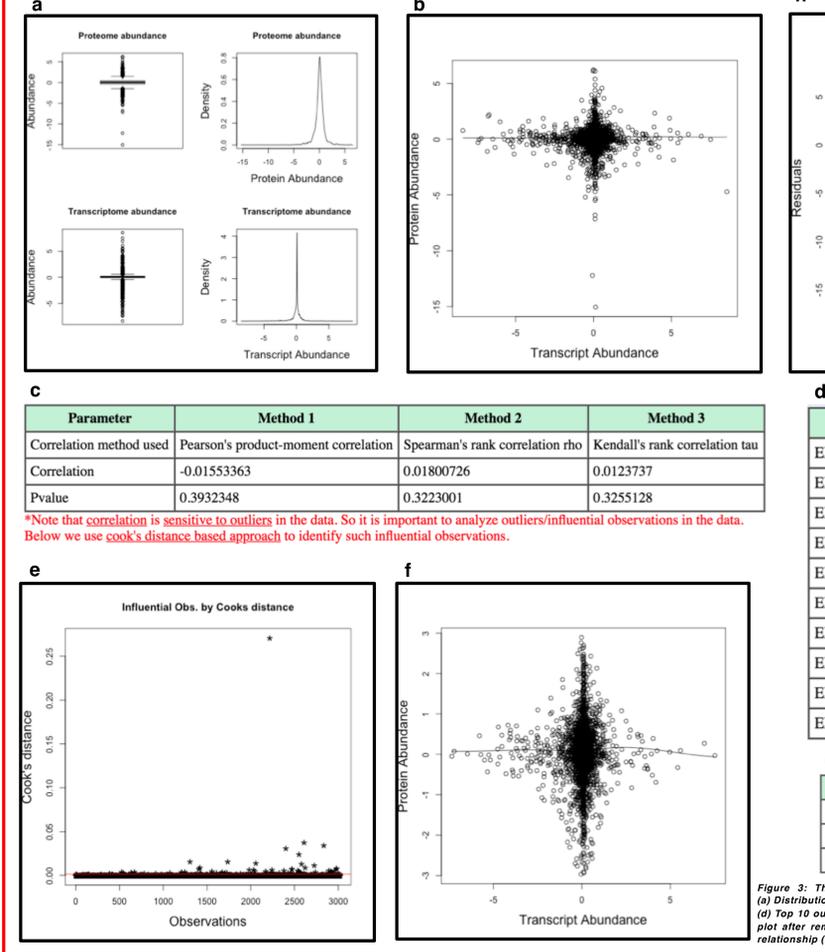


Figure 3: The protein transcript correlation tool generates HTML output with visualizations and statistics. Some of the output figures are (a) Distribution of protein and transcript abundances (b) Scatter plot between protein and transcript abundances (c) Table showing correlation statistics (d) Top 10 outliers or influential points based on the cook's distance (e) Plot showing the cook's distance distribution of each observations (f) Scatter plot after removing the outliers/influential points (g) Correlation statistics after removing the influential points (h) Regression plot checking linear relationship (i) Clustering based on the K-means clustering algorithm.

Figure 4: The workflow and the tool implementation in Galaxy

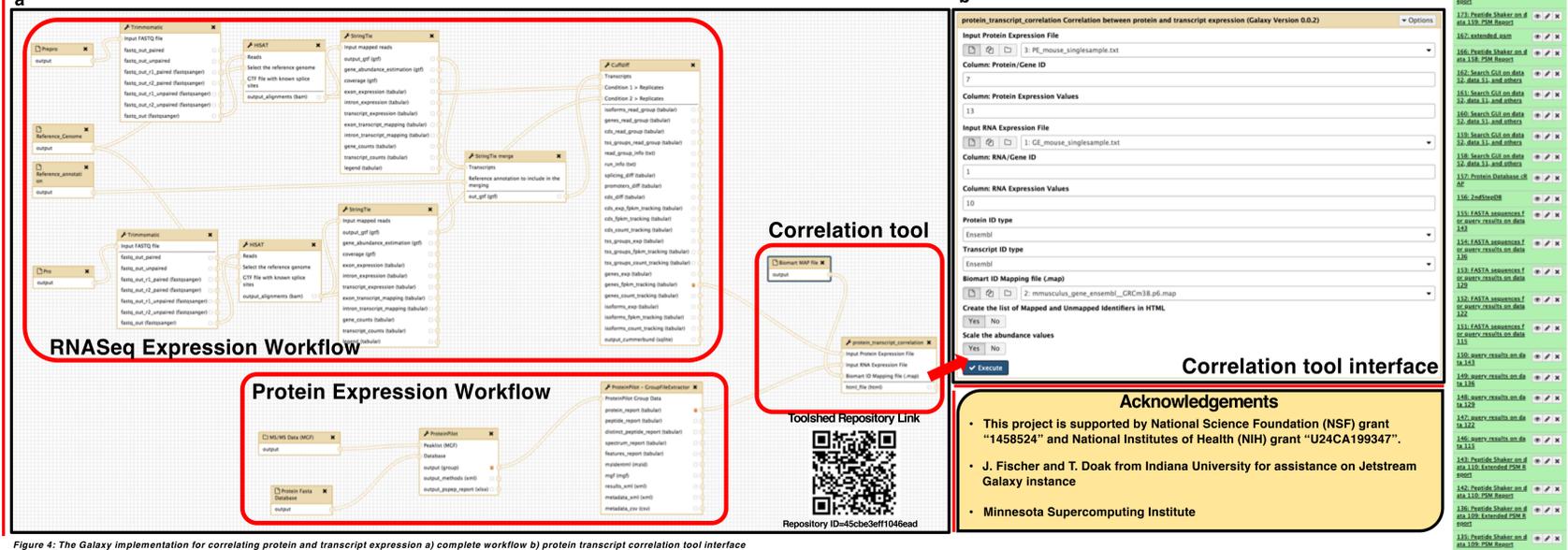


Figure 4: The Galaxy implementation for correlating protein and transcript expression a) complete workflow b) protein transcript correlation tool interface

Acknowledgements

- This project is supported by National Science Foundation (NSF) grant "1458524" and National Institutes of Health (NIH) grant "U24CA199347".
- J. Fischer and T. Doak from Indiana University for assistance on Jetstream Galaxy instance
- Minnesota Supercomputing Institute