

WOP Title

Thesis Proposal for Doctor of Philosophy(PhD) Degree

UBC bioinformatics Graduate Program

Ogan Mancarci, B.Sc

Thesis Supervisor

Dr. Paul Pavlidis

Committee Members

Dr. Clare Beasley

Dr. Robert Holt

Dr. Sara Mostafavi

Chair

Dr. Steven J.M. Jones

Examination Date

June 19, 2015

Contents

Contents	i
List of Figures	ii
1 Motivation and Introduction	1
2 Hypotheses and specific aims	3
2.1 Hypotheses	3
2.1.1 Hypothesis 1:	3
2.2 Specific aims	3
2.2.1 Aim 1: Compilation of cell type specific expression database and make it available to third parties	3
2.2.2 Aim 3: Identification and verification of marker genes	3
2.2.3 Aim 4: Enumeration of cell type proportions	3
3 Background	3
3.1 Expression profiling	3
3.1.1 Microarrays	4
3.1.2 RNA sequencing	4
3.1.3 Analysis of RNA quantification results	5
3.2 Cell type isolation	5
3.3 Cell type deconvolution	6
3.3.1 Reference based deconvolution	6
3.3.2 Reference free deconvolution	6
3.4 Cell types of the central nervous system	6
References	6
4 Figures	9

List of Figures

1	Workflow of the project	9
---	-----------------------------------	---

1 Motivation and Introduction

Brain is a remarkably heterogeneous organ with a wide variety of cell types that can be very closely related (eg. neuronal subtypes) or highly differentiated from each other (eg. neurons and glial cells). Characterization of these cell types is an ongoing challenge in neuroscience. A common problem with characterization attempts stems from the fact that they often focus on a subset of cell types that a research group is interested in. Such studies often include electrophysiological measurements coupled with expression data and attempt to find defining characteristics of a given cell type¹. Major issue with such studies is their limited scope. Such that a characteristic thought to be specific to a cell type in that study might be commonplace in other cell types that were not a part of the study.

Another important problem is that even though the heterogeneity of brain is well known, most large scale expression studies that focus on neurological/neurodegenerative disorders use whole tissue samples to examine the effects of the diseases (**citations**). Even though this approach remains popular due to its relative ease and low cost, it complicates the analysis of the results by obfuscating the source of observed differences (eg. which cell type(s) is effected by the changes) and making harder to detect changes in less abundant cell types due to signal dilution².

We aim to solve the former problem by creating a highly inclusive database of cell type specific expression profiles. This database will be used to detect marker genes that that best represents a given cell type in a particular brain region that is likely to be biologically relevant. The genes that were discovered this way can shed light into characteristics of given cell types and be useful to neuroscientists who want to work on the cell type by being useful biomarkers. Through the markers we found, we partially want to solve the latter problem by using expression of marker genes as a surrogate proportion of given cell type. This will enable us to make previously unkown inferences from the whole tissue data to have a better understanding of the disease. Also we are hoping to use recently established methods² to use the estimated proportions to get more out of differential expression data.

A large amount of cell type specific data gathered by independent groups is available in GEO. Many of these samples are collected independently, using different methods from different mouse strains. Due inherent problems of using a dataset compiled from many different sources, significant effort has to be allocated to make sure the marker genes that were found are reliable. We will do this by analysing independent datasets to make sure marker genes behave as expected. Namely, in a whole tissue dataset, we will show that they are more co-expressed than the background and in single cell datasets they tend to occur in the same cells more frequently than a randomly selected gene set.

We also hope to make the entire dataset more accessible by creating a web app to visualize the data. The users will be able to see expression of genes in the cell types represented in the dataset.

The pipeline for the project can be found in [Figure 1](#).

2 Hypotheses and specific aims

2.1 Hypotheses

2.1.1 Hypothesis 1:

2.2 Specific aims

2.2.1 Aim 1: Compilation of cell type specific expression database and make it available to third parties

1. Gathering high quality gene expression data representing brain cell types
2. Employing quality control measures to minimize intake of flawed data
3. Making the data available in a web application

2.2.2 Aim 3: Identification and verification of marker genes

1. Detecting cell type marker genes in a region dependent based on the localization of their expression
2. Verify marker genes in independent datasets

2.2.3 Aim 4: Enumeration of cell type proportions

1. Estimate cell types
2. later

3 Background

3.1 Expression profiling

Proteins are the main functioning units in a cell. They are the ultimate products of the central dogma the way it is traditionally described. DNA transcribing RNA and RNA is translated into proteins. Quantification of specific proteins is, while possible, much more difficult than quantification of specific RNA molecules and often not scalable to the same degree **citation?**. RNA quantification on the other hand have been historically used much more frequently on a high throughput fashion particularly with the rise of microarrays and more

recently RNA sequencing. A common concern about RNA quantification is that they do not always correlate with protein levels³. Therefore care should be taken when considering their biological significance.

3.1.1 Microarrays

RNA microarray is the most common way to quantify RNA in a high-throughput manner and its development have been a transformative force in many branches of biology⁴. Microarrays are fabricated by planting single strand probes, that are specific to a location on a target genome, in high density, to a known location on a solid surface. These probes will later be hybridized to a labeled cDNA library acquired by amplification of a target transcriptome. The amount of cDNA that hybridized to a probe is later quantified by staining the label attached to the cDNA library⁵. Often, multiple probes target a closeby region on the genome to form a probeset that target a single gene. The standard output for most normalization techniques is a summarization of the probes that make up a probeset⁶.

There are a multitude of microarray chips for researchers to choose from. These chips primarily differ in the probesets they use that can cover a different number of genes and/or cover the same genes using different sequences.

Microarrays are extensively used in neurobiology as a default tool for expression analysis. A wide array of data is available both in tissue⁷⁻⁹ and isolated cell type^{1,10,11} level. Since the level of heterogeneity in the brain is not fully understood, samples aiming to feature single cell types often have a risk of contamination by unintended cell types.

3.1.2 RNA sequencing

RNA sequencing (RNA seq) is a much more recent method of RNA quantification. RNA seq is performed by sequencing of amplified cDNA that is acquired from a target transcriptome. Unlike microarrays they do not target specific locations on the genome hence can give a much comprehensive picture of the transcriptome including gene discovery and splicing variant identification. The quantification done by normalizing the number of reads found from a single transcript to the length of the transcript and to the total number of transcripts¹². In general RNA seq is more prone to technical artifacts due to the stochasticity of the sequencing process **citation**. This effect is particularly powerful for lowly expressed genes which often make up the majority of the data¹³ **add figure (S8) from the paper**. Increasing the sequencing depth, read lengths and using unique labels for molecules before any amplification are ways developed to alleviate such artifacts.

With development of microfluidics it became possible for RNA sequencing to be performed on single cells **citations**. Due to scarcity of the starting product, technical artifacts resulting from amplification and sequencing is more prominent **citation**.

RNA-seq analysis, especially single cell studies are starting to gain popularity in neuroscience [14;]. Due to it's heterogeneous structure brain is a prime target for single cell studies that allows differentiation of individual cell types with much less concern of isolating heterogeneous samples.

3.1.3 Analysis of RNA quantification results

Most common use of RNA quantification is to observe differential expression between two or more groups. This is done in order to find what genes are effected by certain conditions or expression is different between different cell or tissue types to a degree of statistical significance^{6,15,16}. Another common method of analysis is coexpression, where researchers attempt to identify genes that show similar changes in expression accross different samples **citation**. This information is often used to derive functional relationships between the genes. **citation**

A higher level analysis of RNA quantification results uses more complex methods to deconvolute cell type proportions as discussed in a **later section**

3.2 Cell type isolation

Isolation of single cell types is necesarry as a precursor to their proper characterization or to analyze specific cells in different conditions such as their response to diseases or chemicals. There are multiple ways to isolate the cell types of interest with varying degrees of precision and quality. Most commonly, such methods rely on one or more marker genes that is specific to the cell type, selectively isolating cells that express the marker. A well established method is Fluorescence Activated Cell Sorting (FACS) where one or more protein or RNA is labeled to be fluorescently active, either through genetic manipulation or using in situ hybridization or labeled antibodies respectively. Cells are then gated according to specific conditions (eg. expression of a gene, absence of a gene). Another established method of isolation is Immunopanning where antibodies layered to a plate are used to hold the cell that express a specific surface marker. Finally, a relatively recent marker based isolation is Translating Ribosome Affinity Purification. This method combines the promoter of a marker gene with the coding region of L10a ribosomal subunit fused with GFP. The tissue is degraded in mass and after fixation, ribosomes marked with GFP are isolated, ensuring only translating RNAs from the target cell type are isolated **a figure for all methods. add citations**¹⁷. Alternatively cells can be isolated either by

relying on markers as in the previous methods or known physical characteristics by manually picking them, or by Laser Capture Microdissection (LCM).

The resulting samples from each of these methods has varying purity and data quality. Trap stands out with having the worst record in purity while PAN seems to be inducing the most level of stress to the cells¹⁸.

3.3 Cell type deconvolution

Expression levels obtained from whole tissues contain signals from multiple cell types. Expression profile of complex tissues can be modeled as below

$$X_{ij} = \sum_{k=1}^K W_{ik} h_{kj} + e_{ij}$$

(adapded from Shenn-Orr et. al.¹⁹)

where X_{ij} is the expression value from a complex sample for genes $j = 1, 2 \dots p$ and samples $i = 1, 2, \dots n$ and W_{ik} is a matrix containing cell type proportions for samples $i = 1, 2, \dots n$ and cell types $k = 1, 2, \dots K$

3.3.1 Reference based deconvolution

3.3.2 Reference free deconvolution

3.4 Cell types of the central nervous system

References

- 1 Sugino K, Hempel CM, Miller MN, Hattox AM, Shapiro P, Wu C *et al.* Molecular taxonomy of major neuronal classes in the adult mouse forebrain. *Nature Neuroscience* 2006; **9**: 99–107.
- 2 Chikina M, Zaslavsky E, Sealfon SC. CellCODE: A robust latent variable approach to differential expression analysis for heterogeneous cell populations. *Bioinformatics* 2015;: btv015.
- 3 Maier T, Güell M, Serrano L. Correlation of mRNA and protein in complex biological samples. *FEBS Letters* 2009; **583**: 3966–3973.
- 4 Hoheisel JD. Microarray technology: Beyond transcript profiling and genotype analysis. *Nature Reviews*

Genetics 2006; **7**: 200–210.

5 Hegde P, Qi R, Abernathy K, Gay C, Dharap S, Gaspard R *et al.* A concise guide to cDNA microarray analysis. *BioTechniques* 2000; **29**: 548–550, 552–554, 556 *passim*.

6 Gautier L, Cope L, Bolstad BM, Irizarry RA. Affyanalysis of Affymetrix GeneChip data at the probe level. *Bioinformatics* 2004; **20**: 307–315.

7 Iwamoto K, Kakiuchi C, Bundo M, Ikeda K, Kato T. Molecular characterization of bipolar disorder by comparing gene expression profiles of postmortem brains of major mental disorders. *Molecular Psychiatry* 2004; **9**: 406–416.

8 Kang HJ, Kawasawa YI, Cheng F, Zhu Y, Xu X, Li M *et al.* Spatio-temporal transcriptome of the human brain. *Nature* 2011; **478**: 483–489.

9 Torrey EF, Webster M, Knable M, Johnston N, Yolken RH. The Stanley Foundation brain collection and Neuropathology Consortium. *Schizophrenia Research* 2000; **44**: 151–155.

10 Okaty BW, Miller MN, Sugino K, Hempel CM, Nelson SB. Transcriptional and electrophysiological maturation of neocortical fastspiking GABAergic interneurons. *The Journal of neuroscience : the official journal of the Society for Neuroscience* 2009; **29**: 7040–7052.

11 Dougherty JD, Maloney SE, Wozniak DF, Rieger MA, Sonnenblick L, Coppola G *et al.* The Disruption of Celf6, a Gene Identified by Translational Profiling of Serotonergic Neurons, Results in Autism-Related Behaviors. *The Journal of Neuroscience* 2013; **33**: 2732–2753.

12 Wang Z, Gerstein M, Snyder M. RNA-Seq: A revolutionary tool for transcriptomics. *Nature Reviews Genetics* 2009; **10**: 57–63.

13 Łabaj PP, Leparc GG, Linggi BE, Markillie LM, Wiley HS, Kreil DP. Characterization and improvement of RNA-Seq precision in quantitative transcript expression profiling. *Bioinformatics* 2011; **27**: i383–i391.

14 Zeisel A, Muñoz-Manchado AB, Codeluppi S, Lönnerberg P, Manno GL, Jureus A *et al.* Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* 2015; **347**: 1138–1142.

15 Islam S, Zeisel A, Joost S, La Manno G, Zajac P, Kasper M *et al.* Quantitative single-cell RNA-seq with unique molecular identifiers. *Nature Methods* 2014; **11**: 163–166.

16 Darmanis S, Sloan SA, Zhang Y, Enge M, Caneda C, Shuer LM *et al.* A survey of human brain transcriptome diversity at the single cell level. *Proceedings of the National Academy of Sciences* 2015; **112**:

7285–7290.

17 Heiman M, Kulicke R, Fenster RJ, Greengard P, Heintz N. Cell typespecific mRNA purification by translating ribosome affinity purification (TRAP). *Nature Protocols* 2014; **9**: 1282–1291.

18 Okaty BW, Sugino K, Nelson SB. A Quantitative Comparison of Cell-Type-Specific Microarray Gene Expression Profiling Methods in the Mouse Brain. *PLoS ONE* 2011; **6**: e16493.

19 Shen-Orr SS, Tibshirani R, Khatri P, Bodian DL, Staedtler F, Perry NM *et al.* Cell typespecific gene expression differences in complex tissues. *Nature Methods* 2010; **7**: 287–289.

4 Figures



Figure 1: Workflow of the project