

**Mathematical Methods and Models in Biosciences**

June 18-23, 2023, Pomorie, Bulgaria

<https://biomath.math.bas.bg/biomath/index.php/bmcs>

## Identification and analysis of cell-specific expressed genetic variants from scRNA-seq data

Allen Kim<sup>1</sup>, Kai Saito<sup>2</sup>, Zhe Yu<sup>1</sup>, Hovhannes Arestakesyan<sup>1</sup>,  
Evgenia Unianova<sup>1</sup>, Nathan Edwards<sup>3</sup>, Anelia Horvath<sup>1</sup>

<sup>1</sup>McCormick Genomic and Proteomic Center,  
School of Medicine and Health Sciences,  
The George Washington University, Washington, DC 20037, USA  
[horvatha@email.gwu.edu](mailto:horvatha@email.gwu.edu)

<sup>2</sup>Northeastern University, Boston, MA 02115, USA

<sup>3</sup>Department of Biochemistry and Molecular and Cellular Biology,  
Georgetown University, Washington, DC 20057, USA

Low cellular frequency variants may indicate pre- or early-somatic clonality in cancer and normal tissues, or cell-specific RNA variance. Currently, most genetic variation is analyzed from bulk sequencing datasets, where low cellular frequency variants are difficult to distinguish from artifacts.

To address challenges posed by low frequency variation events, we have developed a computational framework for identification and analysis of Single Cell-specific Expressed Single Nucleotide Variants (sceSNVs) from single cell RNA-sequencing (scRNA-seq) data. Central for the framework is our new tool SCEExecute, which enables the execution of various software designed for bulk sequencing data on barcode-stratified, extracted on-the-fly, single-cell alignments. Applying SCEExecute in conjunction with tools for analysis of bulk sequencing data, we explored, for the first-time, expressed genetic variation at cell-level across 28 publicly available cancer and normal datasets, including prostate cancer, non-small cell lung carcinoma, cholangiocarcinoma, neuroblastoma, normal fetal adrenal and normal embryo. This analysis identified over 100,000 previously unreported expressed SNVs, including somatic mutations and RNA-originating variance, such as posttranscriptional modifications and locus-specific transcriptional infidelity. Our analysis shows that over 70% of these variants cannot be identified with the current bulk-based variant callers. Furthermore, approximately 10% of these novel variants show preferential expression in particular cell clusters and pseudo-time stages. Single-cell RNA e-QTL (scReQTL) analysis revealed that the expression of such sceSNVs correlates with increased expression of their harboring gene. Moreover, differential gene expression analysis between cells expressing these sceSNVs and the neighborhood cells expressing the reference allele, showed deregulation of functional gene-networks of the

SNV-harboring gene. Asymmetrically expressed sceSNVs across multiple samples are enriched in genes participating in DNA-repair, replication, and cell cycle pathways. We exemplify our analyses with a novel missense substitution – 6:26104128\_G>T, expressed in a gene encoding one of the core histones (HIST1H4C<sup>V61F</sup>). We demonstrate that HIST1H4C<sup>V61F</sup> is correlated to high expression of HIST1H4C and deregulation of the HIST1H4C-related gene network, the observation being more pronounced in neurons, across multiple cancer samples.

Our findings suggest that there is an unappreciated repertoire of cell-level expressed nucleotide variation, possibly recurrent and common across samples, that participates in transcriptome function and dynamics in both cancer and normal cells. Their appearance and, for some, relationship to certain gene-sets and cell types, suggests novel mechanisms and function for the expressed genetic variation, including in cancer progression and cell fate.