

Multi-omics data integration in the Cloud: Kawther Abdilleh^{1,2}, Boris Aguilar^{1,3}, J. Ross Thomson⁴

Analysis of Statistically Significant Associations Between Clinical and Molecular Features in Breast Cancer Institute for Systems Biology-Cancer Genomics Cloud (ISB-CGC)^{1,} General Dynamics Information Technology², Institute for Systems Biology,³ Scientific Computing, Google Cloud⁴

Abstract

- Breast Cancers are among the most common forms of cancers impacting women with over 1 million diagnoses every year worldwide.
- They are complex cancers characterized by distinct clinical outcomes, morphological and molecular features.
- The availability of large-scale cancer data allows us to study these entities in concert to gain a more holistic picture of Breast Cancer
- This study aimed to:
- Use a cloud-based multi-omics data analysis approach to identify novel associations between clinical outcomes and molecular features in Breast Cancer
- Leverage available open-access National Cancer Institute (NCI) cancer data in ISB-CGC BigQuery tables to identify genes and proteins that are significantly associated with Breast cancer clinical features
- Demonstrate how to compute and implement complex statistical methods in <u>Google BigQuery</u> directly on multi-omics cancer data on the Cloud

Methods

- Employed a novel multi-omics cloud-based approach to analyze statistical associations between available clinical, genomic and proteomic cancer data on the Google Cloud Platform through the <u>ISB-CGC</u>, one of the National Cancer Institute's (NCI) Cloud Resources.
- The Cancer Genome Atlas (TCGA) clinical, genomic, proteomic data stored in <u>Google BigQuery</u> tables hosted by ISB-CGC were analyzed using standard SQL.
- Google BigQuery tables, a managed service consisting of a columnar database backed by a massively parallel analytics engine. Data is stored in a highly distributed manner making it possible to split up SQL queries automatically, resulting in super-fast processing times.
- Statistical associations between TCGA categorical clinical features and protein/RNA expression were computed using the Kruskal Wallis test implemented as BigQuery user-defined functions.
- Only significant associations (p-value < 0.001) were used in the analysis.
- Data visualizations were implemented in cloud-based python notebooks that are publicly available.





Google Cloud