**Breakout Session 5:**

**Cloud-Based Machine Learning and Biomarker Visual Analytics for Salivary Proteomics**

Dr. Marcelo Freire
*Associate Professor, J. Craig Venter Institute*

*"Cloud-Based Machine Learning and Biomarker Visual Analytics for Salivary Proteomics"*

Marcelo Freire, DDS, PhD, DMedSc
Associate Professor, JCVI

Date: 01-18-23

Change view ⬌ | Log in

# Human Salivary Proteome *wiki*

Search this wiki

**Search**

Advanced Search

**Browse** ▾ **Search** ▾ **Contribute** ▾ **Analyze** ▾ **Learn** ▾ **More** ▾
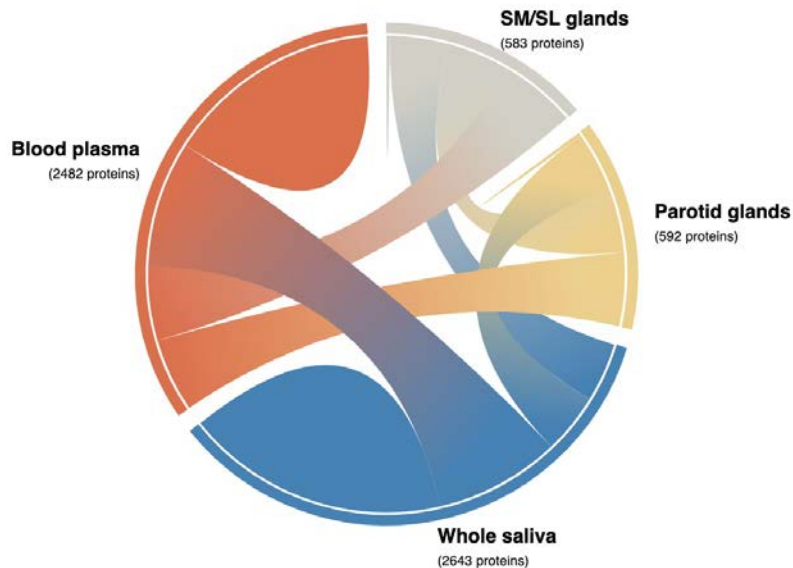
Page | Discussion | History | More ▾

**Get Started**

Click on the chord diagram below to browse the salivary protein catalog or use the search box to find specific entries. Only proteins with at least 2 distinct peptides identified in mass spectrometry experiments are included. Each arc in the diagram represents the set of proteins found in the connected tissue or sample types. Please hover over the individual arcs for more details.

**Salivary Protein Map**

**SM/SL glands**
(583 proteins)

**Blood plasma**
(2482 proteins)

**Parotid glands**
(592 proteins)

**Whole saliva**
(2643 proteins)

## How to Contribute

We truly appreciate your visit and participation. Ways that you can contribute to the wiki include:

- share your mass spec data
- annotate the proteins
- give your thoughts in discussions
- tell us how we can improve the site

**JCVI** J. CRAIG VENTER INSTITUTE™

**HSP WIKI Grant Aims:**

**Aim 1. Transfer the HSP Wiki to a cloud-based server, maintain and develop the database computing services to provide the research community with a reliable, curated, and constantly evolving (up-to-date) data platform for evaluating salivary proteomics datasets.**

**Aim 2. Implement a novel computational infrastructure for salivary proteomic datasets that includes visual analytics, faster processing, and efficient measures of success.**

**Aim 3. Expand the HSP Wiki database to include bacterial and phage proteins.**

*The long-term goal of this project is to establish the resources and develop new tools to facilitate salivary research for both scientific discovery and diagnostic applications.*

**JCVI** J. CRAIG VENTER INSTITUTE™

# Human Salivary Proteome Wiki 2.0

# ANALYSIS PAGE

Cloud Grant Aims:

**Aim 1: Develop a novel proteomics data analytics platform to support explainable and optimal machine learning identification of salivary proteomic biomarkers through visual analytics and parallel computing.**

**Aim 2 : Develop a scalable and on-demand cloud pipeline to identify protein modification, protein structure prediction using AlphaFold2, and 3D models comparison tool using MolStar.**

**Specific Aim 1: Develop a novel proteomics data analytics platform to support explainable and optimal machine learning identification of salivary proteomic biomarkers through visual analytics and parallel computing.**

The goals are:

A) Review and integrate the advantages of the existing different general analytical methods and systems

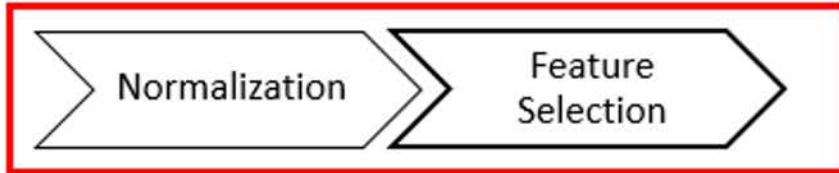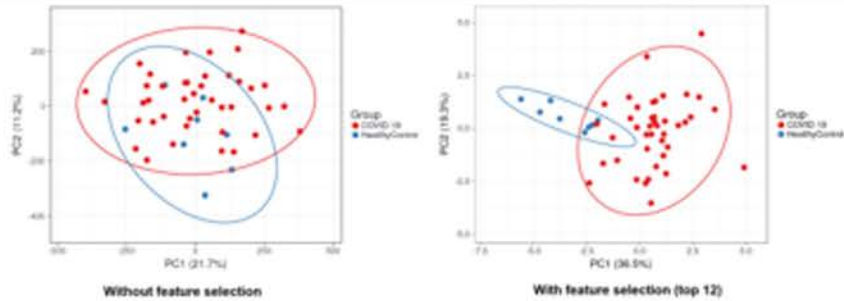B) Design and produce a specific but efficient and user-friendly platform on the Cloud for machine learning salivary proteomics data analysis
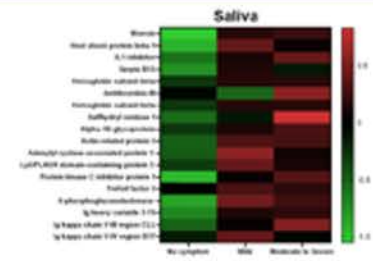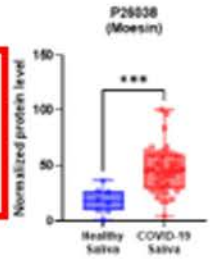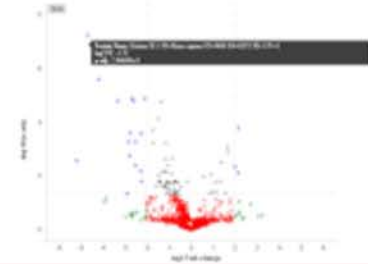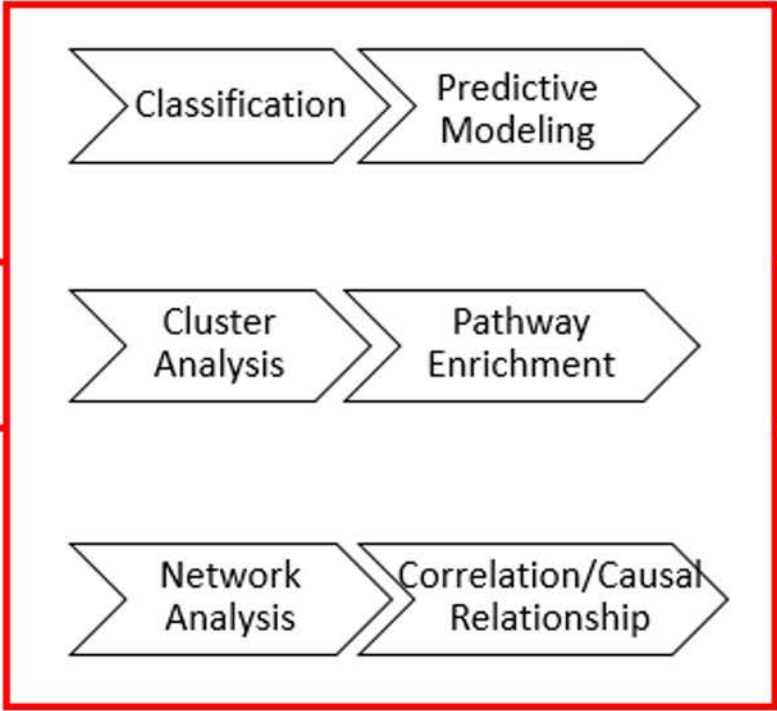
# Goals



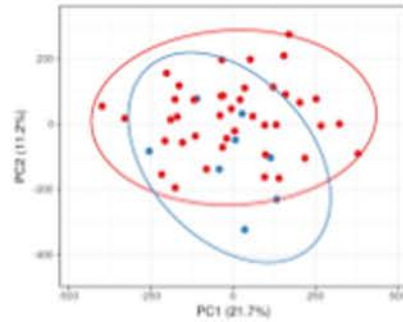A metadata-driven **visual analytics** interface and platform for Web-based exploration and analysis of salivary proteomics data
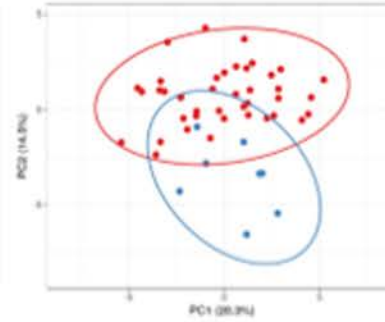
# Preliminary Data



| P-value_smaller_than0.01 | Gene Name |
|---|---|
| P06331 | IGHV4-34 |
| Q9Y3D6 | FIS1 TTC11 CGI-135 |
| P31949 | S100A11 MLN70 S100C |
| P62937 | PPIA CYPA |
| P26599 | PTBP1 PTB |
| P02533 | KRT14 |
| P05109 | S100A8 CAGA CFAG MRP8 |
| P07737 | PFN1 |
| P13646 | KRT13 |
| P29401 | TKT |
| P04083 | ANXA1 ANX1 LPC1 |
| A0A0B4J1X5 | IGHV3-74 |
| P28799 | GRN |
| P61769 | B2M CDABP0092 HDCMA22P |
| O60235 | TMPRSS11D HAT |
| P23528 | CFL1 CFL |
| P01624 | IGKV3-15 |
| P01037 | CST1 |
| P06702 | S100A9 CAGB CFAG MRP14 |
| P05164 | MPO |
| P52209 | PGD PGDH |
| P07195 | LDHB |

Rank sum test p-value < 0.01

Without feature selection   Wald-Wolfowitz test   Wilcoxon rank sum test

ISSN 0006-2979, Biochemistry (Moscow), 2022, Vol. 87, No. 3, pp. 207-214. © Pleiades Publishing, Ltd., 2022.
Published in Russian in Biokhimiya, 2022, Vol. 87, No. 3, pp. 376-385.

## High Serum Progranulin Levels in COVID-19 Patients:
## A Pilot Study

Inflammation Research (2022) 71:369–376
https://doi.org/10.1007/s00011-022-01545-7

**Inflammation Research**

ORIGINAL RESEARCH ARTICLE

**The prognostic value of S100A calcium binding protein family members in predicting severe forms of COVID-19**

JCVI J. CRAIG VENTER INSTITUTE™

# Infrastructure

**Specific Aim 2: Develop a scalable and on-demand cloud pipeline to identify protein modification, protein structure prediction using AlphaFold2, and 3D models comparison tool using MolStar.**

# Tools Needed

**Dimensionality Reduction**
- Principal component analysis (PCA)

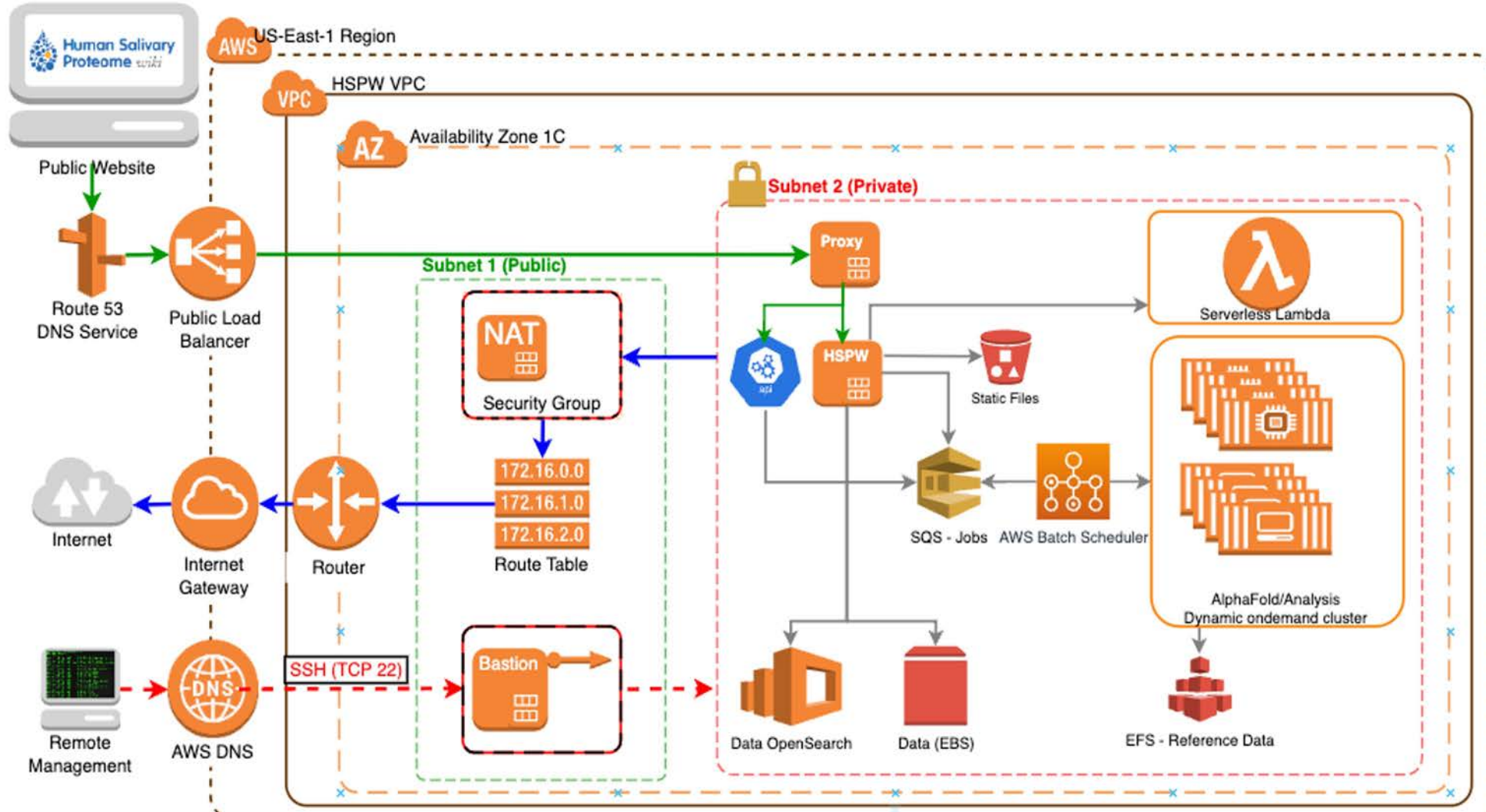**Protein Clustering**
- CD-HIT
- MMSEQ2
- kClust

**Differential Abundance**
- LIMMA
- Fold Change analysis Based on statistical Tests.

**Classification**
- Random Forest (RF)

**Heatmap**
- pheatmap (K-means clustering)

**Statistical Tests**
- t-test (or ANOVA)
- Wilcoxon Rank Sum Test
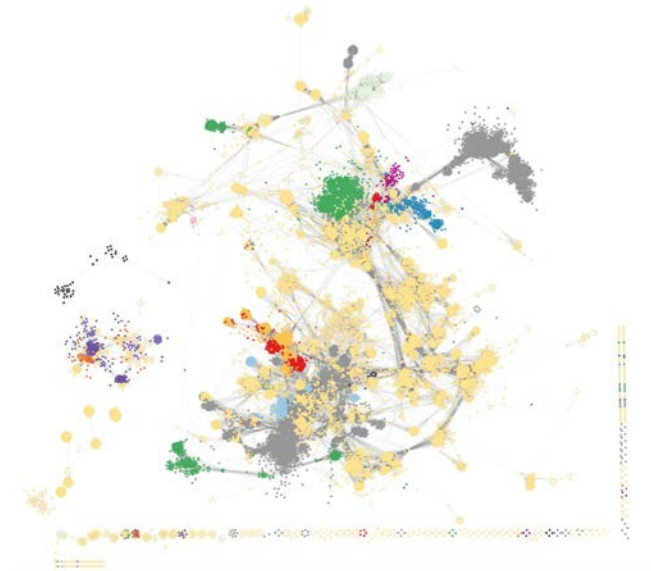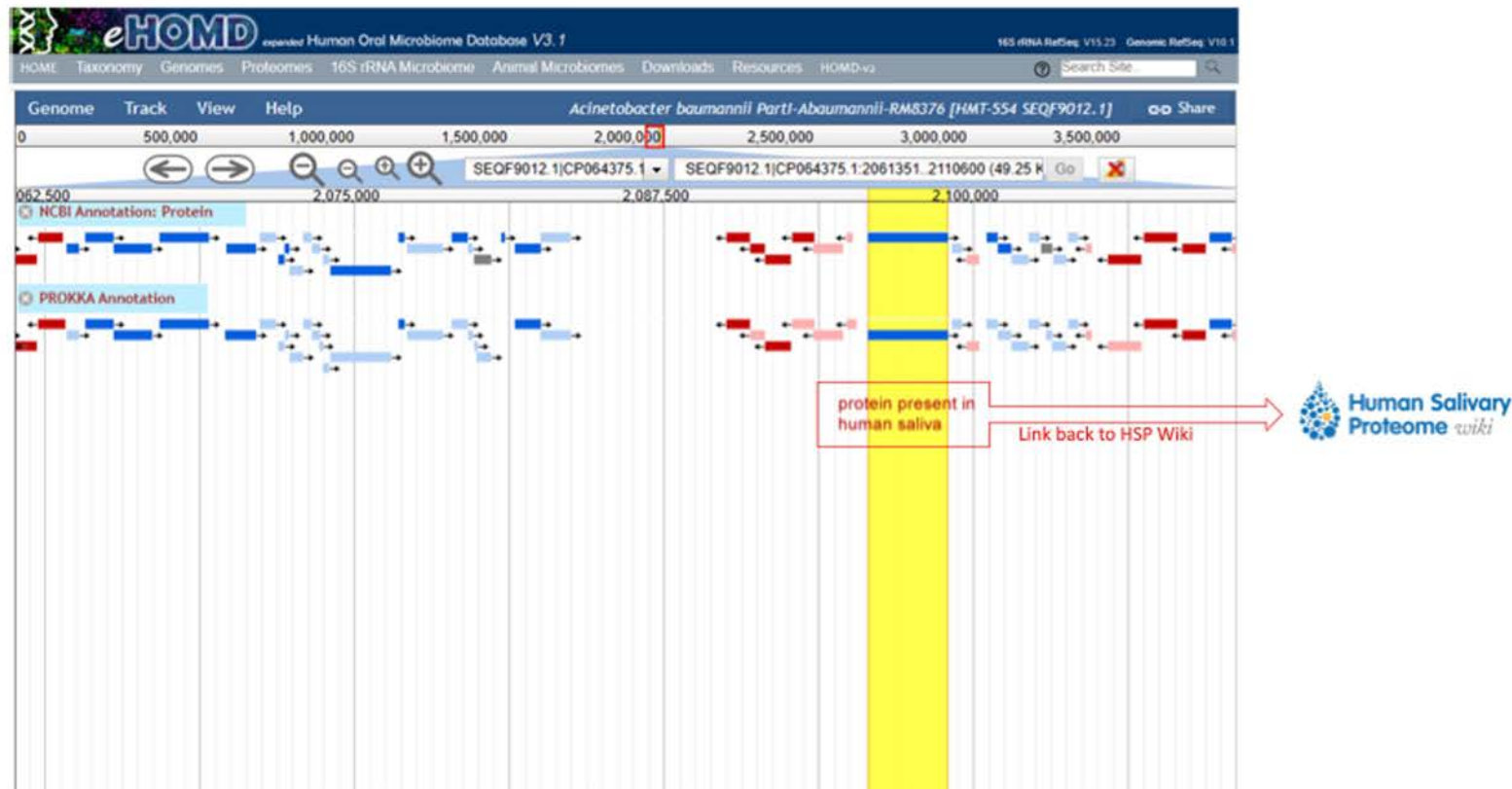- Wilcoxon Signed Rank Test

**Pathway Enrichment**
- TopGO
- KEGG Pathway (KEGGREST)
- GOANA and KEGGA (LIMMA)

**Network Analysis**
- Cytoscape (RCyc3)
- Ggraph, igraph
- STRING (database)

# Future Insights: Host-Microbial Proteins



Genes encoding the protein can be displayed in HOMD Genome Browser, with link back to WiKi

protein present in human saliva

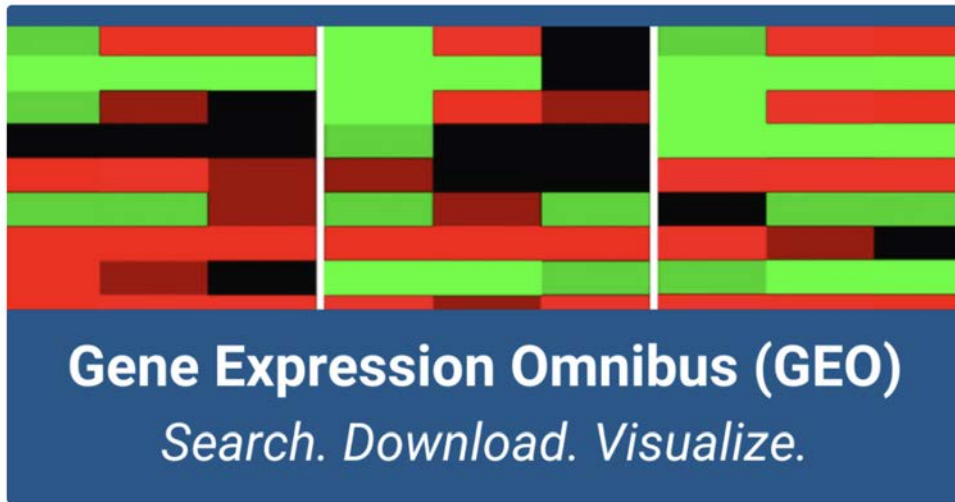Link back to HSP Wiki

Human Salivary Proteome *wiki*

**Overview of 16,745 predicted prophage-like entities in HOMD bacterial genomes (HOMD v0.1 draft).** Nodes represent predicted prophage genomes; selected phages are colored by the host in which they reside (e.g. *Streptococcus*: green; *Tannerella*: orange; *Prevotella*: purple; *Aggregatibacter*: red; *Fusobacterium*: blue).
From: Dr. Kauffmm (U at Buffalo)

From: Drs. Jessica Mark Welch and George Chen (Forsyth Institute)

JCVI J. CRAIG VENTER INSTITUTE™

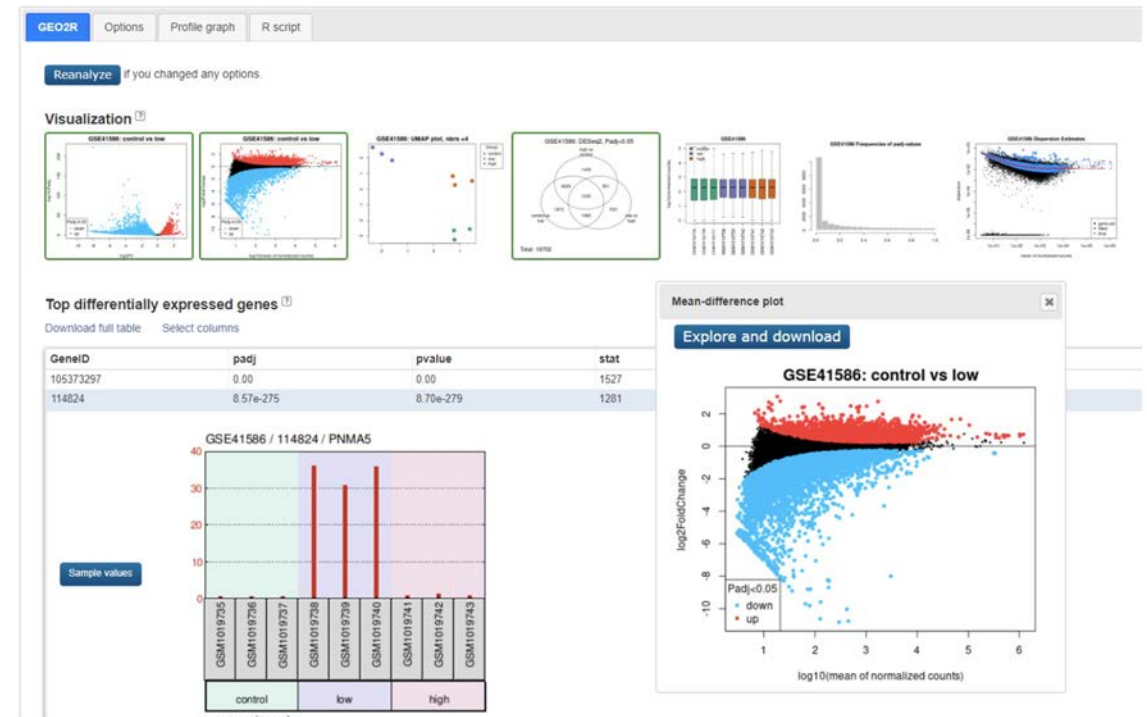# Open Source Visual Analytics (NCBI)



Figure 1: Screenshot of GEO2R differential gene expression analysis results, including Volcano, Mean difference, Mean variance, UMAP, Venn, Boxplot, and Histogram plots.

Source: NCBI

JCVI J. CRAIG VENTER INSTITUTE™

# Acknowledgements
# - Team: HOMD+HSP Wiki

**JCVI**
- Marcelo Freire
- Indresh Singh
- Harinder Singh
- Wan Hin (Marco) Choi
- Max Qian

**The Forsyth Institute**
- Jessica Mark Welch
- Floyd Dewhirst
- Tsute (George) Chen
- Markus Hardt

**University at Buffalo**
- Stefan Ruhl
- Kathryn Kauffman

**University of Tennessee**
- Yanhui Zhang

National Institute of Dental and Craniofacial Research

JCVI J. CRAIG VENTER INSTITUTE™

# Questions



mfreire@jcvi.org

@drmfreire

*"Enhancing human immunity to impact global disparities"*