

CBCB

Center for Bioinformatics and Computational Biology

# Using Secure Computation for Statistical Analysis of Quantitative Genomic Assay Data

**Justin Wagner** 

Ph.D. Candidate University of Maryland, College Park Advisor: Hector Corrada Bravo

## Genomic Assay Analysis Pipeline



# Microbiome Sequencing



Used with Permission from Nathan Olson

# Microbiome Sequencing (Continued)



Used with Permission from Nathan Olson

# Metagenomic Association Studies

- Presence/Absence
  - Tests if presence or absence of specific bacteria is associated phenotype
- Differential Abundance
  - Tests for difference in mean abundance of bacteria between phenotypes
- Alpha Diversity
  - Test for association with diversity between phenotypes



# Integrative Analysis

- Large epidemiological projects are vital to gain insights of microbiome between populations
- Gathered across institutional domains and countries
- Examples
  - Global Enteric Multicenter Study
  - Human Microbiome Project





http://www.medschool.umaryland.edu/CVD/Projects/Global-Enteric-Multicenter-Study-GEMS/

https://www.hmpdacc.org/

# Microbiome-Based Identification

- Fierer *et al*. 2010
  - Forensic identification of objects touched by an individual
- Franzosa et al. 2015
  - Unique identification of ~30% of individuals in Human Microbiome Project at later time points
- Human Read Contamination





Fierer *et al*. and Franzosa *et al*. techniques use count vector for individual sample

Noah Fierer, 6477–6481, doi: 10.1073/pnas.1000162107 ; Eric A. Franzosa et al. PNAS 2015;112:E2930-E2938

### Secure Multi-Party Computation

• Two parties have input **x**, **y**, and want to compute F(x,y)



# Yao's Garbled Circuits Protocol

- Parties agree on representation of function
- Input values are
  - encrypted and exchanged
  - such that, neither party learns anything about the inputs
  - except what can be inferred by the output

#### Garbled Circuits Protocol

Parties  $P_1$  and  $P_2$  agree on a function F(x,y) to compute





P1 garbles the circuit:

- Maps each input bit to a random string which is as an encryption key
- Encrypts each output gate using appropriate keys



P2 receives garbled circuit:

- P1 sends encrypted truth table entries
- P1 also sends encryption keys mapping to its input



P1 and P2 run Oblivious Transfer:

- P2 receives key corresponding to its input
- P1 does not learn P2's input





Ε.

### Implemention

- Used FlexSC library which is part of ObliVM project
- Coded analysis functions
  - Chi-square
  - Odds Ratio
  - Two-sample t-test
  - Alpha diversity
- Ran analysis on three datasets
  - Moderate to Severe Diaherrial Disease (754 Features, 992 Samples)
  - Personal Genome Project (277 Features, 168 Samples)
  - Human Microbiome Project (97 Features, 694 Samples)

# Running Time (Minutes)

- Split datasets up with mix of case and controls for each party
- Ran computation between two EC2 instances
  - r3.2xLarge instances with 2.5 GHz processors and 61 GB RAM
- Sparse matrix implementation Two-sample t-test ran in under one hour for largest dataset



# Network Usage (MB)

- Large data transfer overhead
- Transferred data within an AWS EC2 Region

Sparse From Garbler	НМР	PGP	MSD
Alpha			
Diversity	5862.87	5305.95	41568.3
Chi Square	203.64	543.32	1644.5
Odds Ratio	79.91	194.89	718.15
Differential			
Abundance	5055.84	6593.47	41020.37

Sparse From Evaluator	HMP	PGP	MSD
Alpha			
Diversity	5.87	5.22	43.78
Chi Square	15.78	17.07	145.2
Odds Ratio	15.78	17.07	145.2
Differential Abundance	18.26	19.84	167.9

#### Accuracy

Normalized Mean Squared Error: ||x-y||<sup>2</sup>/||x||<sup>2</sup> with x as the value output by R and y the value from our implementation

	PGP	НМР	MSD
Chi-square statistic	7.84e-07	7.48e-06	7.02e-08
Chi-square P-value	2.00e-07	2.14e-06	9.72e-08
Odds Ratio	1.60e-13	5.42e-13	2.44e-13
Differential abundance			
t-statistic	0.023	0.0017	0.0012
Differential abundance			
Degrees of Freedom	2.7e-4	2.5e-4	0.0028
Differential abundance			
P-value	0.0024	0.0026	0.0011

# Proposed Solutions for Other Pipeline Components



# Visualization

- Community Proportion with Samples Grouped By Phenotype
- Visualized as a Stacked Bar Plot
- Proportion could be calculated using secure computation or differential privacy



# Proposed Solutions for Other Pipeline Components



# Jaccard Index Estimate using MinHash

- Probabilistic technique
- Used in MASH distance measure and MHAP for assembly
- Would work with Private Set Intersection



Ondov et al. Genome Biology (2016) 17:132DOI 10.1186/s13059-016-0997-x

# Acknowledgements

Corrada Bravo Lab

UMD Center for Bioinformatics and Computational Biology

Joseph N. Paulson

Xiao Wang



## Thank You

• Questions?

# Backup Slides

Bacterial community distances between keyboard keys and fingertips.





©2010 by National Academy of Sciences

Metagenomic codes (overview).



Eric A. Franzosa et al. PNAS 2015;112:E2930-E2938



©2015 by National Academy of Sciences

Temporal stability of metagenomic codes.



Eric A. Franzosa et al. PNAS 2015;112:E2930-E2938



©2015 by National Academy of Sciences

• Chi-squared test



Samples



• Two sample t-test



Samples

Mean Abundance

Variance



 $s_i^2 = \frac{1}{n-1} \sum_{1}^{j} (\mathbf{M}_{ij} - \overline{x}_i)^2$ 

t-statistic



• Two sample t-test



Samples

Mean Abundance

$$\overline{x}_{OTU_{i,k}} = \frac{1}{n_{k_1} + n_{k_2}} \left( \sum_{j}^{n_{k_1}} M_{ij} + \sum_{j}^{n_{k_2}} M_{ij} \right)$$

Variance

$$s_{OTU_{i,k}}^2 = \frac{\sum_{j=1}^{n_{k_1}} M_{ij}^2 + \sum_{j=1}^{n_{k_2}} M_{ij}^2 + \frac{\left(\sum_{j=1}^{n_{k_1}} M_{ij} + \sum_{j=1}^{n_{k_2}} M_{ij}\right)^2}{n_{k_1} + n_{k_2}}}{n_{k_1} + n_{k_2}}$$

• Alpha Diversity



Samples

Simpson's index: 
$$D = \frac{\sum n(n-1)}{N(N-1)}$$

*n* is number if feature counts for feature i *N* is total number of counts observed in a sample



Sort by second row

Third row is used as a Counter for Chi-Square Test and Odds Ratio

Add a fourth row to compute sum of squares for t-test