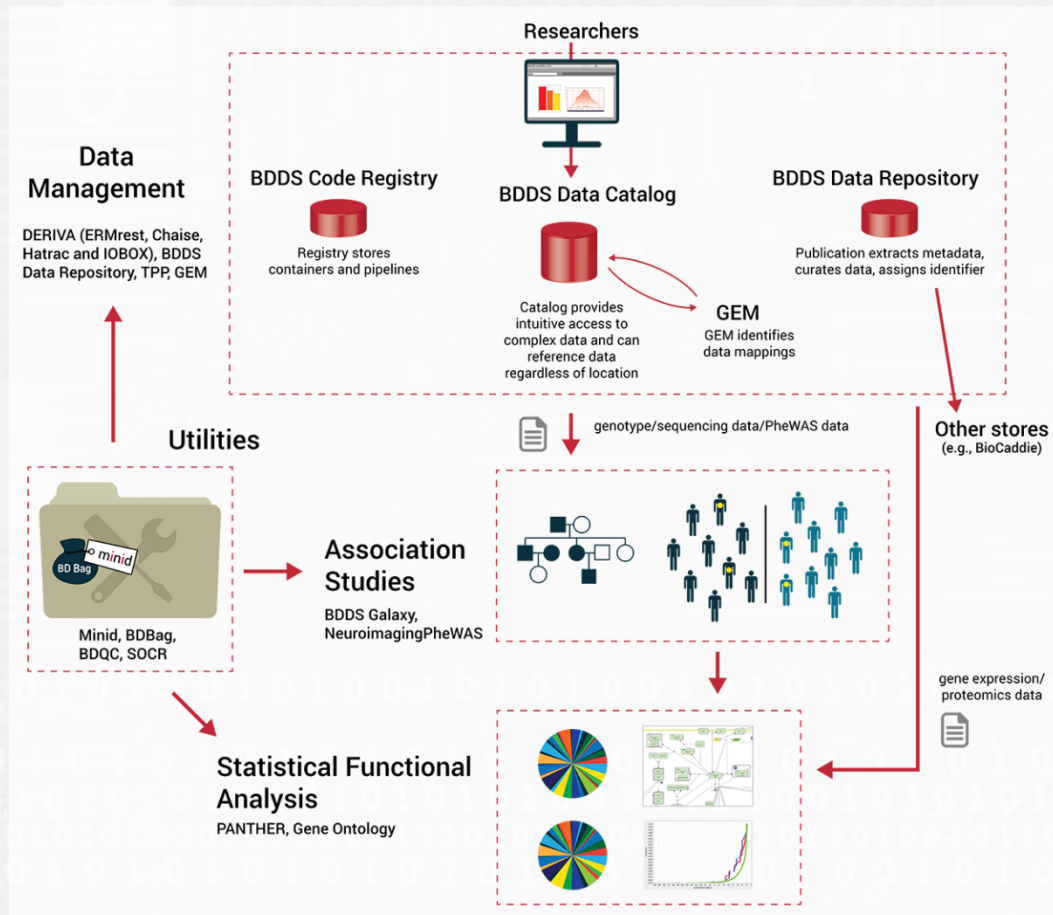# Big Data for Discovery Science (BDDS)

"What can we do today that we couldn't do before."

Arthur Toga, PI

Sept. 30, 2016

Santa Rosa, CA

# BDDS Platform – Integrated Tools for Discovery
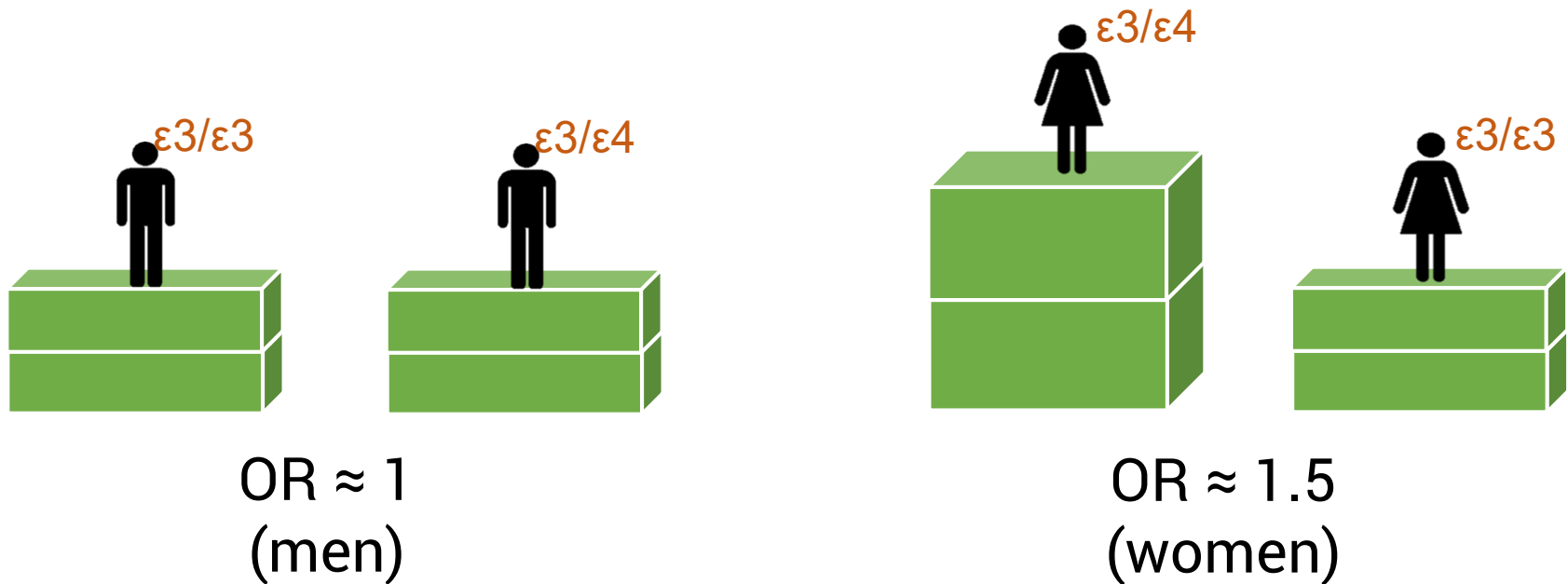
# BDDS Case Studies

- Data Aggregation
  - Sex/AD
  - Statins
- TReNA
  - Ben Heavner, Ravi Madduri
- PheWAS
  - Carl Kesselman, Lu Zhao
- BDDS Demos this afternoon
  - Dry Creek Valley I

# Farrer 1997 Meta-Analysis

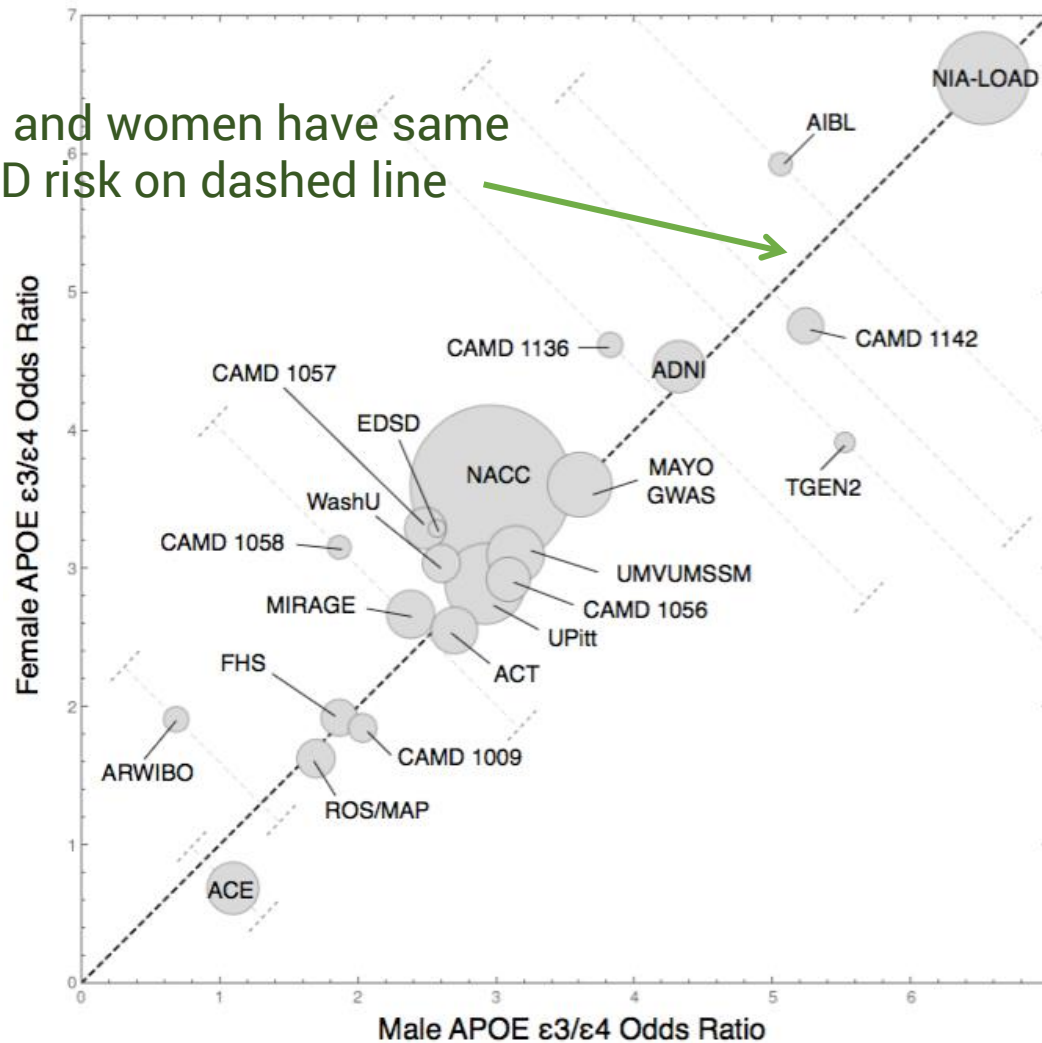Effects of Age, Sex, and Ethnicity on the Association Between Apolipoprotein E Genotype and Alzheimer Disease

A Meta-analysis

Lindsay A. Farrer, PhD; L. Adrienne Cupples, PhD; Jonathan L. Haines, PhD; Bradley Hyman, MD, PhD; Walter A. Kukull, PhD; Richard Mayeux, MD; Richard H. Myers, PhD; Margaret A. Pericak-Vance, PhD; Neil Risch, PhD; Cornelia M. van Duijn, PhD; for the APOE and Alzheimer Disease Meta Analysis Consortium

ε3/ε3          ε3/ε4                    ε3/ε4          ε3/ε3

OR ≈ 1          OR ≈ 1.5
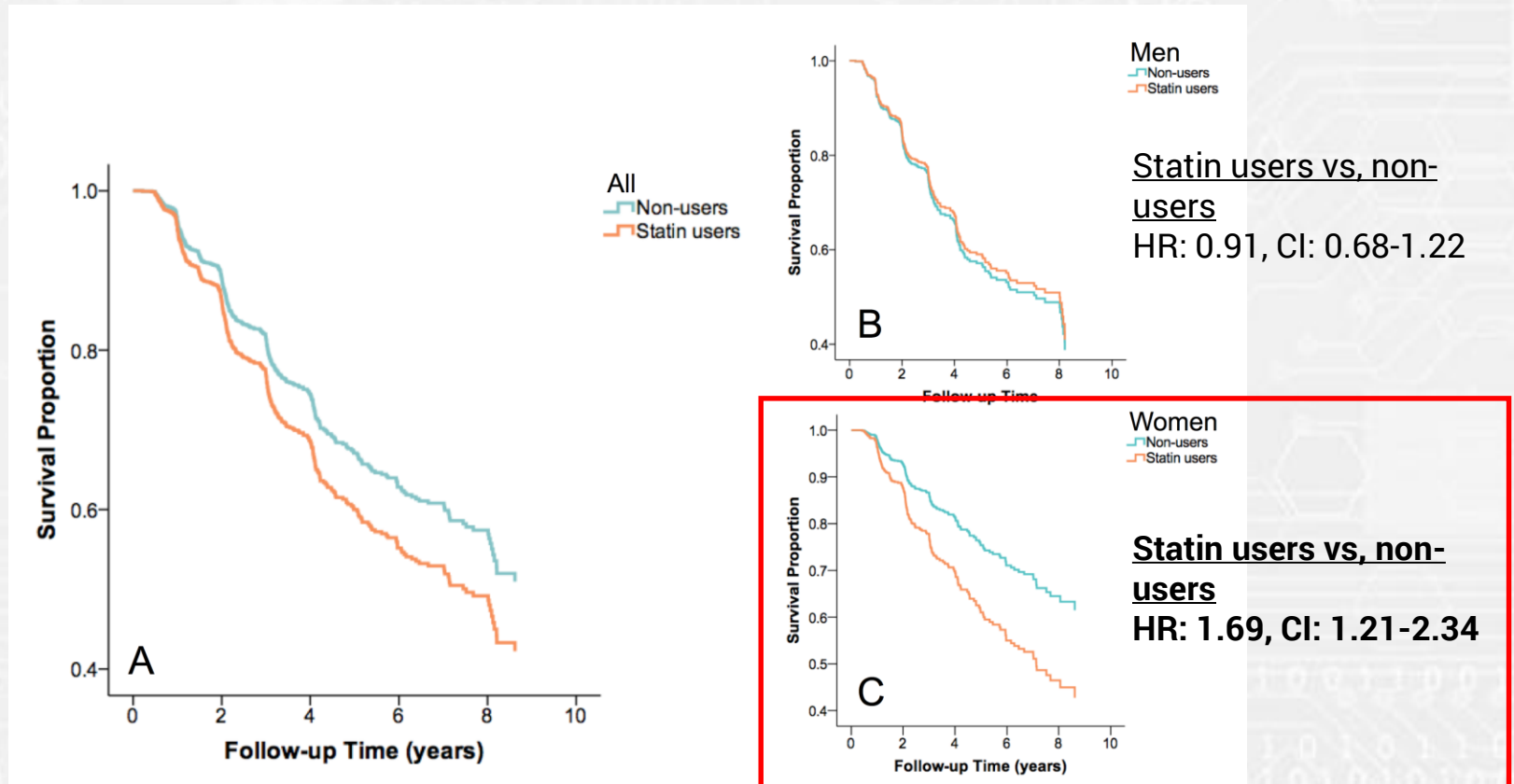(men)          (women)

# Data Set Comparison



Men and women have same AD risk on dashed line

# Background and Objective

- Statins are widely prescribed to treat high cholesterol in men and women

- However, lipophilic statins enter the brain and may impact brain cholesterol that plays a critical role in brain functioning, like estrogen production

- **Objective: to understand the relationship between lipophilic statin use and gender on the brain, cognition, and Alzheimer's disease**

# RESULTS – CLINICAL CONVERSION



All
Non-users
Statin users

Men
Non-users
Statin users

Statin users vs, non-users
HR: 0.91, CI: 0.68-1.22

Women
Non-users
Statin users

**Statin users vs, non-users**
**HR: 1.69, CI: 1.21-2.34**

# Biological Motivation: Two Puzzling Observations

### Genomic Analysis



### Gene Expression Analysis



Image credits:
http://software.broadinstitute.org/software/igv/MutationData
Nature Methods 13, 9–10 (2016) doi:10.1038/nmeth.3692

# Analysis Plan: TReNA



Vernot et al.
doi: 10.1101/gr.134890.111

- Uniform processing of next generation sequencing data
    - Align to reference genome
    - Identify DNase hypersensitve regions
    - Apply multiple footprinting algorithms to locate putative transcription factor binding sites (TFBSs)
- Evaluate confidence in putative TFBSs
- Use TFBSs as features for machine learning approaches applied to disease-specific research

# Primary Data Source: ENCODE



www.encodeproject.org

# BIG Data CHallenges: TReNA

- Gathering raw data from public repository

- Identifying raw data objects

- Transferring data objects

- Scalable data analysis

  - Integrating data from disparate sources

- Providing results for downstream analysis

# BDDS Solutions:
## Enabling Trena — bdbag



https://github.com/ini-bdds/bdbag

# BDDS SOLUTIONS:
## ENABLING TRENA — ENCODE TO BDBAG



http://encode.bdbag.org/

# BDDS SOLUTIONS:
## ENABLING TRENA – MINIDS



https://github.com/ini-bdds/minid

# BDDS Solutions:
# Enabling TReNA – BDDS Galaxy
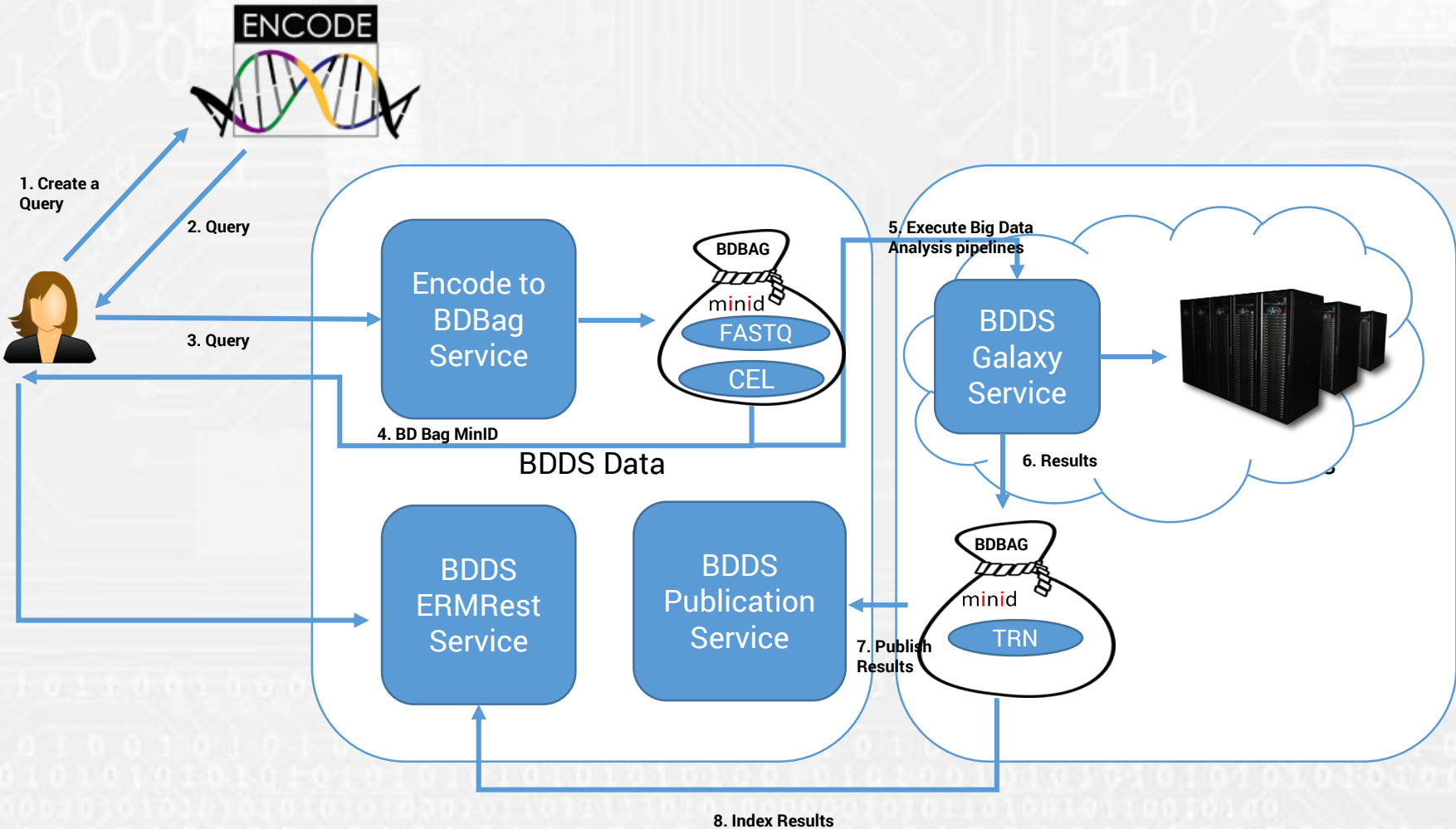


https://bdds.globusgenomics.org/
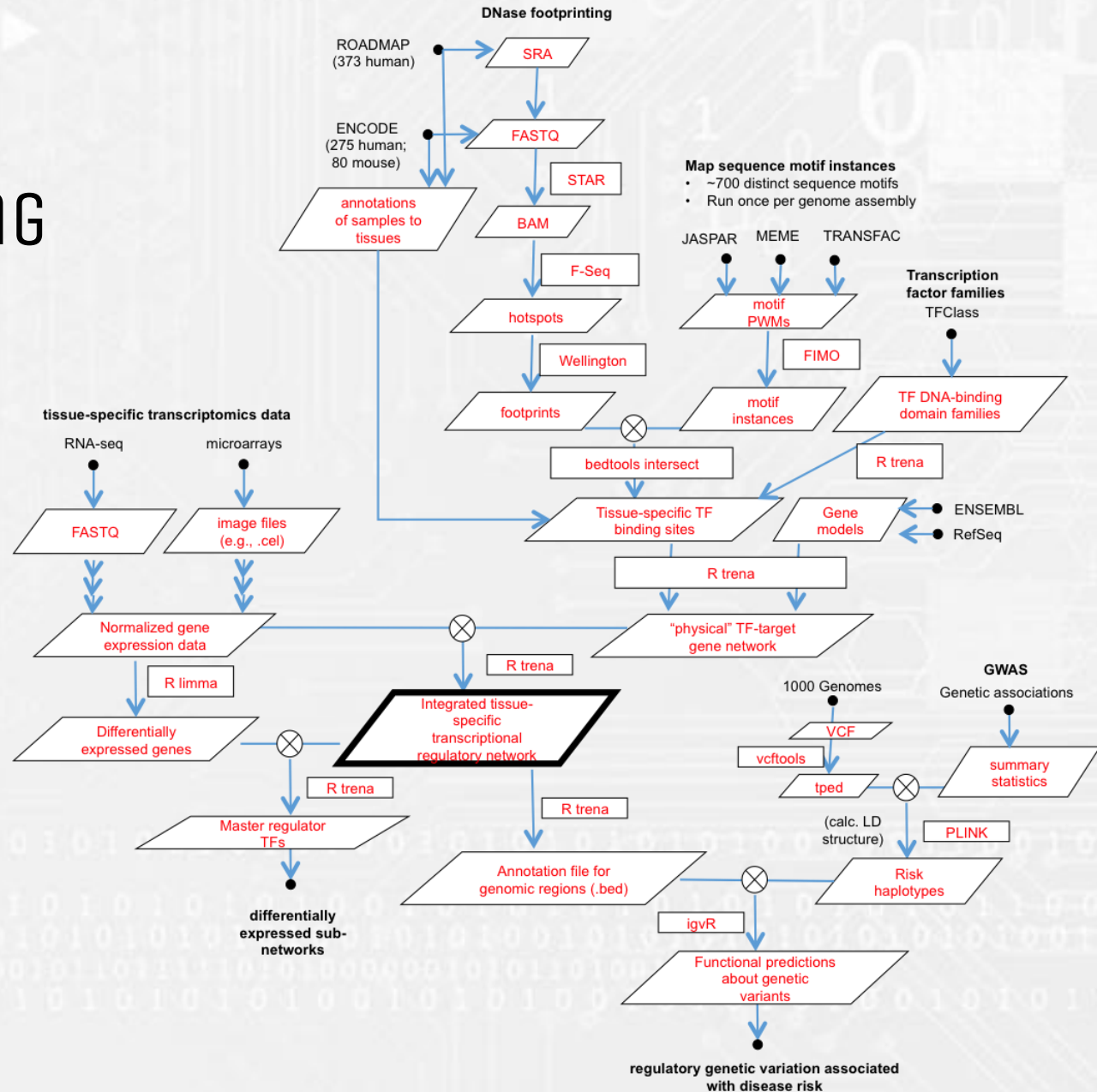
# BDDS SOLUTIONS:
## ENABLING TRENA

- BDBags and BDBag Tooling

- ENCODE to BDBag Web Service

- minids and minid Tooling

- Cloud based Data Transfers

- BDDS Galaxy Analytics Platform

# Implementing the DNase Hypersensitivity Workflow on the BDDS Galaxy Platform

# Plan Overview

# DNase Footprinting Data Processing Workflow

# BDDS Galaxy Implementation

TRN Model is working given tissue-specific TFBS counts from TFBS counts in promoters and corresponding gene expression data.

# PUBLISHING RESULTS

**Identifier:**
ark:/99999/fk4dj5pv64

**Created:**
2016-08-09 02:52:57.562893

**Creator:**
Alex Rodriguez (None)

**Checksum:**
TEST-cb5a81b06e7d138d7357fac5cce80e65330b0ff42f727ca87deae09bed4cb742

**Status:**
ACTIVE

**Locations:**
galaxy#bdds/scratch/madduri/bdds_trena_lymphoblast_bag.zip

**Titles:**
Lymphoblast DNase Footprinting Results (BDDS)

# What can we do today that we couldn't do before?

- Generate BDBags containing researcher-defined subsets of ENCODE data

- Uniquely identify ENCODE data sets using minids

- Copy data directly from ENCODE to BDDS Galaxy platform by specifying minid

- Run complex DNase footprinting analysis workflow on Galaxy platform with full provenance

- Uniquely identify workflow instantiation with minid

- Return analysis results in BDBag with minid

# BDDS Tools Used in the TReNA Approach

- Minids and Minid Tooling

- BDBags and BDBag Tooling

- ENCODE to BDBag Web Service

- Cloud based Data Transfers between ENCODE and BDDS Galaxy

- BDDS Galaxy Cloud-based Analytics Platform

# A Platform for Phenome Wide Association Studies (PheWAS)

# Neuroimaging PheWAS

- ## What is PheWAS?
  - One SNP -> a wide variety of neuroimaging phenotypes (inverse of GWAS)

- ## Why PheWAS?
  - Unbiasedly validates GWAS/single-phenotype studies findings and explores new system-level genetic associations.

- ## Challenges
  - Complexity, heterogeneity, and volume of the data
  - Complex and sophisticated brain image processing
  - Multiple-comparison correction
  - Result visualization

# • PheWAS findings (Zhao,…, Toga, Nat Neurosci, submitted)



Shaw, Molecular Psychiatry (2009) 14, 348–355

Raznahan, Neuroimage (2011) 57, 1517-23

# Image PheWAS

1. Assemble Data Collections
2. Identify subjects with images and extract images
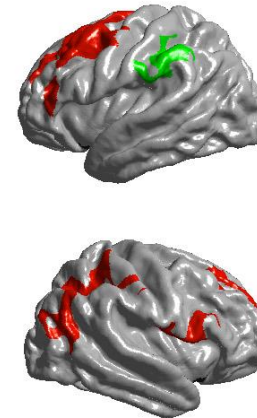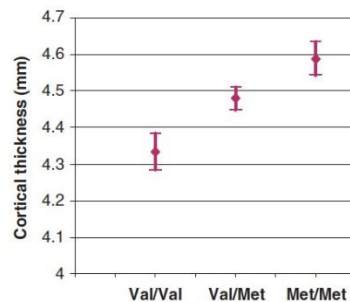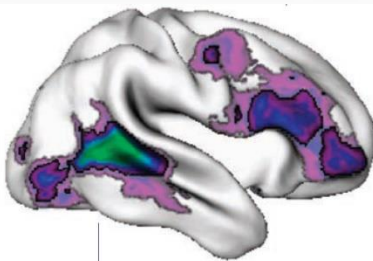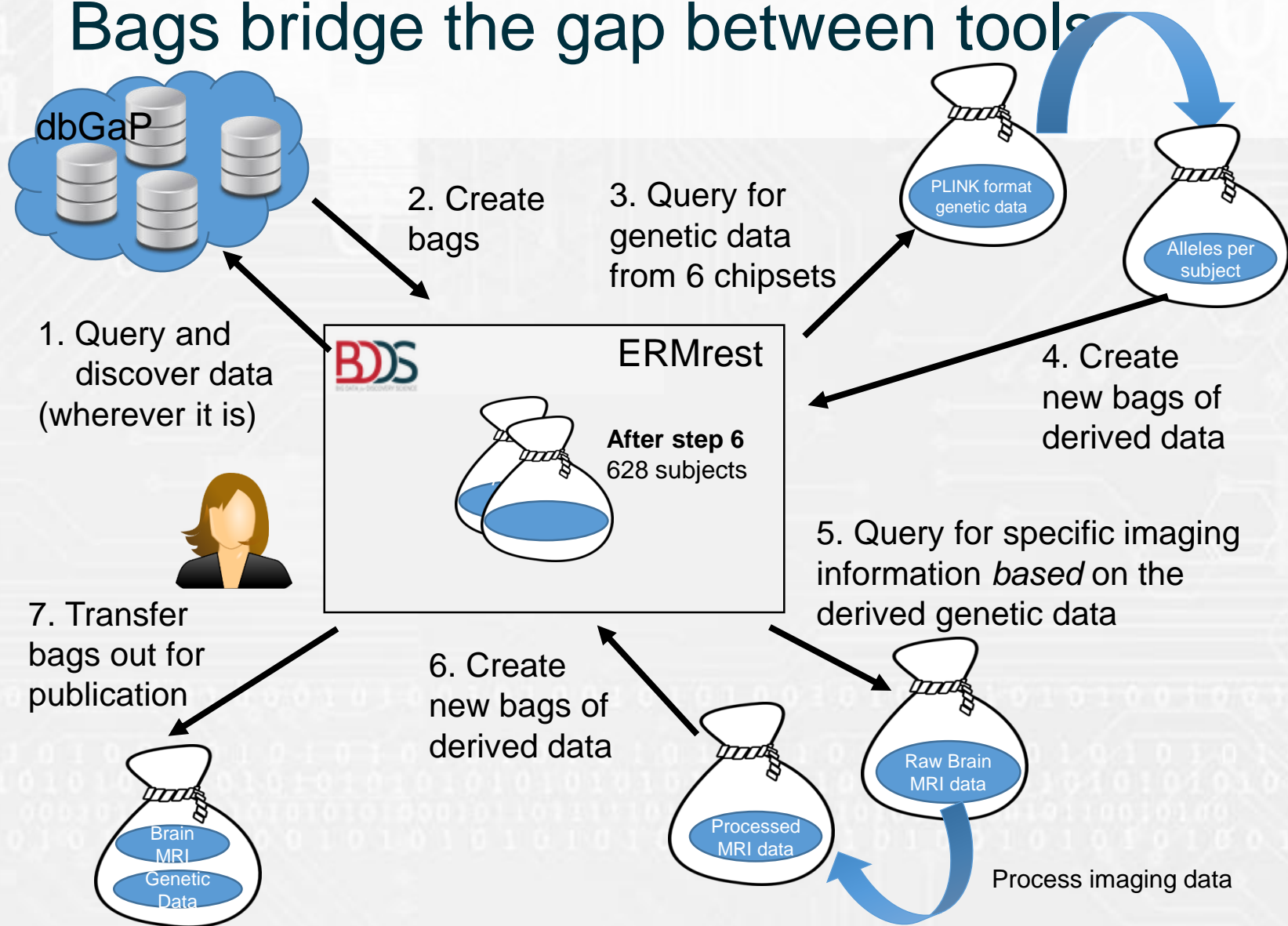3. Compute image phenotypes
   - Use Freesurfer with different atlases and computed measures
4. Associate Freesurfer results with each subject.
5. Quality control on derived data. Rerun on bad results
6. Identify subset of subjects that have variant of interest in SNP being considered
7. Collect up all phenotype data associated with identified subset
8. Do correlation analysis of phenotypes for the SNP to look for predictive correlations.

# Bags bridge the gap between tools



dbGaP

2. Create bags

3. Query for genetic data from 6 chipsets

PLINK format genetic data

Alleles per subject

1. Query and discover data (wherever it is)

ERMrest

After step 6
628 subjects

4. Create new bags of derived data

5. Query for specific imaging information *based* on the derived genetic data

7. Transfer bags out for publication

6. Create new bags of derived data

Brain MRI Genetic Data

Processed MRI data

Raw Brain MRI data

Process imaging data

# Assemble Data Collections

- …create bag with references to dbGap
- Log into dbGap and resolve references
- Assemble results in local directory
- Validate bag
- Ingest bag into catalog

# Philadelphia Neurodevelopmental Consortium

- 8719 subjects in study
  - Baseline clinical elements
- 6 different SNP array chipsets resulting in a combined set of 1,873,486 distinct SNPs (out of a possible 85 million in the human genome).
  - The total combinatorial space of the genomic data is 5,435,533,460 (SNP, subject, allele) tuples across the 8719 subjects
- 997 of the subjects have MRI imaging data

# Managing data collections

# Export Bags for Phenotype calculation

- Bag structure used to collect data sets, along with description of what should be computing

- Collect results of calculation into bag and reingest

- Parcellation process generates 381 distinct regional measurements per subject, for a total of 370,641 regional neuro-anatomical measurements

- Surface-based model generates > 2 millions local measurements per subject, for a total of > 2 billions local neuro-anatomical measurements

- Image data for the 997 subjects consists of 70930 files (including derived images) @ 666GB

# Details on one data element

# QC on derived data

# Complex data relationships…

# NeuroimagingPheWAS Toolbox

# What can we do now we couldn't do before?

- Broad survey for true system-level genetic associations across the whole population
  - All kinds of imaging genome data and processing
  - Not just for PNC, PING
  - Extensible to other phenotypes, not just FreeSurfer
- Build more complex studies the previously possible
  - e.g. TRENA + PheWAS
- Reproducible, complex, multistep big-data analysis

# BDDS Demos
## This afternoon

- TReNA
  - Ravi Madduri, Ben Heavner

- PheWAS
  - Carl Kesselman, Mike D'Arcy, Kristi Clark, Lu Zhao

- Panther
  - Huaiyu Mi, Anushya Muruganujan

- Data Publication
  - Ian Foster

  - Dry Creek Valley I room