# NCBI Data Analyses (aka Pre-computes)

## Emily L. Harris, PhD, MPH

## NHGRI

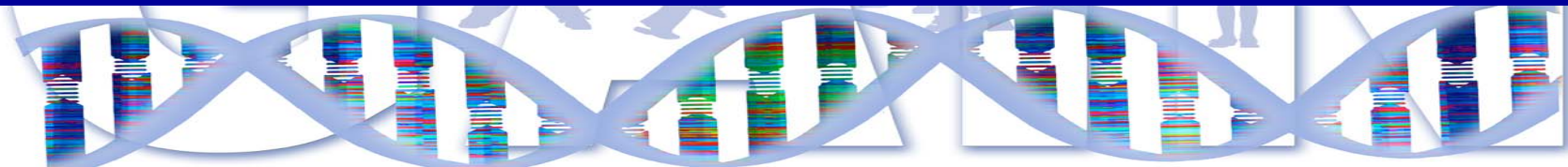GENETIC ASSOCIATION INFORMATION NETWORK

# Purpose

- Provide basic results
  - Preliminary scientific results for investigators
  - Quality check for shared data
- Promote broad data use
  - Protect against claims on genotype/allele frequency data and phenotype-genotype associations
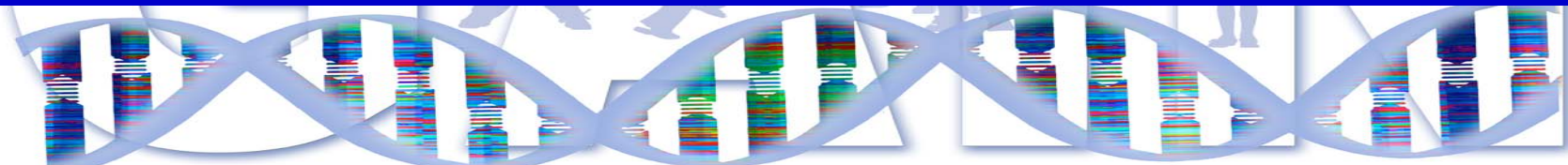
# Analysis Ideas

- Each study
  - For each SNP
  - For each target trait/condition
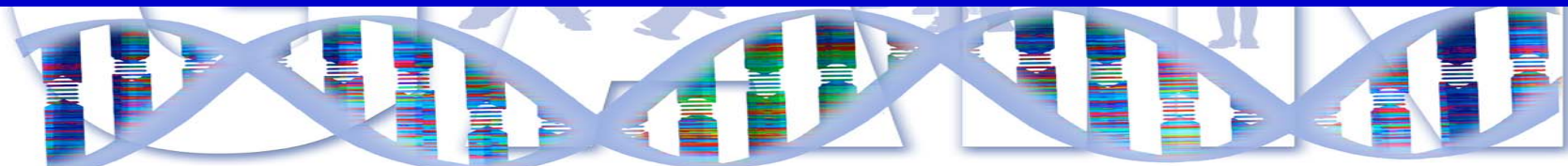  - For multi-SNP analyses
- Cross-study analyses

# Analysis Ideas for each SNP

- Genotype call accuracy (clusters)
- Allele frequencies
- Genotype frequencies
- Tests of Hardy-Weinberg equilibrium

- Stratified on ascertainment status
  - Case-control design: Cases, controls
  - Family designs: Probands, unrelated individuals

# Analysis Ideas
# for each Target Trait/Condition

- Association with each SNP genotype
  - Unadjusted measure of association (e.g., unadjusted odds ratio)
  - Statistical significance, with correction for multiple tests
    - By SNP
    - Genome-wide visual representation

# Visual Representation: Genome-wide

# Visual Representation: A More Detailed Look

# Analysis Ideas
# for Multi-SNP Analyses

- Pairwise linkage disequilibrium
  - By pairs of SNPs
  - Visual representation

- Stratified on ascertainment status
  - Case-control design: Cases, controls
  - Family designs: Probands, unrelated individuals

# Analysis Ideas
# for Multi-SNP Analyses

- Haplotyping for chromosomal regions that surround SNPs with significant associations

# Analysis Ideas

What else would be useful?

⇨   Cross-study analyses?

⇨

⇨

# Analysis Ideas

Questions: Analysis

⇨ How many/which genetic models should we test?

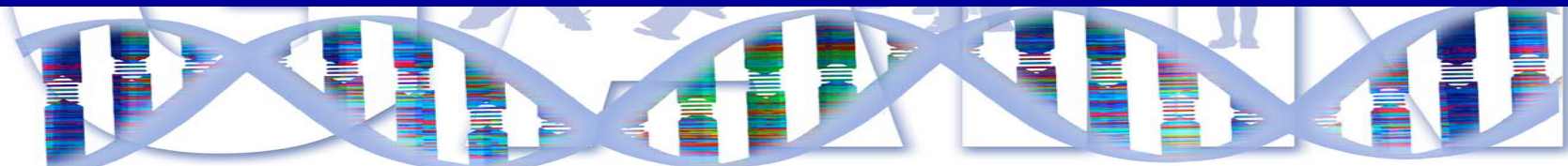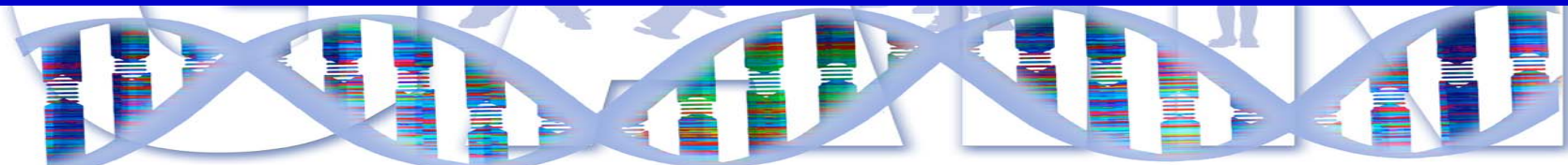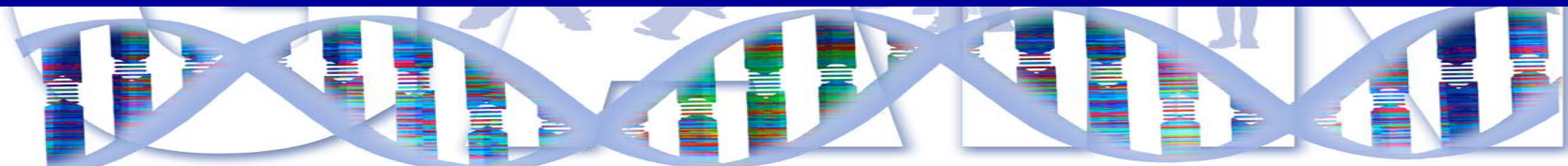⇨ How important to consider multiple testing for p-values in the pre-computed analyses, given that it is a public, unpublished resource?

⇨ What multiple testing corrections should be included for the pre-computed analyses?
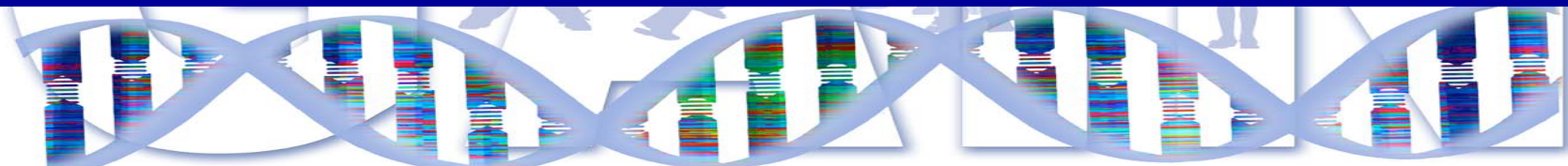
# Analysis Ideas

## Questions: Links to genomic information

For SNPs with a significant association

⇨ Should we include in the report information about location of each SNP within a known gene?

⇨ Should we annotate the genome with these significant associations?

# Availability

- Results initially available only through the controlled access process for approved users to download

- After the 9-month protected period for a specific project, results available on the GAIN public web site

# Thanks!

## NCBI

Al Graeff

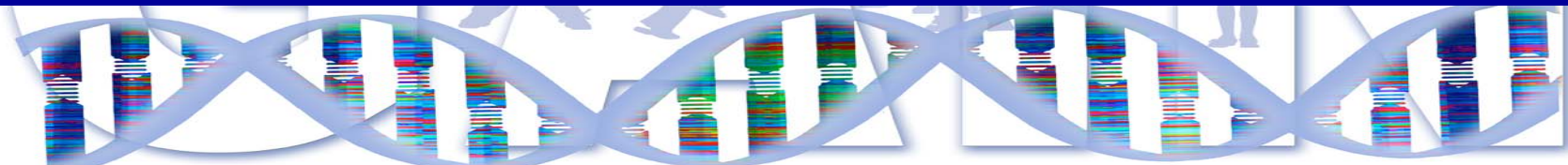Matt Mailman

Jim Ostell

Steve Sherry

## NHGRI

Lisa Brooks

Teri Manolio

Laura Lyman
   Rodriguez

# GWA data availability

- dbGaP – NCBI
- CGEMS – NCI
- NINDS Open Access Repository

# Genome-wide SNP genotyping

- Initial genome wide genotyping in:
  - 276 PD
  - 276 Stroke
  - 276 ALS
  - 276 Controls
  - 200 African Americans

  *FROM NINDS OPEN ACCESS REPOSITORY*

  *M.E. & A.Z. NIA*

- 109k exon-centric assay (phase I)
- 317k HapMap assay (phase II)
- >99.8% call rate, >450,000,000 unique genotypes, >99.9% reproducibility (over 19,000,000 replicate genotypes)

## Genome-wide genotyping in Parkinson's disease and neurologically normal controls: first stage analysis and public release of data

Hon-Chung Fung, Sonja Scholz, Mar Matarin, Javier Simón-Sánchez, Dena Hernandez, Angela Britton, J Raphael Gibbs, Carl Langefeld, Matt L Stiegert, Jennifer Schymick, Michael S Okun, Ronald J Mandel, Hubert H Fernandez, Kelly D Foote, Ramón L Rodríguez, Elizabeth Peckham, Fabienne Wavrant De Vrieze, Katrina Gwinn-Hardy, John A Hardy, Andrew Singleton

### Summary

**Background** Several genes underlying rare monogenic forms of Parkinson's disease have been identified over the past decade. Despite evidence for a role for genetics in sporadic Parkinson's disease, few common genetic variants have been unequivocally linked to this disorder. We sought to identify any common genetic variability exerting a large effect in risk for Parkinson's disease in a population cohort and to produce publicly available genome-wide genotype data that can be openly mined by interested researchers and readily augmented by genotyping of additional repository subjects.

## Genome-wide SNP assay reveals structural genomic variation, extended homozygosity and cell-line induced alterations in normal individuals

Javier Simon-Sanchez[1,2,†], Sonja Scholz[1,†], Hon-Chung Fung[3,†], Mar Matarin[1,†], Dena Hernandez[1], J. Raphael Gibbs[4], Angela Britton[1], Fabienne Wavrant de Vrieze[3], Elizabeth Peckham[5], Katrina Gwinn-Hardy[6], Anthony Crawley[6], Judith C. Keen[7], Josefina Nash[7], Digamber Borgaonkar[3], John Hardy[3] and Andrew Singleton[1,*]

[1]Molecular Genetics Unit, [2]Laboratory of Neurogenetics and [3]Computational Biology Core, National Institute on Aging, [4]Human Motor Control Section and [5]Neurogenetics Branch, National Institute of Neurological Disorders and Stroke, National Institutes of Health, Bethesda, MD, USA, [6]Unidad de Genética Molecular, Departamento de Genómica y Proteómica, Instituto de Biomedicina de Valencia-CSIC, 46010, Valencia, Spain and [7]Coriell Institute for Medical Research, Camden, NJ, USA

# Genome-wide SNP genotyping

PD and control genotype data posted publicly

downloaded by >300 unique visitors.

Stroke and ALS data available Jan 2007.

Allows entire data download or search for specific SNPs



http://ccr.coriell.org/ninds/