# Maximizing Public Data Sources for Sequencing and GWAS
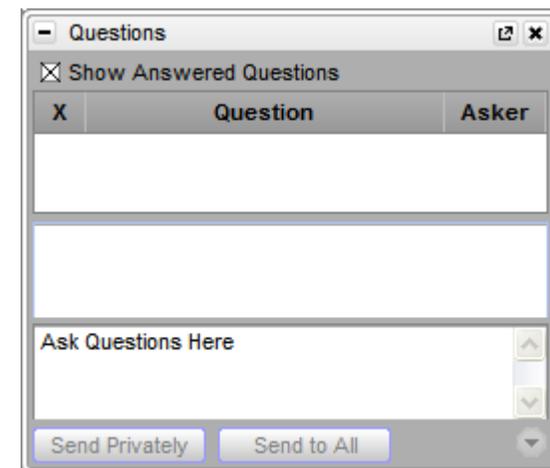
February 4, 2014

G Bryce Christensen
Director of Services

# Agenda

**1** Why Use Public Data?

**2** Where to Find Public Data

**3** Tips for Using Public Data

**4** Manipulating Public Data in SVS

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Why Use Public Data?

- **Reference samples for assessing population structure in GWAS**

- **Replicating results of your own GWAS or other research**

- **Meta-analysis or Mega-analysis**

- **Testing new analytical methods**

- **Reference data for SNP imputation**

- **Increase study size with public controls**

# Sources of Public Data

- **NCBI**
  - dbGaP
  - GEO
  - SRA

- **EGA**

- **HapMap Project**

- **1000 Genomes Project**

- **Hardware vendors**

- **Software vendors**

- **All over the internet…**

# dbGaP



**dbGaP**

The database of Genotypes and Phenotypes (dbGaP) was developed to archive and distribute the results of studies that have investigated the interaction of genotype and phenotype.

- **"Database of Genotypes and Phenotypes"**

- **435 studies in database (as of January 28[th])**

- **Known primarily as a GWAS database, but NGS content is growing**

- **Freely view and download results for many studies**

- **Access to raw phenotype and genotype data requires application process**

# 435 Studies in dbGaP (January 28th)

## GWAS Platforms

- **Affymetrix**
  - SNP-6.0:          51
  - 500k:             15

- **Illumina**
  - HumanHap550:   37
  - HumanHap300:   13
  - HumanCNV370:   11
  - Human610:       35
  - Human660:       26
  - Omni1:           22
  - Omni2.5:         14
  - Human_1M:       12

- **Perlegen**
  - 600k:              4

## NGS Platforms

- **454: 22**

- **GA-II: 49**

- **HiSeq 2000: 72**

- **HiSeq 2500: 3**

# dbGaP Tools

- **GaP Browser**
  - View GWAS study results in context of other genomic annotations

- **GaP Genome Browser**
  - Karyotype views of GWAS study results

- **PheGenI**
  - "Phenotype-Genotype Integrator"
  - Search NHGRI and dbGaP study results by phenotype or by gene
  - Annotated results with links to abstracts and/or dbGaP study pages.

# Applying for dbGaP data

- **Each application is reviewed by a "DAC," or data access committee**
  - I've seen approval time range from 1 to 8+ weeks.

- **Keep proposals relatively simple**
  - Read the instructions and be sure that your application is complete before submitting
  - Contact DAC before submitting if you have special needs or concerns

- **Some datasets require IRB approval to access**
  - Waiver letter is often sufficient

- **Pay attention to data embargoes**

- **External collaborators and contractors must apply separately for access**

- **Pay attention to consent groups**
  - General research use
  - Non-commercial use
  - Disease-specific use

# Using dbGaP Data

- **Know what you are getting—read the documentation!**
  - Original study design
  - Data processing and formats

- **Be patient and thorough as you explore the data--treat it like fresh new data and don't assume that it is "clean."**

- **Phenotype data is usually stored in text files, often with a separate data dictionary.**
  - Read the documentation!

- **Carefully review phenotype data for completeness and consistency.**
  - Data from multi-center projects can be particularly problematic

- **Many studies include three levels of genotype data:**
  - Raw data
    - CEL or iDat files
    - Hardest to use
  - Processed data
    - Genotype calls or Log Ratio values
    - Individual and/or matrix formats
  - QC'ed data
    - As used for the public analysis results
    - Easiest to use (usually in a format supported directly by SVS)

- **Start from the raw or minimally processed data and do your own QC whenever possible.**

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# The "Gotchas"

**A sampling of issues GHI has observed in dbGaP and elsewhere:**

- **Gender discrepancies**

- **Cryptic relatedness**

- **Phenotype data formatted differently between sample groups in a study**

- **Incomplete matching of subjects between raw and processed genotype data.**
  - Example: 500 with raw data, 510 with processed data, 495 with both.

- **Nsp/Sty mismatches in Affy 500k data**

- **Batch effects processed genotypes**

# Example of Batch Effects in a Multi-Center Study



- Caucasian controls from one center have very different allele frequencies than the Caucasian controls from another center…

# GEO – Gene Expression Omnibus



- **"GEO is a public functional genomics data repository… Tools are provided to help users query and download experiments and curated gene expression profiles."**

- **Primarily a gene expression database, but also includes extensive genotype data**

- **Data access:**
  - "Anybody can access and download public GEO data. There are no login requirements."
  - "NCBI places no restrictions on the use or distribution of the GEO data. However, some submitters may claim patent, copyright, or other intellectual property rights in all or a portion of the data they have submitted."

# GEO Data Profile

- **3413 studies, 1335 with human data (1239 mouse, 311 rat, etc.)**

- **Genotype data among the human datasets:**
  - 730 datasets flagged as containing some SNP array data
    - 11,715 samples among 200 data series for Affy 6.0
    - 9689 samples in 152 series for Affy 250k-Nsp
    - 1757 samples in 26 series for Illumina Omni-1
    - 1245 samples in 11 series for Illumina 550k

- **Sample sizes are generally much smaller than with dbGaP**

- **Many studies are based on somatic tissues**

- **GEO database structure is sample oriented, very detailed, and very different from dbGaP**

# GEO: Browsing the Database

- **Browse data by platform to get data for every sample or study to use a particular chip.**
  - 762 samples in 24 studies using Illumina Human1M-Duo.

- **Browse by study design to get data for similar types of studies.**
  - 403 results for "SNP genotyping by SNP array."
  - 654 results for "Genome variation profiling by SNP array."

# Using GEO Data

- **GEO is a good resource for test data and reference data.**

- **There are a few large GWAS studies, but not many.**

- **GEO has several human diversity reference panels available for various genotyping arrays.**
  - Illumina posts HapMap data there for many of their arrays.
  - Other diversity panels from NIA, Mayo, others.

- **Raw and processed data formats are usually available.**

- **"Series Matrix File" is a plain text format that is fairly easy to work with.**

GOLDEN HELIX
*Accelerating the Quest for Significance*™

# SRA: Sequence Read Archive

**Sequence Read Archive**

- **SRA…**
  - "Archives raw oversampling NGS data for various genomes from several platforms"
  - "Shares NGS data with EMBL and DDBJ"
  - "Serves as a starting point for 'secondary analysis'"
  - Provides access to data from human clinical samples to authorized users who agree to the dataset's privacy and usage mandates."

- **SRA primarily stores reads reads (SRA/fastq) and alignments (BAM)**

- **SRA hosts sequence data for some dbGaP and EGA studies**
  - Data not part of public SRA, but searchable summaries do appear on SRA.

- **PubMed abstracts can be linked to research data on SRA**

# Our Team's Experience with SRA



- **A recent Golden Helix webcast featured bison and cattle sequence data from SRA.  Read about it on our blog!**

# EGA: European Genome-Phenome Archive



- **European equivalent of dbGaP**

- **Many EGA datasets are searchable on dbGaP**

- **May be most familiar as the repository for the WTCCC GWAS data**

- **From 2013 IGES talk by Justin Paschall:**
  - Over 450 studies in EGA
  - Extensive sequence data, including 110k BAM files and 35k fastq
  - Current submission rate of about 30TB/month

- **From personal experience: don't forget to request the decryption key…**

# A Few More Sources

- **Illumina provides example data for most of their genotyping chips**
  - Complete HapMap Phase 2 populations for some, subset for others

- **Major imputation software developers have 1000 Genomes reference panels available in their preferred input formats**
  - Beagle
  - Impute2
  - MACH

- **Golden Helix offers several public datasets for download from within SVS**
  - HapMap data for various genotyping chips
  - 1000 Genomes
  - Complete Genomics

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Agenda

**1**    Why Use Public Data?

**2**    Where to Find Public Data

**3**    Tips for Using Public Data

**4**    Manipulating Public Data in SVS

GOLDEN HELIX
*Accelerating the Quest for Significance*™

# Final Tips for Using Public Data

- **Read the documentation BEFORE you download the full archive**

- **Be vigilant with QC**

- **You can't be too careful, especially when combining data from multiple sources**
  - Start from raw data and process each source with a standard protocol. Re-calling genotypes is never a bad idea.
  - Pay special attention to strand orientation
  - Best if all sources were genotyped with the same array, but consider using imputation to combine data from mismatched arrays
  - Always adjust statistical tests for the data source

- **Examine results carefully before reporting or publishing**
  - Give special attention to results involving rare alleles.
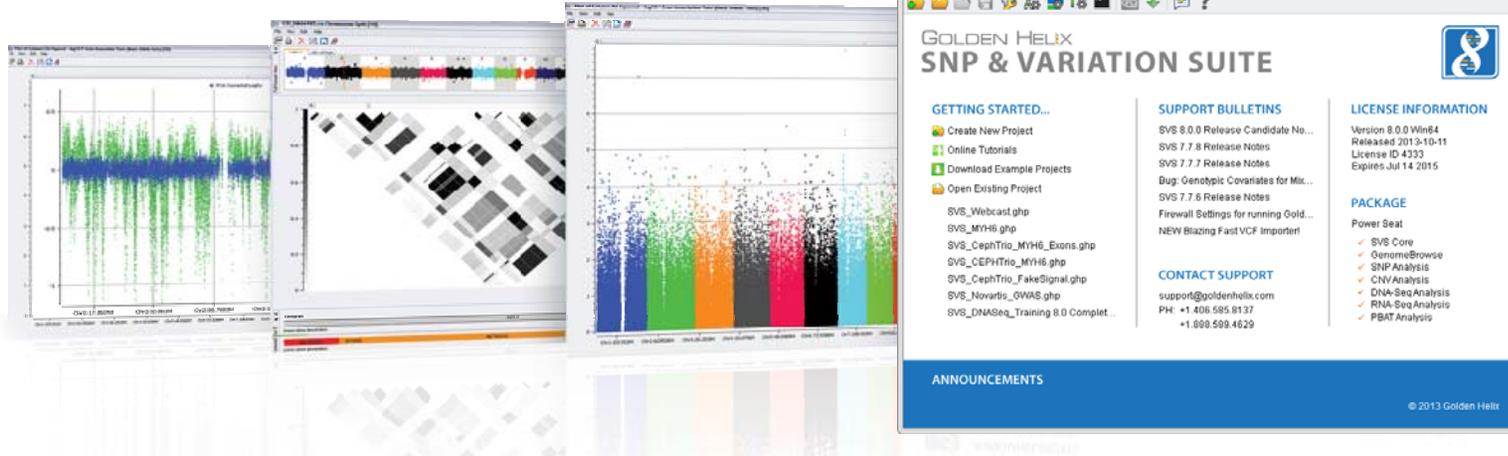  - If something seems fishy, it probably is.

# Challenges of Public Data

**Some of the challenges we hear about at Golden Helix:**

- **"These files are really big!"**
  - Welcome to the world of bioinformatics. Small hard drives need not apply.

- **"Do I need a Linux computer to work with dbGaP data?"**
  - No, but if you're in Windows, you will find that a Linux emulator like CygWin is very useful for manipulating the data. Compression utilities like WinRar and 7-Zip may also be helpful.

- **"There are a bunch of different data formats here…"**
  - Many of the standard formats you find on dbGaP and elsewhere can be read by SVS. Contact us if you're not sure about a particular file—we might already have an import script that will work with it.

- **"I can read the data in text files, but it needs some serious manipulation before I can use it."**
  - Data manipulation? That's one of the most powerful features in SVS…
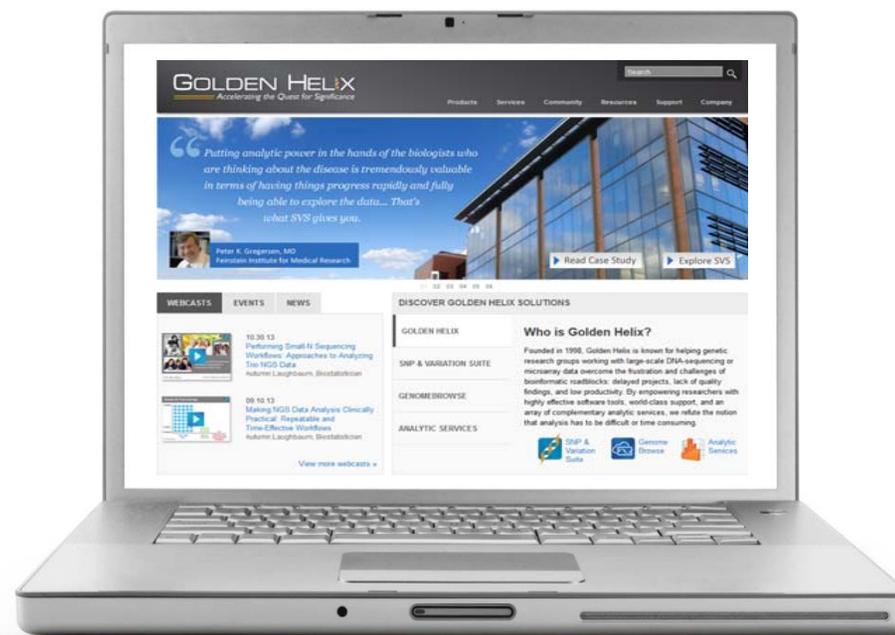
GOLDEN HELIX
*Accelerating the Quest for Significance™*

[Demonstration]

# Questions or more info:

- Email
  [info@goldenhelix.com](mailto:info@goldenhelix.com)

- Request an evaluation of the software at
  [www.goldenhelix.com](http://www.goldenhelix.com)

# Questions?

Use the Questions pane in your GoToWebinar window