# A Walk Through GWAS

January 20th, 2016

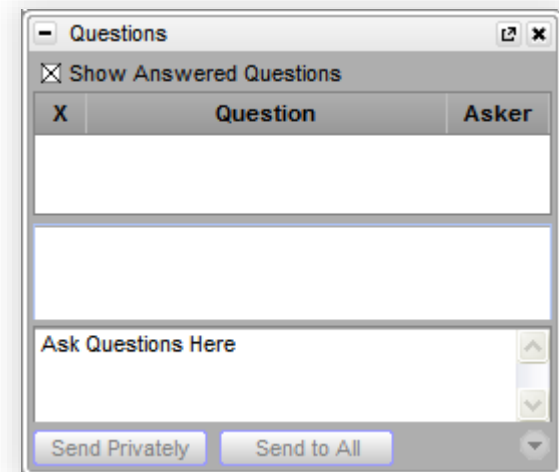Jami Bartole
Senior Field Application Scientist

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Questions during the presentation
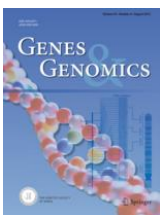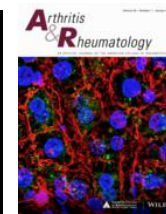
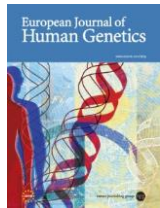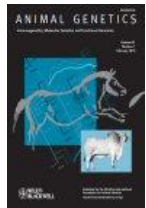Use the Questions pane in your GoToWebinar window

# Golden Helix – Who We Are

**Golden Helix is a global bioinformatics company founded in 1998.**

**We are cited in over 900 peer-reviewed publications**

# Our Customers

**Over 200 organizations world wide, and thousands of users, trust our software.**

# Golden Helix – Who We Are

**When you choose a Golden Helix solution, you get more than just software**

- **REPUTATION**

- **TRUST**

- **EXPERIENCE**



Genetic Testing for Cancer
Dr. Andreas Scherer
Golden Helix, Inc.

- **INDUSTRY FOCUS**

- **THOUGHT LEADERSHIP**

- **COMMUNITY**
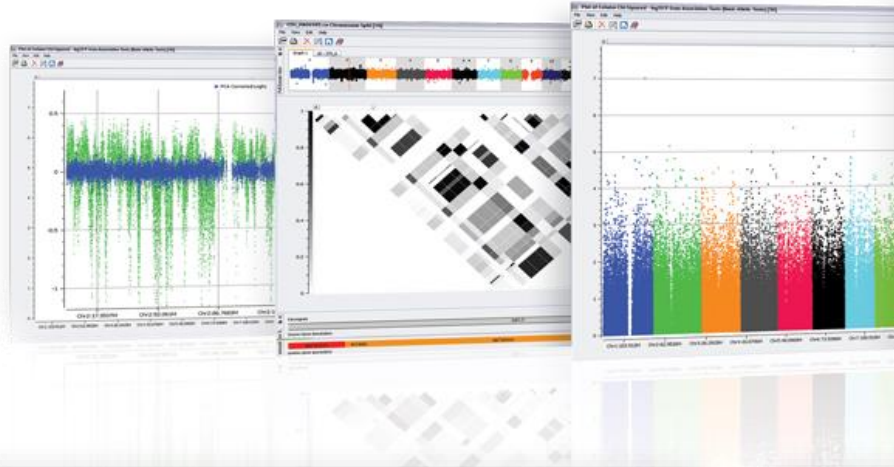
- **TRAINING**

- **SUPPORT**

- **RESPONSIVENESS**

- **TRANSPARENCY**

- **INNOVATION and SPEED**

- **CUSTOMIZATIONS**

# SNP & Variation Suite  (SVS)



## Core Features

- Powerful Data Management
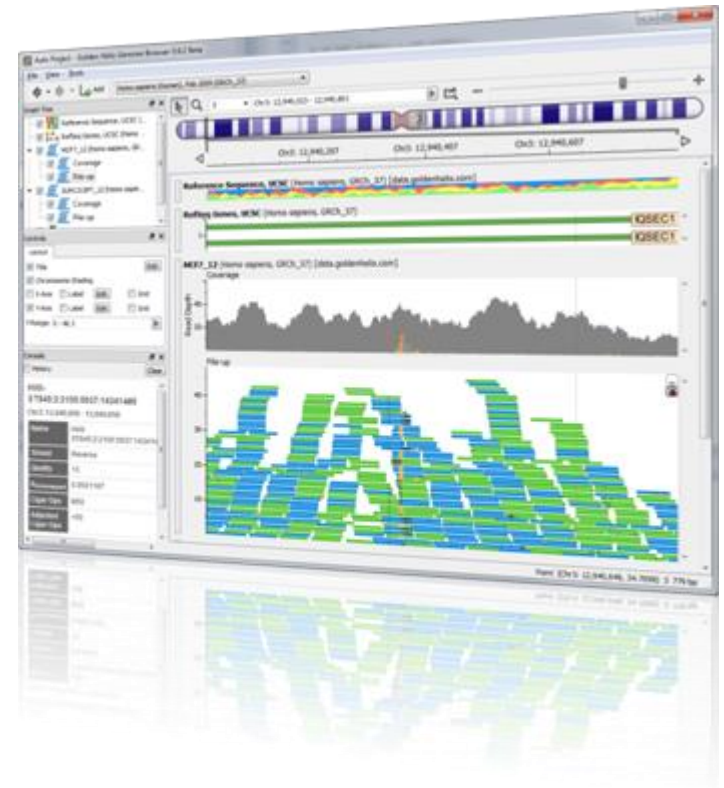- Rich Visualizations
- Robust Statistics
- Flexible

## Applications

- Genotype Analysis
- DNA sequence analysis
- CNV Analysis
- RNA-seq differential expression

# GenomeBrowse

- Powerful visualization software for DNA and RNA sequencing data

- Supports most standard bioinformatics file formats

- Fast and responsive for interactive analysis

- Intuitive controls

- Stream data from the cloud and from your own remote data servers
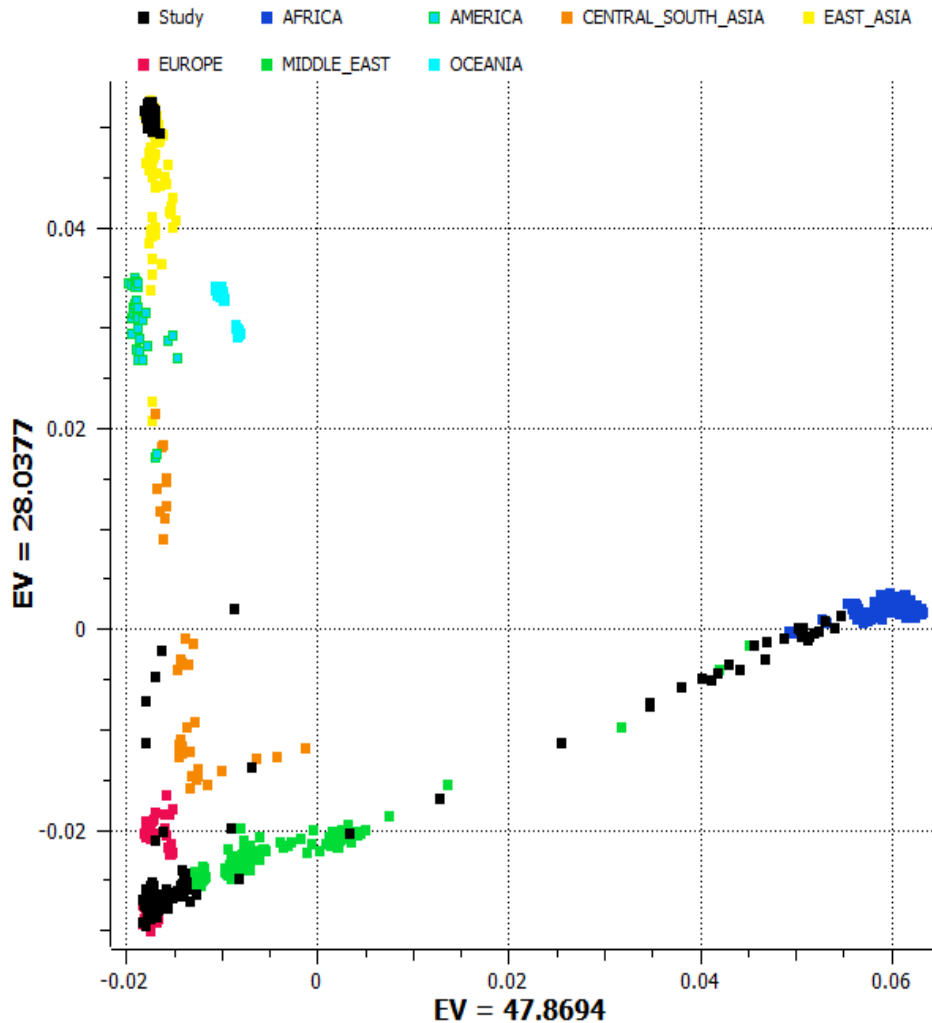
# Approximate Agenda

**1** Background of GWAS

**2** Explore Results from a GWAS Project

**3** Population Stratification in Analysis

**4** Q&A

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# A brief background of GWAS



- First the naïve approaches: Trend Tests, Contingency Tables, Linear/Logistic Regression

- Batch Effects, Population Structure and sharing of controls may violate assumptions of the naïve approaches and result in confounding of results.

- Stratification effects are more pronounced with larger sample sizes.

- Non-independence of samples is especially problematic in agrigenomic applications.

# Summary of GWAS Dataset

- 513 individuals

- 29 populations from the Human Genome Diversity Project (HGDP)

- Illumina Infinium HumanHap550 Genotyping BeadChip

- Simulated Case/Control Phenotype

# Q-Q Plots

# Overview of Methods

| Naïve GWAS | GWAS + Correcting for Population Stratification | Mixed Model Approach |
|---|---|---|

**Corr/Trend Test**

- Quality Control of Samples and Markers

**PCA Correction (Eigenstrat Price 2006)**

- Direct correction of genotype and phenotype data
- Adding PCs as covariates to regression model

**EMMAX (Kang 2010)**

- Using the Genomic Relationship Matrix (IBD) to account for stratification

# Summary

- Performed Basic Association Test
  - Verified contingency table counts
    - http://goldenhelix.com/SNP_Variation/scripts/pages/FrequencyTable.html
  - Q-Q Plot to look for inflation of p-values
    - http://goldenhelix.com/SNP_Variation/scripts/pages/CalculatePseudoLambda.html

- Examined workflow to determine reasons for inflation of p-values
  - Sample/Marker Statistics
    - Call Rate Histograms
  - Cryptic Relatedness through IBD
    - IBD Heat Maps
  - Population Stratification with PCA
    - PCA Plots (2D & 3D)

- Adjusted Analysis for Population Stratification
  - Using PCA from within Association Testing dialog
  - Using PCs as Covariates in Numeric Regression
  - Using Mixed Linear Model Analysis

# Questions or more info:

- Email info@goldenhelix.com

- Request an evaluation of the software at www.goldenhelix.com

- Check out our abstract competition!