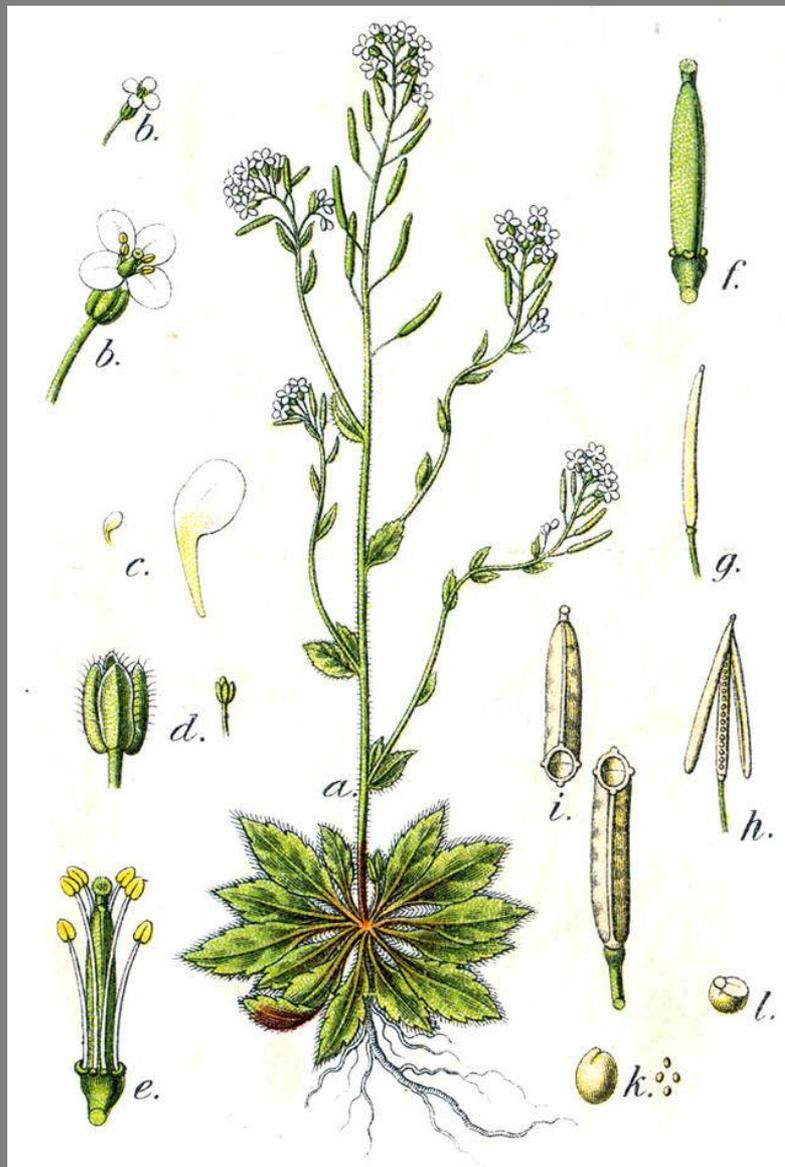# GWAS in a model organism: *Arabidopsis thaliana*

June 9, 2014
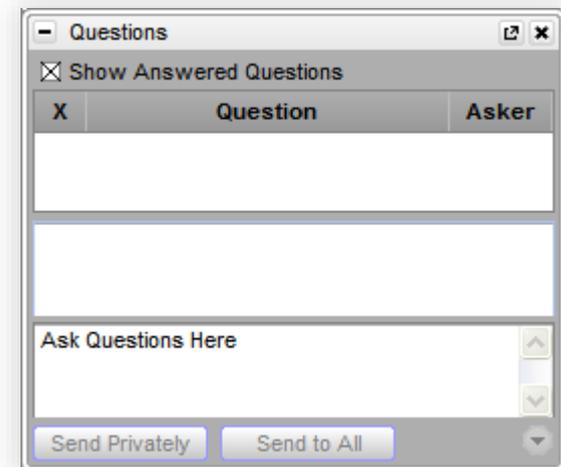
Ashley Hintz
Field Application Scientist
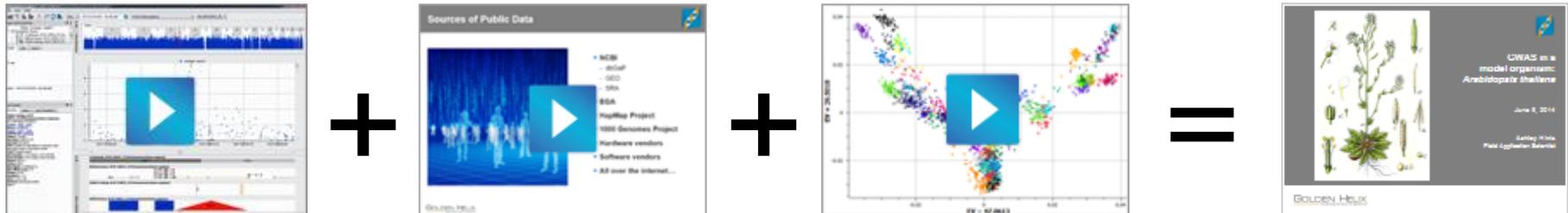
# Questions during the presentation

Use the Questions pane in your GoToWebinar window

# Introduction

- Combining topics of previous webcasts:
  - "Maximizing Public Data Sources for Sequencing and GWAS Studies"
  - "Back to Basics: Using GWAS to Drive Discovery for Complex Diseases"
  - "Mixed Models: How to Effectively Account for Inbreeding and Population Structure in GWAS"

- Additional Dimension: Non-human data from *A. thaliana*
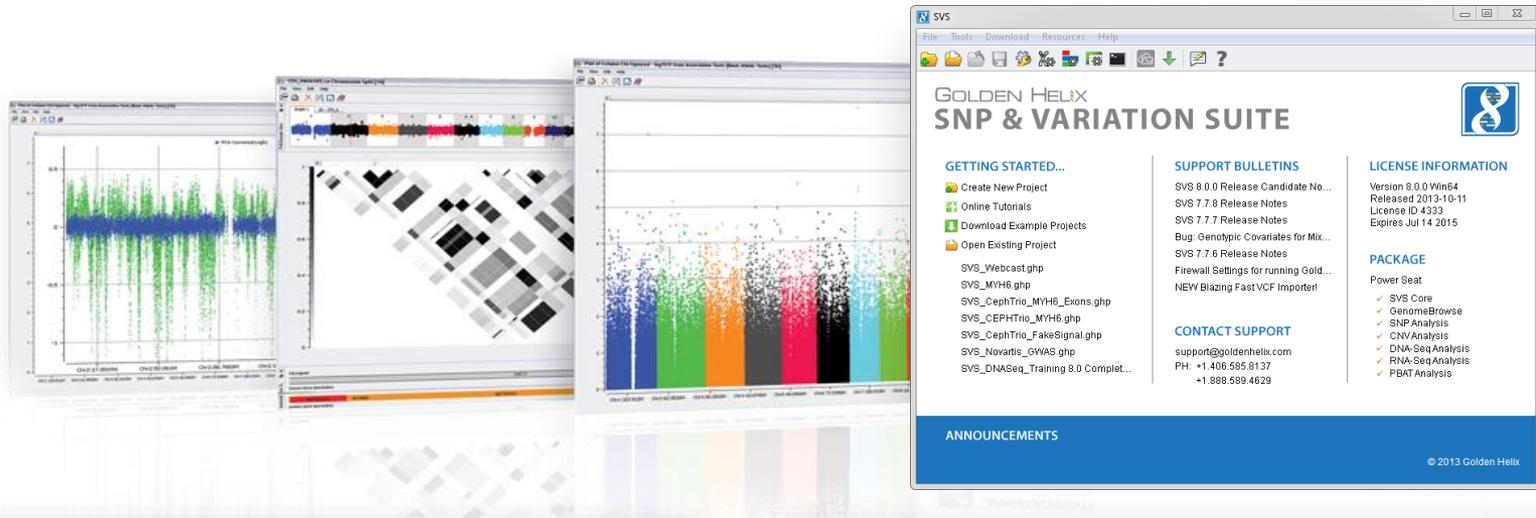
- Combining these topics in a single webcast!

# Agenda

| 1 | *A.* Thaliana Data Overview |
|---|---|

| 2 | Obtaining Public Data |
|---|---|

| 3 | SVS Demonstration |
|---|---|

| 4 | Conclusion and Questions |
|---|---|

# SNP & Variation Suite  (SVS)



## Core Features

- Powerful Data Management
- Rich Visualizations
- Robust Statistics
- Flexible
- Easy-to-use

## Applications

- Genotype Analysis
- DNA sequence analysis
- CNV Analysis
- RNA-seq differential expression
- Family Based Association

# Agenda

GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Dataset Overview

- **Downloaded from Gregor Mendel Institute for Molecular Plant Biology through AtPolyDB**

- **Custom Affymetrix 250K SNP chip**
  - Genotyped 1307 samples covering 214,051 markers

- **107 different phenotypes were recorded from Atwell *et al.* 2010**
  - Flowering, ionomics, defense and development

- **Trait associated with virulence to *Pseudomonas***





www.pseudomonas-syringae.org

# Agenda

**1**    *A. thaliana* Data Overview

**2**    Obtaining Public Data

**3**    SVS Demonstration

**4**    Conclusion and Questions

# Raw Data

- **Genotype Data**

  - Downloaded as .CSV file containing haploid genotypes

  - Imported text file to SVS, used Python script to convert to homozygous diploid genotypes for analysis

- **Sample Data**

  - Text file detailing collection site of all 1307 samples, including GPS coordinates

  - Additional text file with phenotypes for 199 samples, as used in Atwell *et al.* 2010 GWAS paper

  - Both were imported to SVS for further analysis

# TAIR_9 Assembly Annotations

- Downloaded from Arabidopsis.org

- Reference sequence in FASTA format

- Gene annotations in GFF3 format

- Both converted to SVS native TSF format using tools in SVS



The Arabidopsis Information Resource

The Arabidopsis Information Resource (TAIR) maintains a database of genetic and molecular biology data for the model higher plant *Arabidopsis thaliana* . Data available from TAIR includes the complete genome sequence along with gene structure, gene product information, gene expression, DNA and seed stocks, genome maps, genetic and physical markers, publications, and information about the Arabidopsis research community. Gene product function data is updated every week from the latest published research literature and community data submissions. TAIR also provides extensive linkouts from our data pages to other Arabidopsis resources.

# Agenda

**1**    *A. thaliana* Data Overview

**2**    Obtaining Public Data

**3**    SVS Demonstration

**4**    Conclusion and Questions

# Demonstration Agenda

| 1 | Quality Assurance Filters |
|---|---|

| 2 | Principle Components Analysis |
|---|---|

| 3 | Genotype Association Testing |
|---|---|

| 4 | EMMAX |
|---|---|

| 5 | Visualizations |
|---|---|

[Demonstration]

# Agenda

**1**    *A. thaliana* Data Overview

**2**    Obtaining Public Data

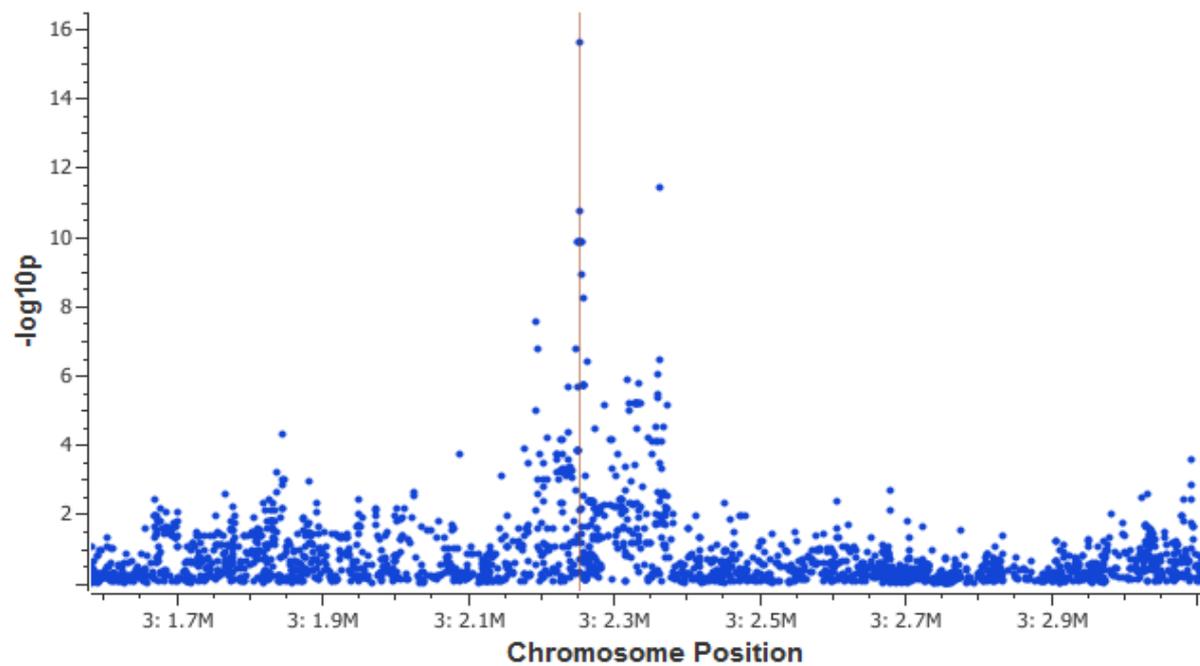**3**    SVS Demonstration

**4**    Conclusion and Questions
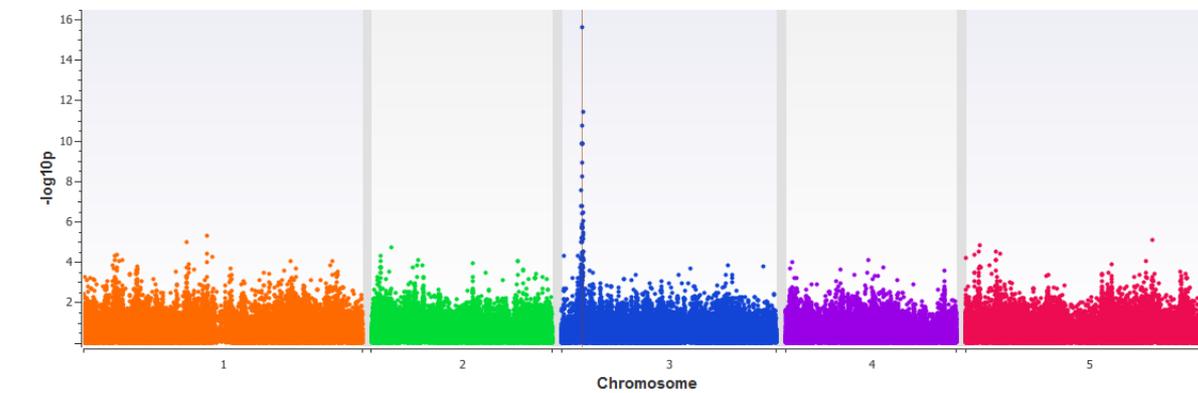
GOLDEN HELIX
*Accelerating the Quest for Significance™*

# Conclusion

Figure 2 | GWA analysis of hypersensitive response to the bacterial elicitor *AvrRpm1*.

# Conclusion

# Questions or more info:

- Email
  info@goldenhelix.com

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

GOLDEN HELIX
Accelerating the Quest for Significance™