# MM-KBAC: Using mixed models to adjust for population structure in a rare-variant burden test
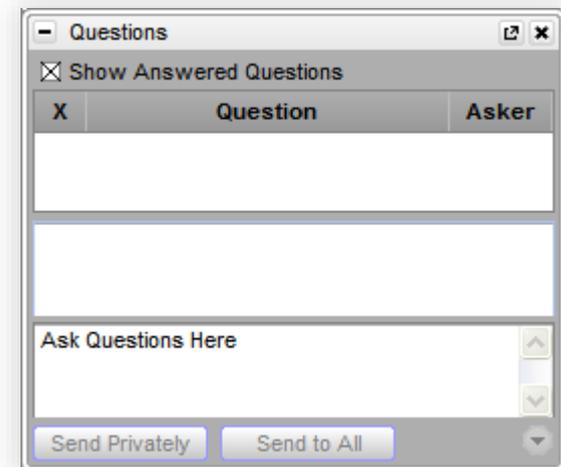
Tuesday, June 10, 2014
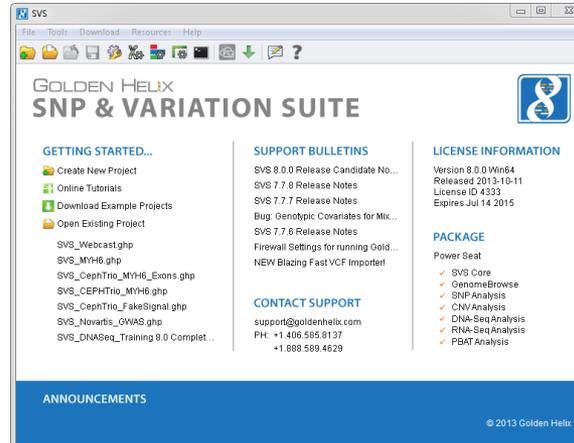
Greta Linse Peterson
Director of Product Management and Quality

GOLDEN HELIX
Accelerating the Quest for Significance™

# Questions during the presentation

Use the Questions pane in your GoToWebinar window

# Golden Helix Offerings



## Services

- Genomic Analytics
- Genotype Imputation
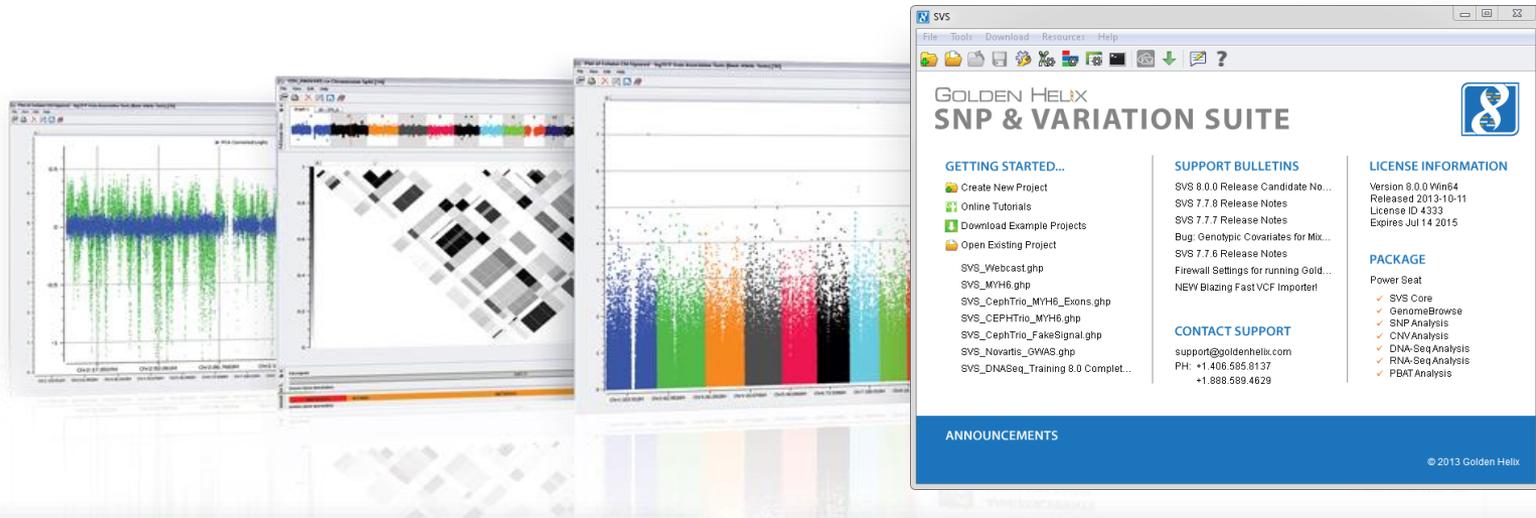- Workflow Automation
- SVS Certification & Training

## Software

- SNP & Variation Suite (SVS) for NGS, SNP, & CNV data
- GenomeBrowse
- New products in development

## Support

- Support comes standard with software
- Customers rave about our support
- Extensive online materials including tutorials and more
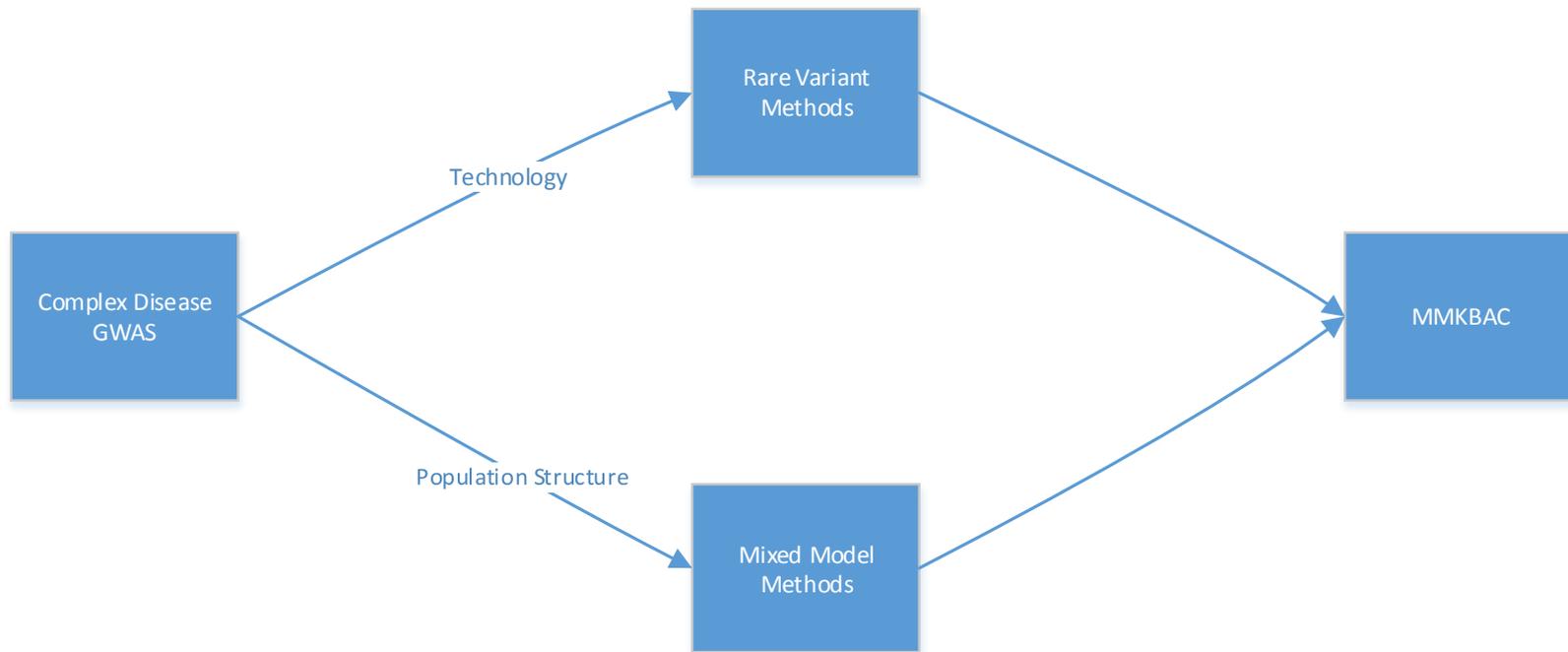
# SNP & Variation Suite  (SVS)



## 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

# Timeline

# Study Design

- **Large cohort population based design (cases with matched controls or quantitative phenotypes and complex traits)**
  - Assumes: independent and well matched samples
  - Can interrogate complex traits

- **Small families (trios, quads, small extended pedigrees)**
  - Can only analyze a single family at a time, looking for de Novo, recessive or compound het variants unique to an affected sample in a single family
  - Looking for highly penetrant variants

# What If????

- **What if we have:**
  - Known population structure
  - Cannot guarantee independence between samples
  - Controls were borrowed from a different study
  - Multiple families with affected offspring all exhibiting the same phenotype
  - Multiple large extended pedigrees of unknown structure

# Just Add Random Effects!

- **Why can't we just add random effects to our regression models for our rare-variant burden testing algorithms?**

- **Existing mixed model algorithms assume a linear model**

- **Kernel-based adaptive clustering (KBAC) uses a logistic regression model**

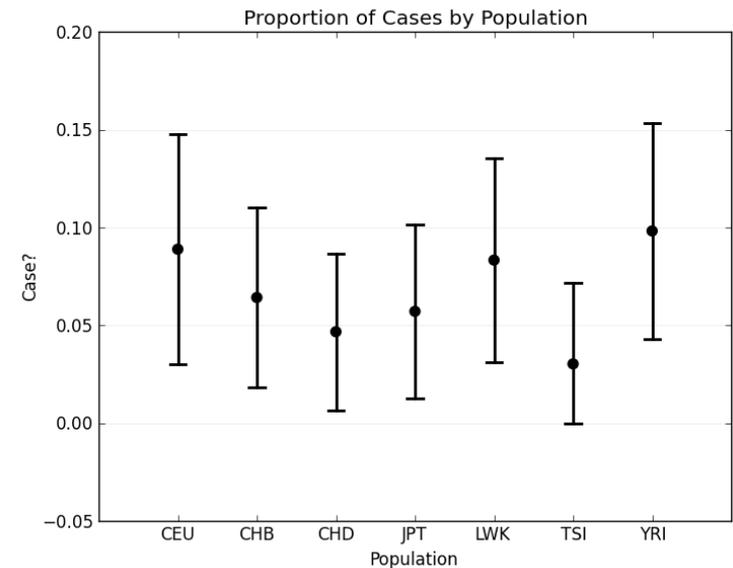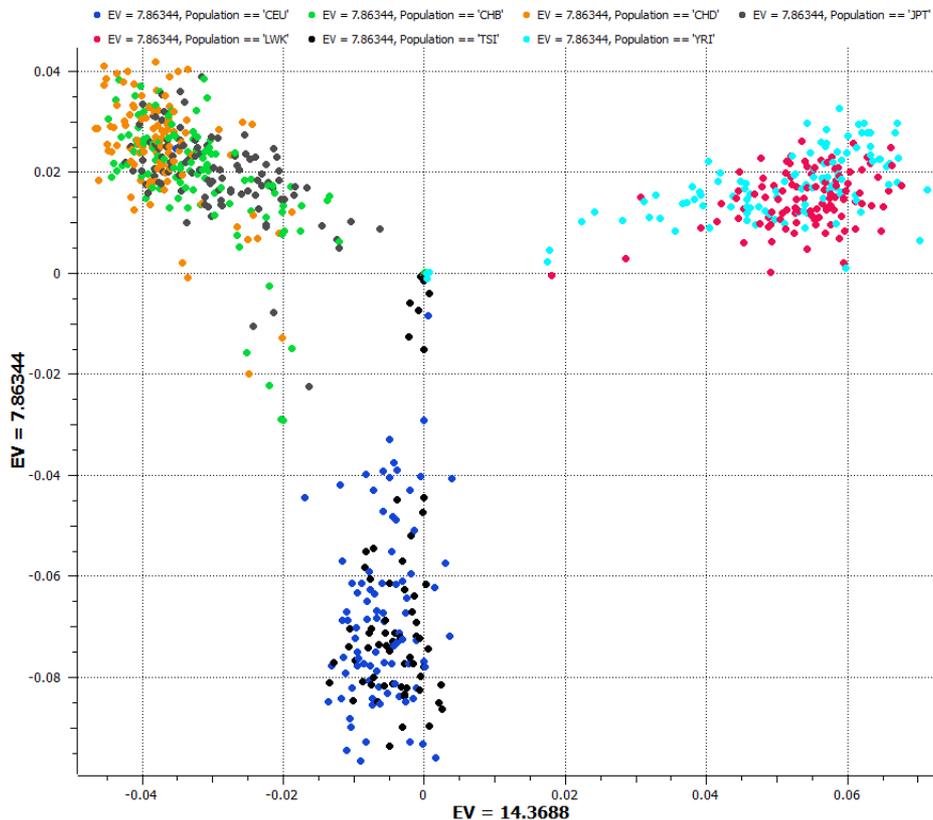- **Hmm what to do....?**

# WARNING!

**What is about to follow are formulas and statistics, specifically matrix algebra…**

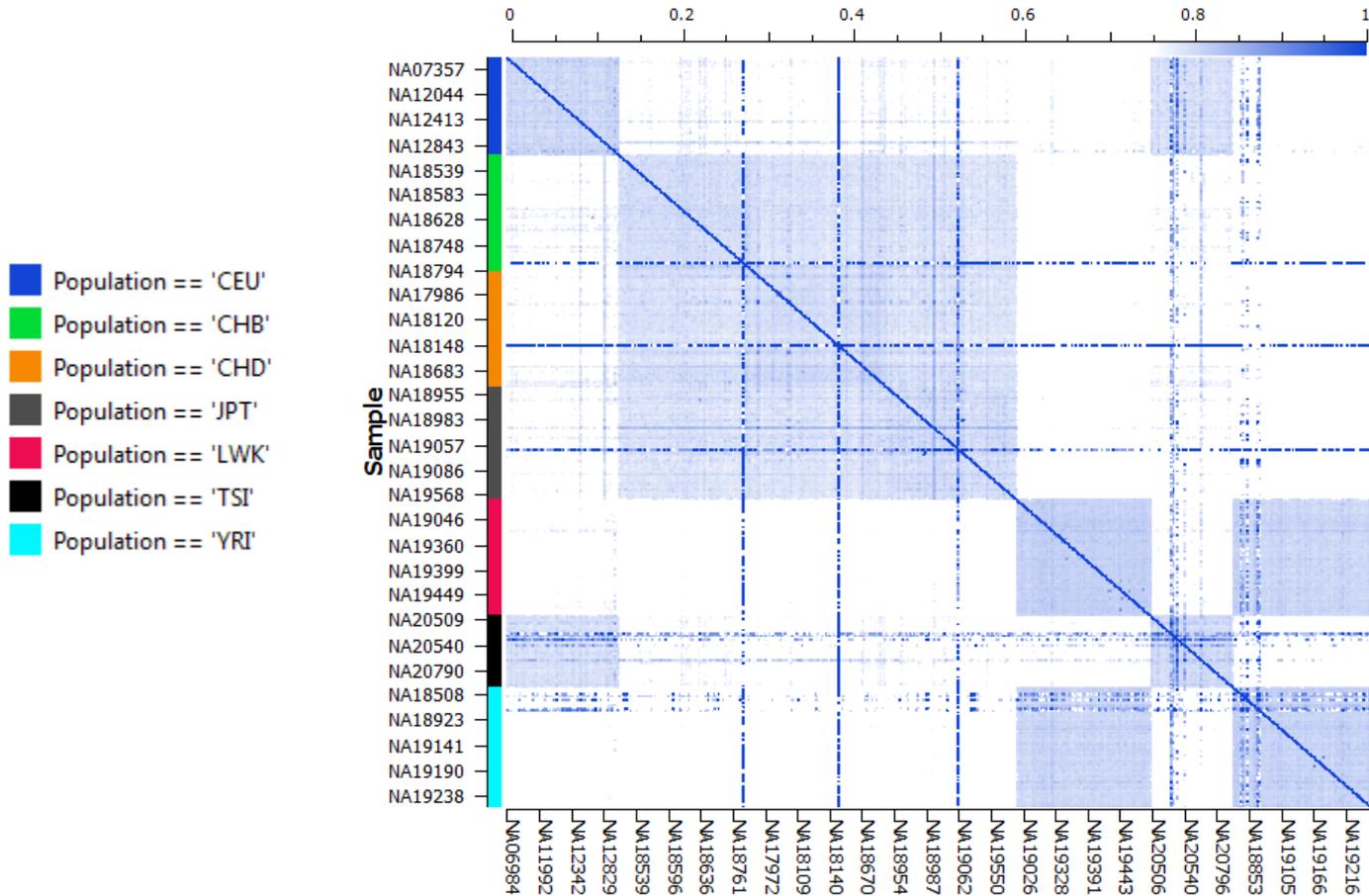**But don't worry we'll end the webcast with a presentation of some preliminary results! So hang in there!**

# But first….

The dataset we have chosen for today is the 1000 Genomes Pilot 3 Exons dataset with a simulated phenotype.

# Relatedness of samples

# Why Mixed Models + KBAC?

- **OK Mixed Models makes sense, but why KBAC?**

- **KBAC was chosen as our proof of concept rare-variant burden test for complex traits**

- **KBAC uses a score test which is trivial to calculate once you compute the reduced model**

- **Mixed models can be added to other burden and kernel tests using the same principles**

# What is KBAC?

- **KBAC = Kernel-based Adaptive Clustering**

- **Catalogs and counts multi-marker genotypes based on variant data**

- **Assumes the data has been filtered to only rare variants**

- **Performs a special case/control test based on the counts of variants per region (aka gene)**

- **Test is weighted based on how often each genotype is expected to occur according to the null hypothesis**

- **Genotypes with higher sample risks are given higher weights**

- **One-sided test primarily, which means it detects higher sample risks**

## A Novel Adaptive Method for the Analysis of Next-Generation Sequencing Data to Detect Complex Trait Associations with Rare Variants Due to Gene Main Effects and Interactions

Dajiang J. Liu[1,2], Suzanne M. Leal[1,2]*

1 Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas, United States of America, 2 Department of Statistics, Rice University, Houston, Texas, United States of America

**Abstract**

There is solid evidence that rare variants contribute to complex disease etiology. Next-generation sequencing technologies make it possible to uncover rare variants within candidate genes, exomes, and genomes. Working in a novel framework, the kernel-based adaptive cluster (KBAC) was developed to perform powerful gene/locus based rare variant association testing. The KBAC combines variant classification and association testing in a coherent framework. Covariates can also be incorporated in the analysis to control for potential confounders including age, sex, and population substructure. To evaluate the power of KBAC: 1) variant genetic data was simulated using rigorous population genetic models for both Europeans and Africans, with parameters estimated from sequence data, and 2) phenotypes were generated using models motivated by complex diseases including breast cancer and Hirschsprung's disease. It is demonstrated that the KBAC has superior power compared to other rare variant analysis methods, such as the combined multivariate and collapsing and weight sum statistic. In the presence of variant misclassification and gene interaction, association testing using KBAC is particularly advantageous. The KBAC method was also applied to test for associations, using sequence data from the Dallas Heart Study, between energy metabolism traits and rare variants in ANGPTL 3,4,5 and 6 genes. A number of novel associations were identified, including the associations of high density lipoprotein and very low density lipoprotein with ANGPTL4. The KBAC method is implemented in a user-friendly R package.
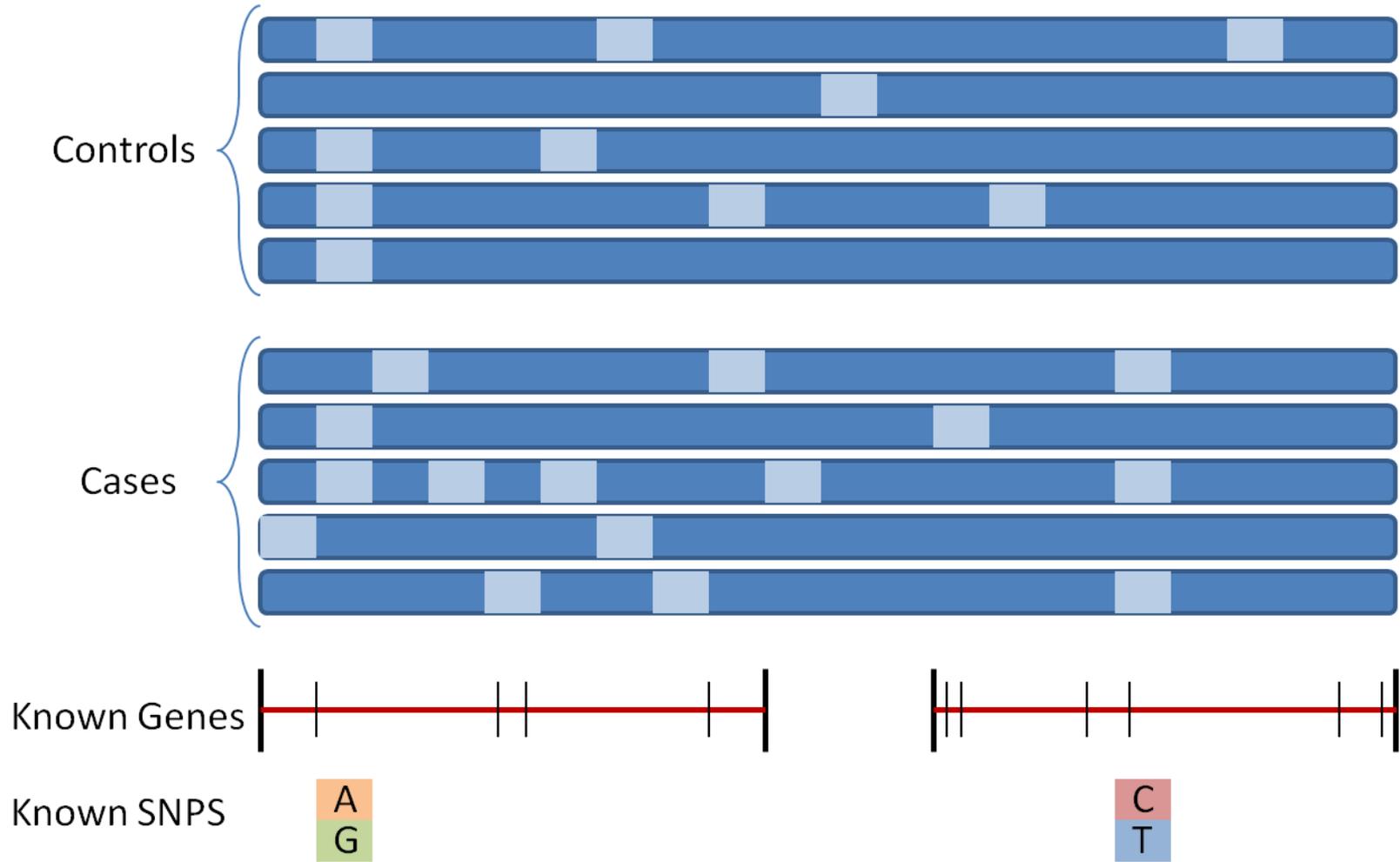
## Introduction

Currently there is great interest in investigating the etiology of complex disease due to rare variants [1–6]. Until recently, indirect mapping of common variants has been the emphasis of complex trait association studies. It has been demonstrated that common variants tend to have modest phenotypic effects while rare variants are likely to have stronger phenotypic effects [7], although not strong enough to cause familial aggregation [8]. For mapping complex diseases due to common variants, instead of genotyping functional variants, tagSNPs are genotyped which act as a proxy for the underlying causal variants. For rare variant association studies, indirect mapping is not an optimal approach due to low correlations ($r^2$) between tagSNPs and rare variants. Instead, direct mapping should be used, where functional variants are analyzed. In order to implement direct mapping, variants must first be identified. Large scale sequencing efforts have begun including the 1000 Genome Project, which will provide a better understanding of the allelic architecture of the genome and a detailed catalog of

human variants. Next-generation sequencing technologies e.g. Roche 454, ABI SOLiD, and Illumina HiSeq, have made it feasible to carry-out rare variant association studies of candidate regions, exomes and genomes.
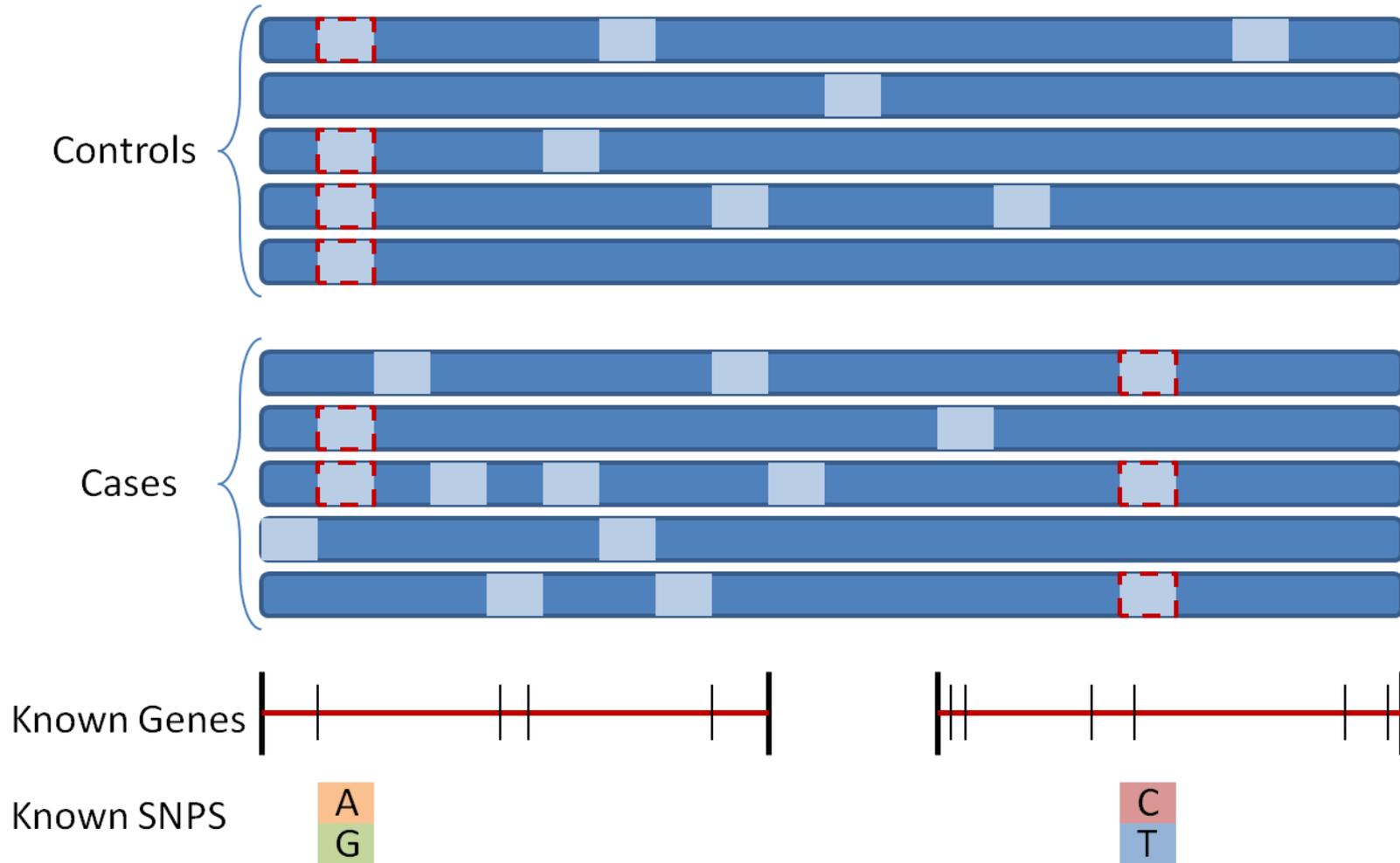
Gene interactions are believed to be involved in a broad spectrum of complex disease etiologies [9]. Although a number of methods have been developed to detect gene interactions between common variants [10–13], their detection has been limited [10]. There is evidence that rare variant interaction also plays a role in disease etiology. In direct association mapping of rare variants, one or more genetic loci are commonly jointly analyzed in order to aggregate information, for example genes with similar functions or residing in the same pathway [3,4]. Therefore it is necessary to account for potential interactions between rare variants in different loci [14] and interactions between common and rare variants [15,16].

Ideally, when carrying out direct mapping, only causal variants should be tested for associations. When DNA samples are sequenced, both causal and non-causal variants are uncovered. Bioinformatics tools [17,18] or filters [1] can be used to predict functionality of
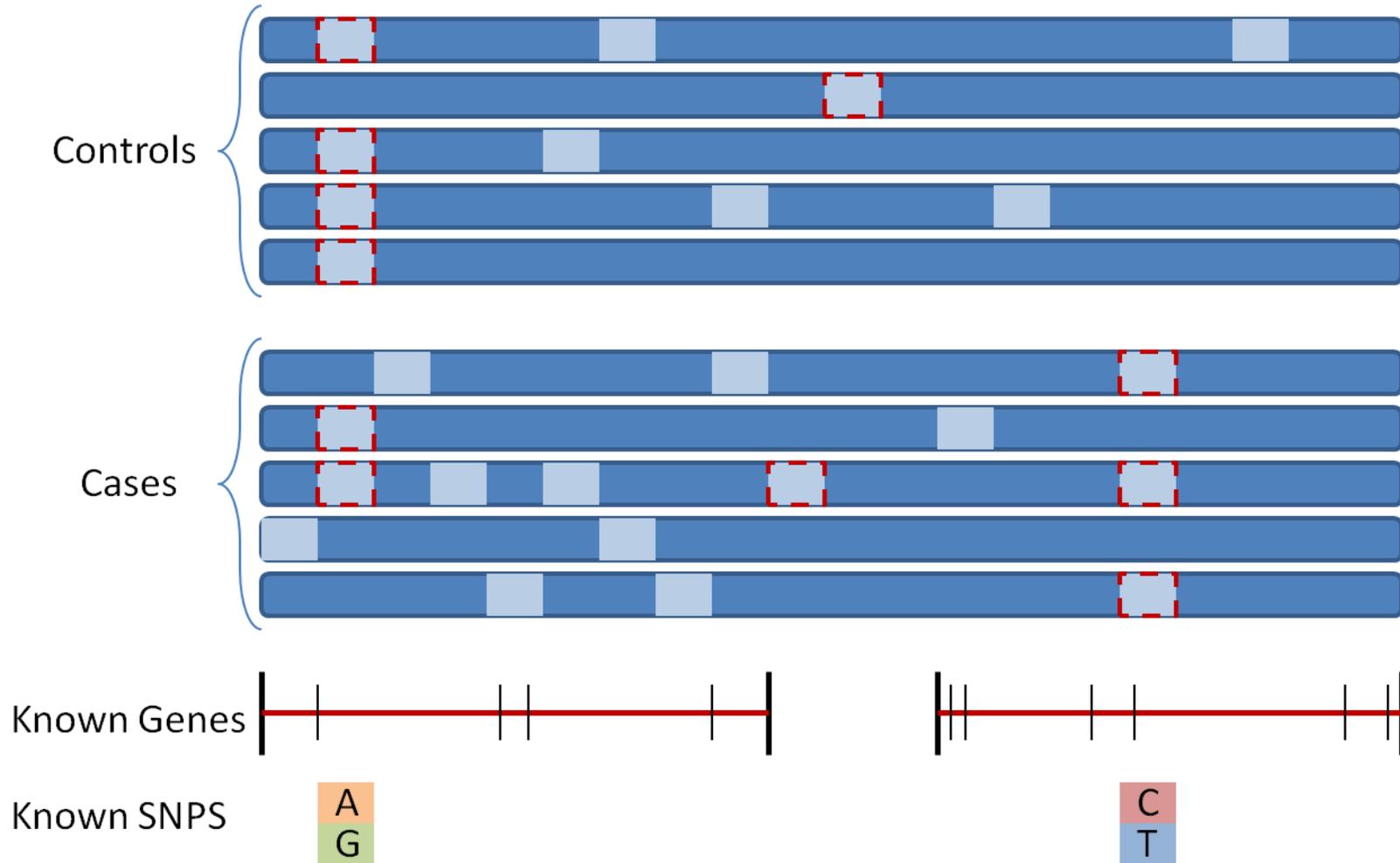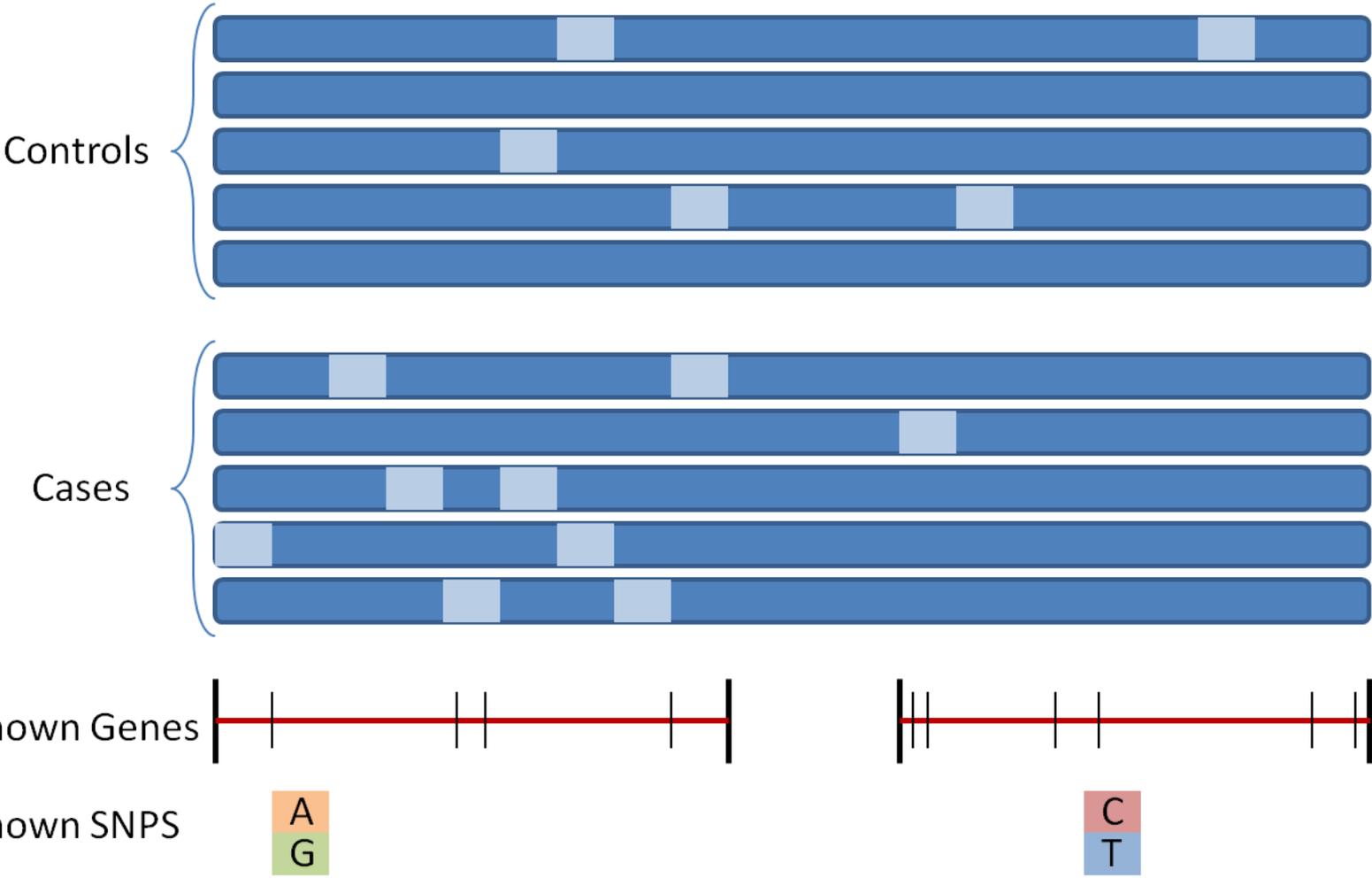
# Rare Sequence Variants

# KBAC Statistic

- $KBAC_1 = \sum_{i=1}^{k} \left( \frac{N_i^A}{N^A} - \frac{N_i^U}{N^U} \right) K_i^0(\hat{R}_i)$

- **Where the weight is defined as:**

$$w_i = K_i^0(\hat{R}_i) = \int_0^{\hat{R}_i} k_i^0(r) dr$$

**The weight can be calculated as a:**

- **Hyper-geometric kernel**

- **Marginal binomial kernel**

- **Asymptotic normal kernel**

- **Monte-Carlo Method is used as an approximation for finding the p-value**

- **The number of cases $n_i^A$ for each genotype $G_i$ approximates a binomial distribution $n_i^A \sim Binom\left(n_i, \dfrac{n^A}{n}\right)$**

- **The case status is permuted among all samples. The covariates and genotypes are held fixed.**

# Logistic Mixed Model Equation

$$\log\left(\frac{P(Y_j = 1 \mid X_j, X_{fjl}, u_j)}{1 - P(Y_j = 1 \mid X_j, X_{fjl}, u_j)}\right) = \beta_0 + \beta_1 X_j + \sum_l (\beta_{fl} X_{fjl}) + u_j$$

**Null hypothesis:** $H_0: \beta_1 = 0$

**The score statistic to test the null of the independence of the model from $X_j$ is:**

$U = \sum_j X_j (Y_j - \mu_j)$, **where**

$\mu_j = h(\eta_j) = \frac{e^{\eta_j}}{1 + e^{\eta_j}}$, **and**

$\eta_j = \beta_0 + \sum_l \beta_{fl} X_{fjl} + u_j$, **and**

$u_j$ **is the random effect for the** $j\wedge(th)$ **sample.**

# Logistic (Reduced) Mixed Model Equation

$$\log\left(\frac{P\big(Y_j = 1 \mid X_{fjl}, u_j\big)}{1 - P\big(Y_j = 1 \mid X_{fjl}, u_j\big)}\right) = \beta_0 + \sum_l \beta_{fl} X_{fjl} + u_j$$

**Which can be rewritten as:**

$$E[Y|u] = h\big(X_f \beta + u\big) = h(\eta) = \mu$$

**And**

$$Var[Y|u] = A = A^{1/2} A^{1/2}$$

**Where $A$ is the variance of the binomial distribution itself, where**

$$\begin{cases} A_{jj} = \mu_j\big(1 - \mu_j\big) \, for \, j = 1..n \\ A_{ij} = 0 \, for \, i \neq j \end{cases}$$

**And the linear predictor for the model is** $\eta = X\beta + u$

**While $h(\blacksquare)$ is the inverse link function for the model**

# Solving the Logistic Mixed Model

**Iterate between creating a linear pseudo-model and solving for the pseudo-model's coefficients**

$$h(\eta) \doteq h(\tilde{\eta}) + \widetilde{\Delta}X(\beta - \tilde{\beta}) + \widetilde{\Delta}(u - \tilde{u})$$

**Where**

$$\widetilde{\Delta} = \left(\frac{\partial h(\eta)}{\partial \eta}\right)_{\tilde{\beta},\tilde{u}}$$

**Rearranging yields**

$$\widetilde{\Delta}^{-1}\big(\mu - h(\tilde{\eta})\big) + X\tilde{\beta} + \tilde{u} \doteq X\beta + u$$

**The left side is the expected value, conditional on $u$, of**

$$P \equiv \widetilde{\Delta}^{-1}\big(Y - h(\tilde{\eta})\big) + X\tilde{\beta} + \tilde{u}$$

**The variance of $P$ given $u$ is**

$$Var[P|u] = \widetilde{\Delta}^{-1}\tilde{A}^{1/2}\tilde{A}^{1/2}\widetilde{\Delta}^{-1}$$

**Where**

$$\tilde{A}_{jj} = \tilde{\mu}_j(1 - \tilde{\mu}_j), \tilde{A}_{ij} = 0 \ (for \ i \neq j)$$

**Pseudo-model:** $P = X\beta + u + \epsilon$ **and** $Var[\epsilon] = Var[P|u]$

NOTE: As an alternative, rather than using the prediction of $u$ from the pseudo-model, we can use the expected value of $u$, which is zero

**Want to solve using EMMA (Kang 2008)**

**Find** $T$ **such that** $Var[T\epsilon] = I$

**So that we can write**

$$TP = TX\beta + Tu + T\epsilon$$

**And use EMMA to solve the mixed model**

$$TP = TX\beta + Tu + \epsilon^*$$

**Where the variance of** $\epsilon^*$**is proportional to** $I$

**It can be shown that this is solved by letting**

$$T = \widetilde{A}^{-1/2}\widetilde{\Delta}$$

First pick starting values of $\tilde{\beta}$ and $\tilde{u}$, such as all zeros. Repeat the following steps until the changes in $\tilde{\beta}$ and $\tilde{u}$ are sufficiently small:

1. Find $\tilde{\eta}$ and $\tilde{u}$ from the original linear predictor equation and the definition of $h(\blacksquare)$

2. Find the (diagonal) $\tilde{\Delta}$ matrix

3. Find the pseudo-model $P$

4. Find the (diagonal) matrix $T$

5. Solve the following for new values of $\tilde{\beta}$ and $\tilde{u}$ using EMMA:
$$TP = TX\beta + Tu$$

    NOTE: The alternative method modifies Step 5 to use EMMA to determine the variance components and to find a new value for $\tilde{\beta}$, while leaving the value of $\tilde{u}$ at its expected value of zero.

After convergence, the alternative method predicts the values of $u$, and computes the final values of $\eta$ and $\mu$ from this prediction

# Computing the Kinship Matrix

# Applying MMKBAC to a real study
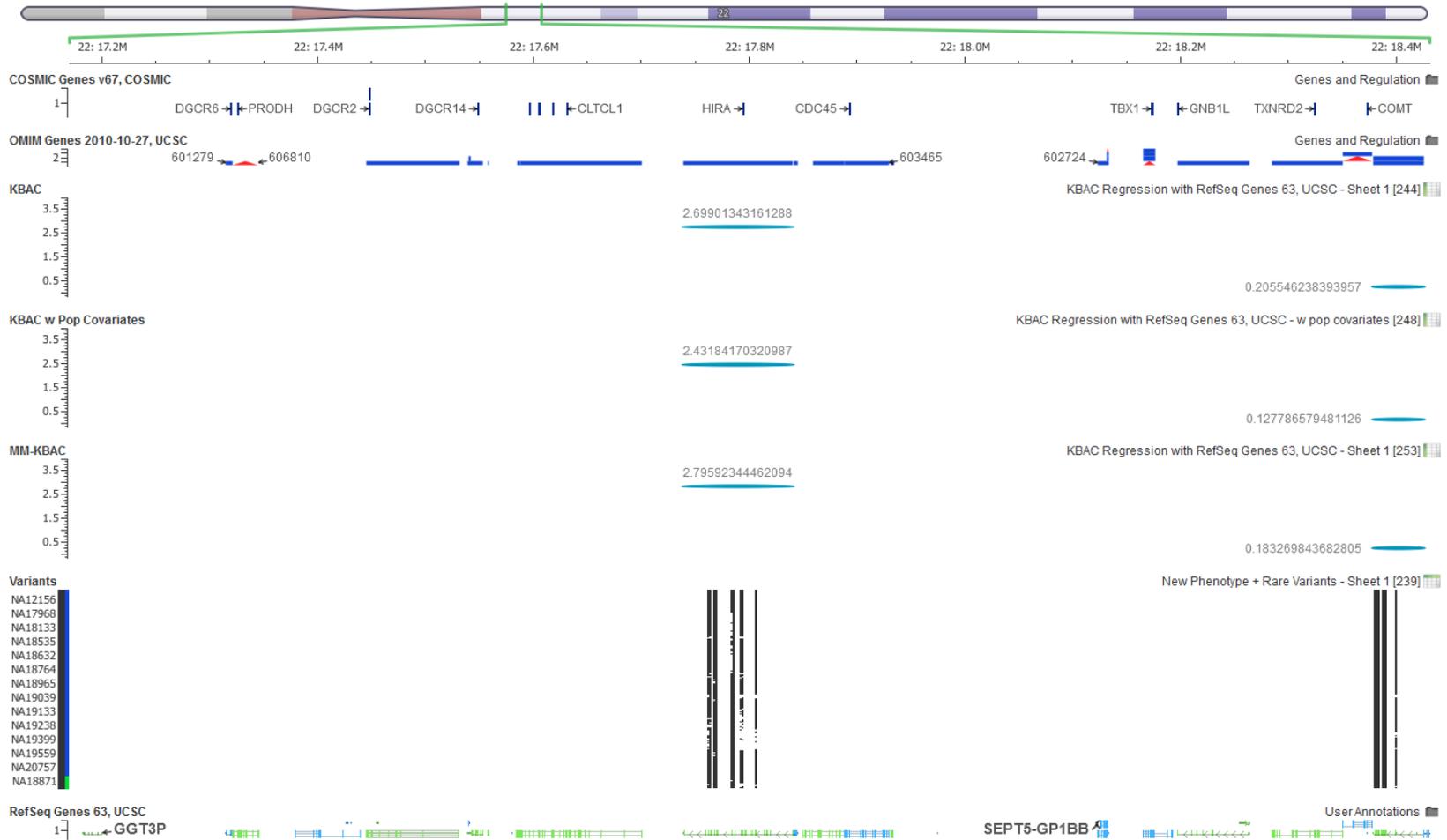
# KBAC vs MM-KBAC QQ Plots



KBAC w Pop. Covariates: $\tilde{\lambda} = 0.902$

# Signal at PSRC1

# Signal at HIRA

# Conclusion

- **This will method will be added into SVS in the near future…**

  **In the meantime…**

- **Like to try it out on your dataset – ask us to be part of our early-access program!**

- **We have submitted an abstract to ASHG, hope to see you there!**

# Announcements

- **Webcast recording and slides will be up on our website tomorrow.**

- **T-shirt Design Contest! Details at www.goldenhelix.com/events/t-shirtcontest.html**

- **Next scheduled webcast is July 22$^{nd}$, but Heather Huson of Cornell University.**

# Questions?

Use the Questions pane in your GoToWebinar window