

# Complex Trait Association Analysis of Rare Variants Obtained from Sequence Data

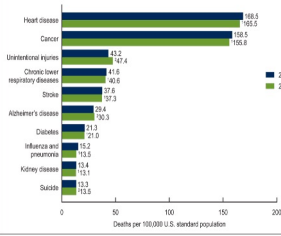
Suzanne M. Leal, Ph.D.  
Sergievsky Family Professor of Neurological Sciences  
Director of the Center for Statistical Genetics  
Columbia University  
sml3@columbia.edu

© 2023 Suzanne M. Leal

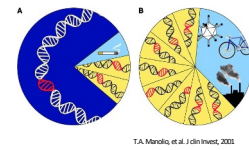
1

## Complex Diseases (Traits)

Top 10 leading causes of death in the United States



Genetic and environmental contribution to complex disorders

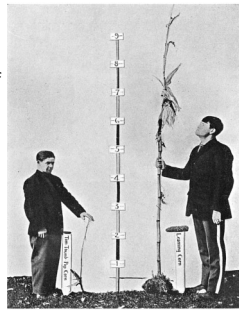


2

## Heritability for Common Traits

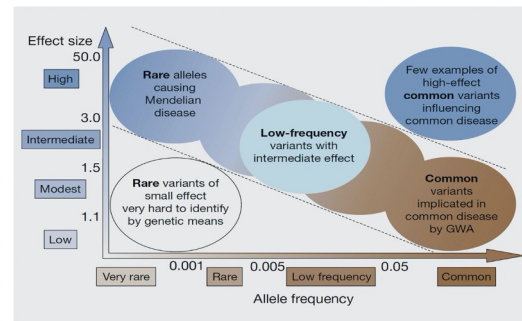
Human height heritability is ~80%

- Strongly associated common variation explain 21–29%
- All common variation explains 60% of height heritability



3

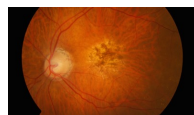
## Allelic Architecture



4

## Complex Disease – Common Variant Associations

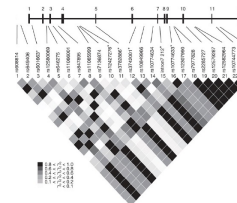
- Disease susceptibility is conferred by variants which are common within populations
  - Variants are old and widespread
- These variants have modest phenotypic effect
- This model is supported by many replicated examples
  - Age Related Macular Degeneration (Klein et al. 2005)
    - Complement factor H (CFH) gene



5

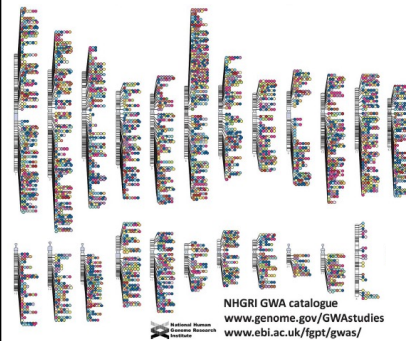
## Studying Complex Traits – Common Variant Associations

- Hundreds of thousands of Single nucleotide polymorphism (SNPs) genotyped and analyzed
  - Indirect mapping
    - Markers usually had a minor allele frequency (MAF) > 0.05
    - Usually not pathogenic – tag SNPs
    - In linkage disequilibrium with disease susceptibility variant



6

## Complex Trait – Common Variant Associations



- Although highly successful in identifying thousands of complex trait loci
- Usually pathogenic susceptibility variant(s) not identified

7

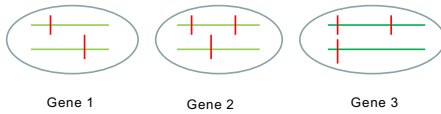
## Complex Disease – Rare Variant Associations

- Complex traits are the result of multiple rare variants
  - Although first thought to large effects, there effect sizes are usually small
- Although these variants are rare, e.g.,  $MAF < 0.005$ 
  - Collectively they may be quite common
- Direct tests of this hypothesis where first reported >15 years ago
  - Dallas Heart Study
    - Small sample ~1,200 individuals
      - Multi-ancestry
      - Used “extreme” sampling
  - Plasma low density lipoprotein levels (Cohen et al. 2004)
    - NPC1L1

8

## Rationale for Rare Variant Aggregate Association Tests

- Testing individual variants with low effect sizes and minor allele frequencies (MAFs)
  - Underpowered to detect associations
- Testing variants in aggregate increases MAFs
  - Improving the power to detect associations



9

## Caveats - Aggregate Rare Variant Association Tests

- Misclassification of variants can reduce power
  - Inclusion of non-causal variants
  - Exclusion of causal variants
- Analysis is limited to
  - Genes
  - Genes within pathways
- Analysis outside of exonic regions is problematic
  - Unlikely a sliding window approach will work
    - Size of window unknown and will differ across the genome
  - A better understanding of functionality outside the coding regions is necessary
    - Predicted functional regions, enhancer regions, transcription factors, DNase I hypersensitivity sites, etc.

10

## Analysis of Rare Variants

- For biobank sized datasets higher frequency rare variants, e.g., 0.5% can be analyzed individually
  - Using same same methods implemented for common variants

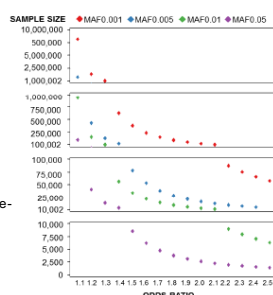
### Example

$$\alpha = 5 \times 10^{-8}$$

Disease prevalence 5%

$$1 - \beta = 0.80$$

\*Note: a more stringent significance criterion may be necessary for genome-wide sequence data. Due to a larger number of effective tests compared to analysis of common variant GWAS panels



11

## A Few Rare Variant Association Tests

- Combined Multivariate Collapsing (CMC)
  - Li and Leal AJHG 2008
- Burden of Rare Variants (BRV)
  - Auer, Wang, Leal Genet Epidemiol 2013
- Weighted Sum Statistic (WSS)
  - Madsen and Browning PLoS Genet 2009
- Kernel based adaptive cluster (KBAC)
  - Liu and Leal PLoS Genet 2010
- Variable Threshold (VT)
  - Price et al. AJHG 2010
- Sequence Kernel Association Test (SKAT)
  - Wu et al. AJHG 2011
- SKAT-O
  - Lee et al. AJHG 2012

Fixed Effect Tests

Random Effect Test

Optimal test

12

### Types of Aggregate Analyses

- Frequency cut offs used to determine which variants to include in the analysis
  - Rare Variants (e.g., MAF<0.05% frequency)
  - Rare and low (MAF=0.05-5%) frequency variants
- Maximization approaches
- Tests developed to detection associations when variants effects are bidirectional
  - e.g., protective and detrimental
- Incorporate weights based upon annotation
  - Frequency
    - e.g., gnomAD
  - Functionality
    - CADD c-scores

13

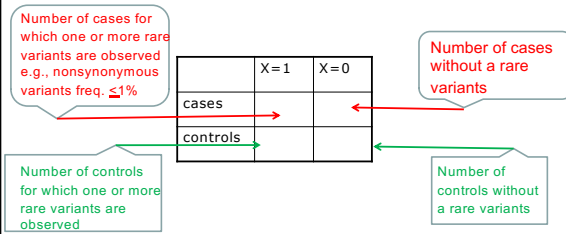
### Methods to Detect Rare Variant Associations Using Variant Frequency Cut-offs

- Combined multivariate & collapsing (CMC)
  - Li & Leal, AJHG 2008
- Collapsing scheme which can be used in the regression framework
  - Can use various criteria to determine which variants to collapse into subgroups
    - Variant frequency
    - Predicted functionality

14

### CMC

- Define covariate  $X_j$  for individual  $j$  as
 
$$X_j = \begin{cases} 1 & \text{if rare variants present} \\ 0 & \text{otherwise} \end{cases}$$
- Compute Fisher exact test for 2x2 table



Can also use same coding in a regression framework

15

### CMC

- Example of coding used in regression framework:
  - Binary coding
 
$$X_j = \begin{cases} 1 & \text{if rare variants present} \\ 0 & \text{otherwise} \end{cases}$$
  - Gene region with 5 variant sites

	Individual	Coding
	1	1
	2	1
	3	0

#### Rare Variant Sites

Green bars: Major allele is observed in the study subject  
Red bars: Minor allele has been observed

16

### Methods to Detect Rare Variant Associations Using Variant Frequency Cut-offs

- Gene-or Region-based Analysis of Variants of Intermediate and Low frequency (GRANVIL)
  - Aggregate number of rare variants used as regressors in a linear regression model
  - Can be extended to case-control studies
    - Morris & Zeggini 2010 Genet. Epidemiol
  - Test also referred to as MZ

17

### GRANVIL

- Example of coding used in regression framework
  - Gene region with 5 variant sites – data available on all sites

	Individual 1	Coded 2/5 (0.4)
	Individual 2	Coded 2/5 (0.4) Note same coding for heterozygous and homozygous genotypes

- Missing data for three of the five variant sites

	Individual 3	Coded 1/2 (0.5)
--	--------------	-----------------

#### Burden Rare Variant (BRV) extension

Individual 1: Coded 2  
Individual 2: Coded 3  
Individual 3: Coded 1

(Auer et al. 2013 Genet Epidemiol)

18

## Methods to Detect Rare Variant Associations Weighted Approaches

- **Group-wise association test for rare variants using the Weighted Sum Statistic (WSS)**
  - Variants are weighted inversely by their frequency in controls (rare variants are up-weighted)
    - Madsen & Browning, PLoS Genet 2009
- **Kernel based adaptive cluster (KBAC)**
  - Adaptive weighting based on multilocus genotype
    - Liu & Leal, PLoS Genet 2010

19

## Methods to Detect Rare Variant Associations Maximization Approaches

- **Variable Threshold (VT) method**
  - Uses variable allele frequency thresholds and maximizes the test statistic
  - Can also incorporate weighting based on functional information
    - Price et al. AJHG 2010
- **RareCover**
  - Maximizes the test statistic over all variants with a region using a greedy heuristic algorithm
    - Bhatia et al. 2010 PLoS Computational Biology

20

## Methods to Detect Associations with Protective & Detrimental Variants within a Region

- **C-alpha**
  - Detects variants counts in cases and controls that deviate from the expected binomial distribution
    - For qualitative traits only
      - Neale et al. 2011 PLoS Genet
- **Sequence Kernel Association Test (SKAT)**
  - Variance components score test performed in a regression framework
    - Can also incorporate weighting
  - Wu et al. 2011 AJHG

21

## Optimal Test

- **SKAT-O**
  - Maximizes power by adaptively using the data to combine a burden test and the sequence kernel association tests
    - Lee et al. 2012 AJHG

22

## Significance Level for Rare Variant Association Tests

- **For exome data where individual genes are analyzed usually a Bonferroni correction for the number of genes tested is used**
  - There is very little to no linkage disequilibrium between genes
- **Bonferroni correction used**
  - e.g.,  $p < 2.5 \times 10^{-6}$  (Correction for testing 20,000 genes)

23

## Determine MAF Cut-offs for Aggregate Rare Variant Association Tests

- **MAF cut-offs are frequently used to determine which variants to analyze in aggregate rare variant association tests**
- **MAF from controls should not be used**
  - Increases in type I error rates
- **Determine variant frequency cut-offs from databases**
  - Using population frequencies for those under study
  - gnomAD
    - <http://gnomad.broadinstitute.org/>

24

### Problem of Missing Genotypes for Aggregate Rare Variant Association Tests

- Same frequency of missing variant calls in cases and controls
  - Decrease in power
- More variant calls missing for either cases or controls
  - Increase in Type I error
  - Decrease in power
- Remove variant sites which are missing genotypes, e.g., >10%
- Can impute missing genotypes using observed allele frequencies
  - For the entire sample
    - Not based on case or control status
- Analyze imputed data using dosages

25

### Dosages

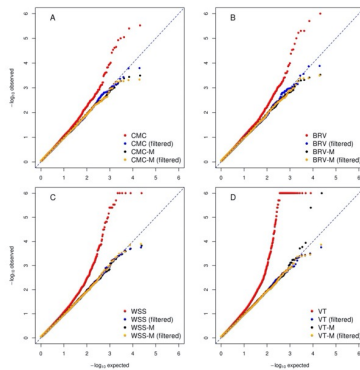
- Genotypes are no longer assigned 0 (1/1), 1 (1/2) or 2 (2/2)
  - Due to uncertainty
- Each genotype is assigned a probability
  - Probabilities sum to 1
- For example
  - Probability of 0 (1/1) genotype is 0.98 and 1 (1/2) genotype is 0.015
- The dosage can be estimated for this example as follows

$$\begin{aligned} 0 \times 0.98 &= 0 \\ 1 \times 0.015 &= 0.015 \\ 2 \times 0.005 &= 0.01 \\ \text{Dosage} &= 0.025 \end{aligned}$$

- Instead of using the most likely genotype the dosage is used

26

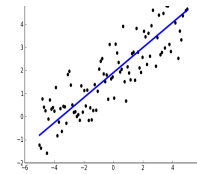
### Results



27

### Rare Variant Aggregate Methods

- Ideally should be performed in a regression framework to adjust for covariates
  - Logistic
  - Linear regression

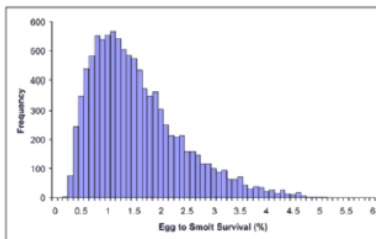


- Almost all rare variant aggregate methods have been extended to be implemented within a regression framework
- Some have also been implemented in a linear mixed model (LMM)/generalized LMM (GLMM)

28

### Analyzing Quantitative Variants

- Most rare variant aggregate analysis methods can be performed on quantitative traits
- If phenotype data includes outliers or deviates from normality
  - Can increase type I errors



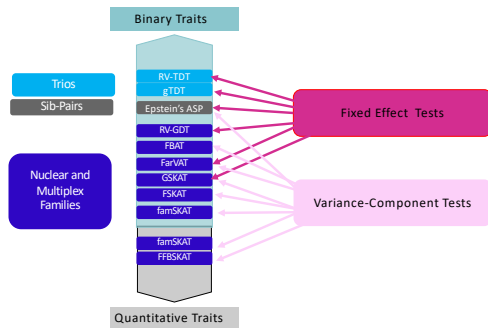
29

### Analyzing Quantitative Variants

- For data that deviates from normality
  - Quantile-quantile normalization
- For data that includes outliers
  - Winsorize
- Don't winsorize and then normalize
- Instead of analyzing quantitative trait values
  - Residual can be generated
    - Adjusting for confounders

30

## Family-based Methods for Rare Variant Aggregate Association Analysis



31

## Linear Mixed Model (LMM) & generalized LMM (GLMM) Analysis of Related & Unrelated Individuals

- LMM is an extension of the linear model to allow for both fixed & random effects and also allows for non-independence of samples
  - Early implementations calculated the kinship matrix  $\Phi$  on the basis of known relationships
  - Amin et al. (2007) proposed to estimate kinships based on genome-wide variant data
    - The generalized relationship matrix (GRM) can be estimated for all individuals using for example identical-by-descent (IBD) sharing
- Extended to binary (case-control) traits - GLMM

32

## LMM and GLMM: Analysis of Related & Unrelated Individuals

- Can be applied to analyze families, cryptically related, & unrelated individuals
  - e.g., UK Biobank
    - 500K study subjects of which 30.3% are  $\leq$  3rd degree relatives & 4.5% sib-pairs
- More recent implementation for large scale data using a variety of methods
  - BOLT-LMM (Loh et al. 2015)
  - FastGWA (Jiang et al. 2019)
  - SAIGE (Zhao et al. 2015)\*
  - REGENIE (Mbatchou et al. 2020) \*
  - SMMAT (Chen et al. 2019)\*\*
- \*Can be used to analyze data where case to control ratio is very unbalanced
  - e.g., 20 cases for every control
- \*\*Cannot be used for UK Biobank Scale data

33

## LMM and GLMM: Analysis of Related & Unrelated Individuals

- To allow for use with biobank sized datasets
- REGENIE does not use the GRM
  - It uses whole genome regression, i.e., the ridge regression
    - In essence, it includes all the SNVs as covariates in the null model
      - Performed by blocks to avoid having to load the entire genome in memory
        - » Using different effect size differences per block
- This large-scale approximation may not control type I error for individuals that are closely related
  - e.g., when only families are being analyzed
  - Can use for example SMMAT
    - Which uses the GRM

34

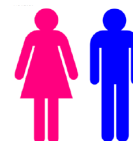
## LMM and GLMM: Analysis of Related & Unrelated Individuals

- A few programs which can perform rare variant aggregate analysis
  - REGENIE - Burden test
  - SMMAT - Burden, SKAT, & SKAT-O tests
  - rvtests (Zhan 2020) implements BOLT-LMM to perform burden association analysis
- An alternative for rare variant aggregate analysis
  - Recode variants within gene regions and then analyze

35

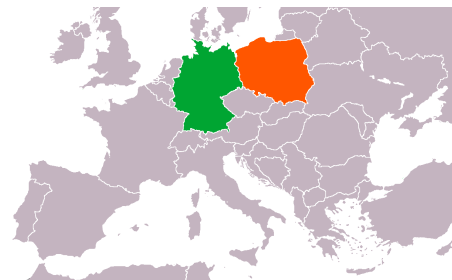
## Rare Variant Association Analysis - Confounders

- Control for covariates in the analysis which are potential confounders
  - Age
  - Sex
  - Batch
  - Body Mass Index (BMI)
  - Smoking pack years
  - Population substructure



36

### Confounder -Population Substructure and Admixis



37

### Population Substructure and Admixture

- If proportion of cases and controls sampled from each population is different
  - Can occur due to
    - Disease frequency is different between populations
    - Sloppy sampling
- Population substructure\admixture can cause detection of differences in variant frequencies within a gene which is due to sampling and not disease status
  - False positive findings can be increased

38

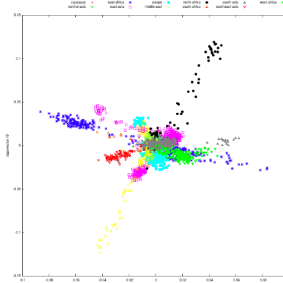
### Example River People



39

### Population Substructure and Admixture

- Currently PCA or MDS are use to control for population substructure\admixture
  - Controls on the global level
  - May not be sufficient
    - For admixed populations
    - Rare variation



40

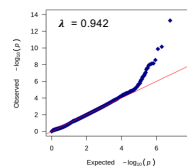
### Rare Variant Aggregate Association Analysis

- When analyzing different populations, e.g.,
  - Africans
  - Europeans
- When analyzing data from different source
  - Analyze each group separately
- Meta-analysis can be used to combine the results from each group

41

### Rare Variant Aggregate Methods

- Best to obtain principal components to include in the regression model (including LMM and GLMM)
  - using variants which are not in LD e.g.,  $r^2 < 0.1$  (pruned)
  - covering a wide range of the allelic frequency spectrum e.g.,  $> 0.1\%$
  - Evaluate how many components need to be included
    - Don't include a fix number of components
      - e.g., 5 or 10 components
- Success of PCA\MDS in controlling for population substructure\admixture can be evaluated through lambda and examining Quantile-Quantile (QQ) plots



42

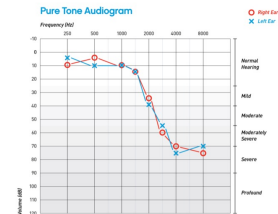
## Part II Example of a Rare Variant Association Study

### Analysis of UK Biobank Exome Data to Study the Etiology of Late-onset Hearing Loss

43

## Age-related Hearing Loss (ARHL) (aka Presbycusis)

- ARHL can impact quality of life and daily functioning
- ARHL is one of the most common adult conditions
  - In the USA
    - ARHL affects 50% of individuals >75 years of age
    - It is estimated that 30-40 million will be affected with significant ARHL by 2030



44

## Goals of the Study

- Using data from the UK Biobank to detect associations between self-reported measures of ARHL and genetic variants
  - **H-aid** self-reported hearing aid use (f.3393: “Do you use a hearing aid most of the time?”)
  - **H-diff** self-reported hearing difficulty (f.2247: “Do you have any difficulty with your hearing?”)
  - **H-noise** self-reported hearing difficulty with background noise (f.2257: “Do you find it difficult to follow a conversation if there is background noise e.g., TV, radio, children playing?”)
  - **H-both** individuals with both H-diff and H-noise
- With an emphasis of understanding the role that rare variation plays in ARHL
  - Current analysis - exome sequence data

45

## UK Biobank

- 500,000 individuals randomly sampled
  - Aged 40-69 at time of enrollment
    - To be followed for at least 20 years
    - Predominantly white Europeans
      - Also includes South Asians and individuals of African Ancestry and smaller number of individuals of a few other ancestries
- Extensive phenotype data
  - Qualitative and quantitative traits
    - ICD-10 and ICD-9 codes
    - Self reports
    - Cognitive test
    - Brain MRIs
    - NMR-metabolomics data
- Genetic Data
  - Genotype and imputed data
  - Exome sequence data
  - Whole genome sequence data
    - 200K currently available
    - Remaining sample - Quarter 1 2023
  - Telomere length data

\*Data showcase can be used to examine phenotypes and sample sizes available

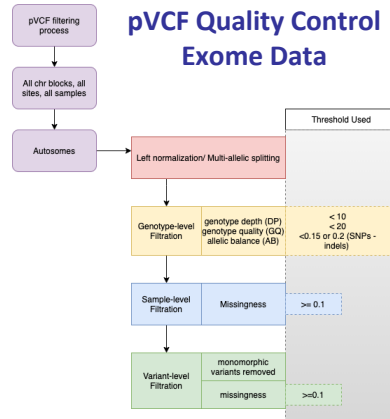
46

## Genetic Data Analyzed

- Exome data
  - ~200,000 participants
- Imputed variant data (secondary replication sample for common variants)
  - ~300,000 participants
    - Did not have exome data at the time of the study

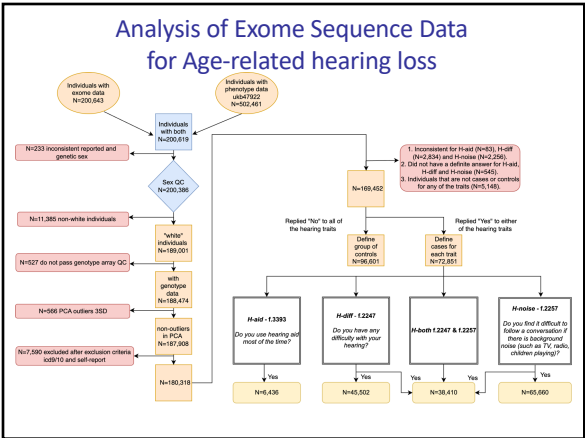
47

## pVCF Quality Control Exome Data

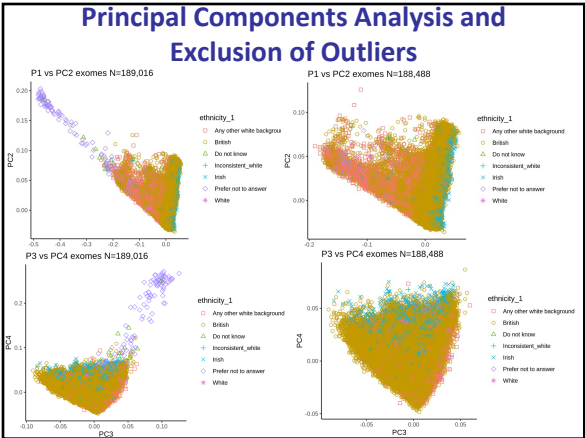


48





49



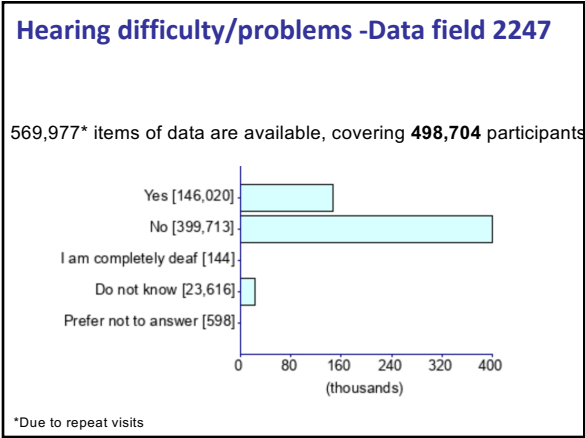
50

- ### Exclusion Criteria Obtained from ICD10, ICD9, & Self Report
- Deafness
  - Early-onset hearing impairment
  - Otosclerosis
  - Meniere's
  - Labyrinthitis
  - Disorders of acoustic nerve
  - Bell's palsy
  - History of chronic suppurative and nonsuppurative otitis media
  - Meningitis
  - Encephalitis, myelitis, and encephalomyelitis
  - Etc.

51

- ### Defining Cases and Controls
- Based on answers obtained from a touch screen
  - Cases - self-reported hearing difficulty
    - f.2247: "Do you have any difficulty with your hearing?"
  - Controls - did not have any self-reported hearing problems
    - **H-aid** hearing aid use (f.3393)
    - **H-diff** self-reported hearing difficulty (f.2247)
    - **H-noise** self-reported hearing difficulty with background noise (f.2257)

52



53

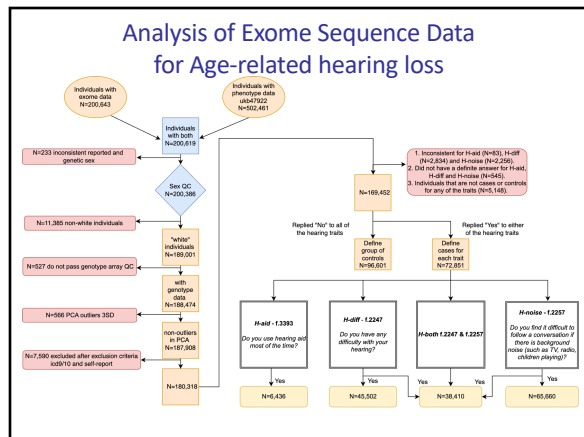
### Repeat measures\*

- Individuals with inconsistent answers removed

	Visit 1	Visit 2	Visit 3	Visit 4	
Study subject A	Problems Hearing	No Hearing Problems	No Hearing Problems	No Hearing Problems	Inconsistent Remove
Study subject B	No Hearing Problems	No Hearing Problems	Problems Hearing	Problems Hearing	Consistent (Case)
Study subject C	No Hearing Problems	No Hearing Problems	No Hearing Problems	No Hearing Problems	Consistent (Control)

\*Majority of study subjects currently have data from only one visit

54

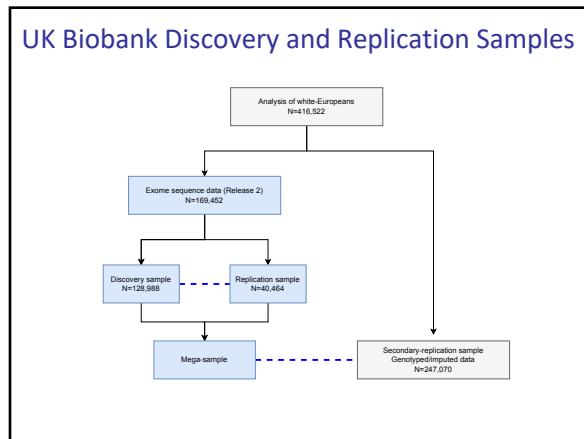


55

### Genetic Data Analyzed

- **Exome data**
  - ~200,000 participants
- **Imputed variant data (secondary replication sample for common variants)**
  - ~300,000 participants
    - Did not have exome data at the time of the study

56



57

### Analysis of Exome Data

- **Analysis performed using generalized linear mixed models (GLMM) (REGENIE)**
  - To control for inclusion of related individuals
    - For the UK Biobank data 30.3% of participants are  $\leq 3$ rd degree relatives & 4.5% sib-pairs
  - Genotype array data (~800K) were used for the ridge regression
    - Data pruned to remove variants with a  $r^2 > 0.1$ 
      - Using exome data for the ridge regression led to an inflated lambda value

QQ Plot using exome data for ridge regression

QQ Plot using genotype data for ridge regression

58

### Analysis of Exome Data

- **Analysis limited to individuals of white European Ancestry**
- **Sex, age, and two PCAs included as covariates**
  - Age for cases first report of hearing difficulty & controls age at last visit
  - The PCAs were recalculated for only individuals included in the analysis
    - Using the pruned genotypes array data ( $r^2 < 0.1$ )

59

### Analysis of Exome data – Single Variant

- **All variants with four or more alternative alleles observed in the sample analyzed**
  - A very low minor allele frequency was used since it was hypothesized some of the variants may have large effect sizes

60

## Analysis of Exome data – Single Variant

- Discovery sample
  - Second release of 150K exome
- Replication sample
  - First release of 50K exomes
- Entire exome sample (200K)
- Secondary Replication Sample\*
  - To replicate findings from the entire exome sample
  - Genotype and Imputed data (Haplotype Reference Consortium Panel)
    - 300K individuals who were not included in the exome data
    - Imputed variants with an INFO score > 0.3 were analyzed

\*Only used for replication of common variants

61

## Significance Levels

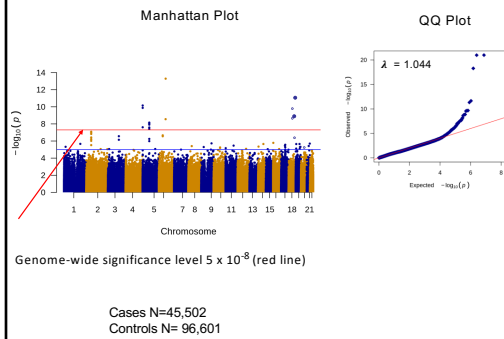
- Discovery sample
  - A genome-wide significance level was used to reject the null hypothesis of no association
    - $p \leq 5.0 \times 10^{-8}$
- Replication sample
  - Permutation was used to obtain empirical p-values
    - Adjusting for the phenotypes and variants brought to replication
      - $p \leq 0.05$

For the replication it is not necessary to use a genome-wide significance level of  $5 \times 10^{-8}$  for single variant tests or  $2.5 \times 10^{-4}$  for gene-based rare variant aggregate analysis. Significance level is adjusted for the number of variants/genes tested in the replication sample

- Bonferroni correction
- Estimate empirical p-values

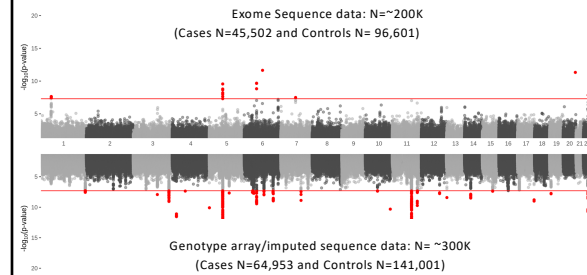
62

## Hearing Difficulty - Data Field 2247



63

## Hudson Plot Discovery and Replication Hearing Difficulty Data Field 2247



64

## Analysis of the Discovery Sample & Replication Single Variant Analysis

### Discovery sample single-variant associations analysis for age-related hearing loss traits

CHR	SNP	EA	EAF	Gene	H-aid			H-diff			H-noise			H-both		
					Beta[OR]	SE	P	Beta[OR]	SE	P	Beta[OR]	SE	P	Beta[OR]	SE	P
5	rs157688232	G	6.65e-04	PDCC6	1.997(1)	0.28	2.25e-07	1.323(7)	0.17	1.12e-07	1.042(8)	0.16	5.50e-07	1.273(4)	0.18	1.62e-07
5	rs49893074	C	5.58e-04	PDCC6	1.997(3)	0.32	1.95e-04	1.393(8)	0.18	7.05e-07	1.072(8)	0.18	6.69e-07	1.283(4)	0.19	5.52e-07
5	rs157370281	G	7.04e-04	PDCC6	1.926(8)	0.28	6.02e-07	1.393(8)	0.16	1.49e-07	1.032(8)	0.16	2.26e-07	1.293(4)	0.17	9.66e-07
6	rs1574480	C	6.09e-04	SLC22A7												
6	rs2242416	G	6.09e-04	CHRF3												
6	rs22423460	G	7.63e-04	MYO6	5.462(19.8)	3.12	1.79e-07	3.94(14.1)	0.90	8.62e-07				3.794(13)	0.90	1.76e-07

65

Mega analysis single variant associations analysis with age-related hearing impairment traits																
CHR	SNP	EA	EAF	Gene	H-aid			H-diff			H-noise			H-both		
					Beta[OR]	SE	P	Beta[OR]	SE	P	Beta[OR]	SE	P	Beta[OR]	SE	P
1	rs15809662	C	0.424	ANKK2				-0.050(0.95)	0.01	2.25e-07						
1	rs2277426	A	0.431	ANKK2				-0.050(0.95)	0.01	3.89e-07						
1	rs1707336	G	0.415	ANKK2				-0.040(0.95)	0.01	3.63e-07						
5	rs57688122	G	7e-04	PDCC6	1.796(8)	0.25	7.06e-07	1.393(8)	0.14	1.04e-07	1.113(3)	0.14	4.96e-07	1.323(8)	0.15	1.11e-07
5	rs49893074	C	6e-04	PDCC6	1.705(5)	0.28	2.48e-05	1.373(8)	0.16	5.19e-07	1.081(8)	0.15	2.19e-07	1.323(8)	0.16	6.63e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1.715(5)	0.24	1.34e-05	1.312(7)	0.14	1.00e-07	1.042(8)	0.14	1.83e-07	1.283(8)	0.15	8.00e-07
5	rs157370281	G	7e-04	PDCC6	1											

Analysis of Exome Data  
Rare Variant Aggregate Analysis

- Genes with at least two variants were analyzed, e.g., predicated loss of function (pLoF) variants
- Max coding was used
- Two masks were used
  - Mask 1 – pLoF variants
  - Mask 2 – pLoF and missense variants
- Minor allele frequency cut-off of <0.01 was used
  - The frequencies for each variant site were obtained from gnomAD non-Finnish Europeans

67

REGENIE Rare Variant Aggregate Analysis

- Three different codes can be used
  - Max
  - Sum
  - Comphet
    - This term is not correct because the phase is unknown
      - Variants may be on the same haplotype

Single variant sites	max	sum	comphet
0000000000000000 →	0	0	0
00000100010000 →	1	2	2
00201011010100 →	2	7	2

<https://raaihub.aithub.io/regenie/options/>

68

Selection of Variants to Include in Rare Variant Aggregate Association Tests

Annotation File	Mask File	AAF file
1:55039839:T:C PCSK9 LoF 1:55039842:G:A PCSK9 missense	+ Mask1 LoF Mask2 LoF,missense	+ 1:55039839:T:C 1.53e-05 1:55039842:G:A 2.19e-06
1:55039839:T:C PCSK9 CADD30 1:55039842:G:A PCSK9 CADD20	+ Mask1 CADD score > 30 Mask2 CADD score > 20	+ 1:55039839:T:C 1.53e-05 1:55039842:G:A 2.19e-06

REGENIE will use information from the annotation and alternative allele frequency (AAF) files to build the Masks (variants to be included in the association testing)

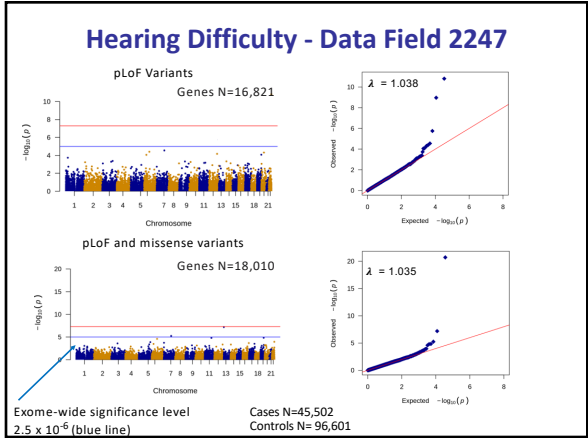
69

Analysis of Exome Data  
Rare Variant Aggregate Analysis

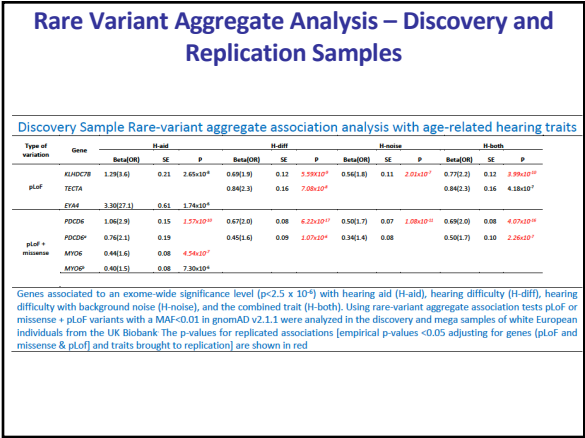
- Exome sample was split
  - Second release of 150K exome were used as the discovery sample.
  - First release of 50K exome were used as the replication sample
- Entire exome sample (200K) was also analyzed\*
- Discovery sample significance level
  - $p < 2.5 \times 10^{-6}$ 
    - 0.05/20,000 Bonferroni correction for testing 20,000 genes
- Replication sample significant level
  - $p < 0.05$
  - Empirical p-values generated
    - Permutation used to adjust for the number of phenotypes and genes brought to replication (pLoF and pLoF & missense)

\*No replication sample available for these findings

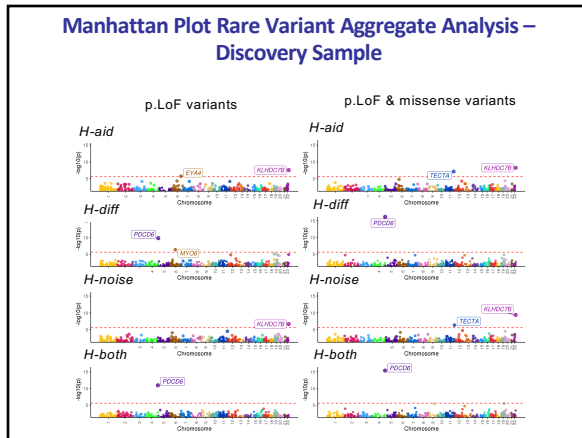
70



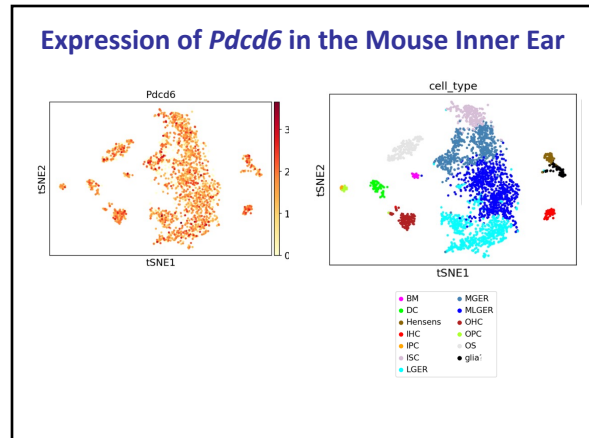
71



72



73



74

- ### Conclusions – Part II
- Replicated some previously reported ARHL genes
    - Some which had not been previously replicated
      - e.g., *BAIAP2L2*, *CRIP3*, *KLHDC7B*, *MAST2*, and *SLC22A7*
  - Identified and replicated a new HL gene, *PDCD6* which has not been previously reported
    - Inner ear expression in humans and mice supports the involvement of gene in HL etiology
    - PDCD6* is a cytoplasmic Ca<sup>2+</sup> binding protein with an important role in apoptotic cell death
  - Rare-variant aggregate analysis demonstrated the important contribution of Mendelian HL genes, i.e. *MYO6*, *TECTA*, and *EYA4* the genetics of ARHL
  - Rare variants for ARHL tend to have larger effect sizes than those for common variants
    - Rare variants should play an important role in risk prediction by increasing accuracy
  - For additional information see
    - Cornejo-Sanchez et al. (2023) Eur J Hum Genet in press PMID: 36788145

75