

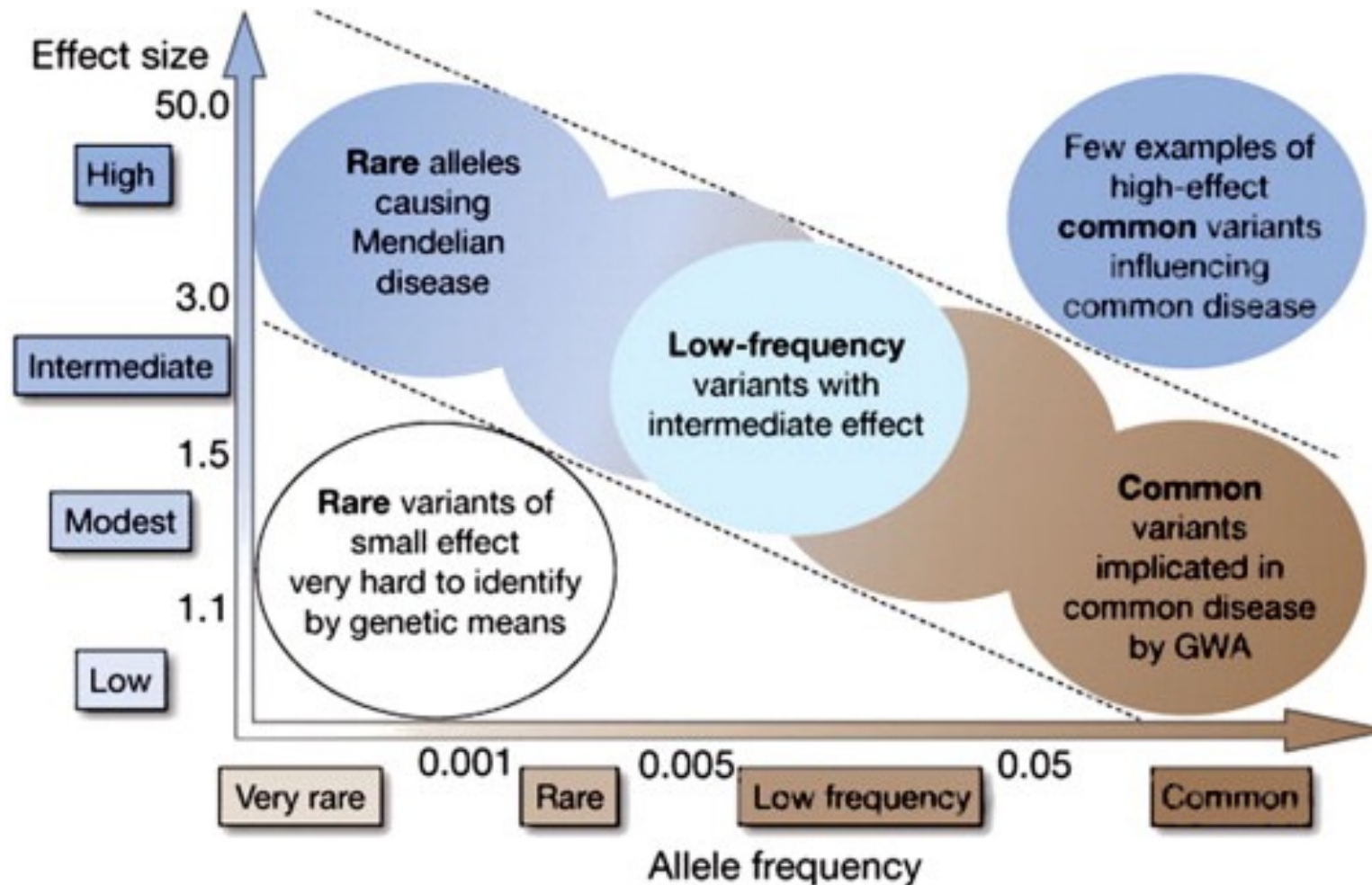
# Session 10:

# Rare variant association studies

---



# Identifying genetic variation associated with disease



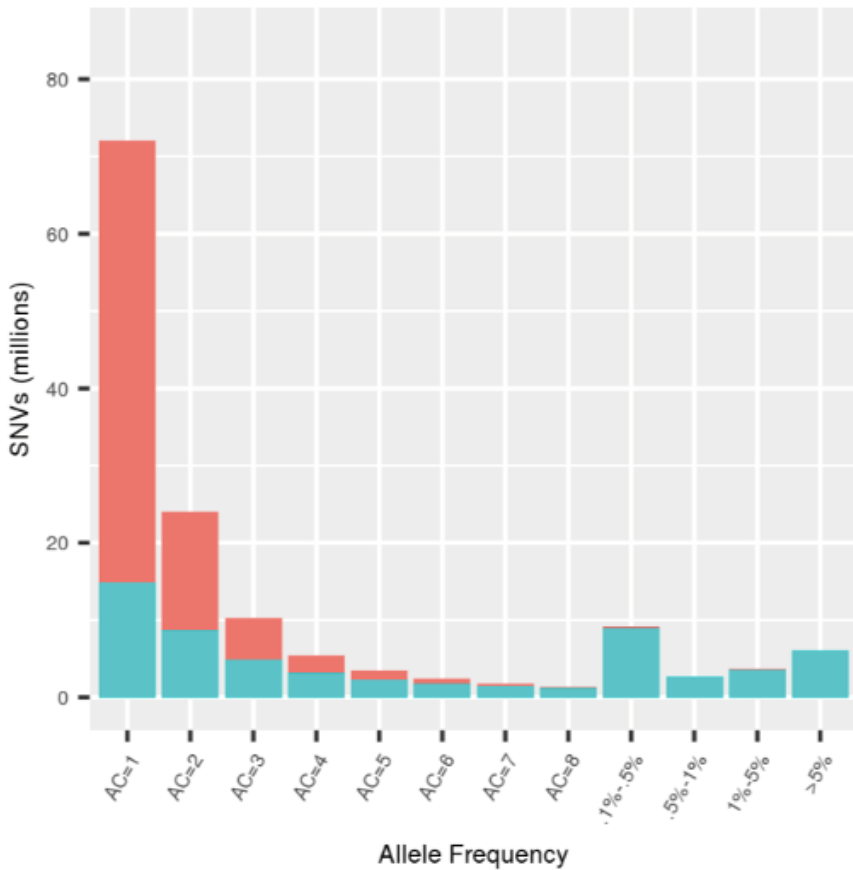
# Introduction – Rare variants

---

- > Usually less than 1% (depending on who you ask)
- > Traditional single variant association analysis have low statistical power and/or are not valid
  - MAF=0.5% in 1,000 cases and 1,000 controls implies 20 minor alleles total
  - Low cell counts lead to invalid statistical tests/low power
- > Because the genome has many more rare variants than common variants, more stringent significance levels might be required, further reducing power

# Most of the human genetic variation is rare

N=10,545 genomes, 150 million variants



N=40,722 genomes, 384 million variants

	All unrelated individuals ( <i>n</i> = 40,722)		Per individual			
	Total	Singletons (%)	Average	5th percentile	Median	95th percentile
<b>Total variants</b>	<b>384,127,954</b>	<b>203,994,740 (53)</b>	<b>3,748,599</b>	<b>3,516,166</b>	<b>3,563,978</b>	<b>4,359,661</b>
SNVs	357,043,141	189,429,596 (53)	3,553,423	3,335,442	3,380,462	4,125,740
Indels	27,084,813	14,565,144 (54)	195,176	180,616	183,503	233,928
<b>Novel variants</b>	<b>298,373,330</b>	<b>191,557,469 (64)</b>	<b>29,202</b>	<b>20,312</b>	<b>24,106</b>	<b>44,336</b>
SNVs	275,141,134	177,410,620 (64)	25,027	17,520	20,975	36,861
Indels	23,232,196	14,146,849 (61)	4,175	2,747	3,145	7,359
<b>Coding variation</b>	<b>4,651,453</b>	<b>2,523,257 (54)</b>	<b>23,909</b>	<b>22,158</b>	<b>22,557</b>	<b>27,716</b>
Synonymous	1,435,058	715,254 (50)	11,651	10,841	11,056	13,678
Nonsynonymous	2,965,093	1,648,672 (56)	11,384	10,632	10,856	13,221
Stop/essential splice	97,217	60,347 (62)	474	425	454	566
Frameshift	104,704	71,577 (68)	132	112	127	165
In-frame	51,997	29,110 (56)	102	85	99	128

# Poll: Why study rare variants?

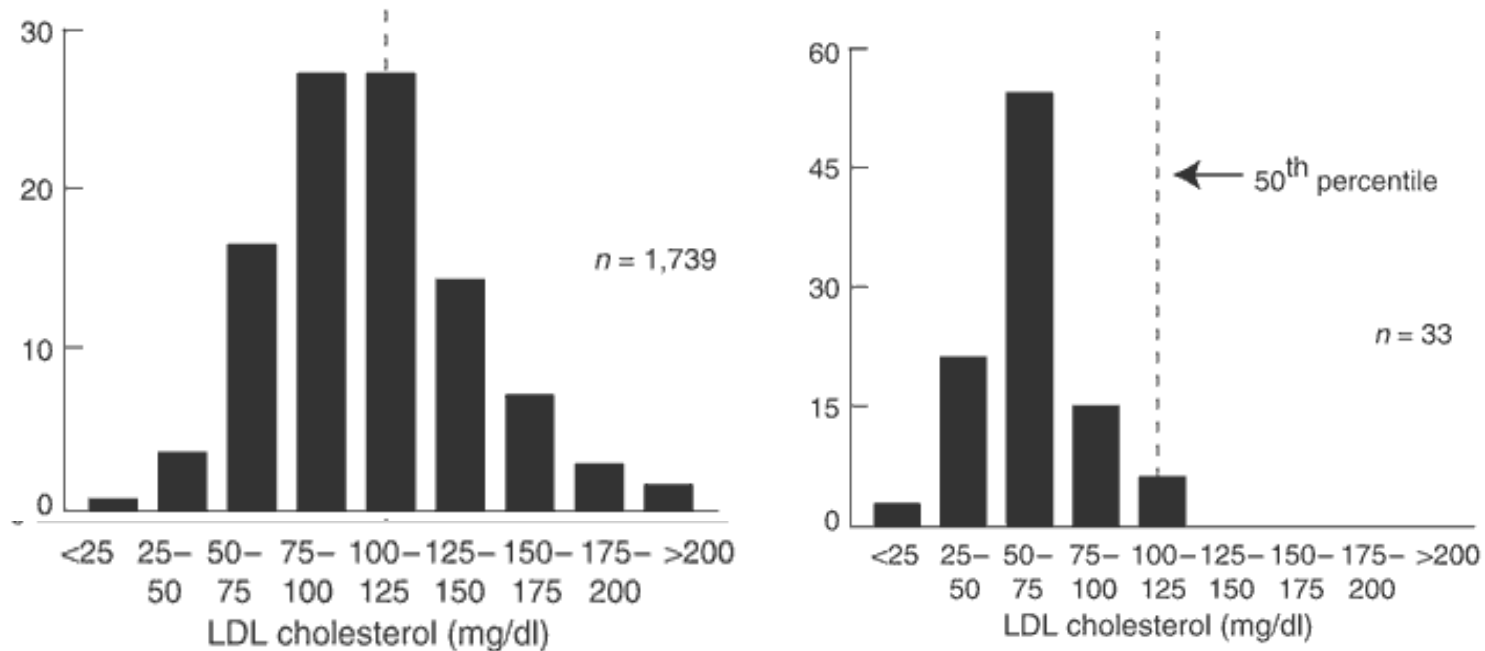
---



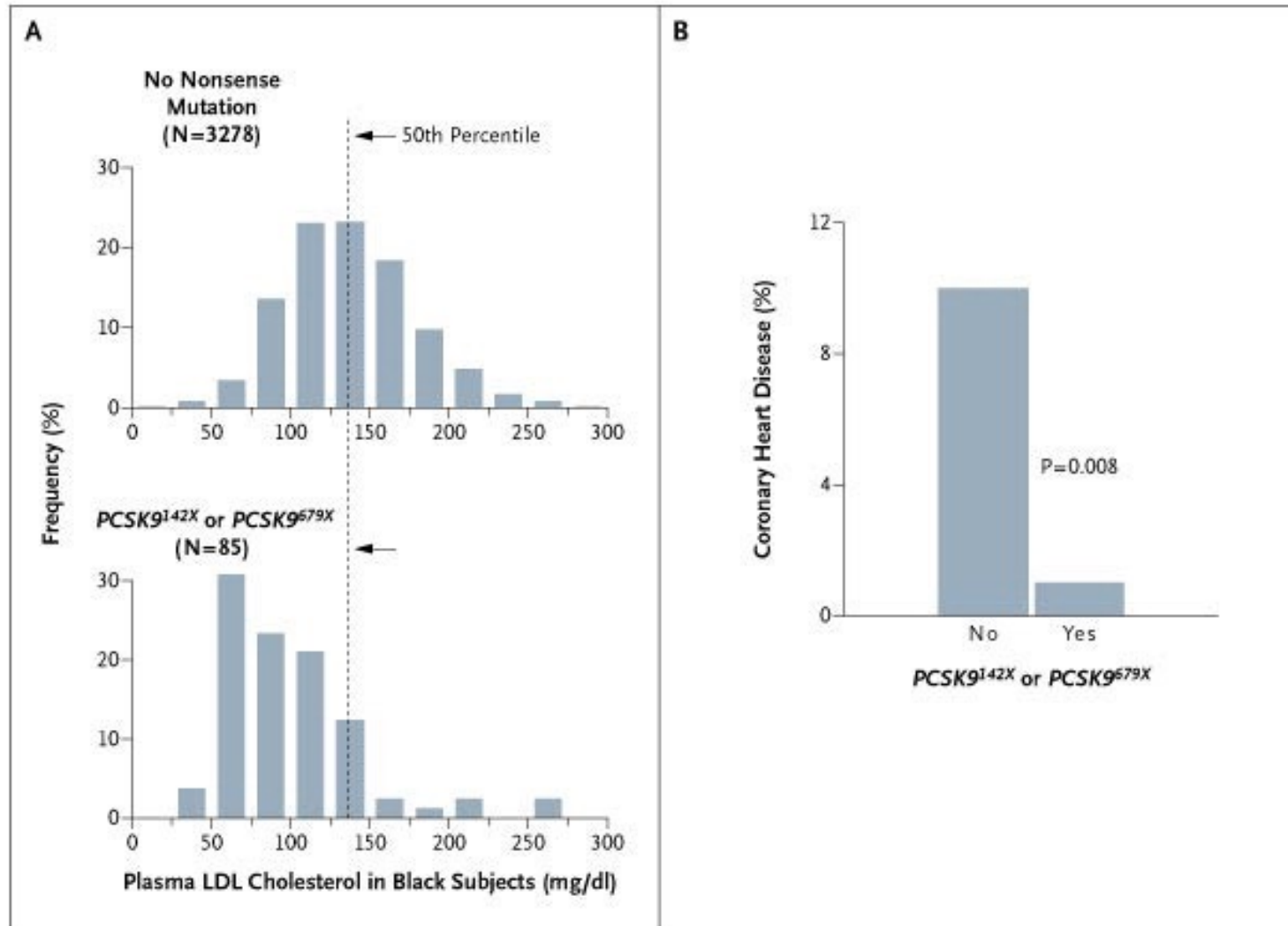
# Why do we care about rare variants when they only affect a small proportion of the population?

## *PCSK9* and LDL cholesterol

Plasma LDL-C levels in African American individuals without (left) and with (right) a nonsense mutation in *PCSK9*.

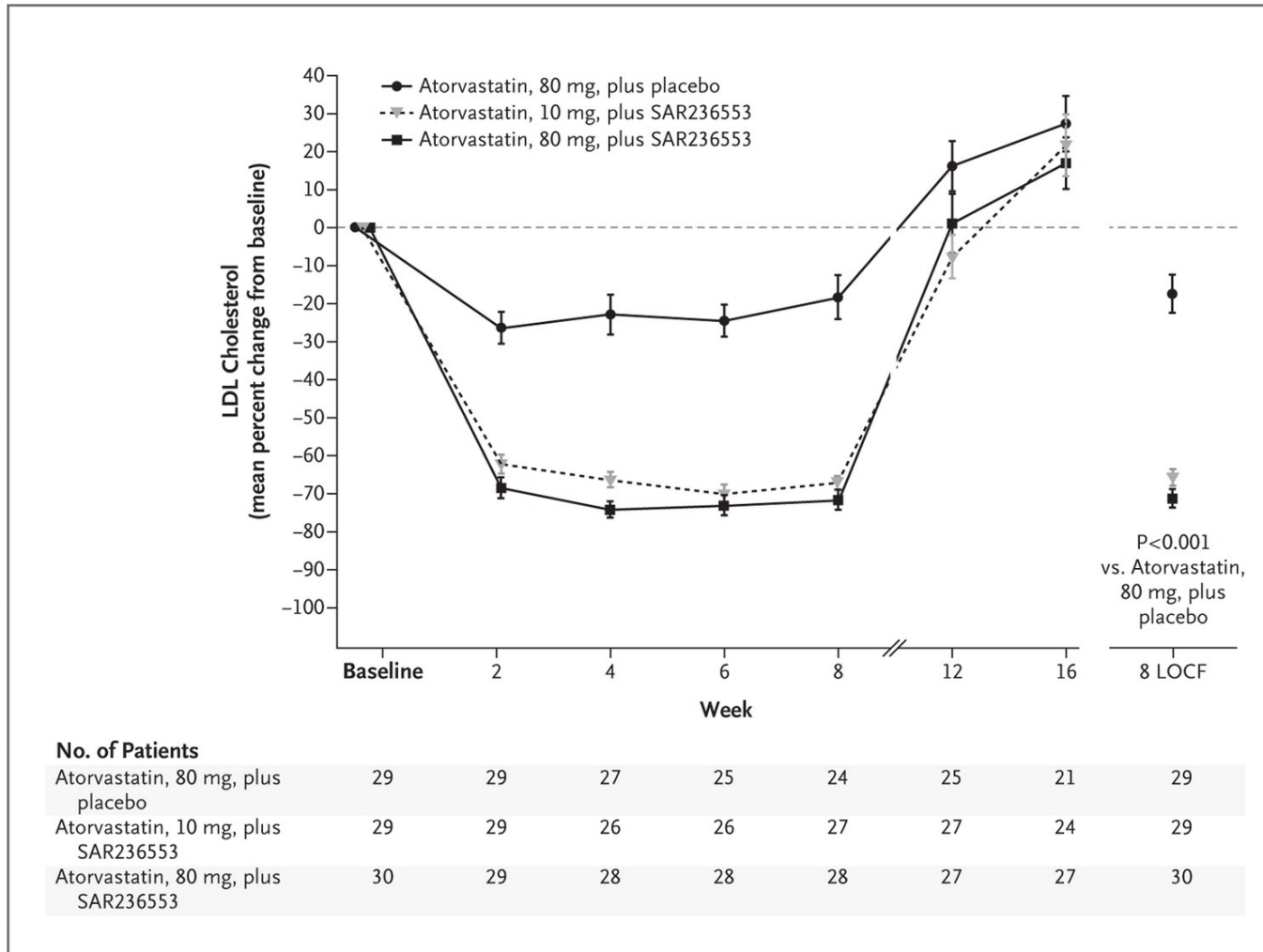


# PCSK9 mutations and coronary heart disease



Cohen, NEJM 2005

# A PCSK9 antibody decreases LDL (8-week trial)





# Study design for rare variant analysis

---

	Advantage	Disadvantage
<b>High-depth WGS</b>	can identify nearly all variants in the genome with high confidence	very expensive
<b>Low-depth WGS</b>	cost-effective and useful approach for association mapping	has limited accuracy for rare-variant identification and genotype calling; compared to deep sequencing, is subject to power loss if the same number of subjects is sequenced
<b>Whole-exome sequencing</b>	can identify all exonic variants; is less expensive than WGS	is limited to the exome
<b>GWAS chip and imputation</b>	inexpensive	has lower accuracy for imputed rare variants; will miss any variants unique to your sample
<b>Exome chip (custom array)</b>	much cheaper than exome sequencing	provides limited coverage for very rare variants and for non-European populations; is limited to target regions

# Breakout room discussion

---

- > If you were to design a study to identify rare (allele frequency <1%) variants associated with ovarian cancer, what approach would you take and why?
- High-depth whole genome sequencing
  - Low-depth whole genome sequencing
  - Whole exome sequencing
  - GWAS chip and imputation
  - Exome chip (custom array)

	Advantage	Disadvantage
<b>High-depth WGS</b>	can identify nearly all variants in the genome with high confidence	very expensive
<b>Low-depth WGS</b>	cost-effective and useful approach for association mapping	has limited accuracy for rare-variant identification and genotype calling; compared to deep sequencing, is subject to power loss if the same number of subjects is sequenced
<b>Whole-exome sequencing</b>	can identify all exonic variants; is less expensive than WGS	is limited to the exome
<b>GWAS chip and imputation</b>	inexpensive	has lower accuracy for imputed rare variants; will miss any variants unique to your sample
<b>Exome chip (custom array)</b>	much cheaper than exome sequencing	provides limited coverage for very rare variants and for non-European populations; is limited to target regions

# Analyses of rare variants

---

- > Many different rare variant tests are available, but most fall into one of two major categories
  - Some are based on aggregating variants (“burden” tests)
    - > CMC (Li and Leal, 2008)
    - > WSS (Madsen and Browning, 2009)
    - > Variable Threshold approach (Price, 2010)
  - Some are based on studying the distribution of variants
    - > C-alpha (Neale, 2011)
    - > SKAT (Wu, 2011)

# Burden tests

---

- > Collapse many variants into a single risk score
  - Combine minor allele counts into one variable
- > Collapsing approach
  - Gene, pathways, functional annotations, etc
  - Much more straight-forward for coding regions
- > Weighing
  - Variant type (predicted function)
  - Variant frequency

# The Cohort Allelic Sums Test - CAST

---

Main Idea: Combine rare variants according to some (arbitrary) feature (gene, genetic region, functional category) and assess the new variable

Step 1: Create an indicator variable  $X$  for individual  $j$ :

$$X_j = \begin{cases} 1 & \text{if rare variants are present} \\ 0 & \text{otherwise} \end{cases}$$

Step 2:  $y = \alpha + \beta X$  (logistic/linear regression)

# Variant Collapsing – 2 approaches

i)

Subject	V1	V2	V3	V4	X
1	1	0	0	0	1
2	0	1	0	0	1
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0
6	0	0	0	0	0
7	0	0	1	1	1
8	0	0	0	1	1

ii)

Subject	V1	V2	V3	V4	X
1	1	0	0	0	1
2	0	1	0	0	1
3	0	0	0	0	0
4	0	0	0	0	0
5	0	0	0	0	0
6	0	0	0	0	0
7	0	0	1	1	2
8	0	0	0	1	1

# The Weighted Sum Statistic (WSS) – often called “Madsen-Browning”

---

- > Main idea: Variants are grouped according to function (e.g., gene), and each individual is scored by a weighted sum of the variant counts.
- > Use permutation to test for an excess of variants in affected individuals.
- > Variants of all frequencies can be included, but variants are weighted according to their frequency in unaffected individuals.

$$\hat{w}_i = 1/\sqrt{q_i(1 - q_i)}$$

$q_i$  is the estimated MAF in controls

# Disadvantages of burden tests

- > Burden tests assume that all variants in a set are causal and associated with a trait in the same direction. If this is not true, power is lost.
- > Solution: Tests that utilize the distribution of rare variants

*APOB* variant counts in individuals with high/low triglyceride levels.

Position	Annotation	High Lipid Level	Low Lipid Level
21078358	Ala4481Thr	2	5
21078359	Ile4314Val	3	0
21078990	Arg4270Thr	6	3
21079417	Val4128Met	1	7
21083082	Thr3388Lys	2	1
21083637	Ser3203Tyr	6	0
21086035	Leu2404Ile	2	3
21086072	Glu2391Asp	2	2
21086127	Thr2373Asn	2	2
21086308	Val2313Ile	2	1
21087477	His1923Arg	6	12
21087504	Asn1914Ser	0	5
21087634	Asp1871Asn	2	0
21091828	Pro1143Ser	0	6
21091872	Arg1128His	0	3
21091918	Asp1113His	1	3
21106140	Thr498Asn	2	0
<b>Singletons</b>		6	4



# SKAT: sequence kernel association test

---

- > In contrast to the C-alpha test, SKAT is regression-based and thereby allows for adjustment of covariates.
- > Uses a variance-component score test in a mixed-model framework to assess regression coefficients for rare variants.

$$\text{logit } P(y_i = 1) = \alpha_0 + \alpha' X_i + \beta' G_i$$

$y_i$ : case-control status;  $\alpha_0$ : intercept;  $\alpha = [\alpha_1, \dots, \alpha_m]'$  is the vector of regression coefficients for the  $m$  covariates;  $X_i$ : fixed effects of covariates;  $\beta = [\beta_1, \dots, \beta_p]'$  is the vector of regression coefficients for the  $p$  observed gene variants in the region;  $G_i$ :  $(G_{i1}, G_{i2}, \dots, G_{ip})$  genotypes for the  $p$  variants within the region

$$H_0: \beta = \mathbf{0} \text{ or } \beta_1 = \beta_2 = \dots = \beta_p = 0$$

# Combined test: SKAT-O

---

- > Picks the best combination of SKAT and a burden test, and then corrects for the flexibility afforded by this choice.
- > Specifically, if the SKAT statistic is  $Q_1$ , and the squared score for a burden test is  $Q_2$ , SKAT-O considers tests of the form

$$(1-\rho) \times Q_1 + \rho \times Q_2, \text{ where } \rho \text{ is between } 0 \text{ and } 1$$

- >  $\rho$  is selected to maximize the power of the test for each variant set
- > When  $\rho = 1$ , SKAT-O is a burden test
- > When  $\rho = 0$ , SKAT-O is a SKAT test
- > When  $0 < \rho < 1$ , SKAT-O is a linear combination of a burden and SKAT test

**Table 2. Summary of Statistical Methods for Rare-Variant Association Testing**

	<b>Description</b>	<b>Methods</b>	<b>Advantage</b>	<b>Disadvantage</b>	<b>Software Packages<sup>a</sup></b>
Burden tests	collapse rare variants into genetic scores	ARIEL test, <sup>50</sup> CAST, <sup>51</sup> CMC method, <sup>52</sup> MZ test, <sup>53</sup> WSS <sup>54</sup>	are powerful when a large proportion of variants are causal and effects are in the same direction	lose power in the presence of both trait-increasing and trait-decreasing variants or a small fraction of causal variants	EPACTS, GRANVIL, PLINK/SEQ, Rvtests, SCORE-Seq, SKAT, VAT
Adaptive burden tests	use data-adaptive weights or thresholds	aSum, <sup>55</sup> Step-up, <sup>56</sup> EREC test, <sup>57</sup> VT, <sup>58</sup> KBAC method, <sup>59</sup> RBT <sup>60</sup>	are more robust than burden tests using fixed weights or thresholds; some tests can improve result interpretation	are often computationally intensive; VT requires the same assumptions as burden tests	EPACTS, KBAC, PLINK/SEQ, Rvtests, SCORE-Seq, VAT
Variance-component tests	test variance of genetic effects	SKAT, <sup>61</sup> SSU test, <sup>62</sup> C-alpha test <sup>63</sup>	are powerful in the presence of both trait-increasing and trait-decreasing variants or a small fraction of causal variants	are less powerful than burden tests when most variants are causal and effects are in the same direction	EPACTS, PLINK/SEQ, SCORE-Seq, SKAT, VAT
Combined tests	combine burden and variance-component tests	SKAT-O, <sup>64</sup> Fisher method, <sup>65</sup> MiST <sup>66</sup>	are more robust with respect to the percentage of causal variants and the presence of both trait-increasing and trait-decreasing variants	can be slightly less powerful than burden or variance-component tests if their assumptions are largely held; some methods (e.g., the Fisher method) are computationally intensive	EPACTS, PLINK/SEQ, MiST, SKAT
EC test	exponentially combines score statistics	EC test <sup>67</sup>	is powerful when a very small proportion of variants are causal	is computationally intensive; is less powerful when a moderate or large proportion of variants are causal	no software is available yet

# Rare variant analyses software

---

## > Rvtests

- <http://zhanxw.github.io/rvtests/>

## > SKAT

- <https://cran.r-project.org/web/packages/SKAT/index.html>

## > SAIGE-GENE

- <https://github.com/weizhouUMICH/SAIGE>

# Issues in rare variant analysis (i)

---

## > Which variants do we include?

1. All variants
  - Most variants likely have no effect on our outcome
2. Only those we think are deleterious
  - How do we determine/predict deleteriousness?
  - What if we get rid of some variants that have effects on our outcome?

## > How should we group variants?

- Rare variants are often grouped by their functional unit such as by gene. This makes variant grouping straight-forward in exome studies
- For whole-genome analysis, alternative approaches such as sliding window or additional functional annotations (conserved regions, regulatory regions etc.) can be used.

# Issues in rare variant analysis (ii)

---

## > Which association test to use?

- If there are multiple variants with risk-increasing effects, burden tests are most powerful
- If there is a mixture of risk increasing and risk decreasing variants and/or most variants do not have an effect, variance-component methods are most powerful
- If no prior information is available, we can conduct both burden and variance component tests. We could also conduct combined tests like SKAT-O. We still to consider multiple testing.

## Issues in rare variant analysis (iii)

---

- > In general, rare variants are more difficult to impute
- > Adjusting for population stratification and cryptic relatedness may be more critical and more complicated for rare variant analyses
- > Rare variants tend to be more recent mutational events and tend to be more geographically localized than common variants