



The Elizabeth H.
and James S. McDonnell III

**McDONNELL
GENOME INSTITUTE**
at Washington University

You've got the data...now what?

Karyn Meltz Steinberg
November 17, 2016
HLA course

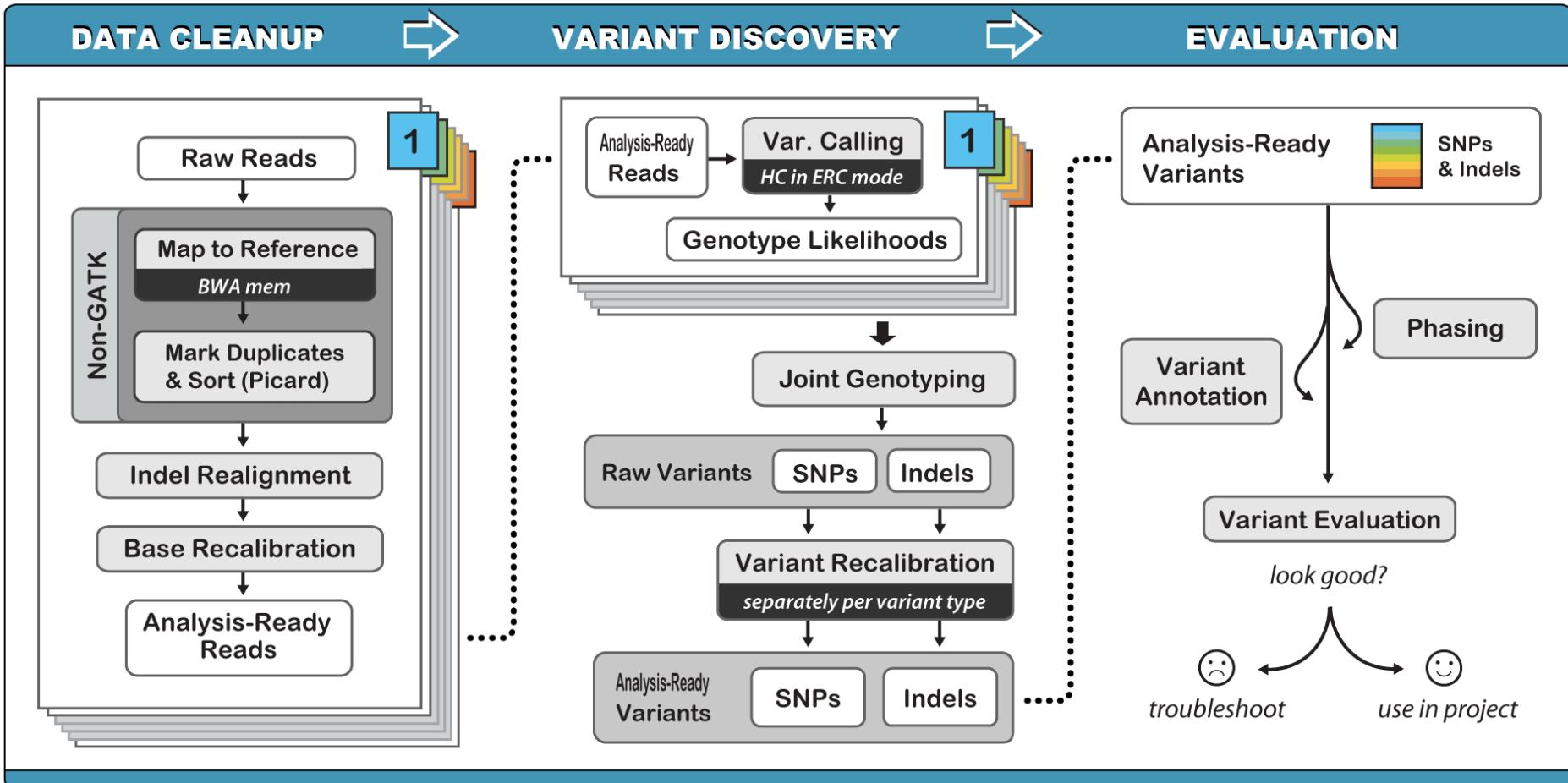
Outline

- Variant calling and filtering
- Variant Evaluation
- Annotation and GEMINI
- Identifying rare variants in large cohorts



Variant calling and filtering



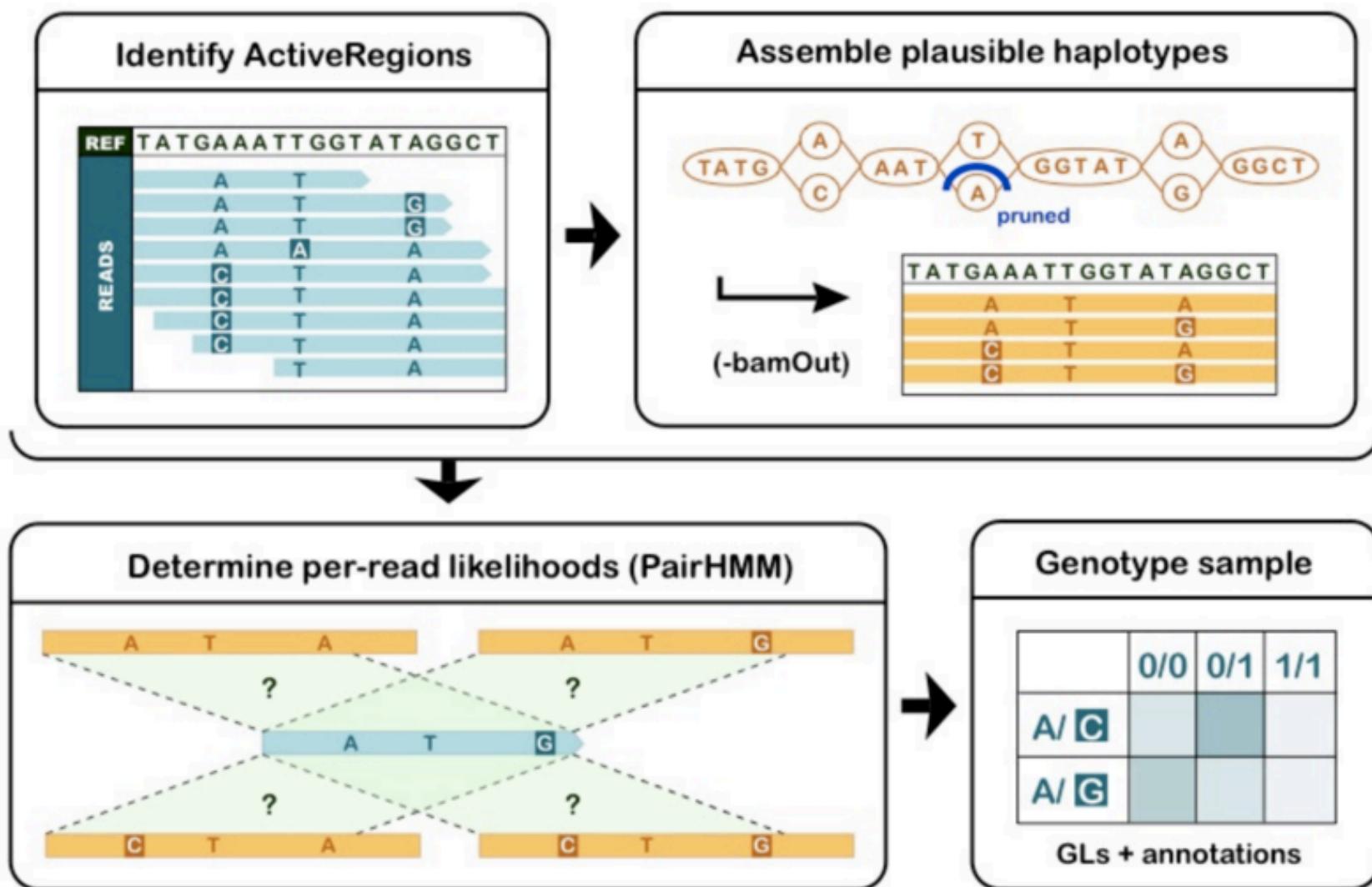


Different methods for calling variants

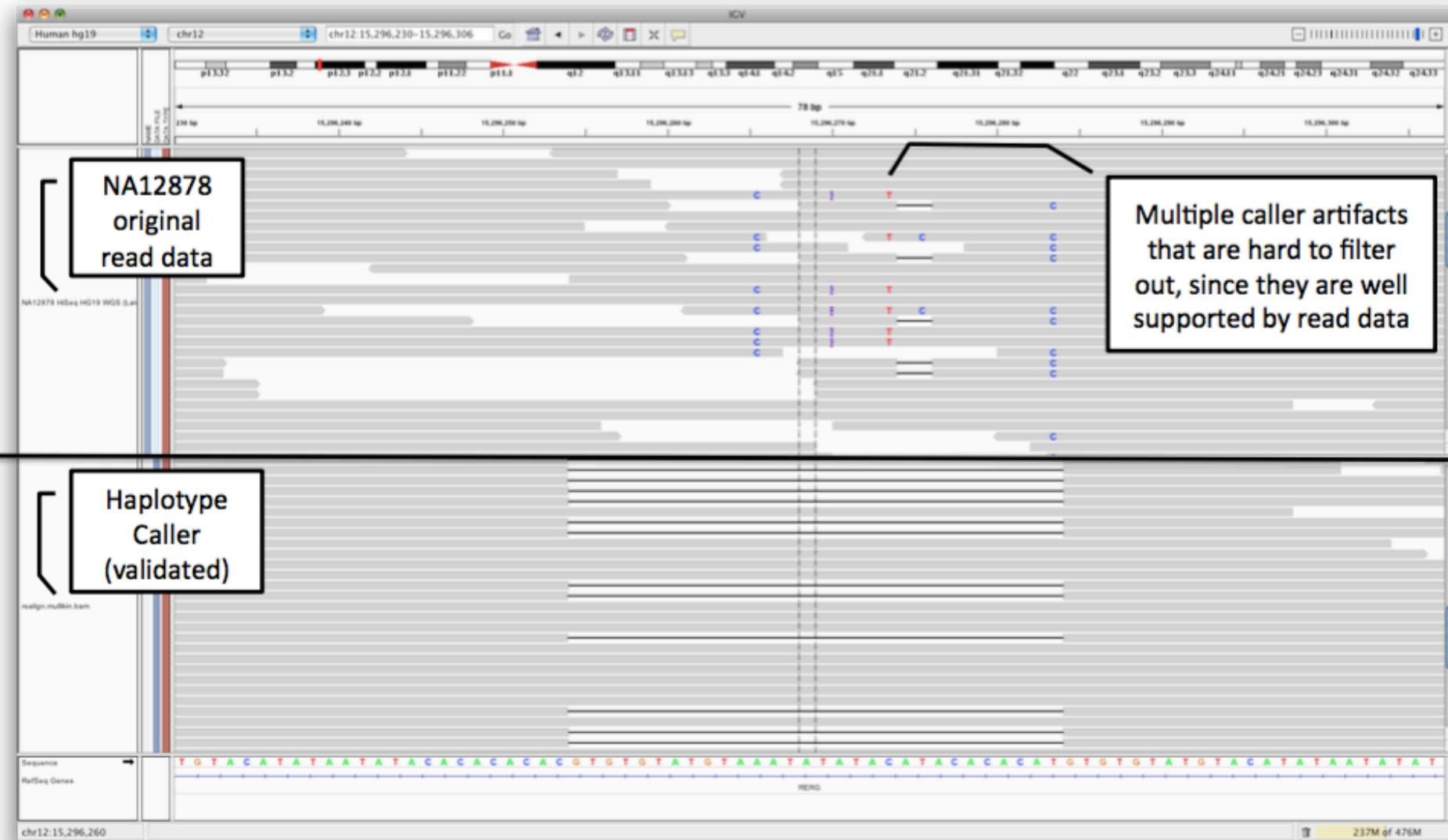
- Call SNVs and indels separately by considering each variant locus
 - Very fast
 - Assumes bases are independent
- Call SNVs and indels simultaneously via Bayesian genotype likelihood model
 - More computationally intensive
 - GATK UnifiedGenotyper
- Call SNVs, indels and SVs simultaneously by performing a local de novo assembly
 - More computationally intensive
 - More accurate—gets rid of many false positives especially indels
 - GATK HaplotypeCaller



HC method illustrated



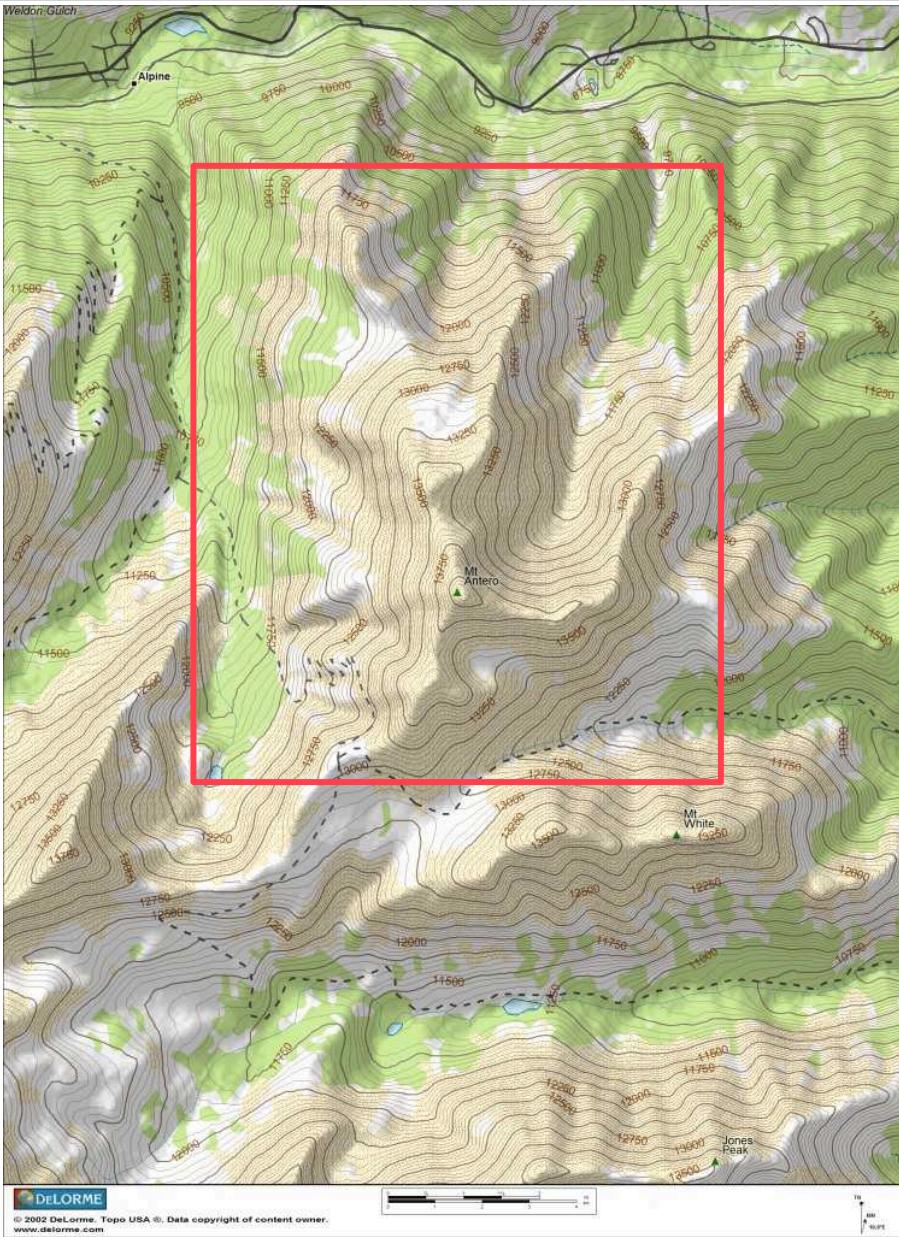
Why all the steps?



GATK recommended filters

- SNPs
 - QD < 2.0 (variant quality/depth of non-ref samples)
 - MQ < 40.0 (Mapping quality)
 - FS > 60.0 (Phred score Fisher's test pvalue for strand bias)
 - SOR > 3.0 (Strand odds ratio, aims to evaluate whether there is strand bias in the data—updated form of FET)
 - MQRankSum < -12.5 (mapping quality of reference reads vs alt reads)
 - ReadPosRankSum < -8.0 (distance of alt reads from end of the read)
- Indels
 - QD < 2.0
 - ReadPosRankSum < -20.0
 - InbreedingCoeff < -0.8
 - FS > 200.0
 - SOR > 10.0



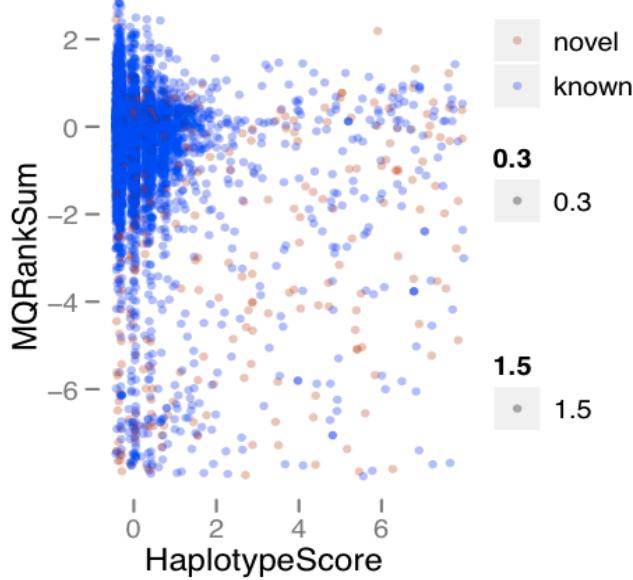
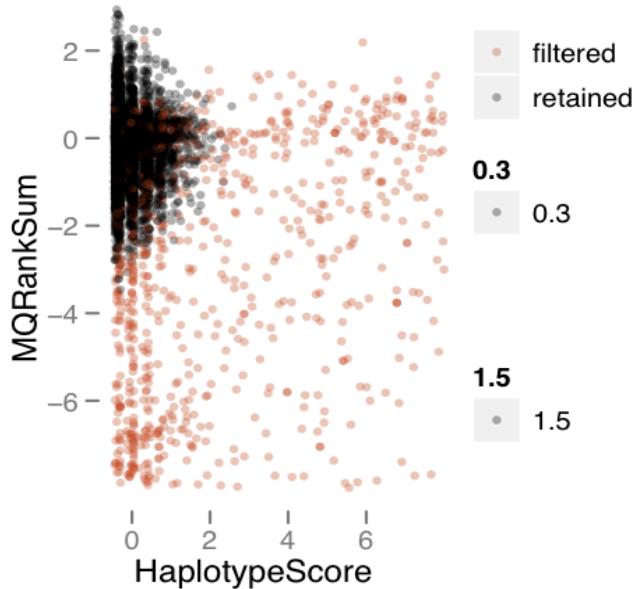
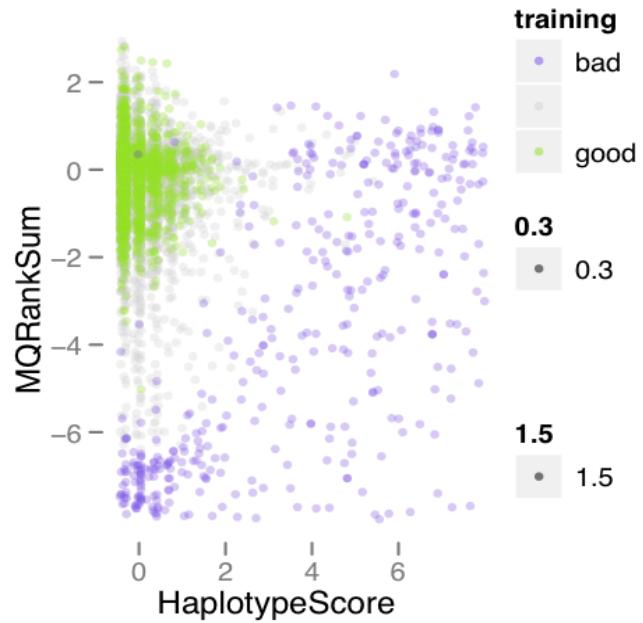
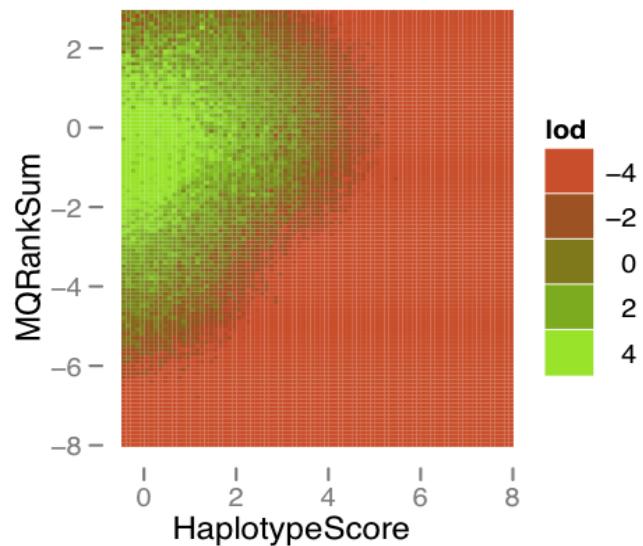


VQSR

- VariantRecalibrator
 - Gaussian mixture model by looking at the annotation values over a high quality subset of input call set and then evaluate all variants
- ApplyRecalibration
 - Apply model parameters to each variant producing a recalibrated VCF file
- Run these separately for SNPs and indels
- Need to have some resource files
 - resource:hapmap,known=false,training=true,truth=true,prior=15.0 HAPMAP
 - resource:omni,known=false,training=true,truth=true,prior=12.0 OMNI
 - resource:1000G,known=false,training=true,truth=false,prior=10.0 G1000
 - resource:dbsnp,known=true,training=false,truth=false,prior=2.0 DBSNP

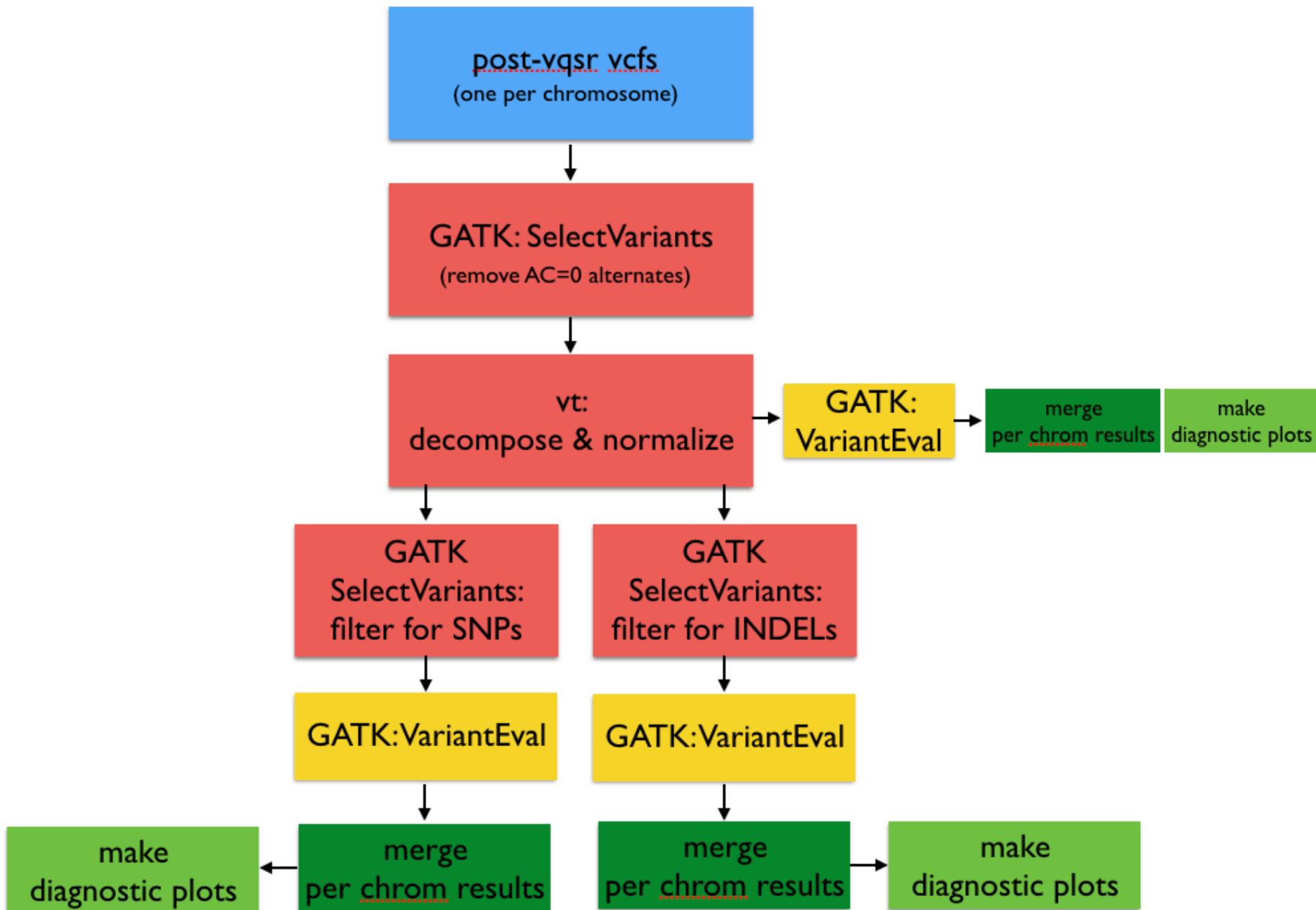


model PDF



Variant Evaluation





Decomposing and normalizing indels

Variant:	Reference Sequence	Alternate Sequence	GGGCACACACAGGG	GGGCACACACAGGG	
	Genome Reference		Variant Call Format		
	GGGCACACACAGGG		POS	REF	ALT
(A)	REF	CAC	6	CAC	C
	ALT	C			
(B)	REF	GCACA	3	GCACA	GCA
	ALT	GCA			
(C)	REF	GGCA	2	GGCA	GG
	ALT	GG			
(D)	REF	GCA	3	GCA	G
	ALT	G			



Normalize indels

Algorithm 1 Normalize a VCF entry

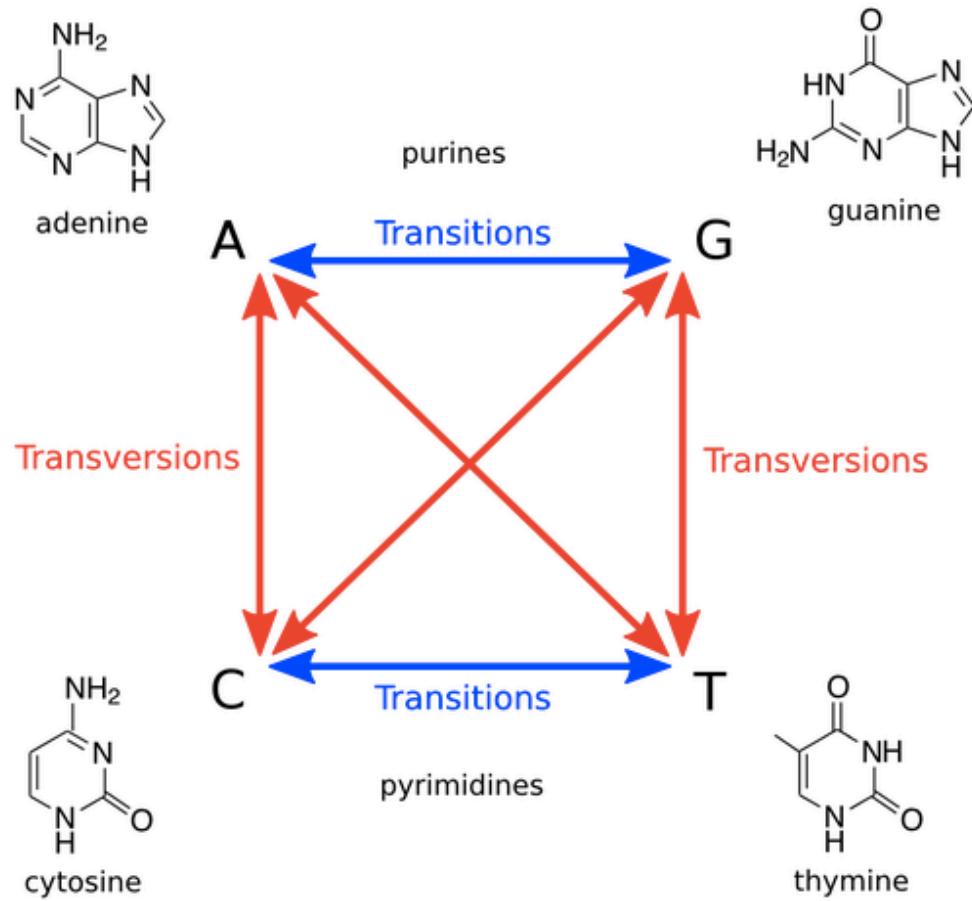
Input: A VCF entry and the reference genome sequence.

Output: A normalized VCF entry

- 1 : do
 - 2 : if all alleles end with same nucleotide then
 - 3 : truncate the rightmost nucleotide of each allele
 - 4 : if any allele is length zero then
 - 5 : extend all alleles by 1 nucleotide to the left
 - 6 : while changes made in the VCF entry in the loop
 - 7 : while all alleles start with same nucleotide and length ≥ 2
 - do
 - 8 : truncate the leftmost nucleotide of each allele
 - 9 : end while
 - 10 : return the VCF entry
-



Transition/transversion ratio (Ti/Tv)



Random = 0.5

WGS = 2.0-2.1

Exome = 3-3.5



Hardy Weinberg Equilibrium

$$(p + q)^2 = p^2 + 2pq + q^2 = 1$$

Where:

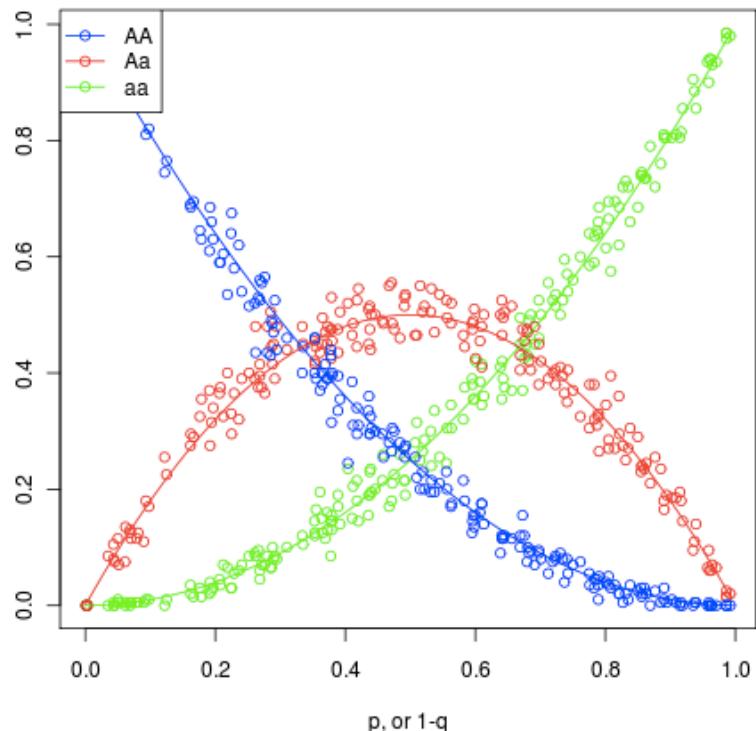
p = the frequency of allele A

q = the frequency of allele a

p^2 = the frequency of individual AA

q^2 = the frequency of individual aa

$2pq$ = the frequency of individual Aa

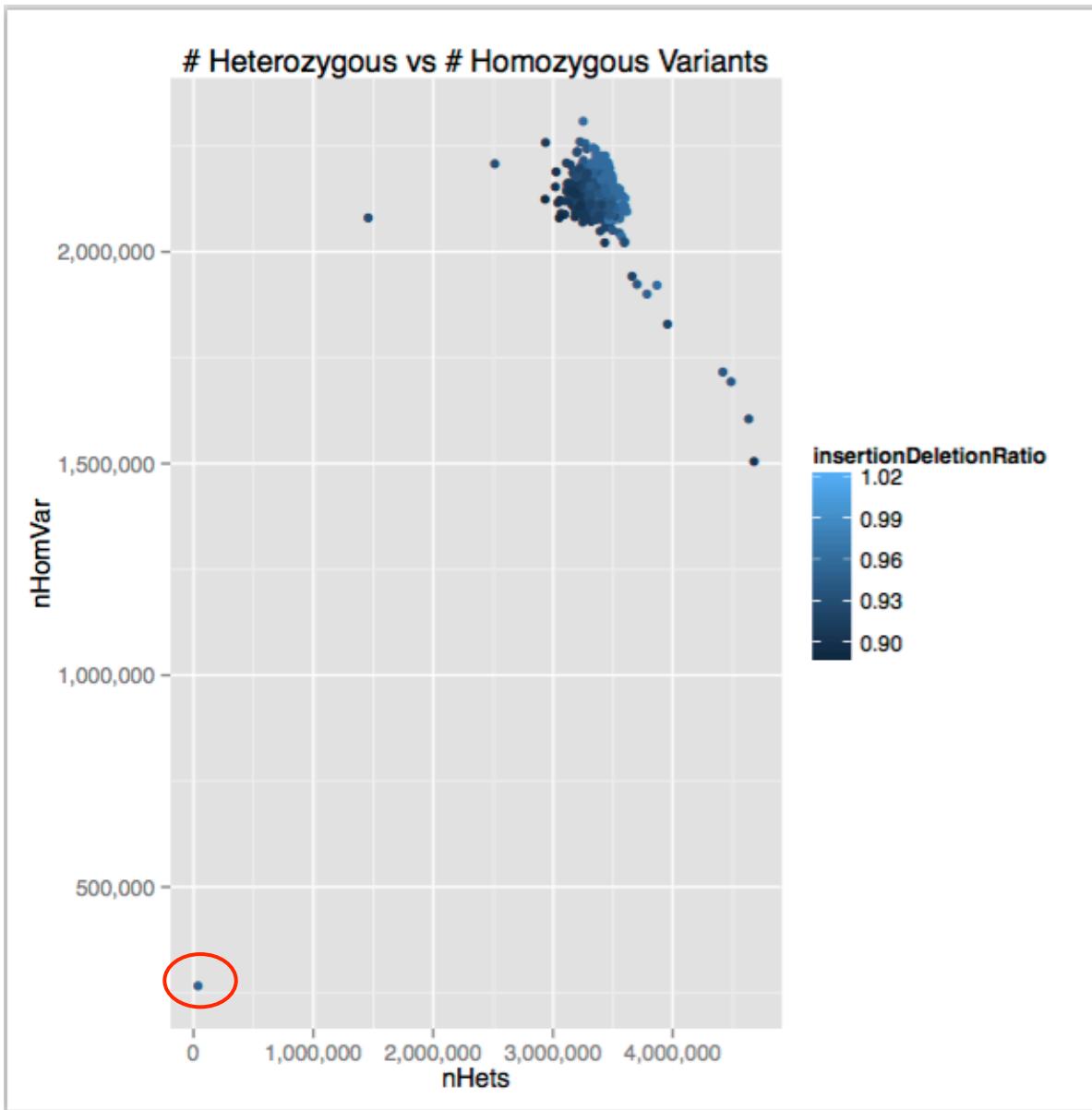


Assumptions of Hardy Weinberg Equilibrium

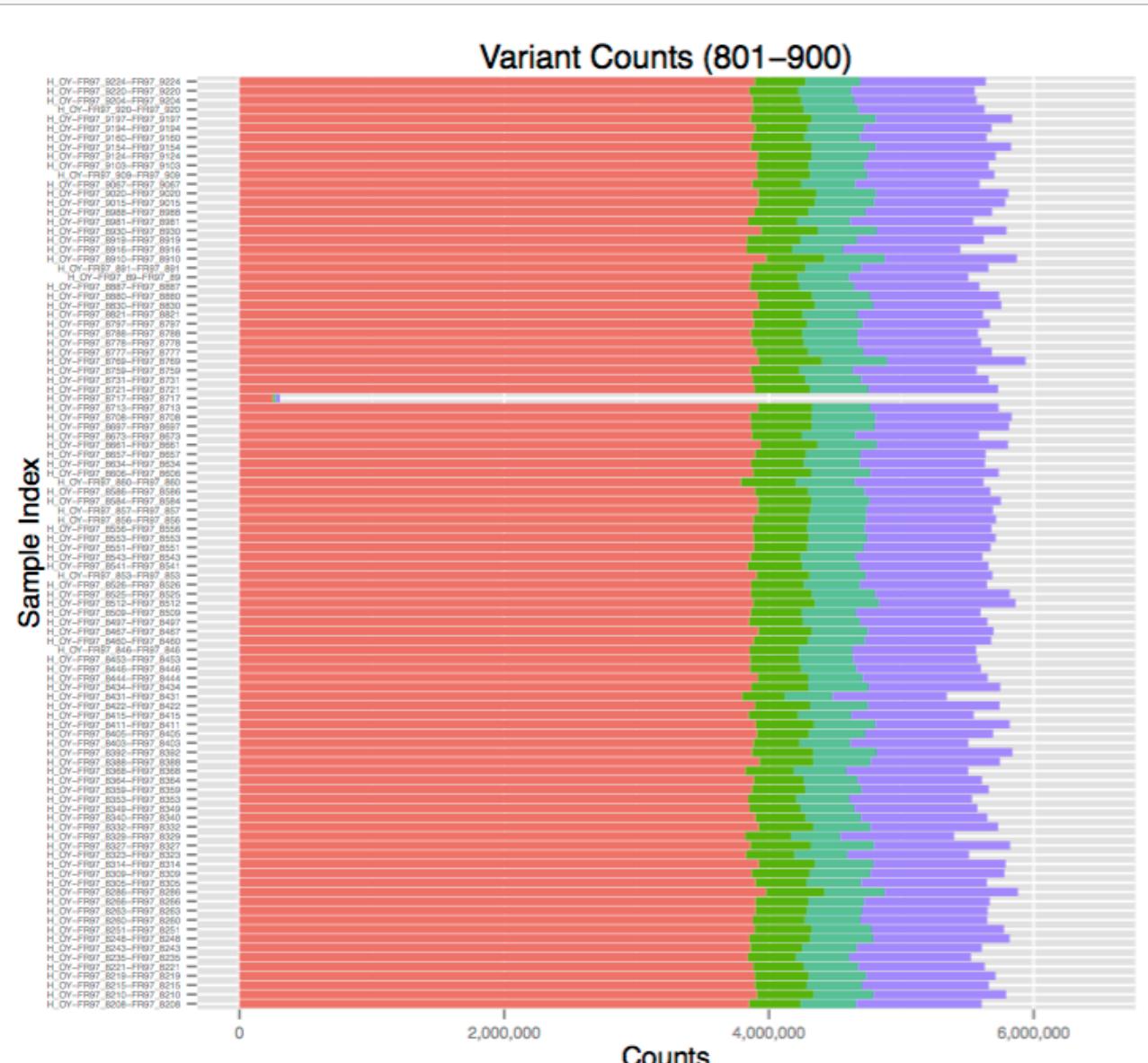
- organisms are diploid
- only sexual reproduction occurs
- generations are non overlapping
- mating is random
- population size is infinitely large
- allele frequencies are equal in the sexes
- there is no migration, mutation or selection



Heterozygous and homozygous variants



Missingness per sample and variant site



CountType nSNPs nMNPs nInsertions nDeletions nComplex nSymbolic nMixed



Other things to evaluate

- Sex check
- Novel vs Known sites
 - dbSNP/1KG concordance (should be >98%)
- Number of singletons per sample
 - excess singletons may indicate a bad sample
- PCA
 - Is there significant population stratification?
 - Do you expect this?
 - What could be causing? Batch effects?



Annotation and GEMINI



Annotation

- Based on conservation
 - SIFT
 - PolyPhen
 - GERP
 - Grantham
 - phyloP
 - CADD
 - ...
- Based on protein function
 - SNPEff
 - VEP
 - ...



SIFT

- Sorting Intolerant From Tolerant
- SIFT predicts whether an amino acid substitution affects protein function
- Given a protein sequence, SIFT chooses related proteins and obtains an alignment of these proteins with the query. Based on the amino acids appearing at each position in the alignment, SIFT calculates the probability that an amino acid at a position is tolerated conditional on the most frequent amino acid being tolerated. If this normalized value is less than a cutoff, the substitution is predicted to be deleterious
 - Scores range from 0-1
 - <0.05 → deleterious

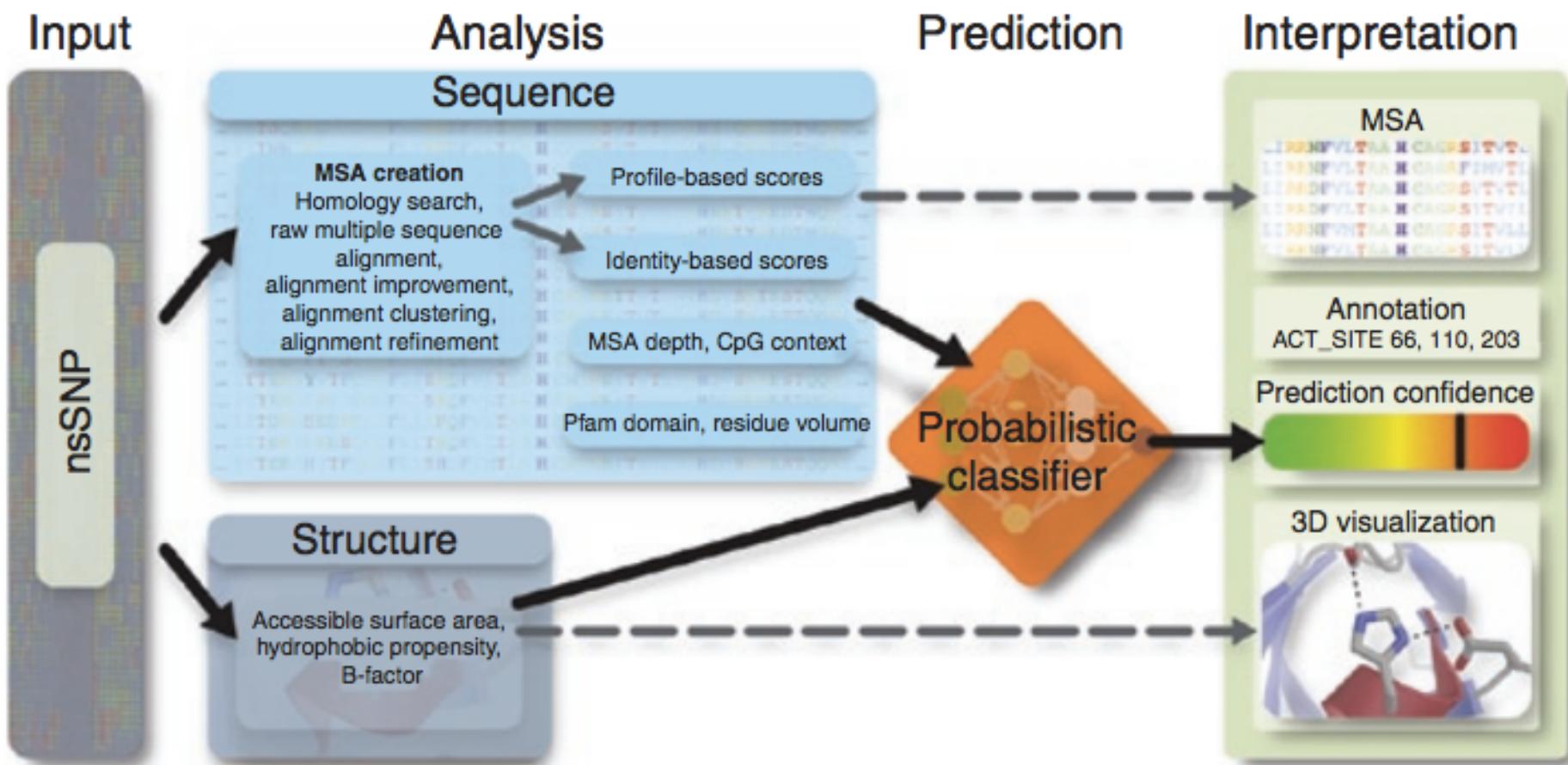


PolyPhen

- **Polymorphism Phenotyping**
 - Now on v2
- Predicts possible impact of amino acid substitution on structure and function of human protein using physical and comparative considerations
 - 8 sequence based predictive features
 - Multiple sequence alignment
 - Comparison of ancestral and mutant allele
 - 3 structure based predictive features



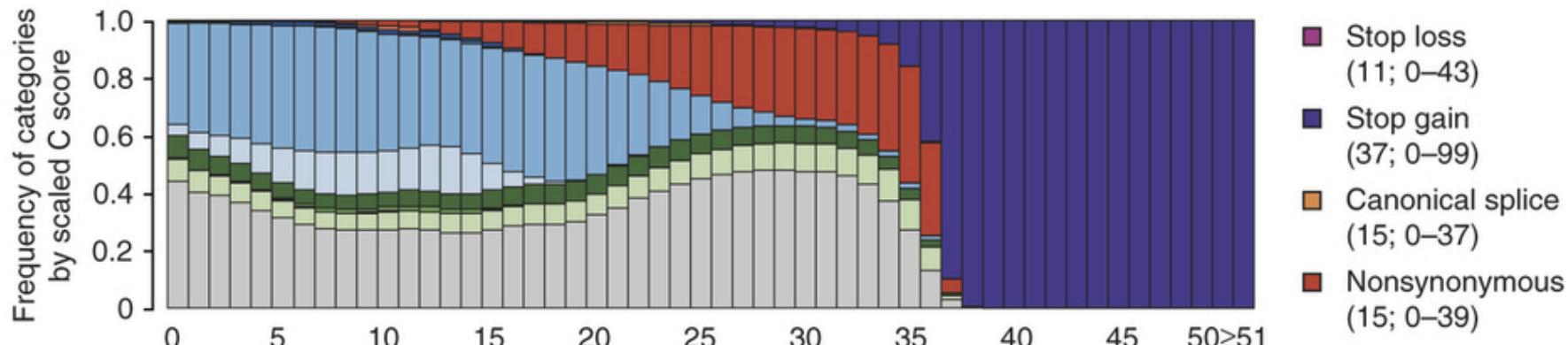
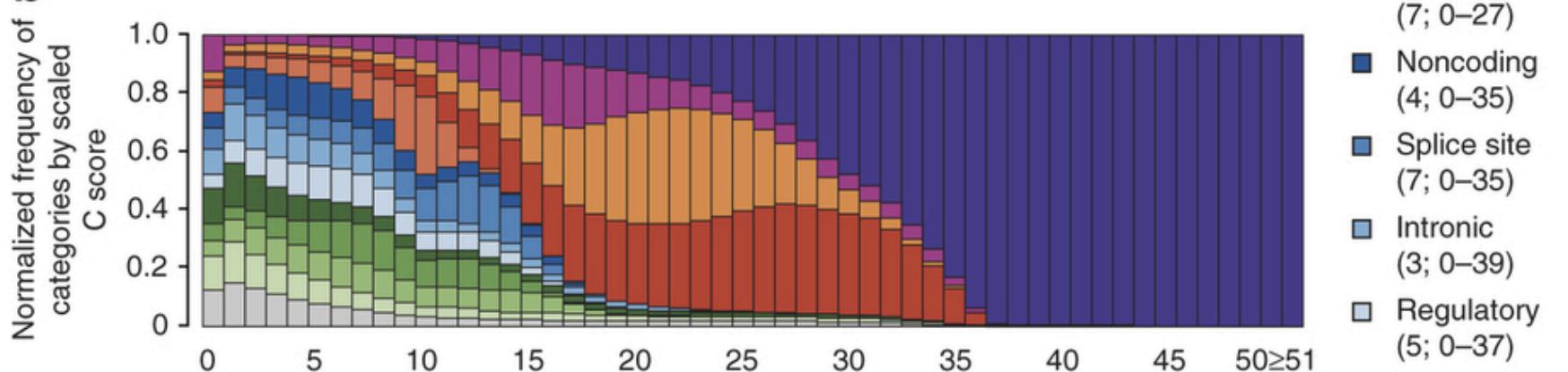
PolyPhen



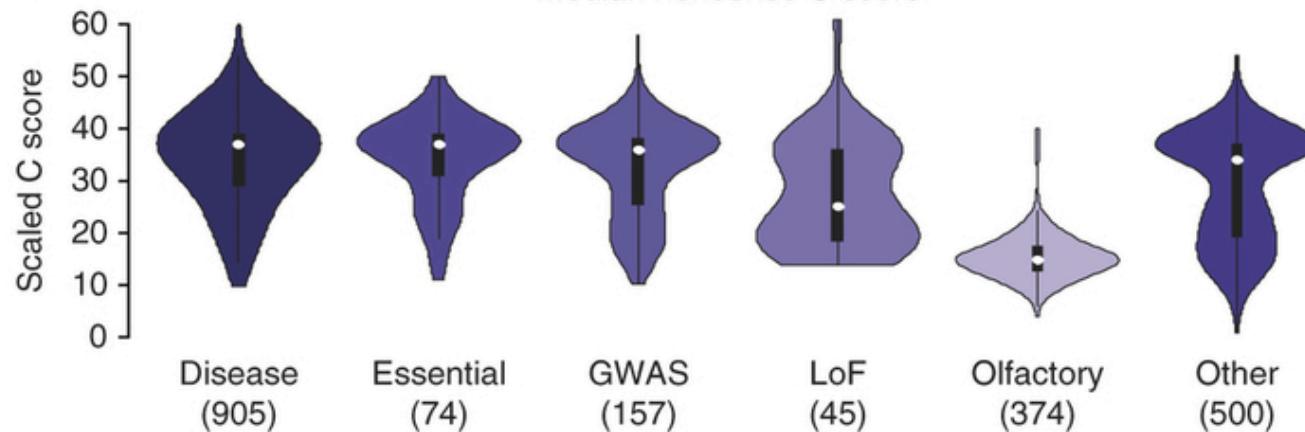
CADD

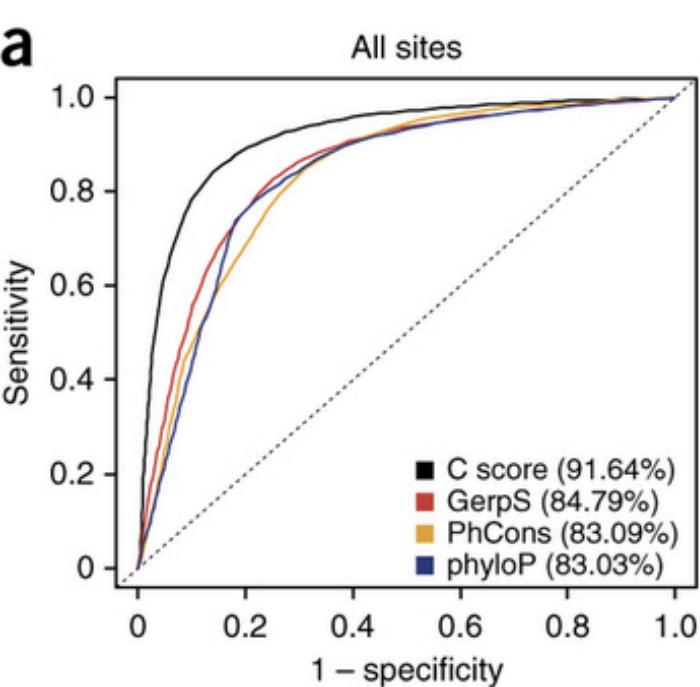
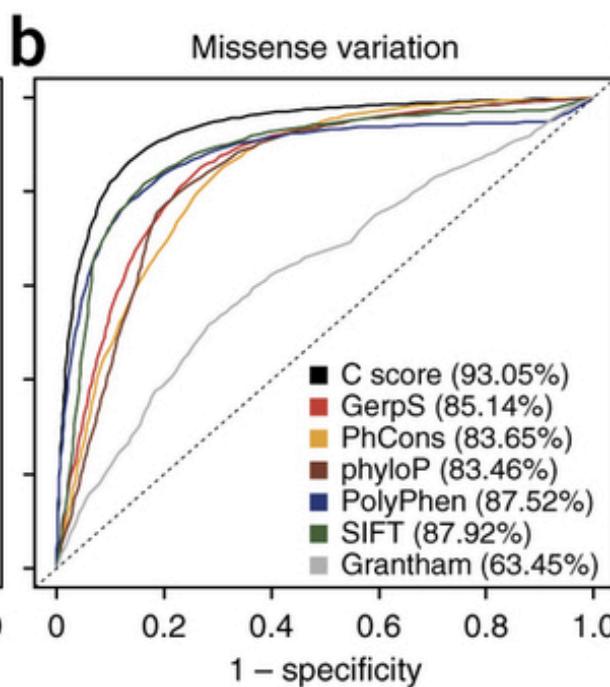
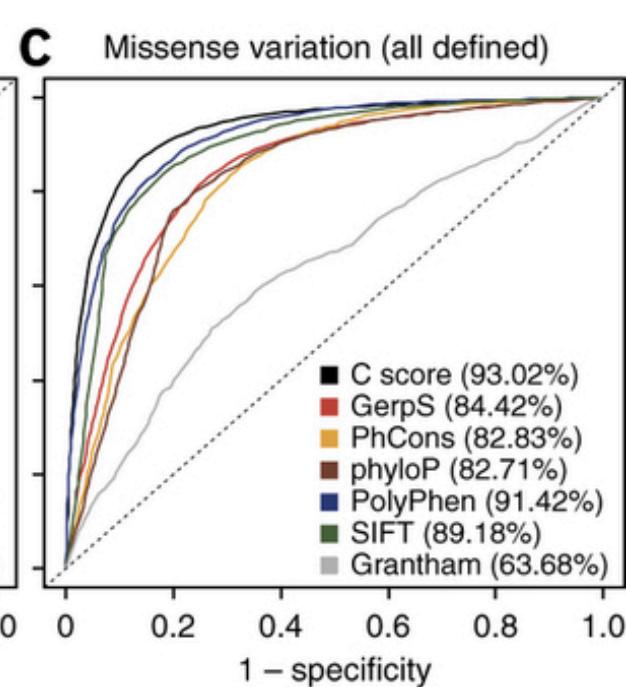
- Combined Annotation Dependent Depletion
- Score deleteriousness of single nucleotide variants as well as indels in human genome
- Integrates multiple annotations into one metric by contrasting variants that survived natural selection with simulated mutations
- C-scores strongly correlate with allelic diversity, pathogenicity of coding and non-coding variants
- Ranks causal variants within individual genome sequences



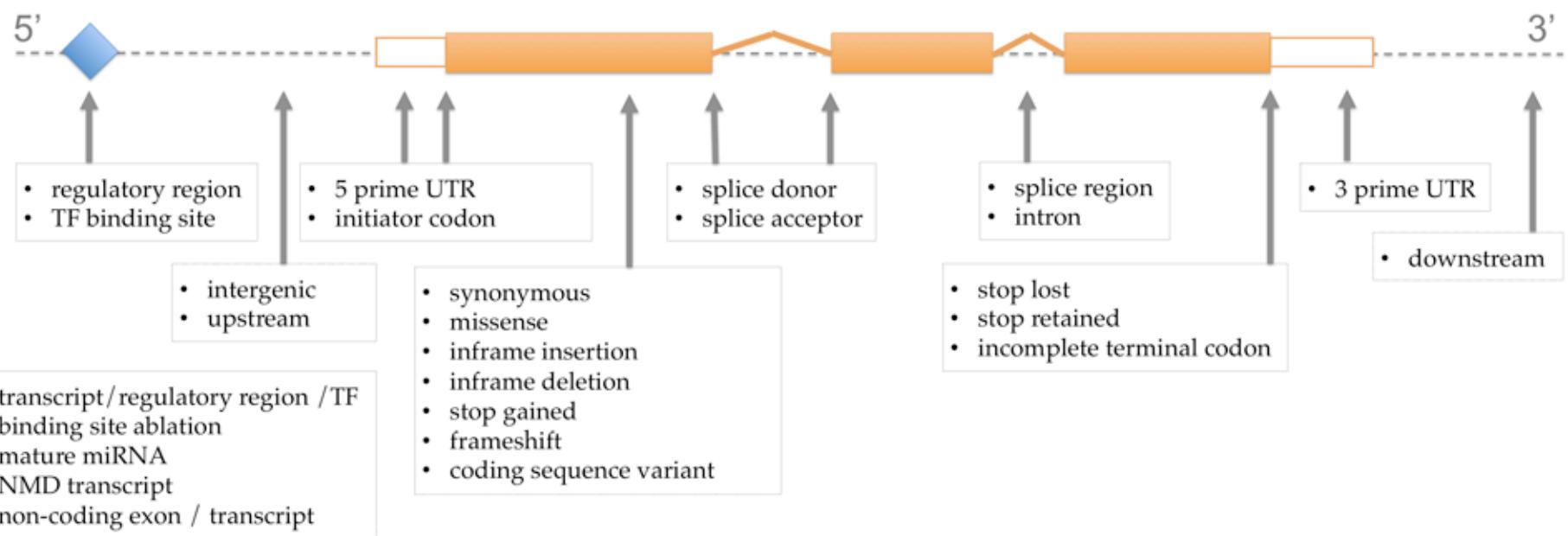
a**b****c**

Median nonsense C score



a**b****c**

Functional Categorization



SNPEff

- Categorizes effects of single nucleotide polymorphisms, multiple nucleotide polymorphisms and indels based on effects on annotated genes
 - Synonymous or non-synonymous
 - 3' or 5'UTR
 - Intronic
 - Upstream, downstream of genes
 - Intergenic regions
- Will output effects for every transcript
- Integrated into Galaxy, GATK and other sequence data analysis pipelines





- Variant Effect Predictor
- Standalone perl script or Web interface
- Output includes
 - Genes and transcripts affected by variants
 - Location of variants (e.g. upstream of transcript, in coding sequence, in non-coding RNA)
 - Consequence of variants on protein sequence (e.g. stop gained, missense, etc)
 - Known variants that match yours and associated minor allele frequencies from 1000 genomes project
 - SIFT and PolyPhen scores



VEP output

Category	Count
Variants processed	498955
Variants remaining after filtering	498955
Novel / known variants	-
Overlapped genes	825
Overlapped transcripts	2888
Overlapped regulatory features	7309

Consequences (all)



- synonymous_variant: 20%
- missense_variant: 18%
- downstream_gene_variant: 15%
- nc_transcript_variant: 11%
- upstream_gene_variant: 10%
- non_coding_exon_variant: 8%
- intron_variant: 6%
- NMD_transcript_variant: 5%
- 3_prime_UTR_variant: 3%
- Others



VEP output

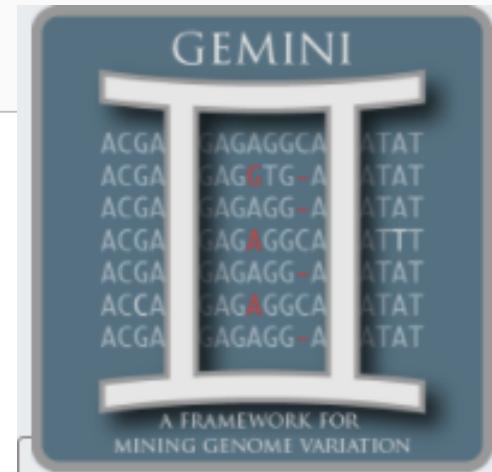
Show/hide columns								
Uploaded variation	Location	Feature	Feature type	Consequence	CDS position	Protein position	Amino acids	
rs116383664	1:1115461	ENSR00000528923	RegulatoryFeature	regulatory_region_variant	-	-	-	-
rs116383664	1:1115461	ENST00000379317	Transcript	upstream_gene_variant	-	-	-	-
rs116383664	1:1115461	ENST00000486379	Transcript	upstream_gene_variant	-	-	-	-
rs116383664	1:1115461	ENST00000379289	Transcript	missense_variant	247	83	R/W	
rs116383664	1:1115461	ENST00000460998	Transcript	upstream_gene_variant	-	-	-	-
rs116383664	1:1115461	ENST00000514695	Transcript	upstream_gene_variant	-	-	-	-
rs116383664	1:1115461	ENST00000379290	Transcript	missense_variant	247	83	R/W	
rs116383664	1:1115461	ENST00000379288	Transcript	missense_variant	28	10	R/W	

 Download
All [VCF VEP TXT](#)
Filtered [VCF VEP TXT](#)

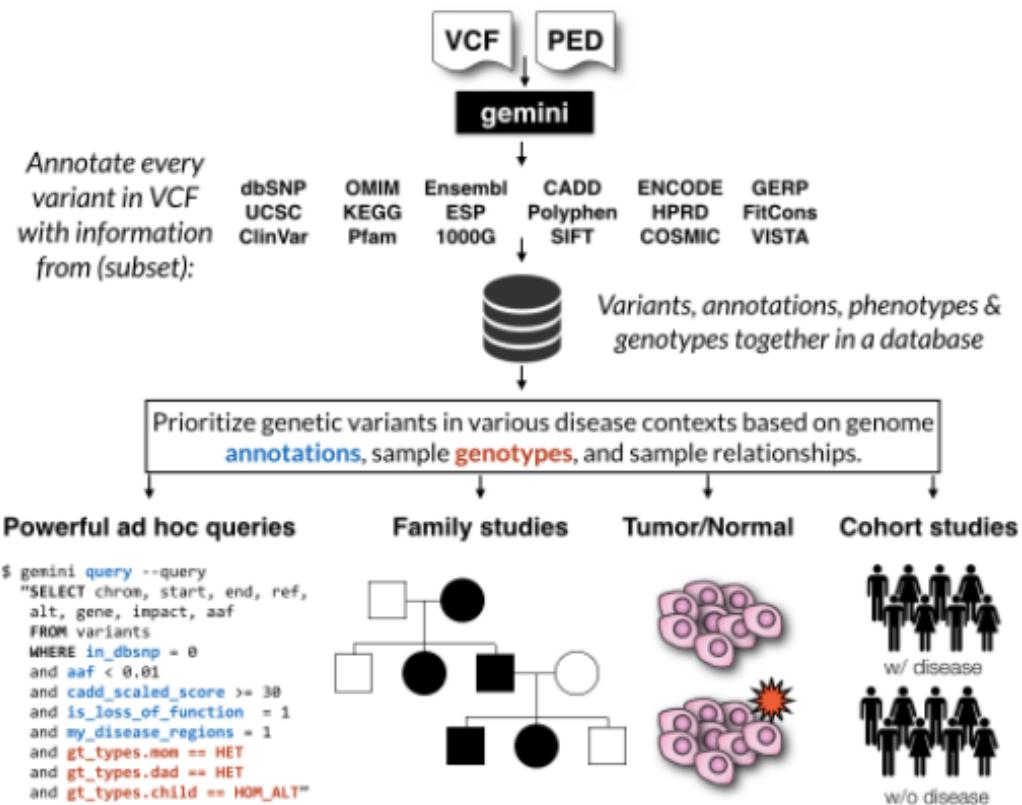
Download filtered results in TXT format (best for Excel)

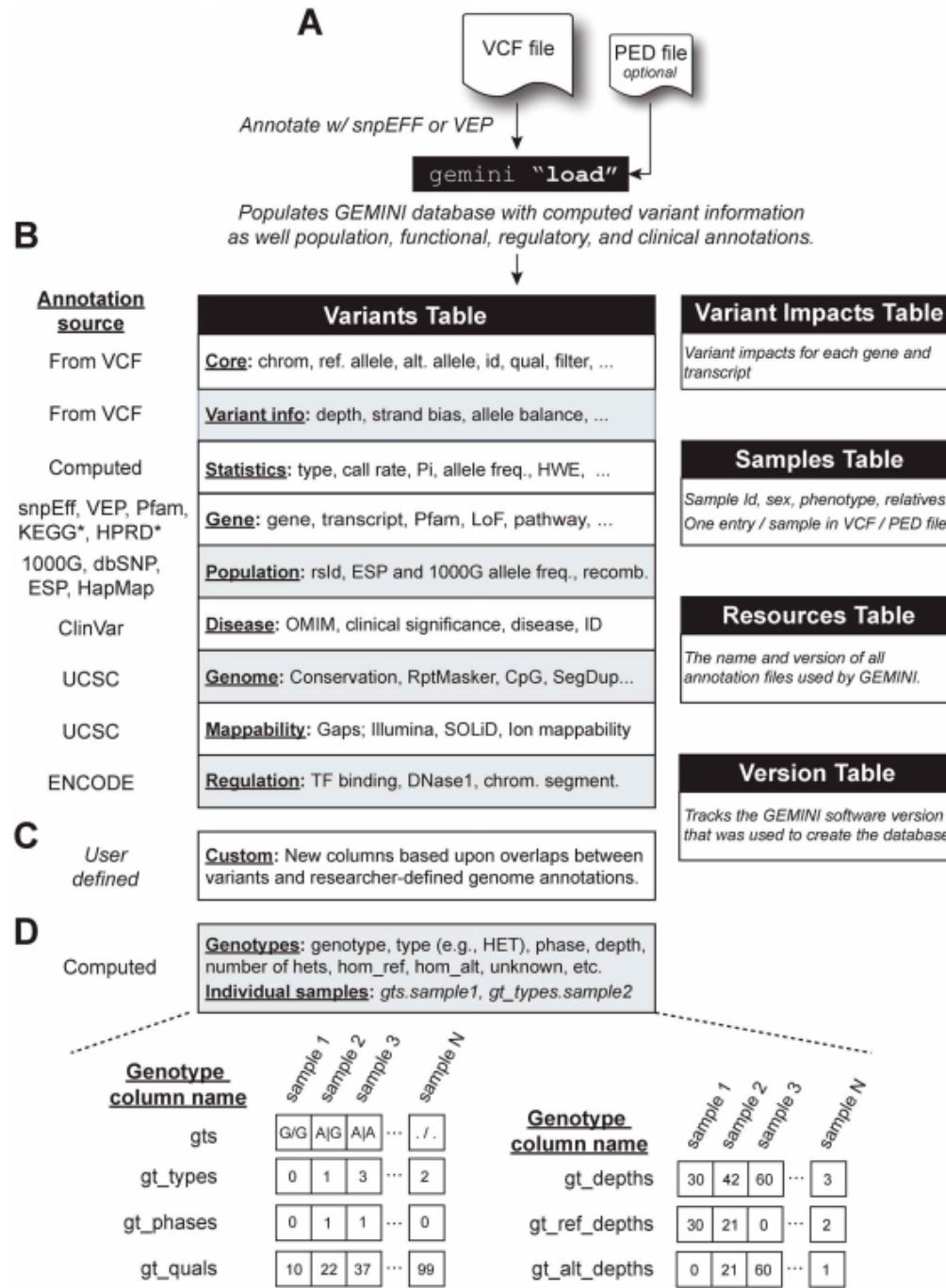


<http://gemini.readthedocs.io/en/latest/index.html>



GEMINI: a flexible framework for exploring genome variation





- Decompose original VCF so that variants with multiple alleles are expanded into distinct variant records (one record for each REF/ALT combination)
- Normalize the decomposed VCF so that variants are left aligned and represented using the most parsimonious alleles
- Annotate with VEP or snpEff
- bgzip and tabix
- Load VCF
 - If family based study be sure to use a ped file!



GEMINI

1. Extract all transitions with a call rate > 95%

```
$ gemini query -q "select * from variants \
    where sub_type = 'ts' \
    and call_rate >= 0.95" my.db
```

2. Extract all loss-of-function variants with an alternate allele frequency < 1%:

```
$ gemini query -q "select * from variants \
    where is_lof = 1 \
    and aaf >= 0.01" my.db
```

3. Extract the nucleotide diversity for each variant:

```
$ gemini query -q "select chrom, start, end, pi from variants" my.db
```

4. Combine GEMINI with bedtools to compute nucleotide diversity estimates across 100kb windows:

```
$ gemini query -q "select chrom, start, end, pi from variants \
    order by chrom, start, end" my.db | \
bedtools map -a hg19.windows.bed -b - -c 4 -o mean
```



- Built in analysis tools
 - de_novo
 - comp_hets
 - autosomal_recessive
 - autosomal_dominant
 - x_linked_recessive
 - x_linked_dominant
 - x_linked_de_novo



- Limit to confidently called genotypes
 - gt-pl-max 10
- Limit to those predicted to disrupt the protein
 - --filter “impact_severity = ‘HIGH’”
- Limit to rare alleles
 - --filter “max_aaf_all < 0.005”
- Limit to high genotype quality
 - --min-gq 20



```
$ gemini de_novo \  
  --columns "chrom, start, end, ref, alt" \  
  --filter "impact_severity = 'HIGH'" \  
  --min-kindreds 2 \  
  test.de_novo.db
```

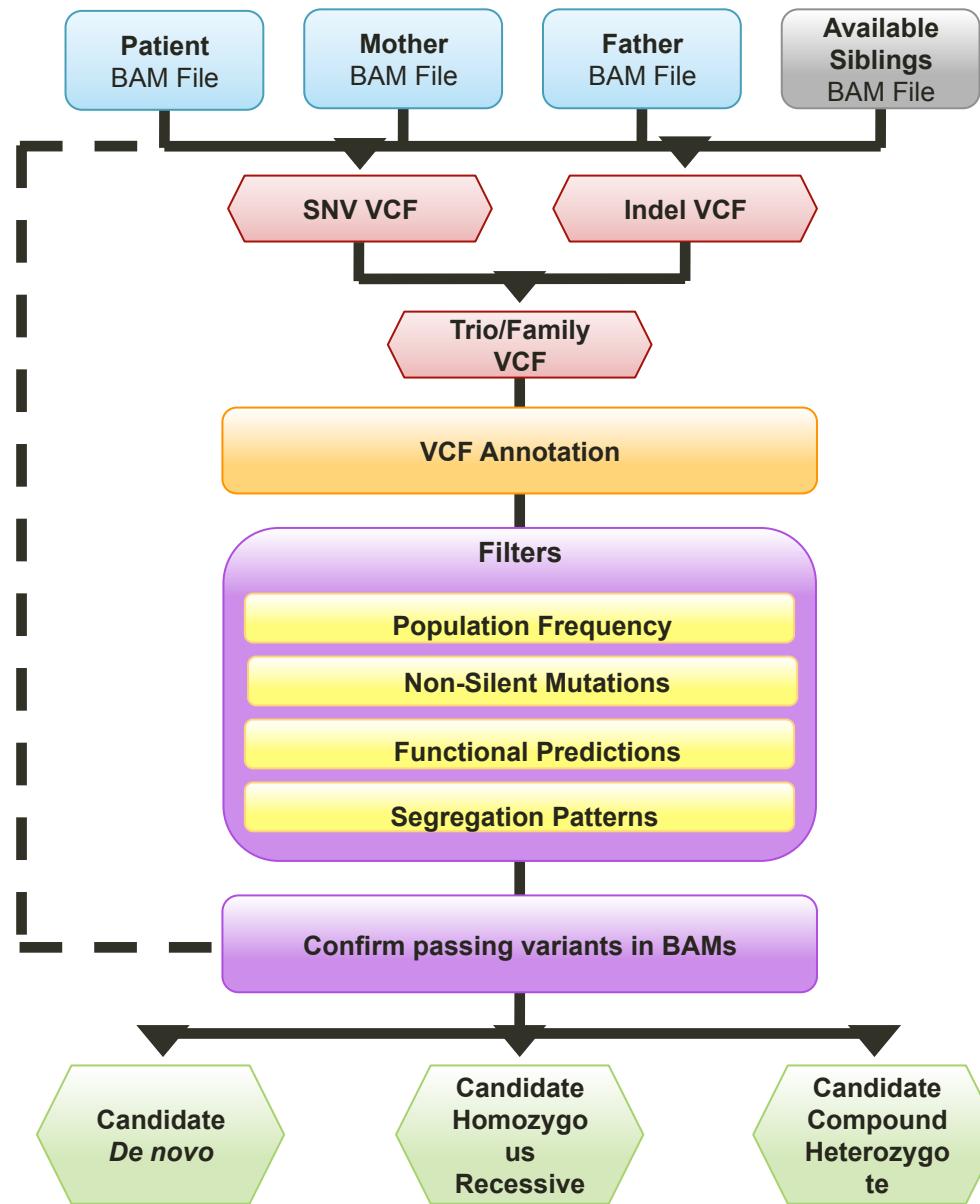
Restrict candidate genes to those genes with a de novo variant
in at least 2 families



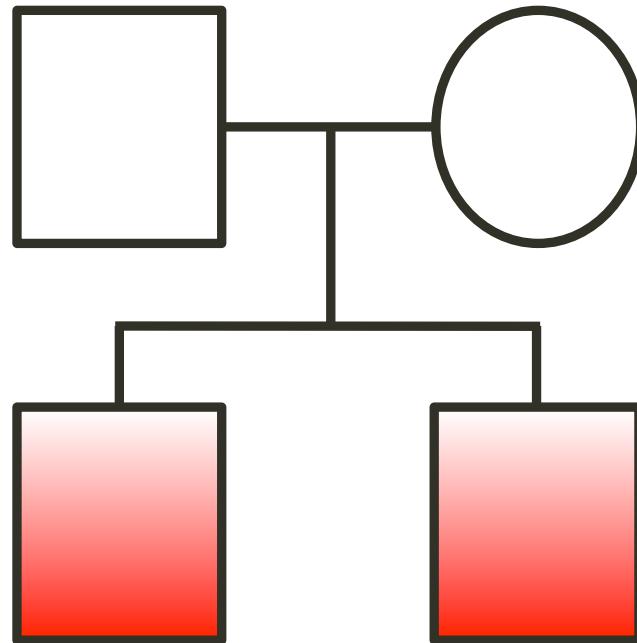
Practical Example: Genomics Pediatric Board

Washington University
Daniel Wegner, Jennifer Waumbach, Thomas Nicholas, F.
Sessions Cole

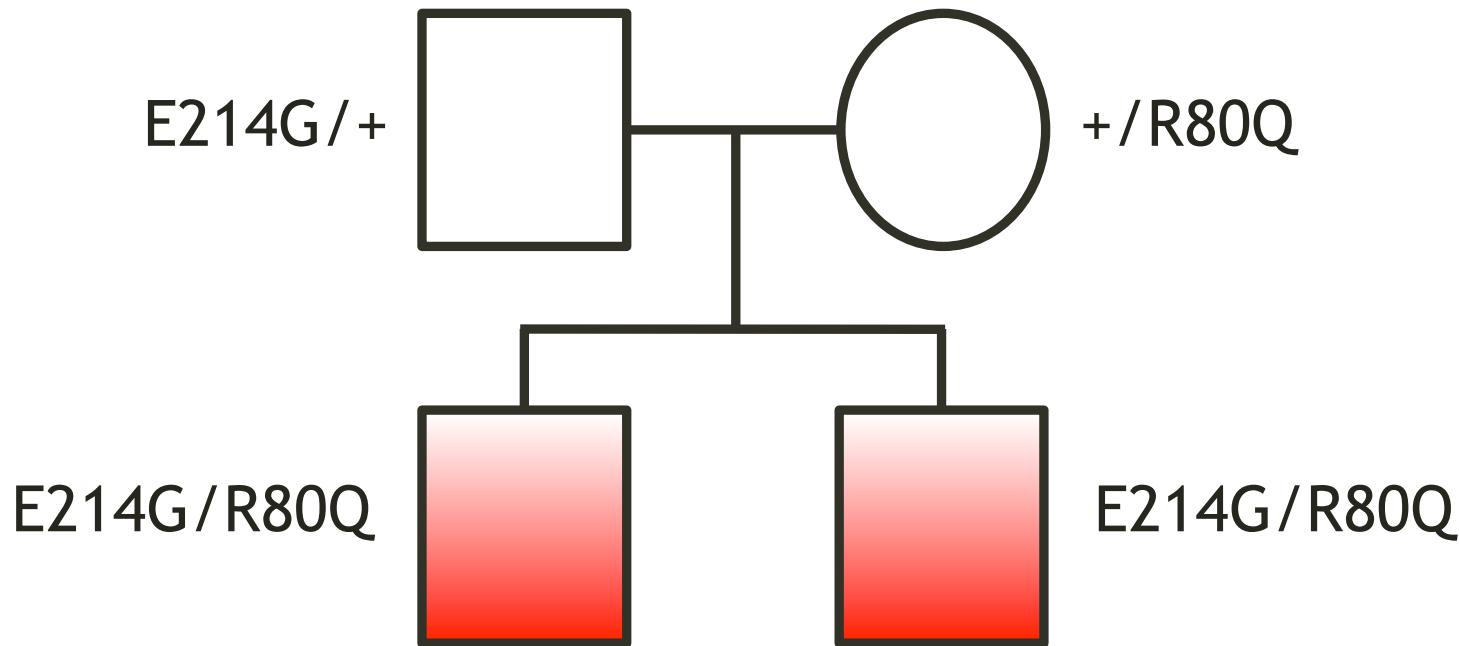


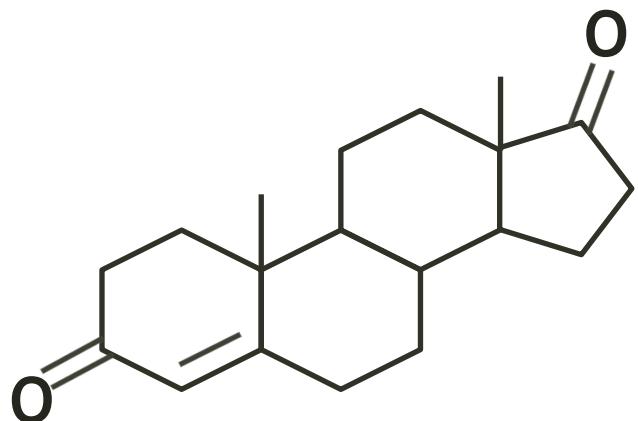


GBD-012 presented 2 male siblings with ambiguous genitalia



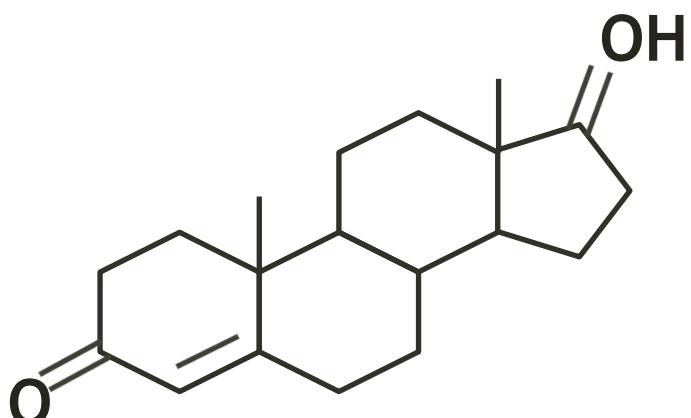
Each sibling is compound heterozygous for *HSD17B3*





Androstenedione

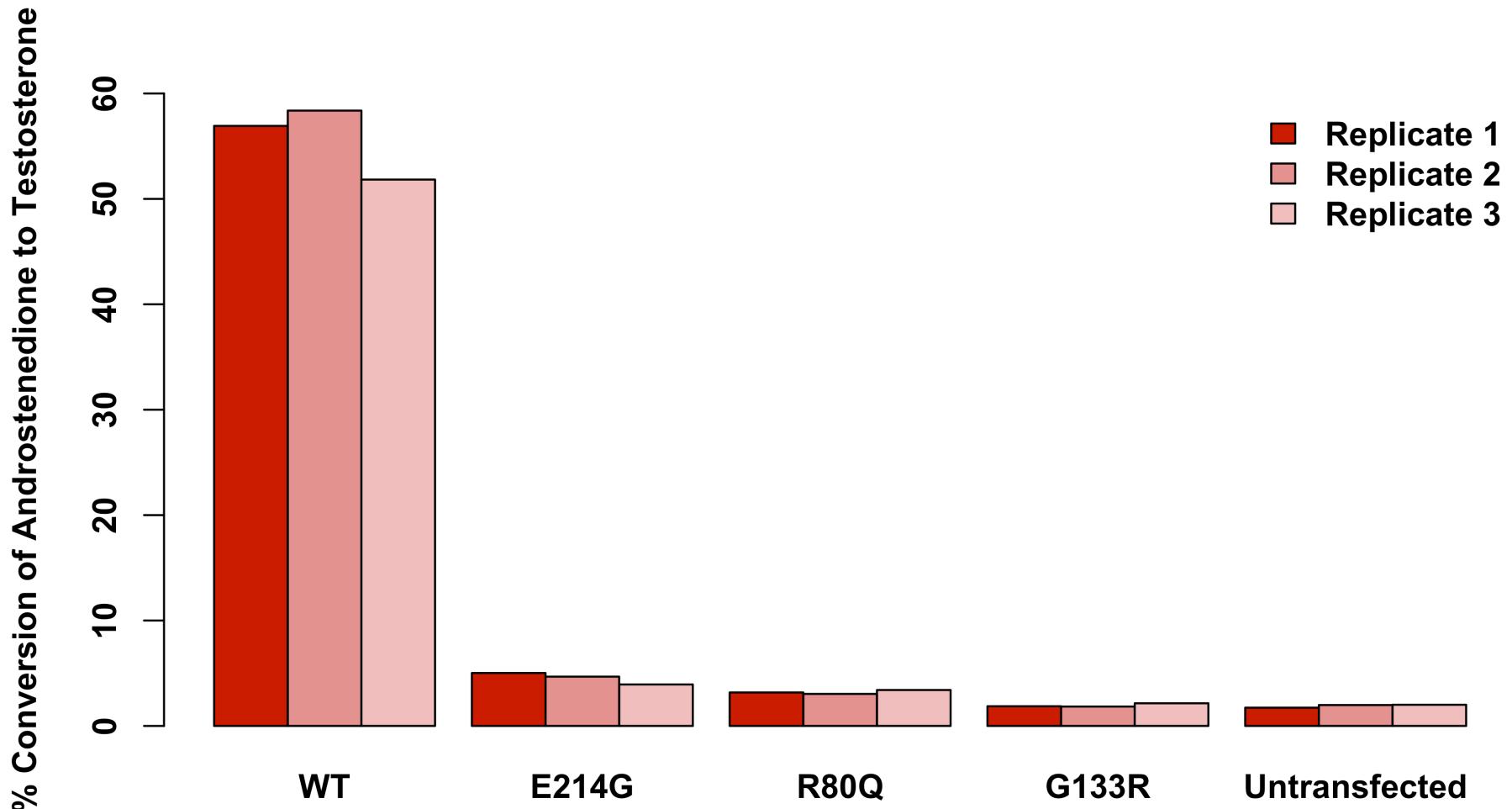
HSD17B3
↔



Testosterone



Decreases in androstenedione conversion in cells with mutations



Identifying rare variants in large cohorts



Why study rare variants?

- Identify the complete genetic architecture of traits
 - Rare variants should identify new susceptibility loci at known trait loci
-
- May have higher penetrance
 - May also have higher effect size
 - Most variants that impact function are expected to be deleterious and removed by natural selection

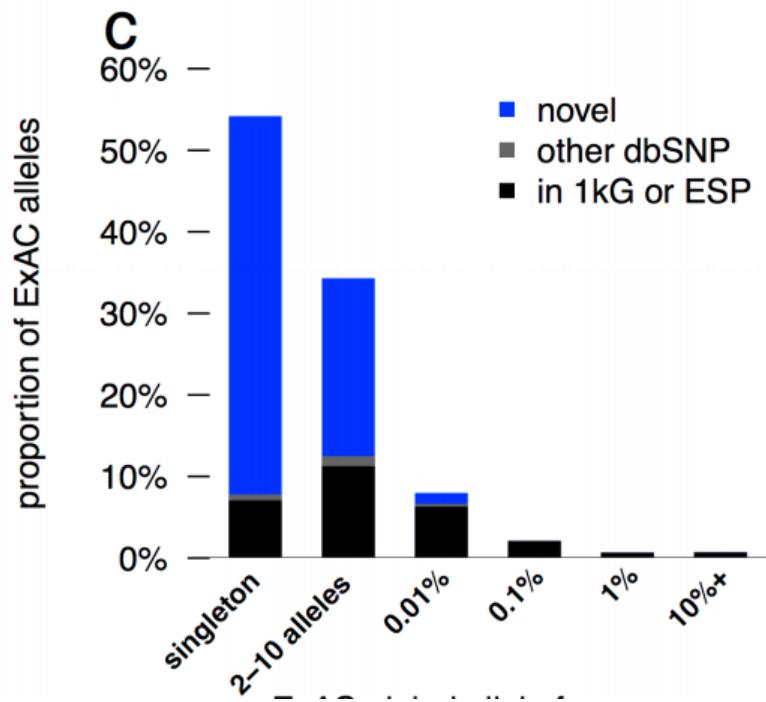
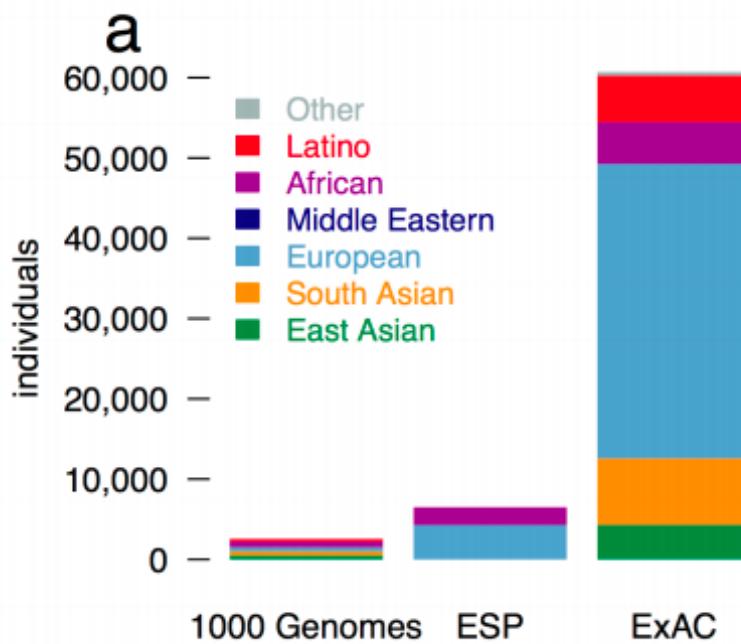


Many rare, likely functional coding variants

	# SNVs	Singletons	Doubletons	Tripletons	> 3 occurrences
Synonymous	270,263	128,319 (47%)	29,340 (11%)	13,129 (5%)	99,475 (37%)
Nonsynonymous	410,956	234,633 (57%)	46,740 (11%)	19,274 (5%)	110,309 (27%)
Nonsense	8,913	6,196 (70%)	926 (10%)	326 (4%)	1,465 (16%)
Ratio (NS/S)		1.8 to 1	1.6 to 1	1.4 to 1	1.1 to 1



Increasing sample size to 60K identifies even more rare variants



MOST VARIATION IS RARE



Single Variant Association Testing

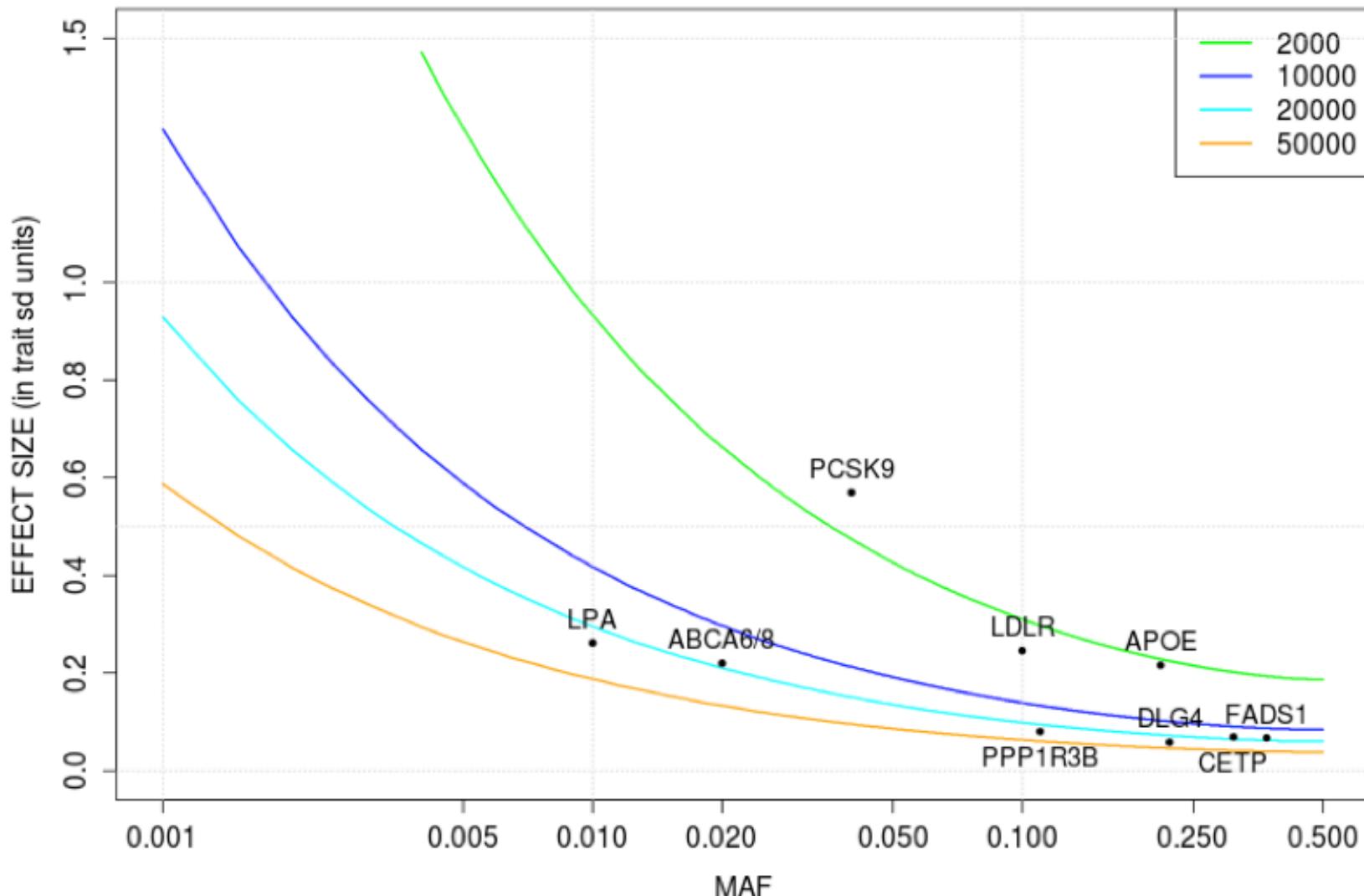
Sample sizes (cases and controls) required
for 80% power to detect
Effect Size=2.0 at $\alpha=5\times10^{-8}$

MAF	Sample size
0.05	2,500
0.01	12,000
0.001	117,000



With increased sample sizes we have greater power to detect associations at lower effect sizes and MAF

Power = 0.8 at $\alpha = 5e-7$



Power depends on:

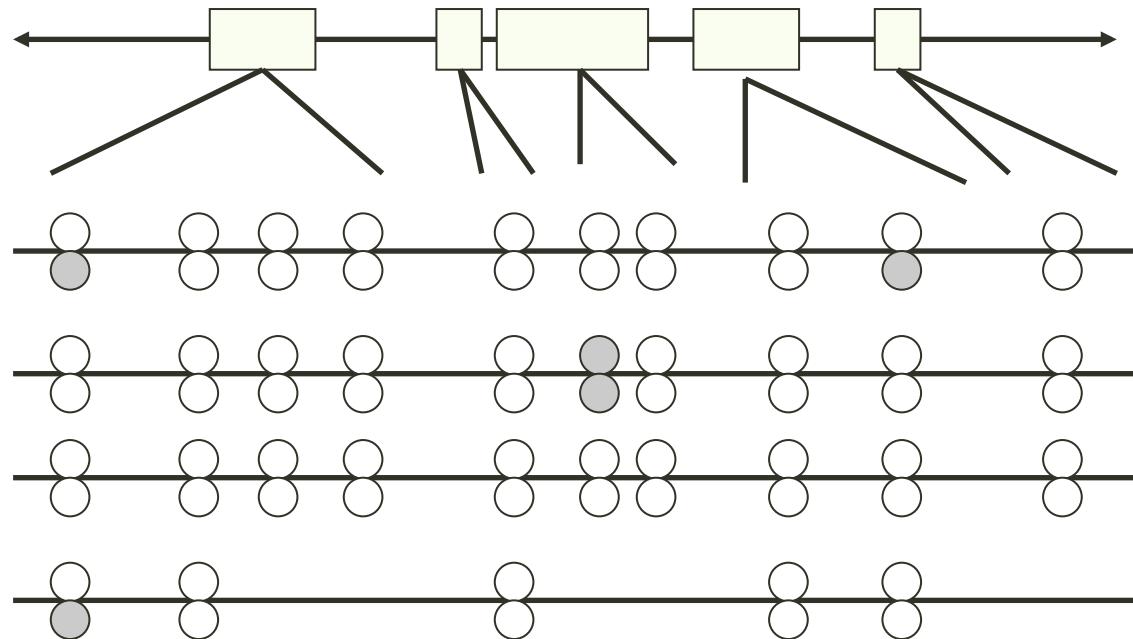
Frequency

Effect Size

Even with large effects, rare variants can only be detected in large samples



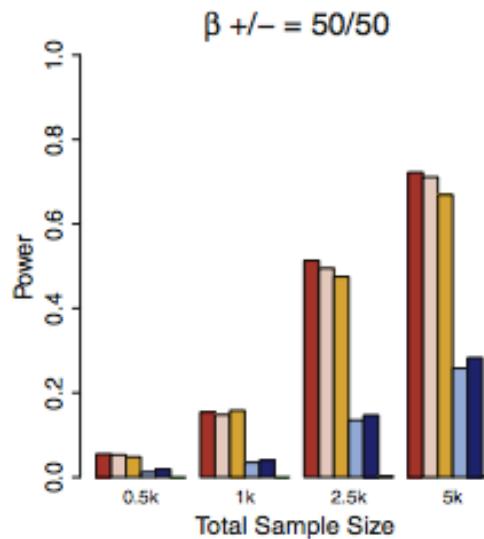
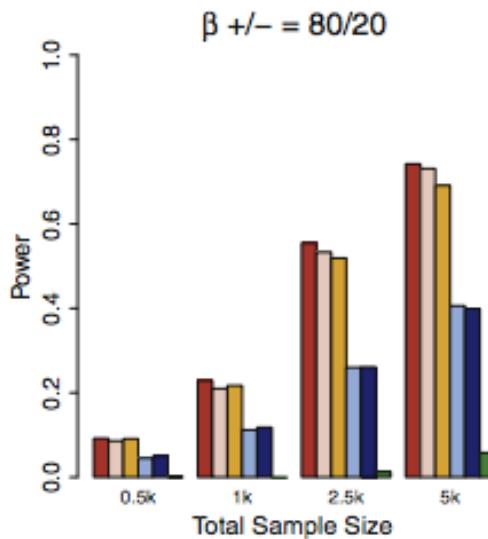
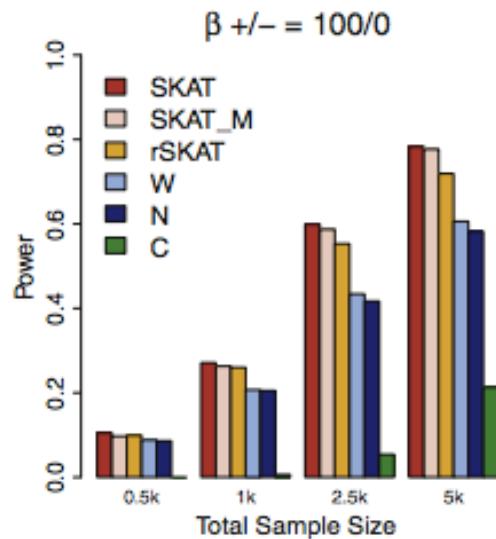
What if we could group variants together?



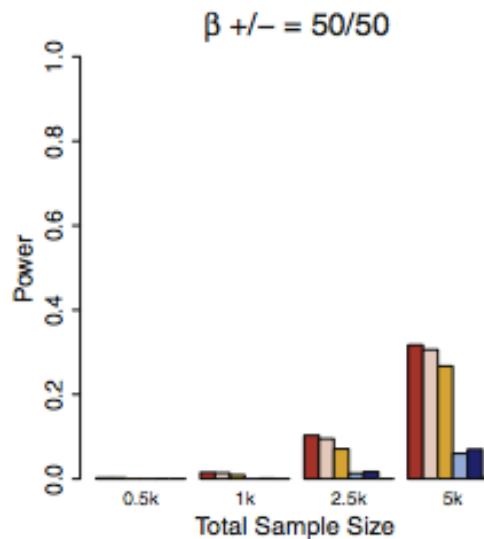
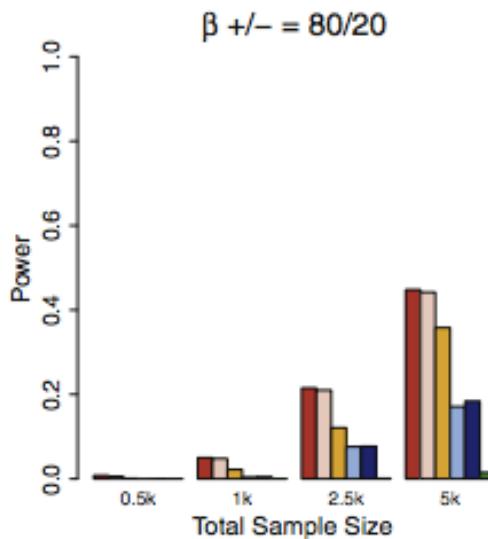
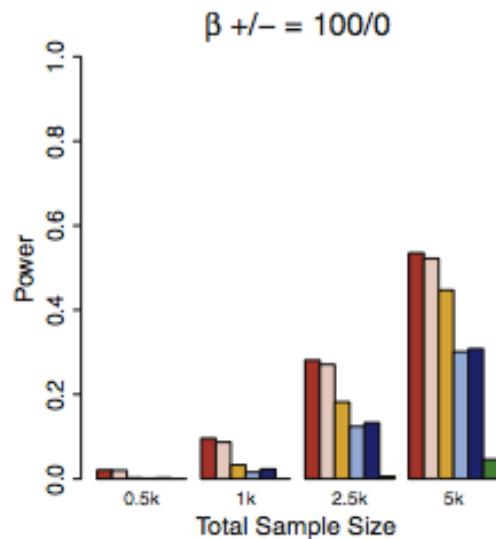
- Sequence kernel association test
 - Regression approach
 - Directly performs multiple regression of a phenotype on genotypes for all variants
 - Tests for association between common and rare variants in a region
 - Dichotomous or continuous phenotype
 - Adjusts for covariates like population stratification
 - Incorporates flexible weight functions to boost power
 - e.g. increase weight of variants with lower MAF



Continuous Trait



Dichotomous Trait

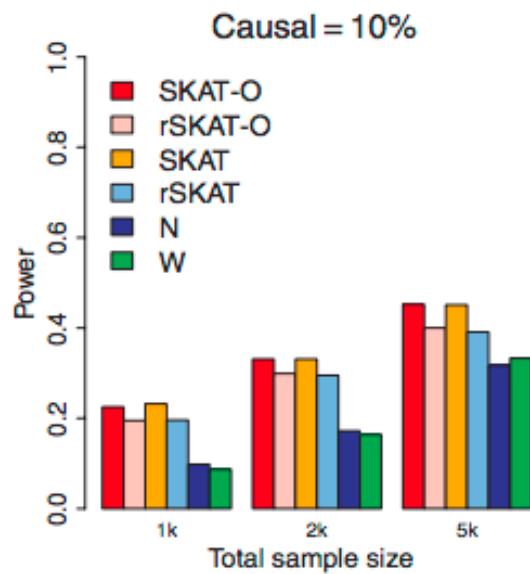


SKAT-O

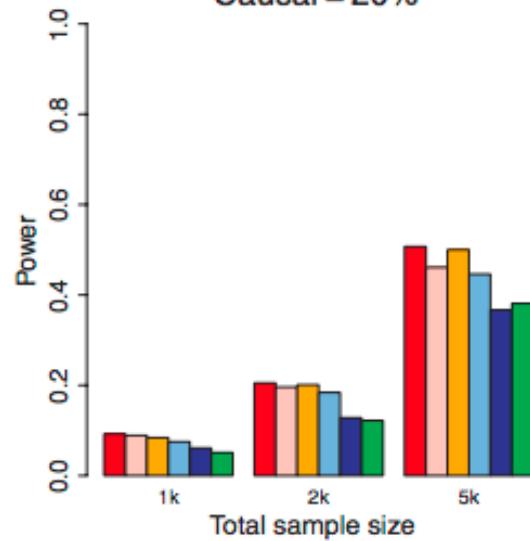
- Use both burden-based tests and SKAT
- Incorporate a correlation structure of variant effects through family of kernels
 - When effects of variants are perfectly correlated, it is essentially just a burden test
- Optimal test derived by estimating correlation parameter in kernel matrix to maximize power



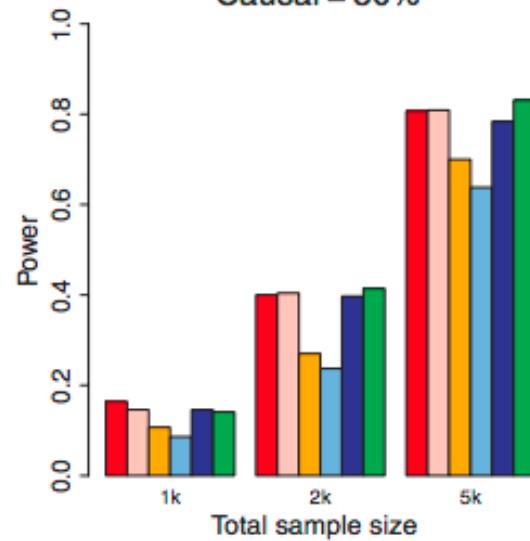
Continuous trait



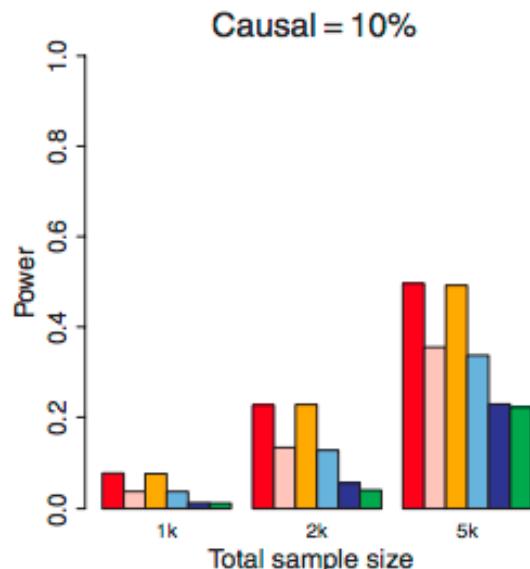
Causal = 20%



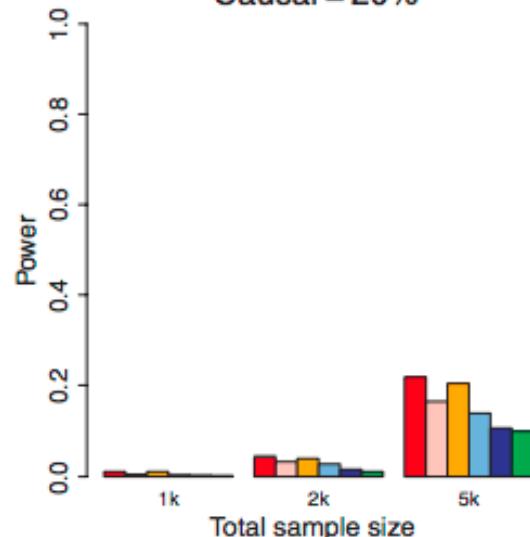
Causal = 50%



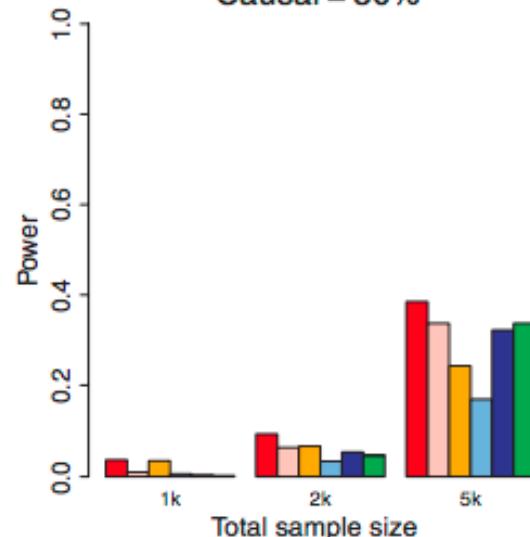
Dichotomous trait



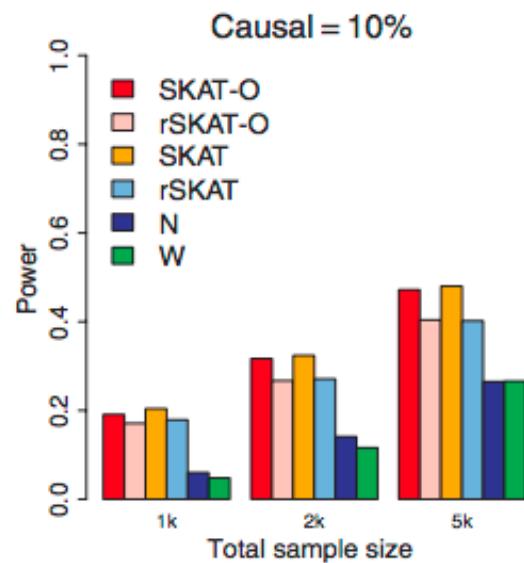
Causal = 20%



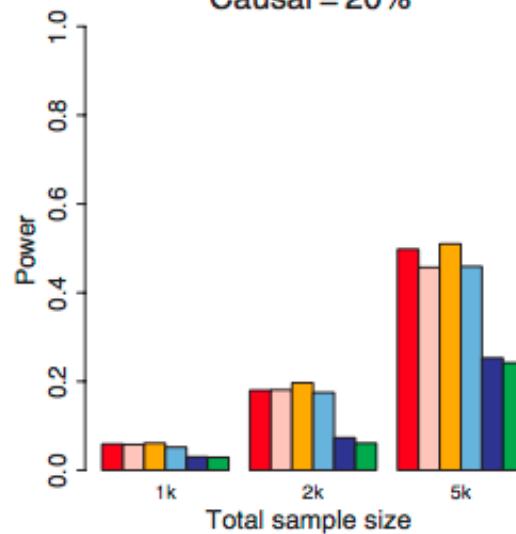
Causal = 50%



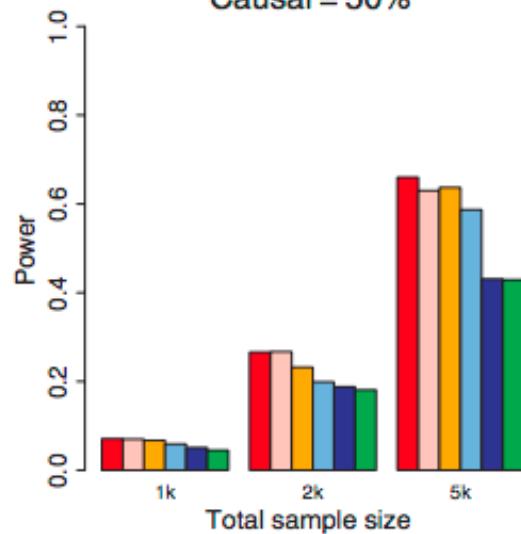
Continuous trait



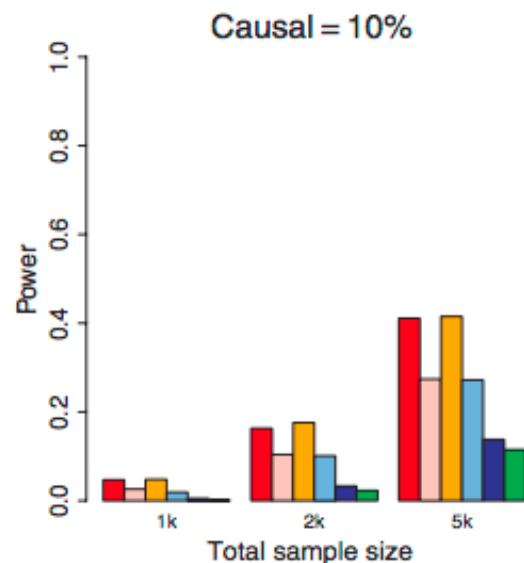
Causal = 20%



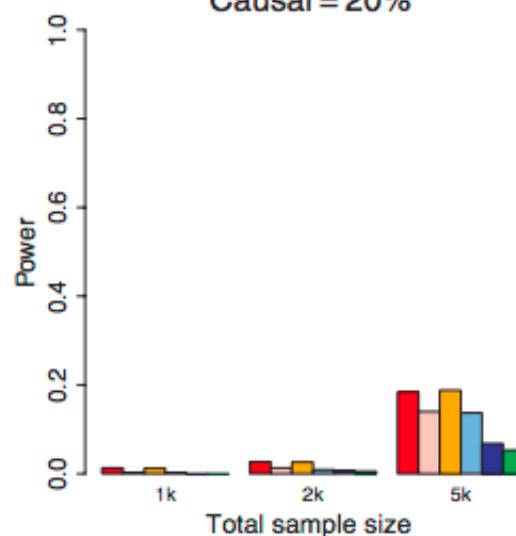
Causal = 50%



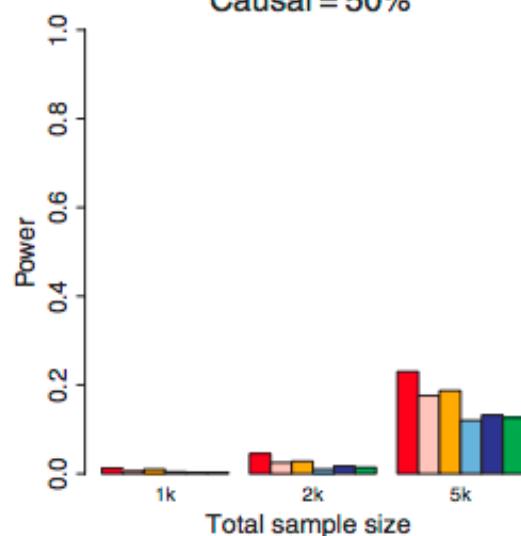
Dichotomous trait



Causal = 20%



Causal = 50%



EMMAX

- Statistical test for large scale association studies accounting for sample structure (e.g. relatedness)
- Avoids repetitive variance component estimation
 - Each loci only explains a small fraction of complex traits
- Very computationally efficient
- <http://genome.sph.umich.edu/wiki/EMMAX>



Maximize the power

- Choose threshold for defining rare carefully
- Focus on loss of function (LOF) variants only
- Focus on individuals with extreme trait values for continuous phenotypes
- Focus on individuals with strong family history of disease for dichotomous traits
- Population isolates



Practical Example: FinMetSeq Project

Washington University, University of Michigan, UCLA, FIMM,
University of Kuopio

Karyn Meltz Steinberg and Adam Locke



Finland Population History

Early Settlement

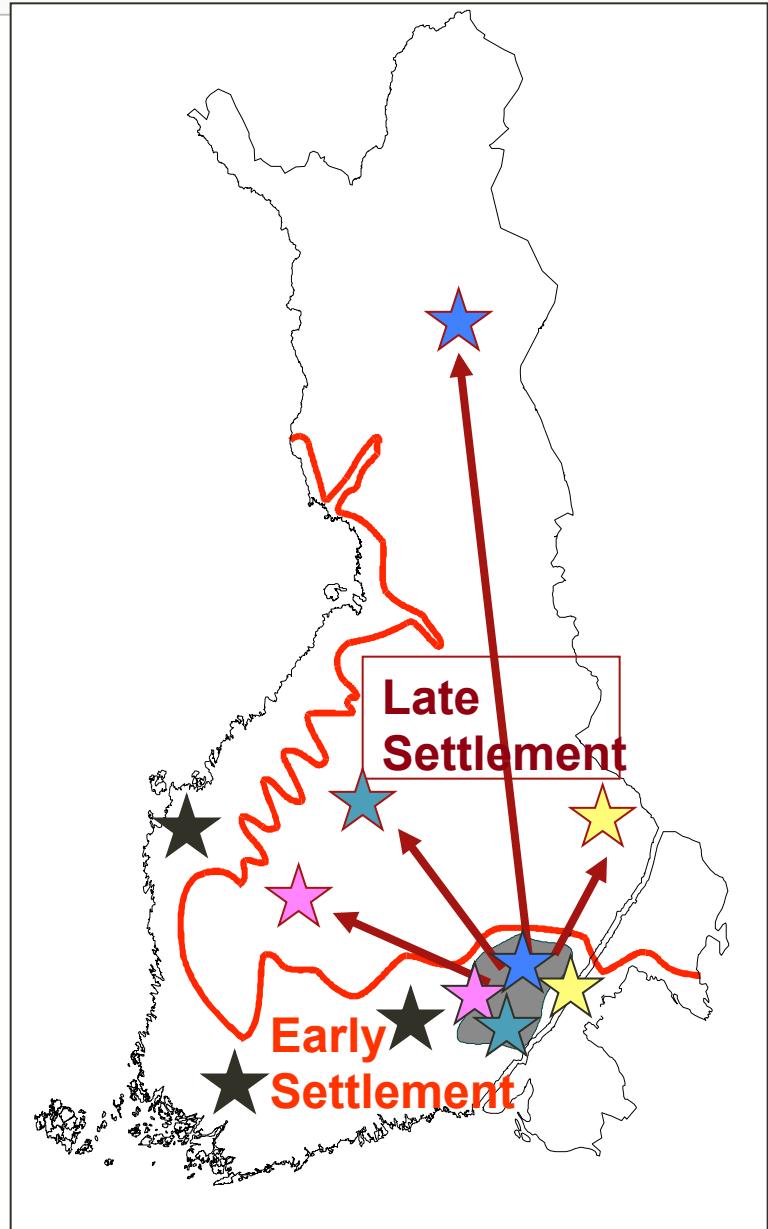
- 2,000-10,000 years ago
- South and Coast

Late Settlement

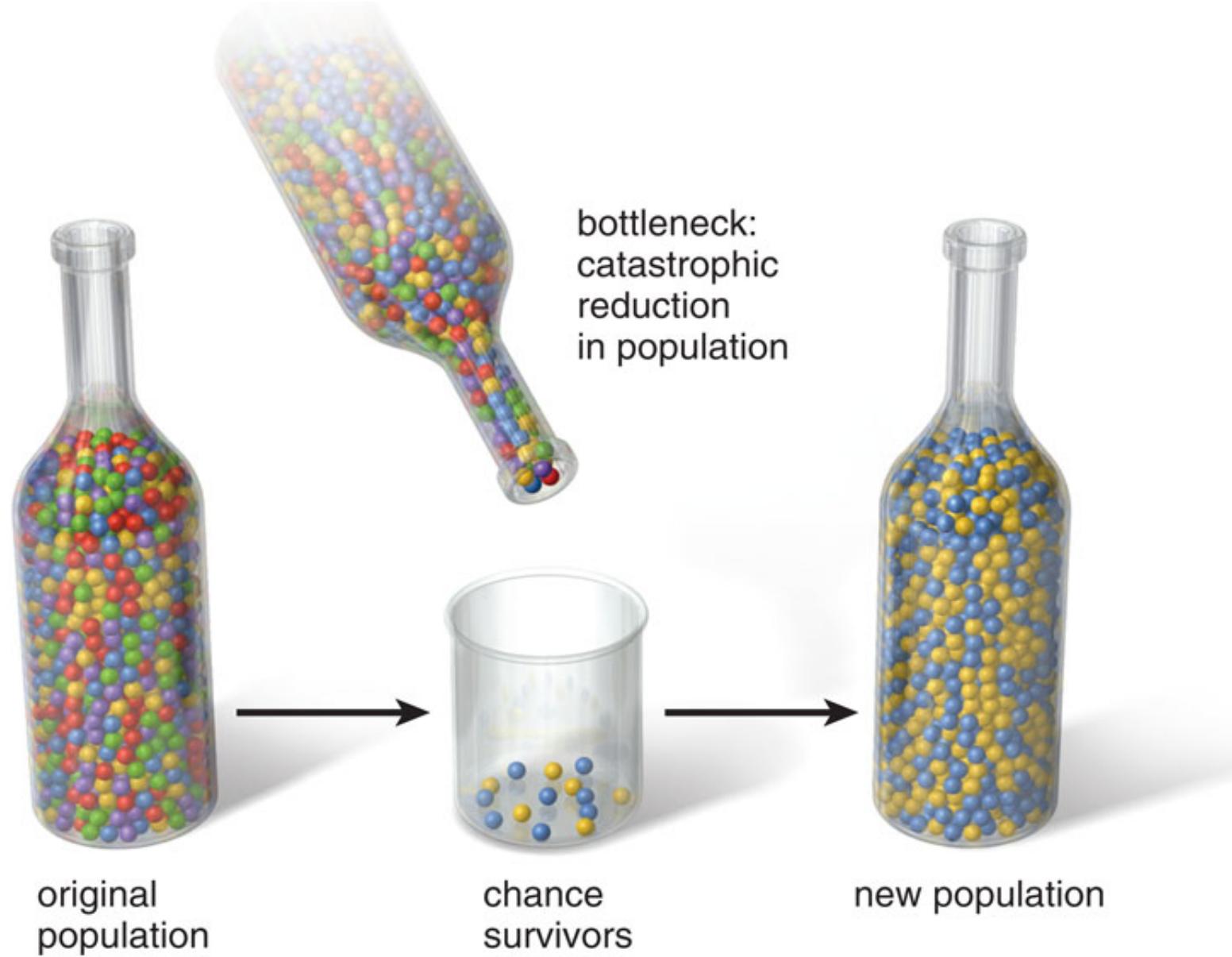
- 16th century
- Multiple bottlenecks

Expansion

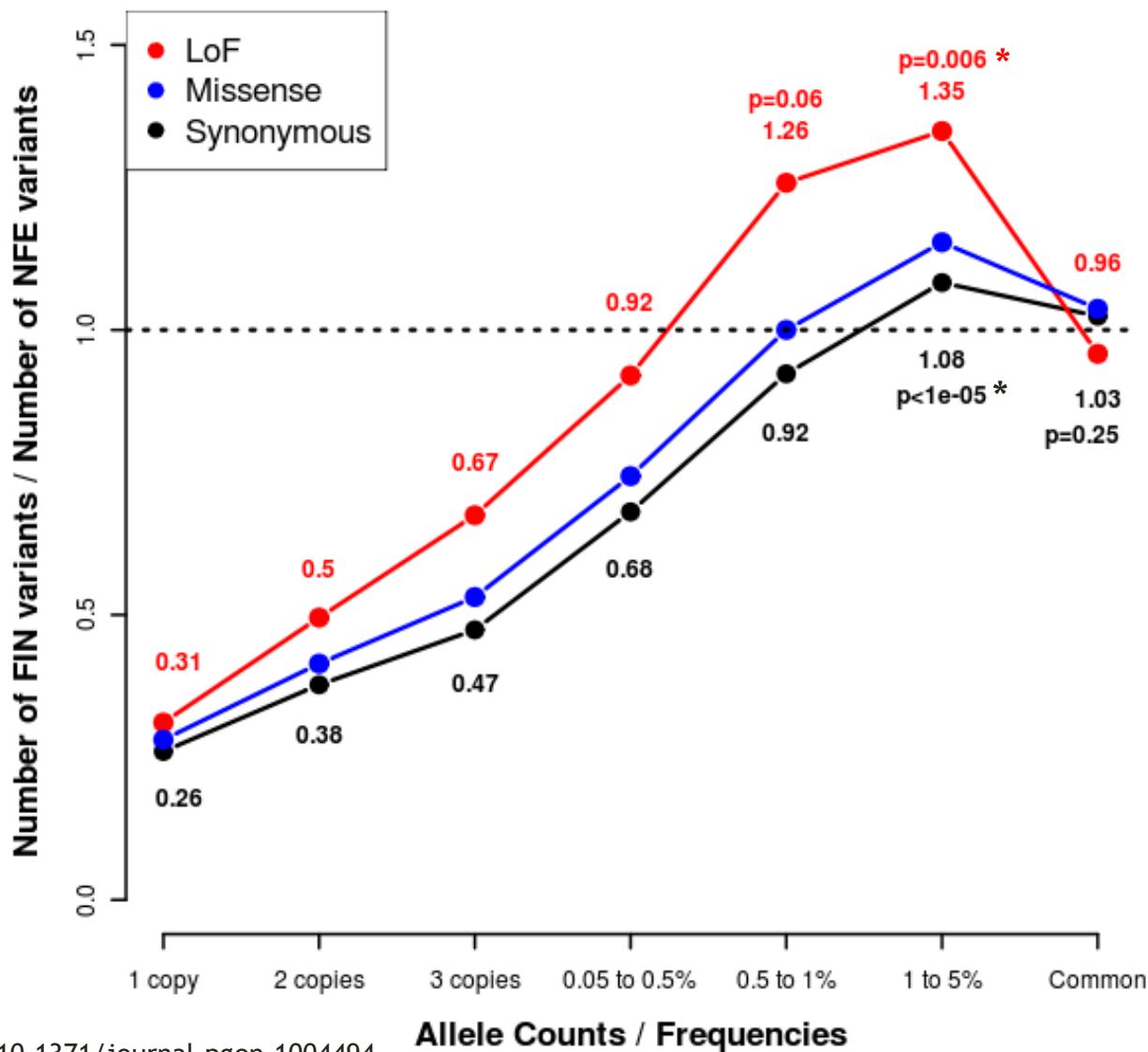
- 18th century (pop 250K)
- Today (pop 5.3M)



Population isolates: Bottleneck effect and genetic drift



There are proportionally more LoF variants in Finnish individuals compared to Non-Finnish Europeans



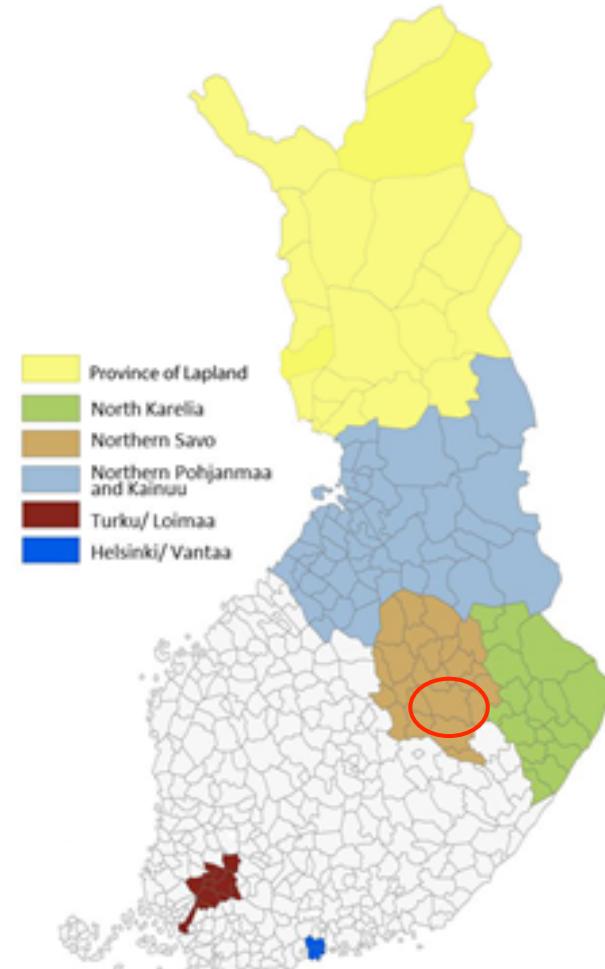
Major Questions

- Do currently unknown infrequent/rare variants contribute substantially to common disease risk and/or important quantitative traits relevant to cardiometabolic disorders?
 - Novel loci
 - Fine mapping known loci and identifying additional signals
- Can we relate sequence variants that we identify to specific clinical outcomes?
 - Finnish disease heritage



Cohorts for exome sequencing

- METSIM
 - Population based study of 10,197 men from Kuopio
 - Ages 45-70 years at baseline
 - Baseline study conducted in 2005-2010
- FINRISK
 - Health examination surveys every five years since 1972
 - Sample sizes 6,000-12,000 individuals
 - Ages 25-74
 - Stratified by 10 year age group, gender and study area
 - Included all participants from eastern and northern areas except those that overlap METSIM



Kuopio



Association analysis methodology

- Generate phenotype residuals and transform separately for METSIM and FINRISK
- Inverse normalize phenotype residuals for all quantitative traits (except height)
- Use EMMAX to account for relatedness within and across studies
- Five major trait groups
 - Lipids
 - Glucose and insulin
 - Anthropometrics
 - Blood pressure
 - Other (inflammatory markers, etc)



Overview

- Protein truncating variants + Non-synonymous variants MAF<0.1, deleterious by all 5 algorithms in Purcell et al 2015
- 12,880 genes
- Bonferroni corrected p-value 3.88e-6



ABCA1 and ApoA1 ($p=3.26e-06$)

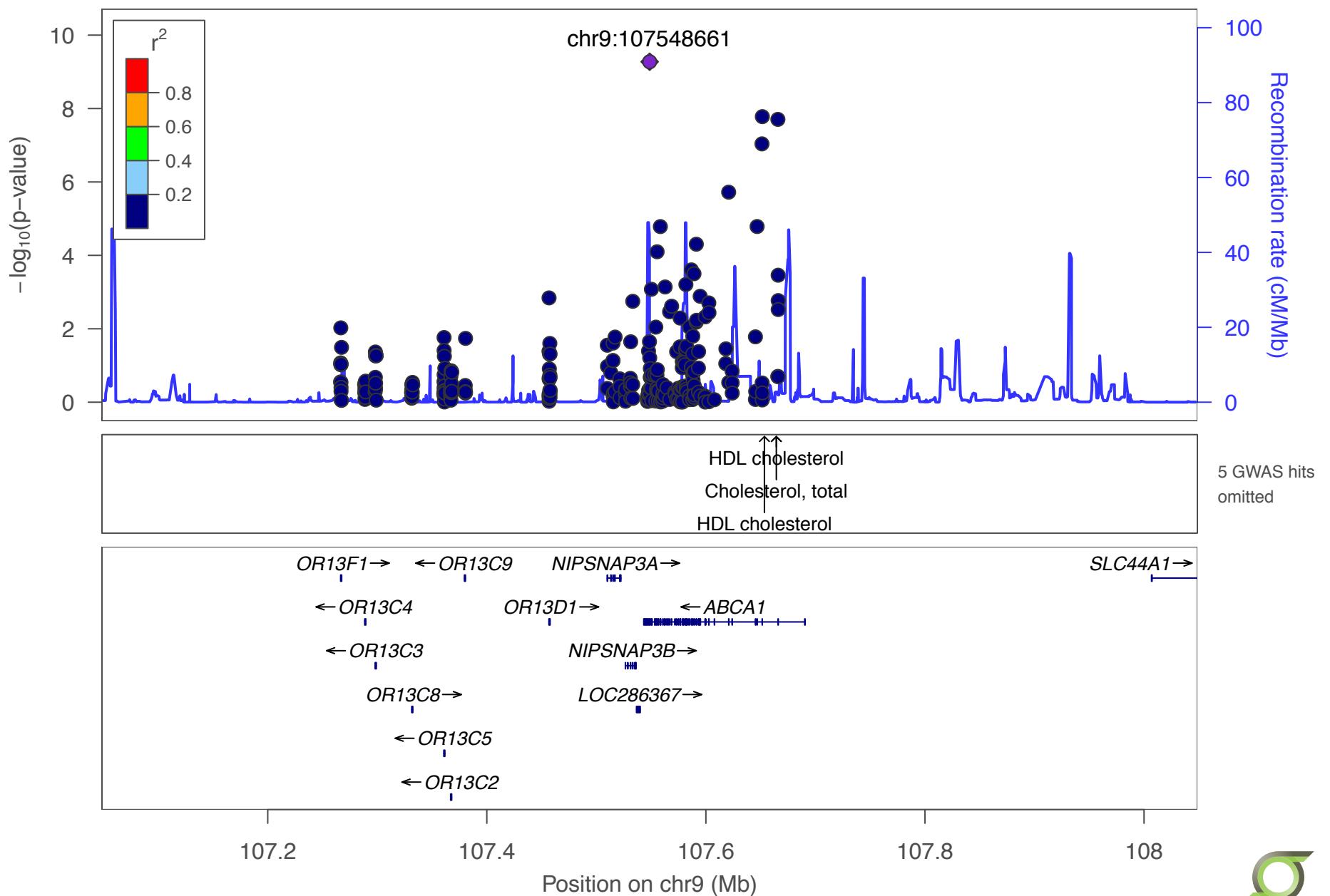
Variants in gene based test	MAC	MAF	ExAC_AF_FIN	ExAC_AF_NFE	P-value	Beta	CADD	CSQ
9:107548661_A/G (Lead SNP)	7	2.32E-04	4.50E-04	0.00E+00	4.42E-05	-1.531	32	Missense
9:107550254_G/A	3	7.72E-05	1.50E-04	1.50E-05	4.14E-03	-1.644	34	Missense
9:107550798_T/G	1	2.57E-05	NA	NA	NA	NA	31	Missense
9:107558416_T/C	15	8.24E-04	NA	NA	7.90E-02	-0.4499	32	Missense
9:107560784_C/T	6	3.35E-04	9.00E-04	3.00E-05	1.20E-01	-0.6298	35	Missense
9:107568659_G/A	1	2.57E-05	0.00E+00	0.00E+00	NA	NA	39	Stop gain
9:107589238_C/G	7	3.60E-04	1.50E-04	3.30E-03	3.12E-01	-0.3791	26.4	Missense
9:107593948_C/T	1	2.57E-05	NA	NA	NA	NA	21.1	Missense
9:107599789_A/G	1	2.57E-05	NA	NA	NA	NA	29.6	Missense



ABCA1 and ApoA1 ($p=3.26e-06$)

- The membrane-associated protein encoded by this gene is a member of the superfamily of ATP-binding cassette (ABC) transporters. ABC proteins transport various molecules across extra- and intracellular membranes.
- With cholesterol as its substrate, this protein functions as a cholesterol efflux pump in the cellular lipid removal pathway. Mutations in this gene have been associated with Tangier's disease and familial high-density lipoprotein deficiency.
[provided by RefSeq, Jul 2008]
- Gene previously associated with Cholesterol (Willer et al, 2013), Lipid metabolism phenotypes (Kettunen et al, 2012), Metabolite levels (Inouye et al, 2012), Metabolic syndrome (Kristiansson et al, 2012), Coronary heart disease (Lettre et al, 2011), HDL cholesterol (Waterworth et al, 2010; Willer et al, 2008)





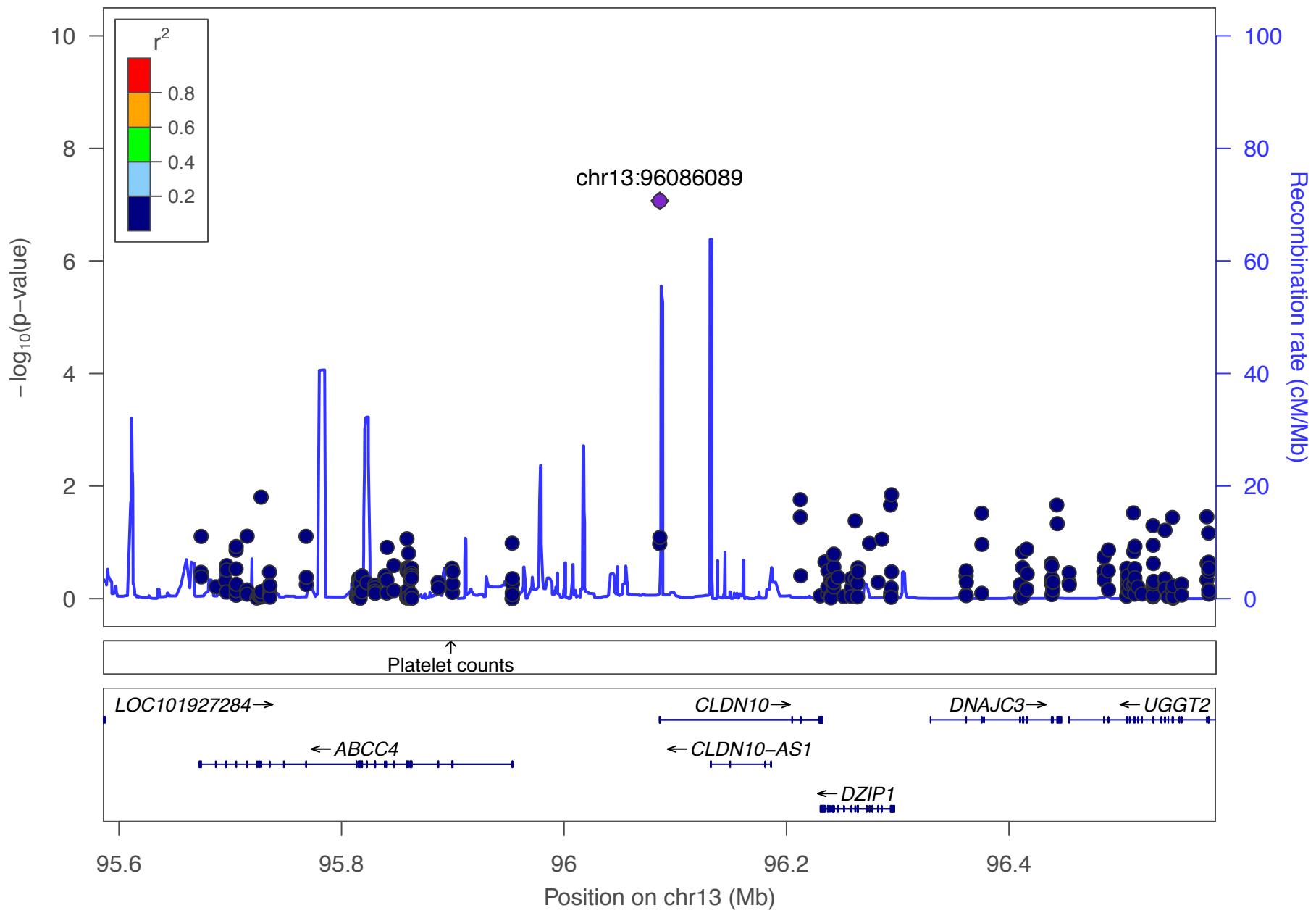
CLDN10 and Creatinine levels ($p=2.90E-08$)

Variants in gene based test	MAC	MAF	ExAC_AF_FIN	ExAC_AF_NFE	P-value	Beta	CADD	CSQ
13:96086089:T/A (Lead SNP)	40	1.80E-03	1.50E-03	1.50E-05	8.55E-08	0.8487	24.2	Start loss
13:96212394:C/T	4	2.00E-04	1.50E-04	0.00E+00	3.55E-02	1.046	40	Stop gain
13:96229547:G/T	1	5.00E-05	0.00E+00	1.50E-05	NA	NA	22.1	missense

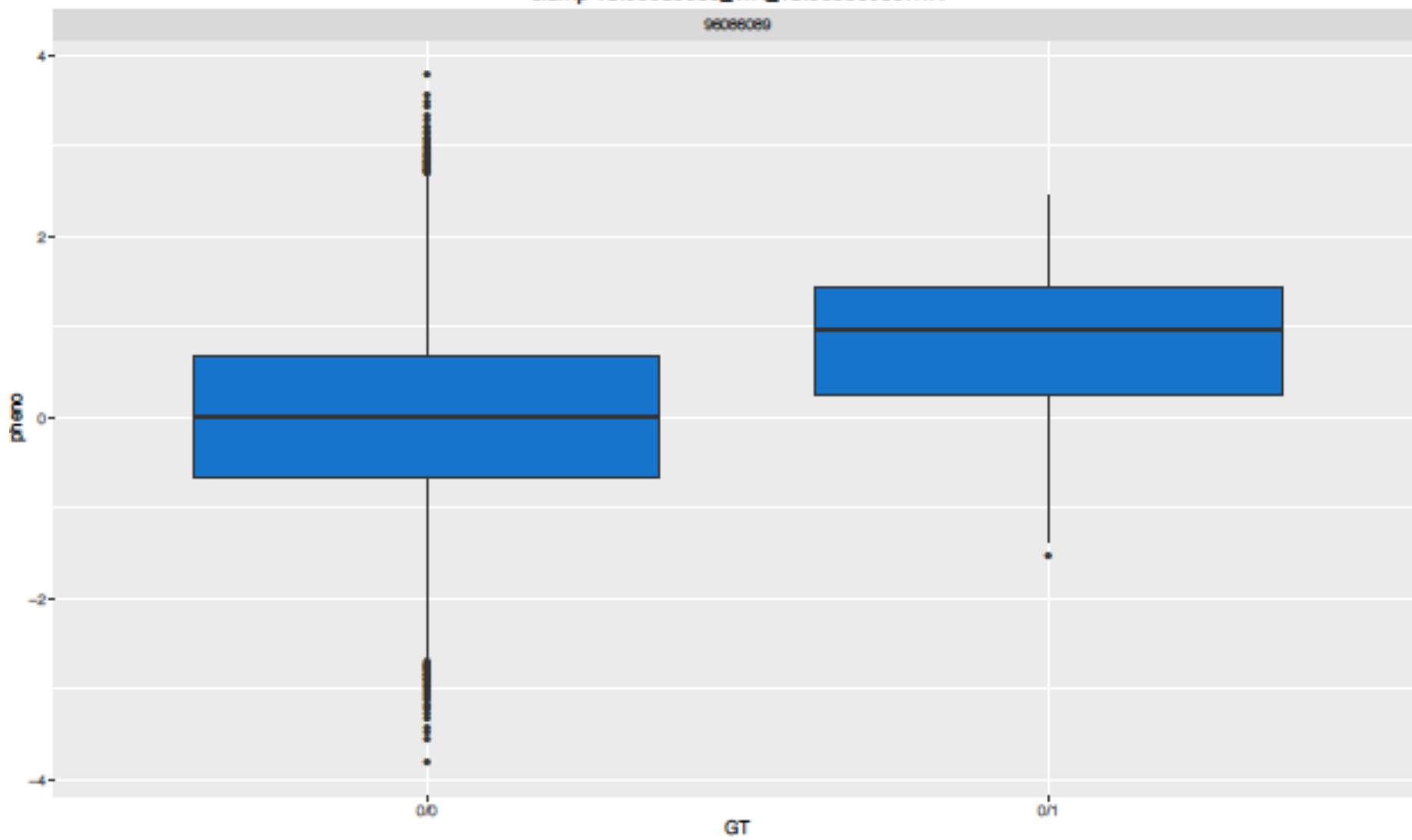
Claudins

- integral membrane proteins and components of tight junction strands
- Plays a major role in tight junction-specific obliteration of the intercellular space, through calcium-independent cell-adhesion activity
- May form permselective paracellular pores
 - isoform 1 appears to create pores preferentially permeable to cations
 - isoform 2 for anions.
- Plays a key role in controlling cation selectivity and transport in the thick ascending limb (TAL) of Henles loop in kidney.

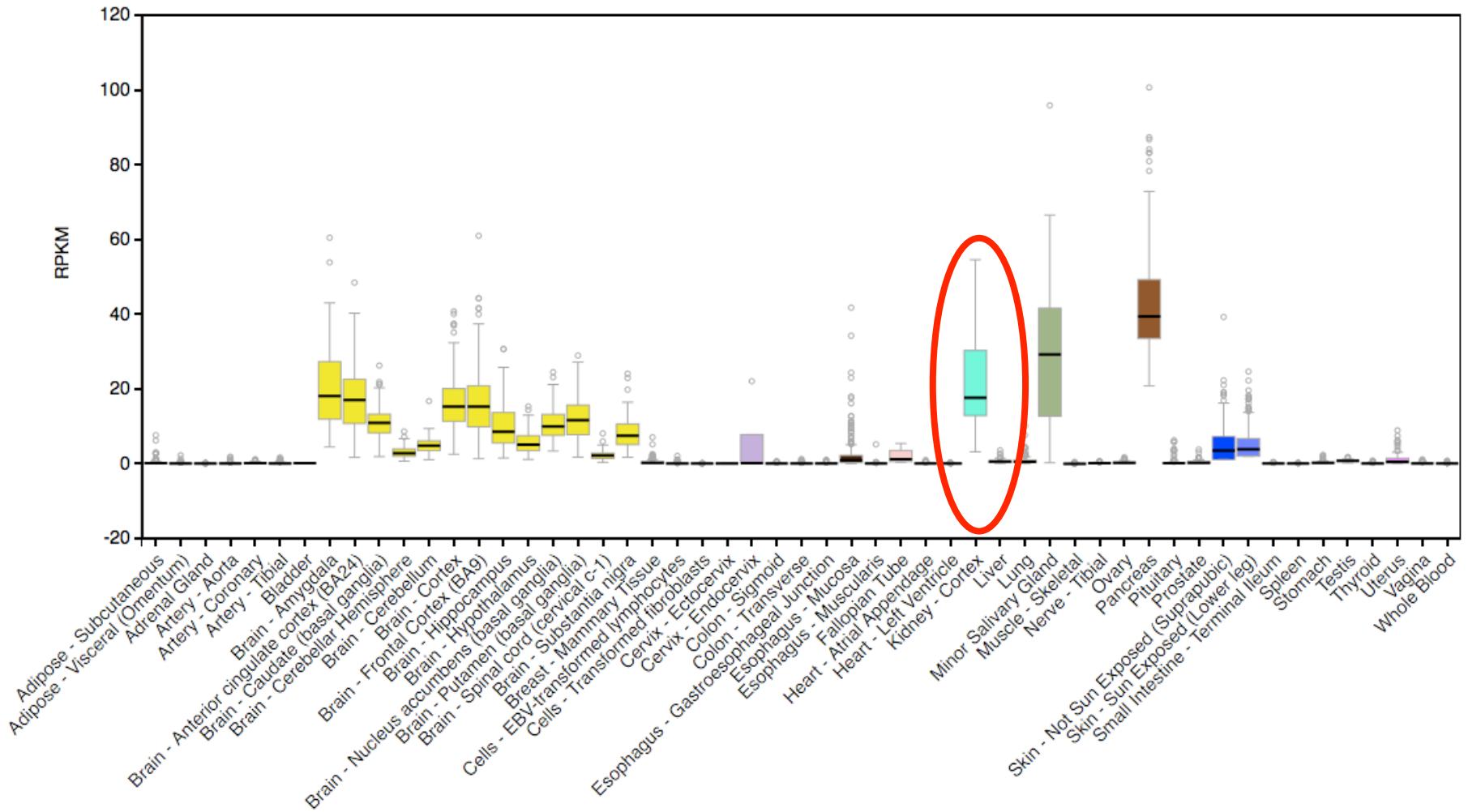




Crea
clump 13:96086089_T/A_13:96086089:T/A
96086089



Gene expression from GTEx



Summary

- Call variants using local de novo assembly (GATK HaplotypeCaller)
- Filter variants using hard filters or VQSR
- Decompose and normalize indels
- Filter variants based on many metrics
 - Ti/Tv, HWE, het/hom, missingness
- Annotate variants
- For small, family based studies GEMINI can be used to analyze inheritance patterns of rare variants
- For large, case-control cohort studies rare variants can be identified using gene based association tests

