

Study designs for genetic association studies

Genotyping vs. sequencing

- Genotyping: Target a particular genetic variant and "measure" it
- Sequencing: Target a region (could be the whole genome) and "measure" the entire region (all base-pairs)
- From an bioinformatic/analysis point of view, genotyping data is much easier to handle.

Genotyping technologies (low-throughput)

Illumina



1500 - 300 SNPs

SNPlex



400 - 40 SNPs

Sequenom



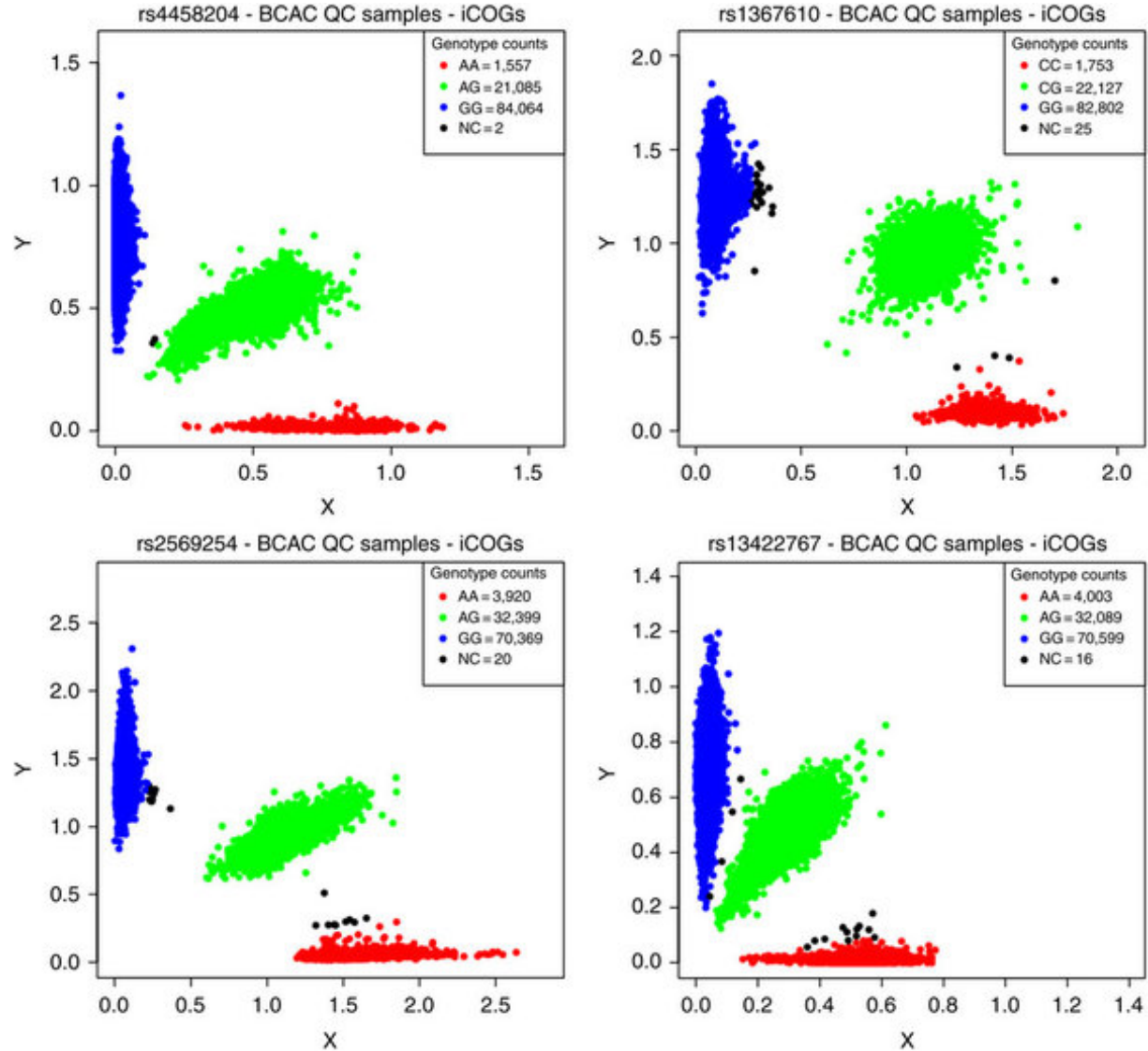
40 - 5 SNPs

TaqMan



10 - 1 SNPs

Genotyping Output



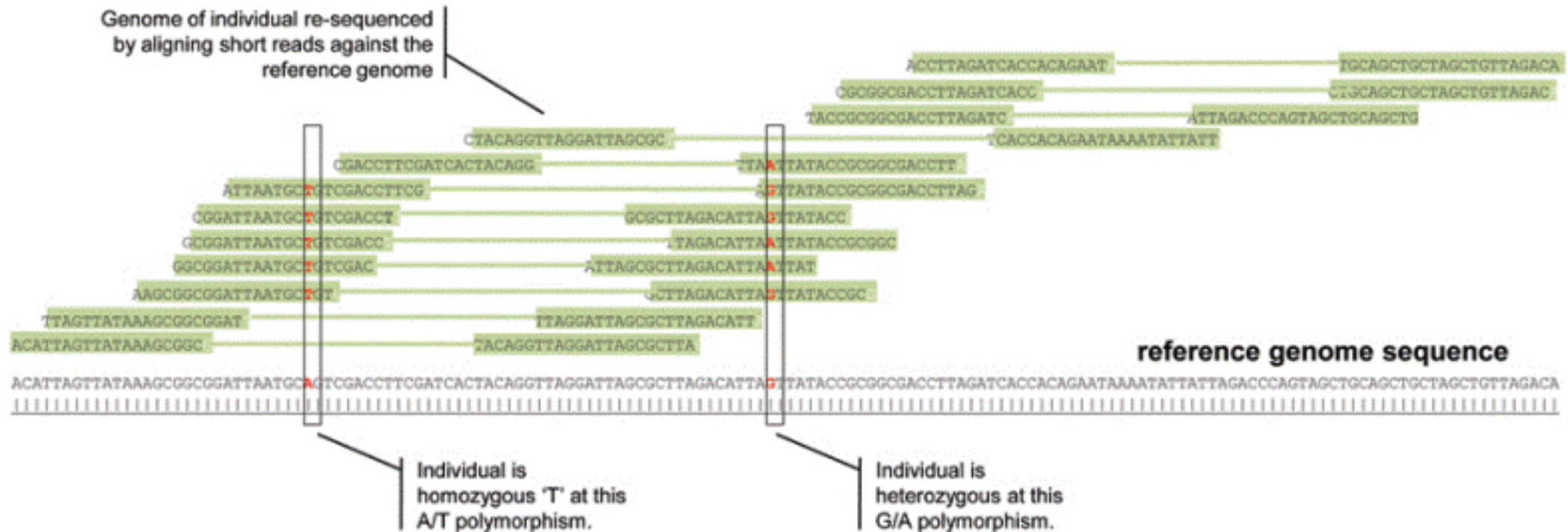
Next-generation sequencing

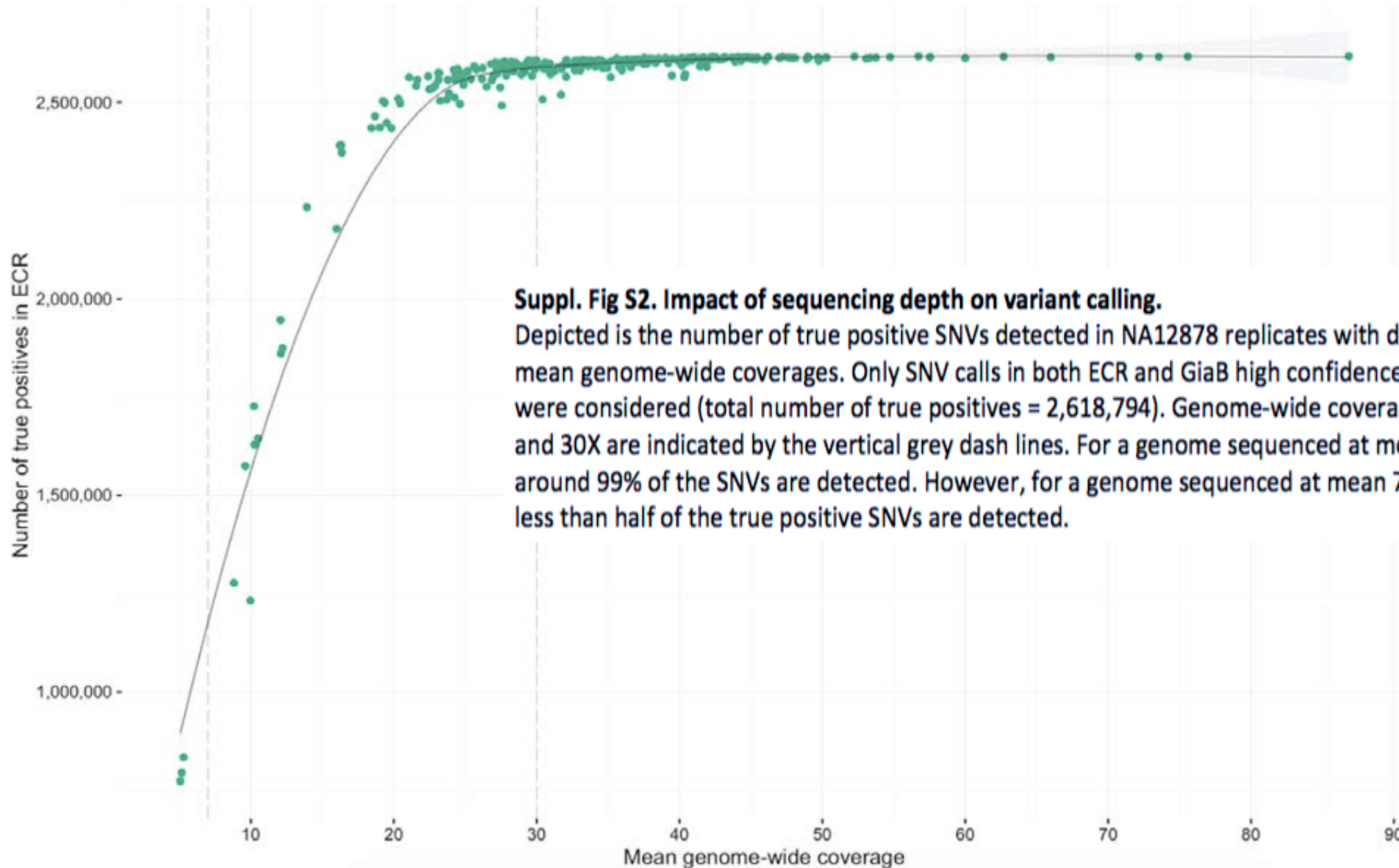
- Capture ALL base-pairs in our region of interest
 - Whole genome sequencing
 - Whole exome sequencing
 - Targeted sequencing
- Allows you to identify any variants that might be unique to your samples
- More expensive than genotyping
- Need more expansive IT support than genotyping



Sequencing alignment and depth

- Depth: The number of times a base-pair is sequenced

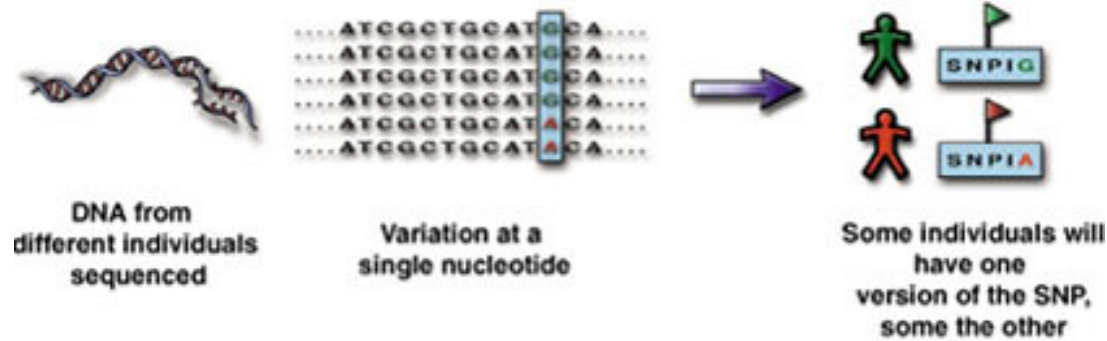




Suppl. Fig S2. Impact of sequencing depth on variant calling.

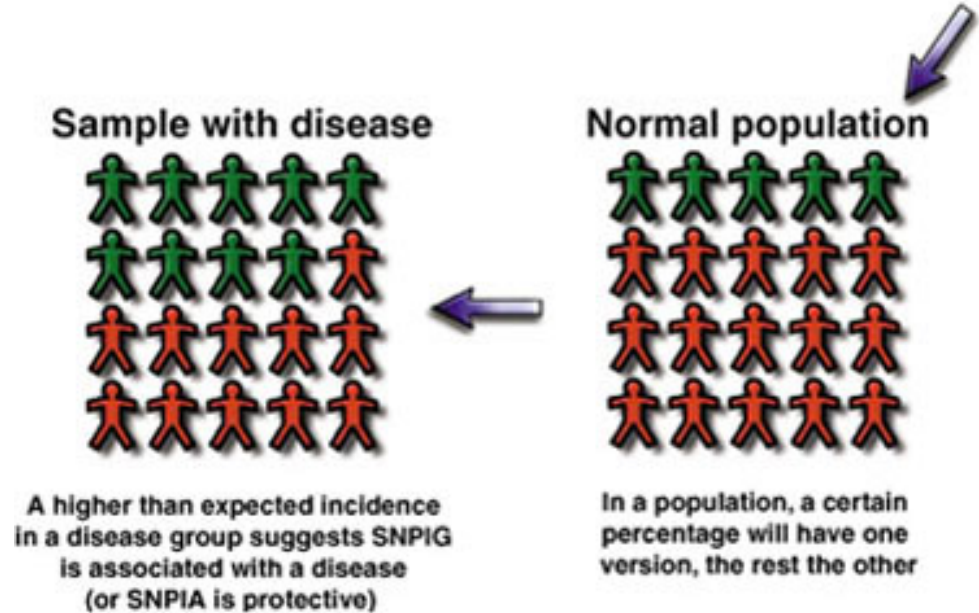
Depicted is the number of true positive SNVs detected in NA12878 replicates with different mean genome-wide coverages. Only SNV calls in both ECR and GiaB high confidence regions were considered (total number of true positives = 2,618,794). Genome-wide coverages of 7X and 30X are indicated by the vertical grey dash lines. For a genome sequenced at mean 30X, around 99% of the SNVs are detected. However, for a genome sequenced at mean 7X coverage, less than half of the true positive SNVs are detected.

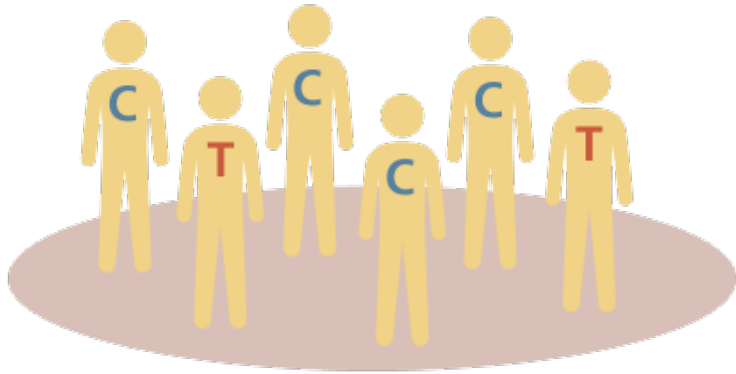
Genetic association studies using SNPs



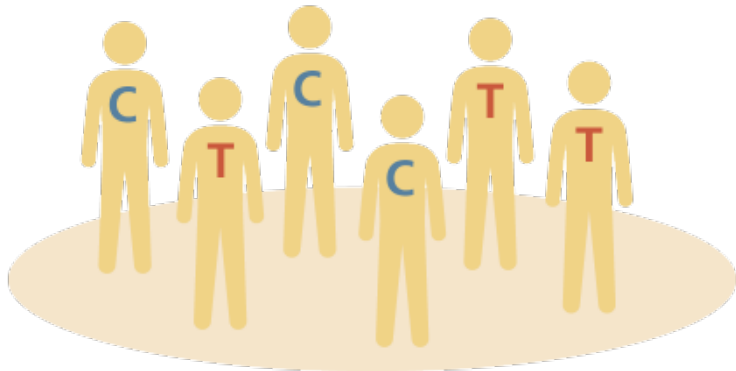
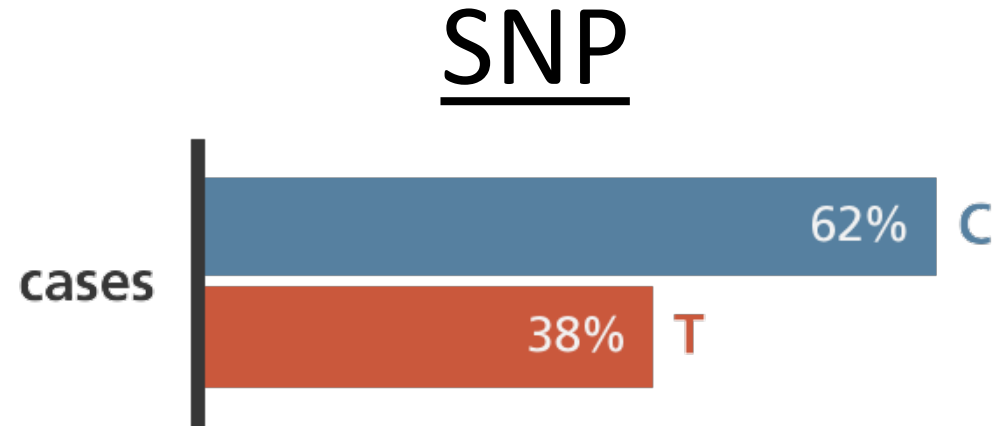
Why we like SNPs:

- Abundant in the genome
- Easy to measure





cases (n=1,000)
people with heart disease



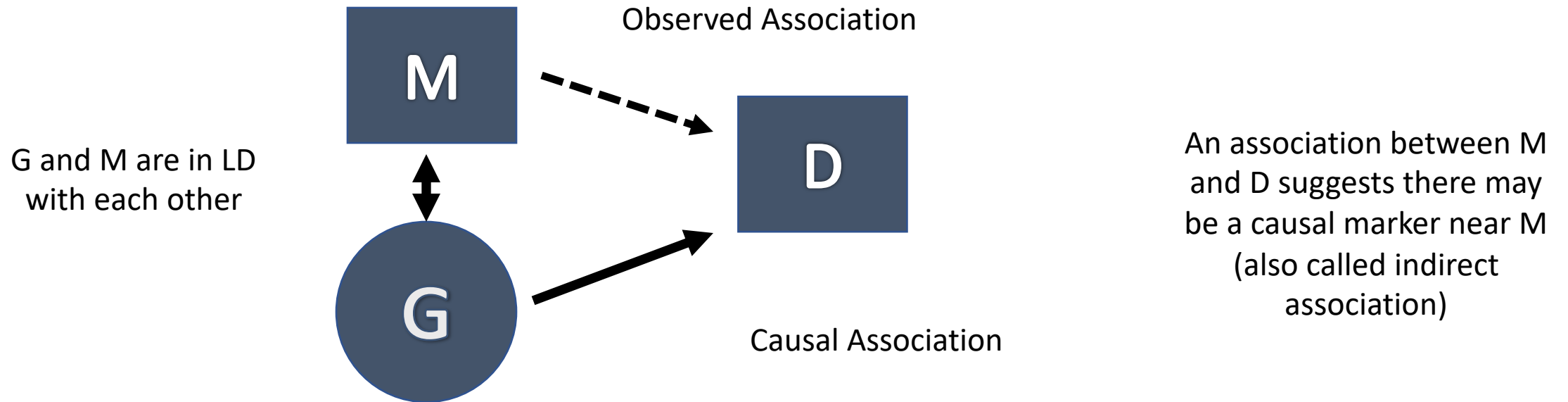
controls (n=1,000)
people without heart disease



In the early days: Candidate gene studies

- Pick your favorite gene
- Map the genetic variation in the gene region
 - Sequence a small population (<100 subjects)
 - The HapMap project provided a map of common genetic variants
- Chose which SNPs to genotype in the entire population
 - If you choose your SNPs carefully, they can explain the majority of genetic variation in the gene (LD!) – also known as “Tagging”
 - Caveat: Rare variation will not be captured (here “rare” often means below 5%)

The use of “tags” (proxy markers)



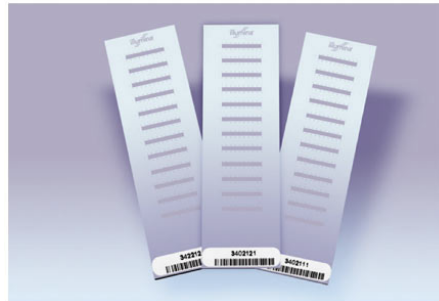
If the r^2 between M and G is 0.5 you need to double your sample size to obtain the same power as if you had measured G directly

Tagging – main idea

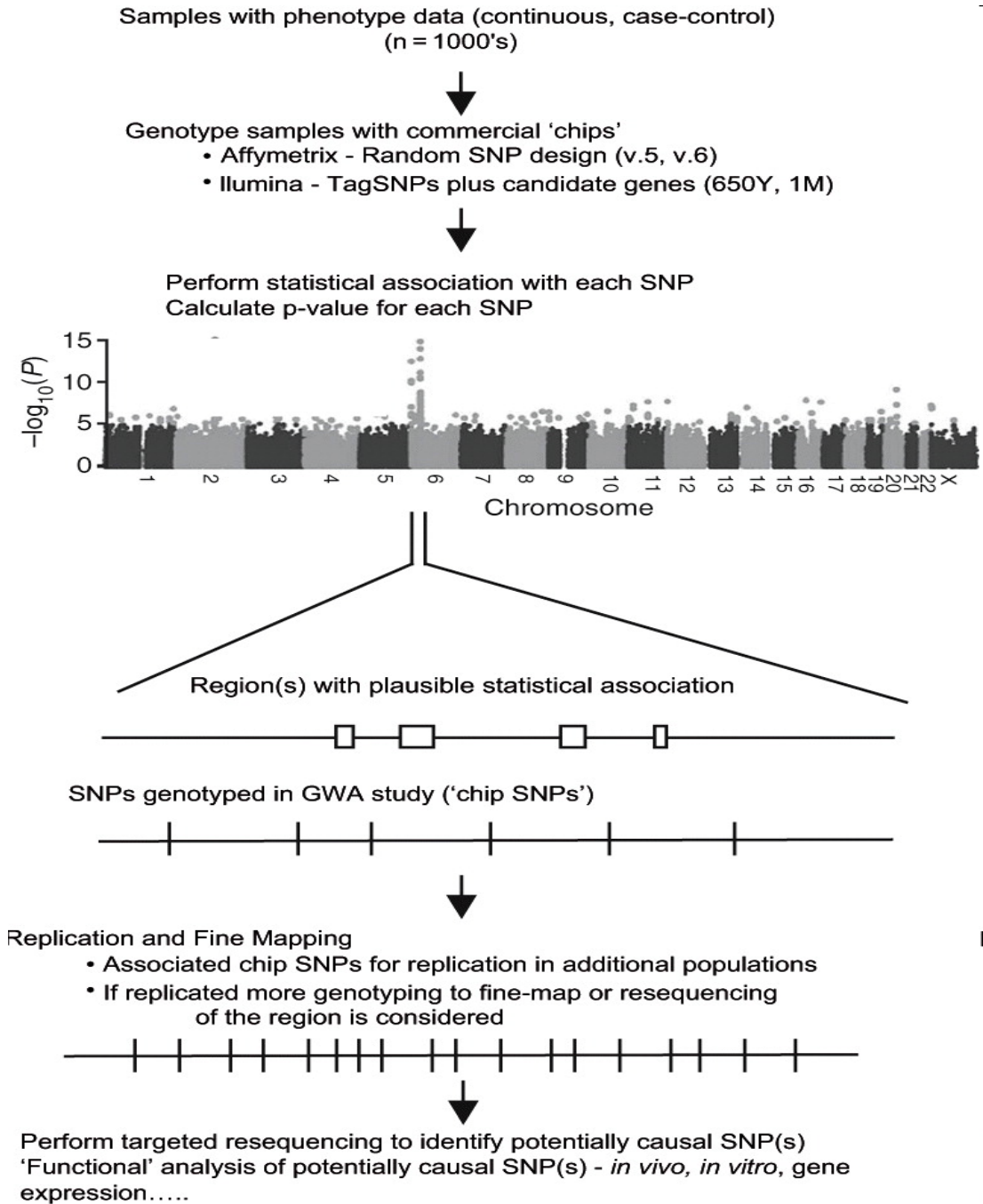
- We can leverage LD and choose non-redundant sets of SNPs that will explain the majority of genetic variation in our region of interest.
- When there is strong LD in a region, we will have very limited loss of power in our association studies even though we are only genotyping a few SNPs.

Genome-wide association studies (GWAS)

Screen the genome for SNPs that are associated with your trait (agnostic approach)

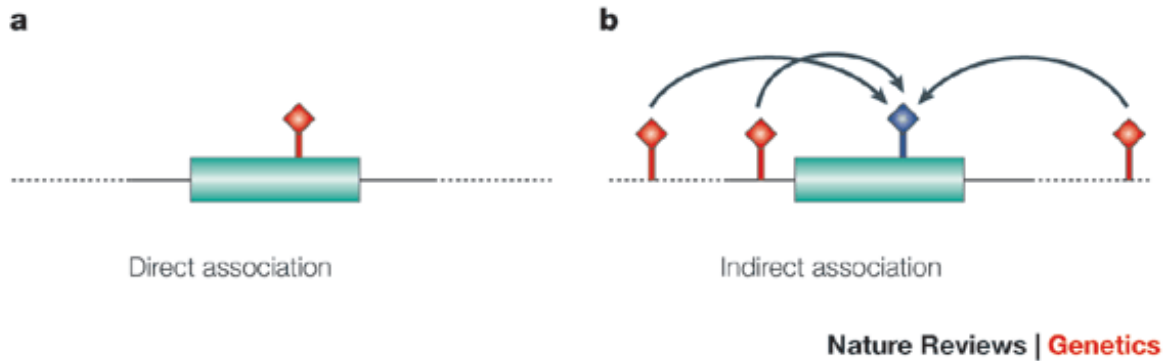


Rieder et al. 2008



GWAS relies heavily on LD

1) Indirect association



2) Imputation

Typical imputation scenario

HapMap or 1,000 Genomes	0	0	1	1	1	0	0	1	1	0	0	1	1	1	Reference haplotypes	
	0	0	0	0	0	1	1	1	0	1	1	1	0	0		1
	1	1	1	1	1	0	0	0	1	0	0	0	0	0		0
	1	0	1	1	0	0	0	1	1	1	1	1	0	0		1
Cases and controls typed on SNP chip	1	?	?	?	2	?	0	?	?	?	?	0	1	?	1	Study genotypes
	1	?	?	?	1	?	0	?	?	?	?	?	0	?	0	
	0	?	?	?	1	?	1	?	?	?	?	1	0	?	1	
	1	?	?	?	2	?	0	?	?	?	?	0	1	?	1	
	?	?	?	?	2	?	0	?	?	?	?	0	0	?	0	
	1	?	?	?	1	?	1	?	?	?	?	1	0	?	?	
	0	?	?	?	2	?	0	?	?	?	?	0	1	?	1	
	1	?	?	?	1	?	1	?	?	?	?	1	1	?	2	

Imputation (I)

- Cost efficient
 - Can assess more SNPs than we genotyped
- Maximize our sample size
 - Fill in missings for already genotyped SNPs
- Allow us to combine data from existing platforms that genotype different SNPs

Imputation (II)

- We can infer genotypes for SNPs we did not genotype (or failed in the lab)
 - **Input:** 550,000 SNPs in 10,000 individuals
 - **Reference panel:** 2,504 individuals from the 1,000 Genomes project (>80M markers excluding singletons)
 - **Output:** Imputed data for >80M markers for your 10,000 individuals

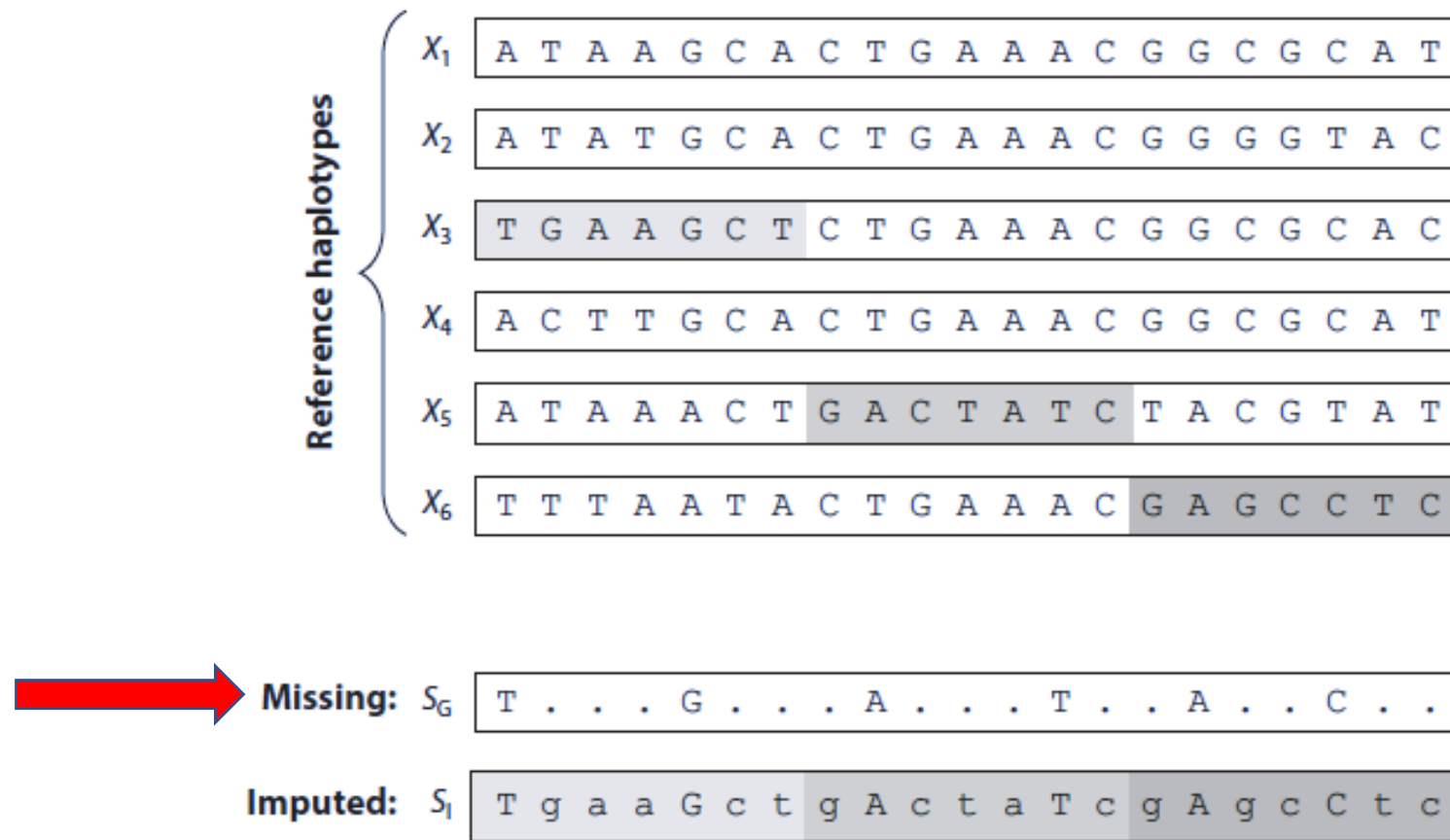


Figure 2

An illustration of genotype imputation, showing the process of imputation for a study haplotype (S_G) genotyped at 6 markers using a reference panel of sequenced haplotypes at 21 markers. The alleles in S_G are used to match short segments from the reference panel. For example, in the first genomic segment, the alleles T and G imply that the corresponding segment might have been copied from haplotype X_3 . In the second segment, the alleles A and T imply that haplotype X_5 might have been copied. Proceeding similarly, the study haplotype can be represented as a mosaic of DNA segments from haplotypes X_3 , X_5 , and X_6 . Consequently, the missing sites can be imputed to obtain the final imputed haplotype, S_I .

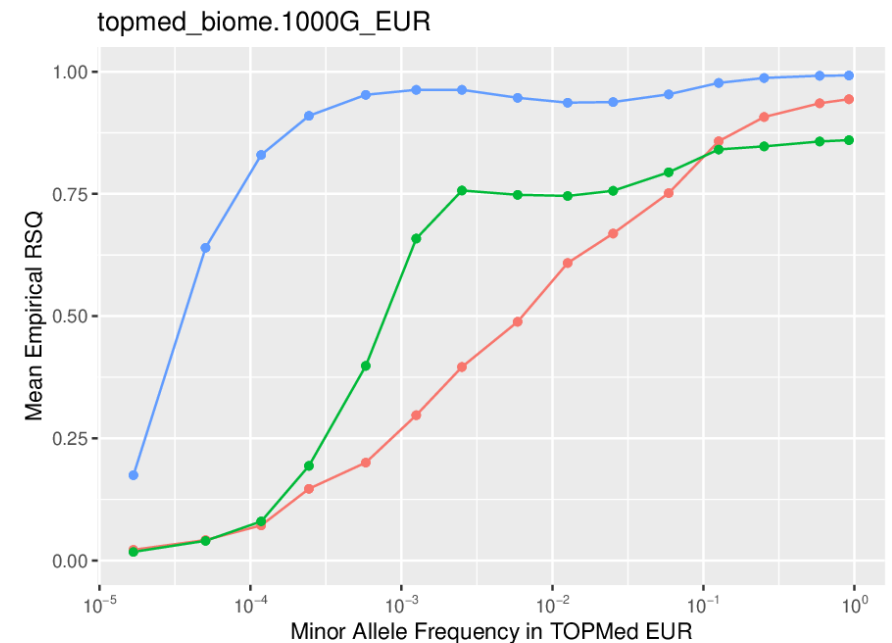
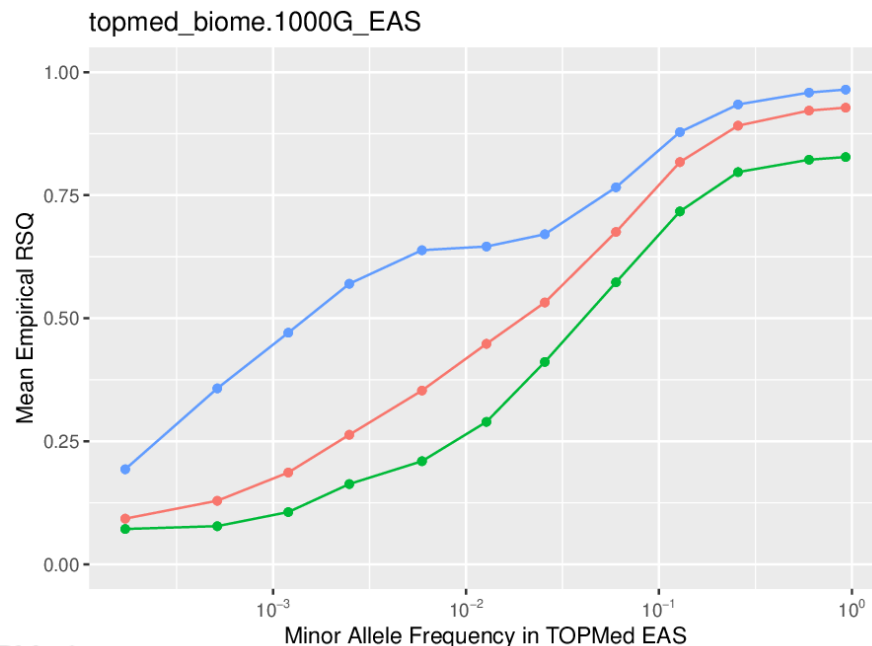
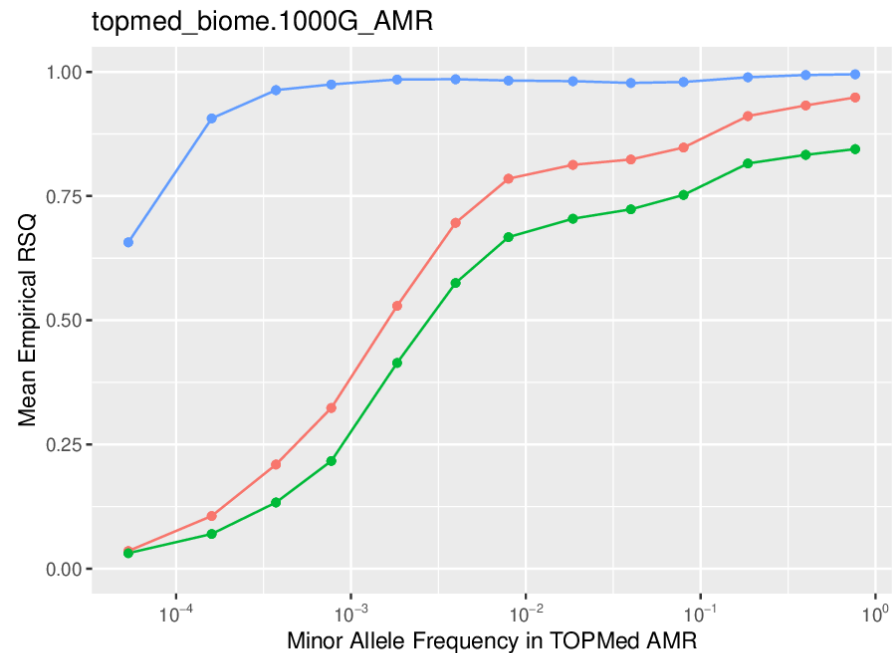
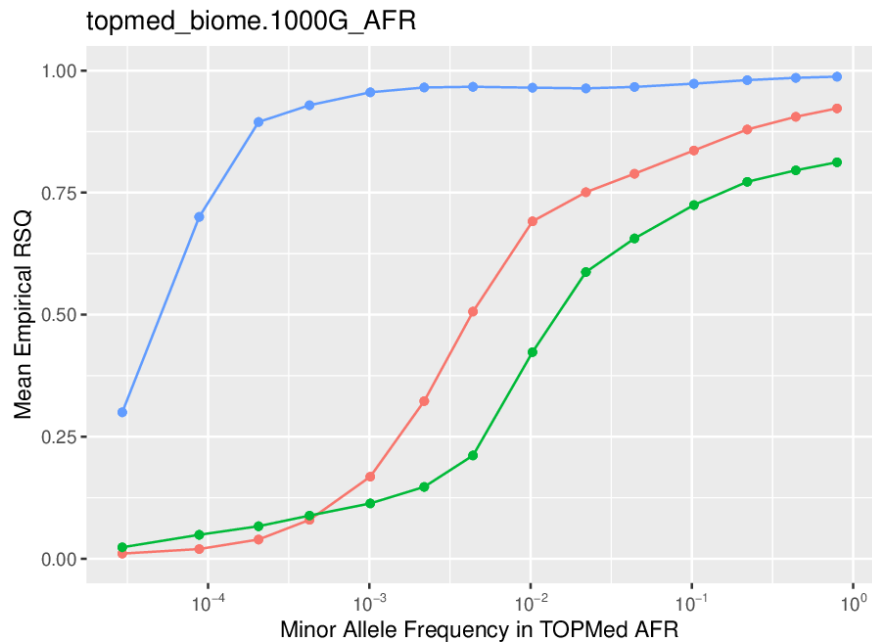
Imputation (III)

- Many imputation algorithms employ a Hidden Markov Model (HMM) method to compare the set of genotypes for each individual in your study to the haplotypes in the reference panel in order to resolve the untyped SNPs.
- Software: MACH, minimac, IMPUTE2, Beagle, PLINK
- Outputs:
 - Posterior probabilities for each potential genotype with three data points per SNP/individual [IMPUTE and BEAGLE]
 - “Dosage” of each imputed genotype ranging between 0-2, representing copies of the reference allele (continuous number) [MACH and BEAGLE].

Imputation (IV)

- The imputation quality score r^2 measures how well a SNP was imputed.
 - Ranges between 0 and 1.
 - Typically, a cut-off of 0.30 or so will flag most of the poorly imputed SNPs, but only a small number (<1%) of well imputed SNPs.
- Factors that affect imputation quality:
 - Number of genotyped SNPs in your data
 - Size of reference panel
 - Similarity in genetic ancestry between reference and study samples
 - Allele frequency

Reference Panels	N	Ancestry
HapMap	60	EUR
1000 Genomes Phase 3	2,504	Mixed
CAAPA	883	African American
HRC	32,470	EUR
TopMed	97,256	Mixed



Panel ● 1000G ● HRC ● TOPMed

COMMENT

CITIES To inform policy, urban scholarship must get organized and funded [p.165](#)

HISTORY A biography of Enrico Fermi, Italy's fallible atomic physicist [p.168](#)



POLITICS The causes Einstein championed offer a window on his time [p.170](#)

OBITUARY Roger Yonchien Tsien, fluorescent-biology pioneer, remembered [p.172](#)



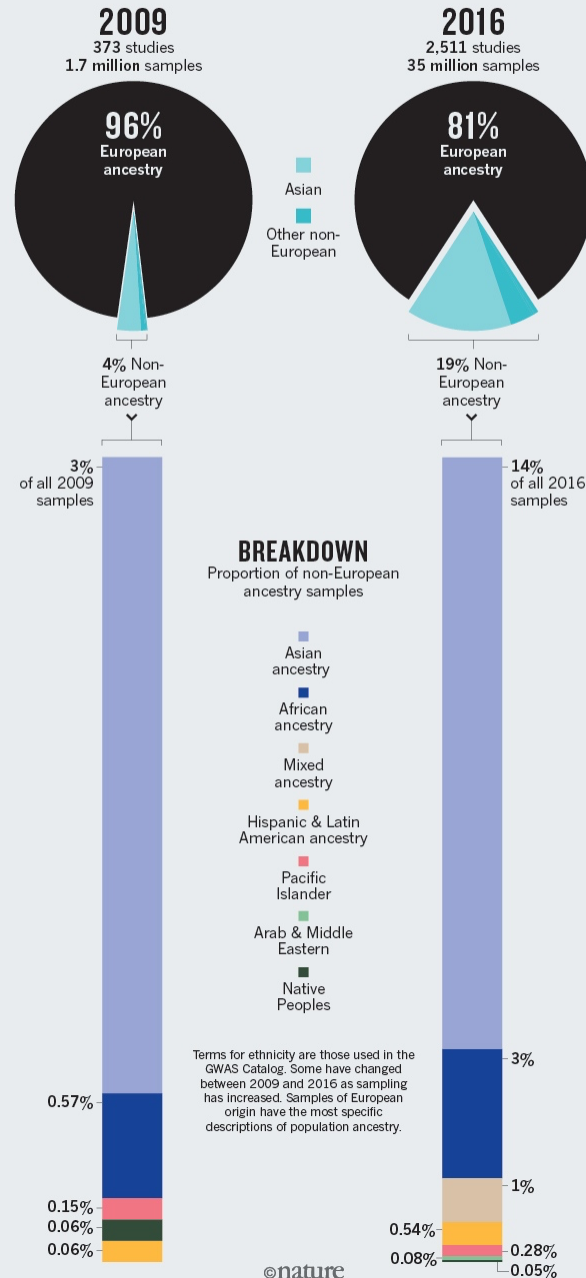
Certain drugs may be less effective, or even unsafe, in some populations because of genetic differences.

Genomics is failing on diversity

An analysis by Alice B. Popejoy and Stephanie M. Fullerton indicates that some populations are still being left behind on the road to precision medicine.

PERSISTENT BIAS

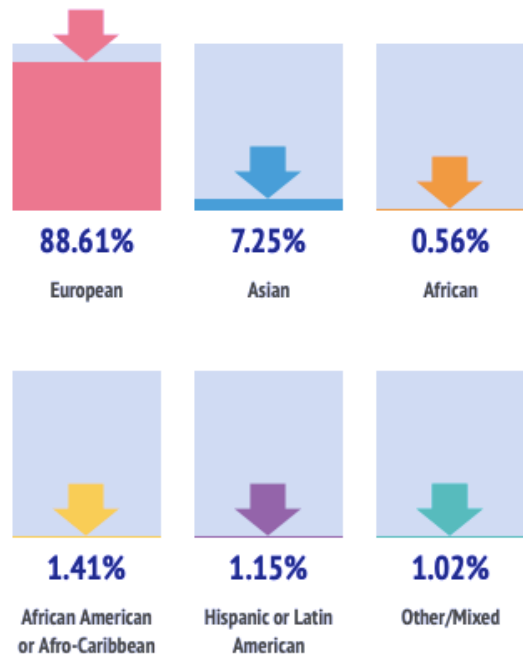
Over the past seven years, the proportion of participants in genome-wide association studies (GWAS) that are of Asian ancestry has increased. Groups of other ancestries continue to be very poorly represented.



Popejoy and Fullerton, Nature 2016

Total GWAS participants diversity

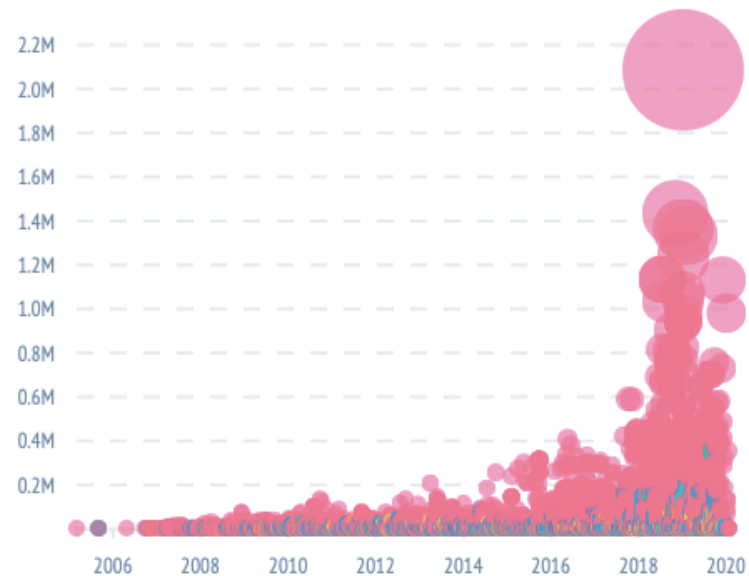
Version 1.0.0. Last check for data: 2020-04-02 06:03:43 .



Ancestry over time by parent term

Discovery Stage

All parent terms OR Search for one or more traits



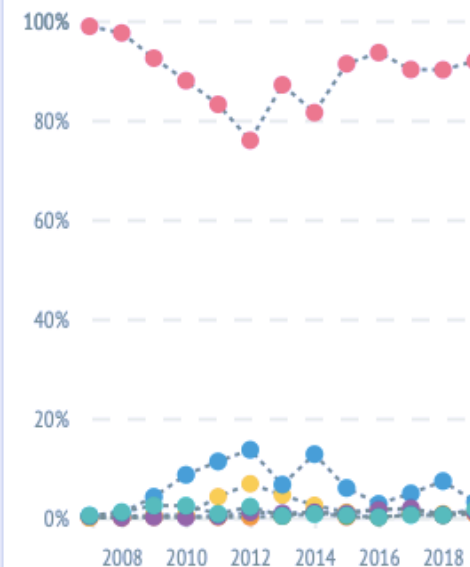
VIEW ALL

- European
- Asian
- African
- African American or Afro-Caribbean
- Hispanic or Latin American
- Other/Mixed

Participants across all parent terms

Discovery Stage

All ancestries Include not recorded

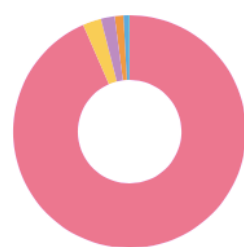


Participants by ancestry

Discovery Stage

Click to show associations discovered

Cardiovascular disease



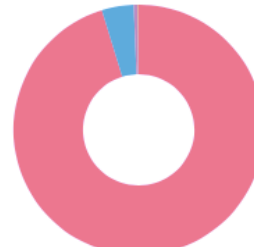
- | | |
|-----------------|--|
| European 93.46% | African American or Afro-Caribbean 2.57% |
| Asian 0.76% | Hispanic or Latin American 1.91% |
| African 1.21% | Other/Mixed 0.1% |

Participants by ancestry

Discovery Stage

Click to show associations discovered

Cancer



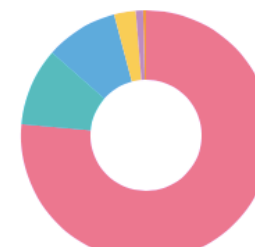
- | | |
|-----------------|---------------------------------------|
| European 95.25% | African American or Afro-Caribbean 0% |
| Asian 4.11% | Hispanic or Latin American 0.52% |
| African 0.11% | Other/Mixed 0% |

Participants by ancestry

Discovery Stage

Click to show associations discovered

Metabolic disorder



- | | |
|-----------------|--|
| European 76.43% | African American or Afro-Caribbean 2.81% |
| Asian 9.4% | Hispanic or Latin American 0.94% |
| African 0.38% | Other/Mixed 10.05% |

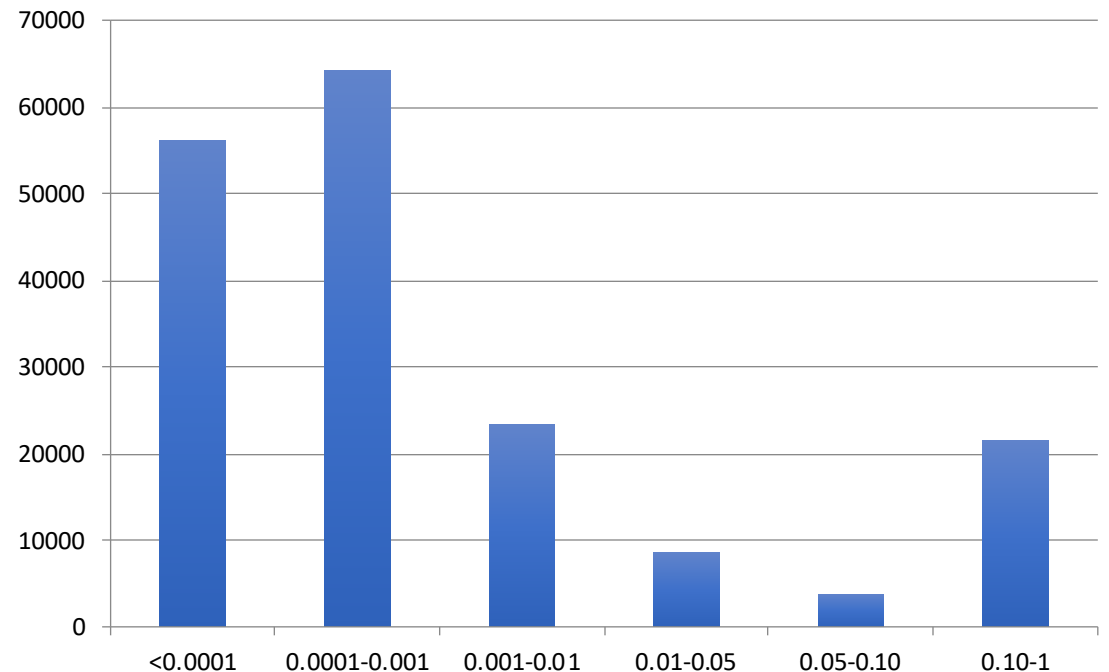
Breakout Room Discussion:

- Explore the racial/ethnic breakdown of GWAS as reported on the website <https://gwasdiversitymonitor.com>. What do you notice about recent trends? What populations seem over- and under-represented in genetic studies?
- What are your ideas for how we can we increase the diversity of study participants in genetic epidemiology?

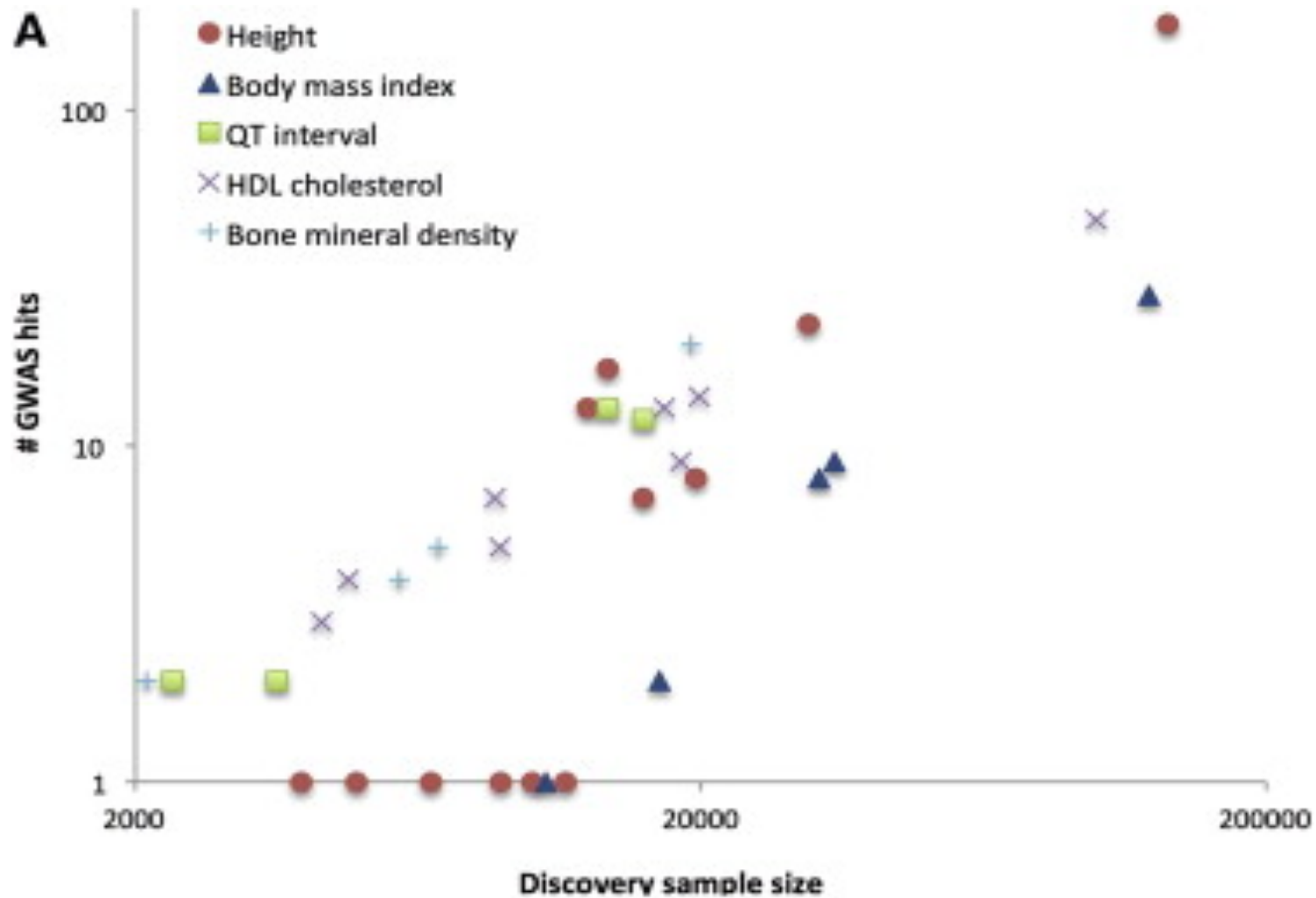
The exome array (~240,000 genetic variants)

- Design based on exome and whole-genome sequencing data from across the world (at the time mostly unpublished data)
 - 9000 samples of European ancestry, 2000 samples of African ancestry, 500 samples each of Hispanic and Asian ancestry

MAF distribution of Exome array data in the Women's Genomic Health Study, n=22,618 (~58,000 variants were monomorphic)

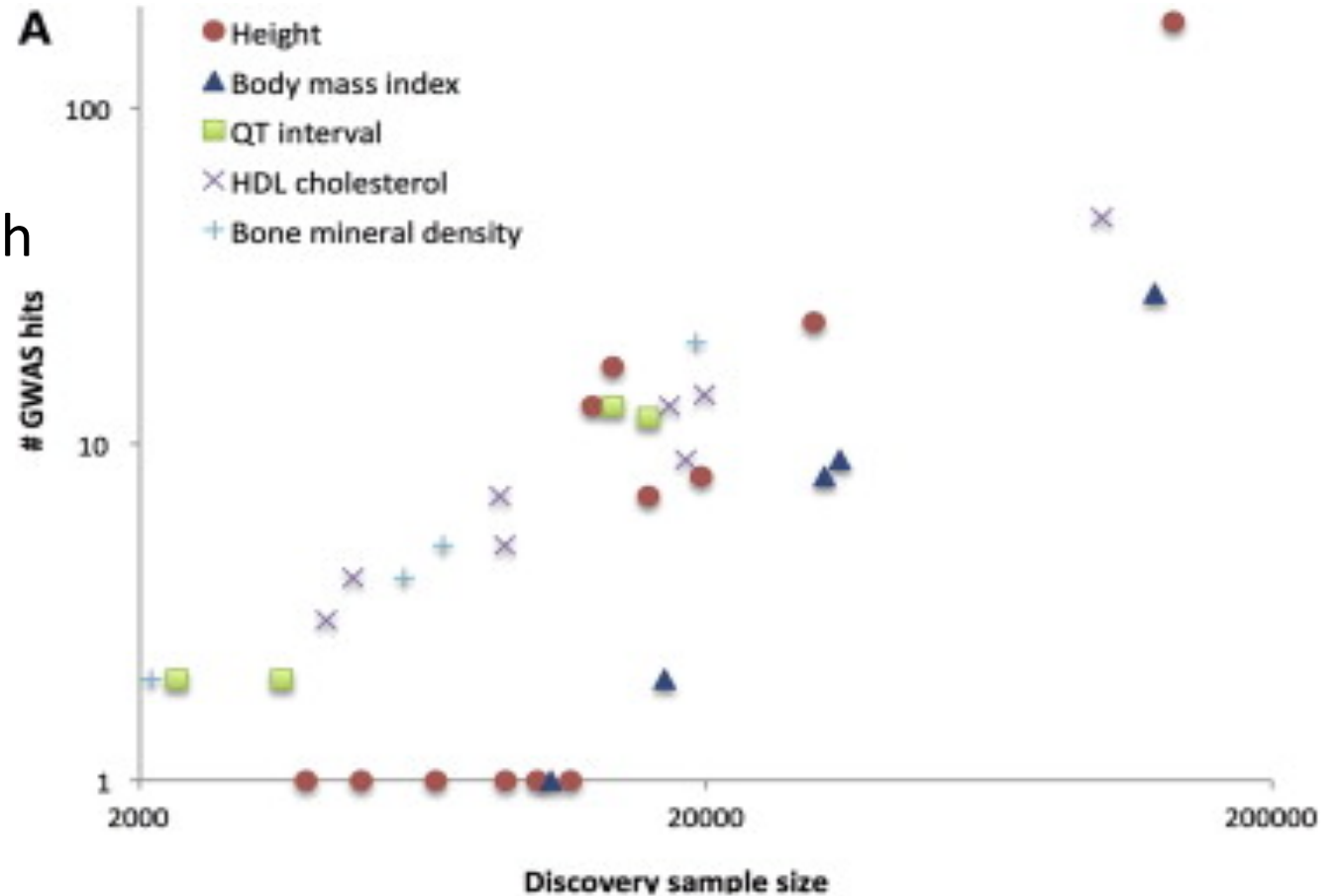


Customized large-scale genotyping arrays



Customized large-scale genotyping arrays

- Idea: Can we design a custom array with 100,000s of SNPs and reduce the price if we commit to genotyping MANY subjects?
- Cost of these arrays are approximately 20% of GWAS arrays, thus enabling far more subjects to be genotyped. Genotyping using a uniform array has also enabled direct comparison across phenotypes.



Customized large-scale genotyping arrays

- **MetaboChip**
 - Custom array designed to test ~200,000 SNPs of interest for metabolic and cardiovascular disease traits.
 - Genotyped in > 100,000 subjects
- **ImmunoChip**
 - Custom array designed to test 195,806 SNPs for immune-mediated diseases.
 - Genotyped in > 150,000 subjects
- **Cardiochip**
 - Custom array that contains 50,000 SNPs across 2,000 genes associated with cardiovascular disease.
 - Genotyped in > 210,000 subjects
- **OncoArray**
 - Custom array designed to test ~500,000 SNPs related to multiple cancers: breast, colorectal, lung ovary and prostate.
 - Genotyped in > 400,000 subjects

Combination arrays

- Emerged over the last few years
 - Includes both GWAS and exome array SNPs
 - Often allows for custom content
 - Target biobanks (e.g. UK Biobank)

Pricing (CIDR, Apr 2020)

Affymetrix Genotyping - GWAS and Custom	
UK Biobank 821K Axiom Array	~\$150 - \$210 Inquire for pricing, sample number dependent
Custom Array (up to 750K SNPs)	~\$180 - \$240 Inquire for pricing
Custom Array (up to 50K SNPs)	~\$120 - \$170 Inquire for pricing

Illumina Genotyping – GWAS					
Omni5		Omni 2.5		Multi_Ethnic	
Omni 5/plus exome	\$450/\$470	Omni 2.5/plus exome	\$310/\$330	Global	\$130-\$165
OmniExpress		Global/Asian Screening Array		Core	
OmniExpress/plus Exome	\$170/\$190	Global/Asian Screening	\$75-\$100	Core/plus Exome	\$95/\$115
PLUS OPTIONS: Custom content can be added to most GWAS and Consortium arrays. Please Inquire for pricing. If FFPE DNA Source cost increases by \$100 per sample.					

Pricing Sequencing (CIDR, Apr 2020)

Illumina Sequencing		
Whole Genome, low pass 4X*		Inquire for pricing
Whole Genome (30X)	>25 samples	\$1,200-\$1,600 sample number dependent
Whole Exome (51 Mb)	>90% @ 10X to >95% @ 20X	~\$475 - \$800 coverage and sample number dependent
Whole Exome (71 Mb) Includes UTRs	>90% @ 10X to >95% @ 20X	\$540 - \$1,000 coverage and sample number dependent
Whole Exome PLUS (plus 6.8-24 Mb custom regions)		Inquire for pricing
Custom Targeted (500 kb – 34 Mb options)		~\$200 - \$1,500
Custom Targeted (amplicon; 10-250kb)		~\$80-~\$200
* Please Inquire for other options. If FFPE DNA Source, costs increase ~ 25%. Somatic variation calling is supported for tumors.		