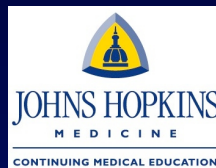


# **Genome-wide association studies**

**Karen Mohlke, PhD  
Department of Genetics  
University of North Carolina**



*Current Topics in Genome Analysis 2012*

*Karen Mohlke, PhD*

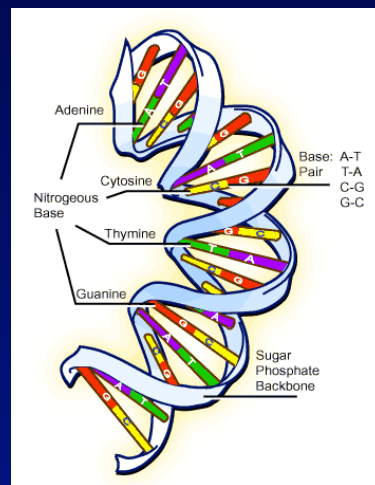
*No Relevant Financial Relationships with  
Commercial Interests*

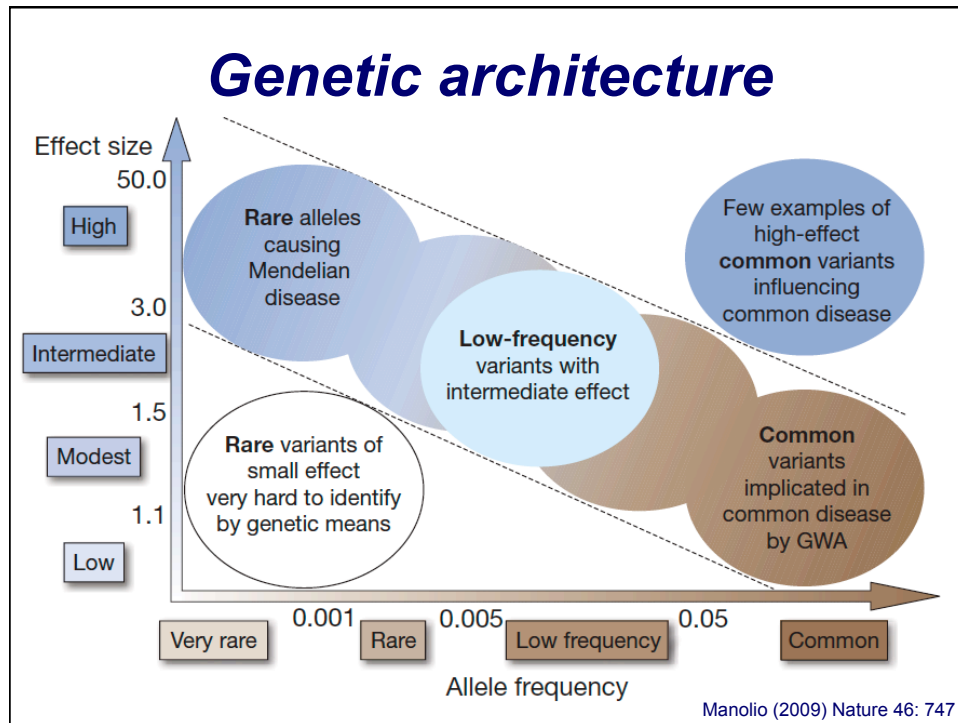
## Complex traits



## Common and rare variants


GGATTCAC**T**GCAAAATCG  
GGATTCAC**T**GCAAAATCG  
GGATTCACAGCAAAATCG  
GGATTCAC**T**GCAAAATCG  
GGATTCAC**T**GCAAAATCG  
GGATTCAC**T**GCAAAATCG  
GGATTCAC**T**GCAAAAT**G**G  
GGATTCACAGCAAAATCG  
GGATTCACAGCAAAATCG  
GGATTCAC**T**GCAAAATCG





## Genome-wide association (GWA)

- What is the goal?
- How are studies performed?
- What can we learn from the associated regions?
- What do the findings tell us about disease?



## **GWA Studies**

- **Benefits of GWA vs classical mapping**
  - More powerful vs linkage for common, low penetrance variants
  - Better resolution than linkage
  - No need to select candidate genes
- **Requirements of GWA**
  - Catalog of human genetic variants
  - Low cost, accurate method for genotyping
  - Large number of informative samples
  - Efficient statistical design and analysis

## **Goals of a GWA study**

- **Test a large portion of the common single nucleotide genetic variation in the genome for association with a disease or variation in a quantitative trait**
- **Find disease/quantitative trait-related variants without a prior hypothesis of gene function**

## ***Steps in a GWA study***

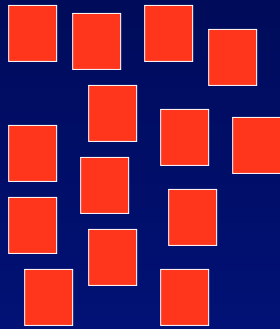
- **Samples**
- **Genotyping**
- **Quality control**
- **Statistical analysis**
- **Replication**

## ***Phenotype***

- **Disease (case/control)**
  - Rare
  - Common
- **Quantitative trait**
  - Easy to measure: Weight, height
  - Requires testing: Coronary artery thickness
  - Requires experiment: Gene expression

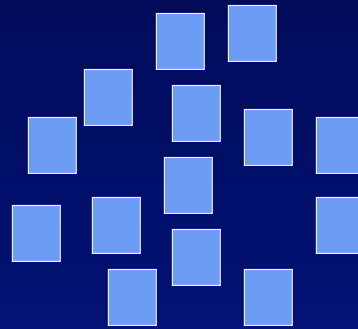
## ***Selection of cases and controls***

### **Cases**



**Definition of case?**

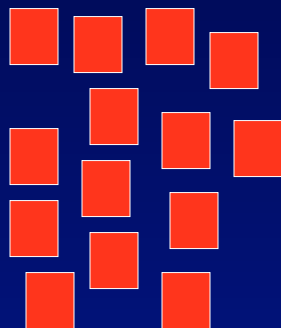
### **Controls**



**Definition of control?**

## ***Selection of cases***

### **Cases**

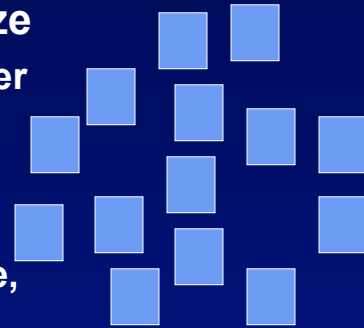


- **Potential criteria to enrich genetic effect size**
  - More severely affected individuals
  - Require other family member to have disease
  - Younger age-of-disease onset

## Selection of controls

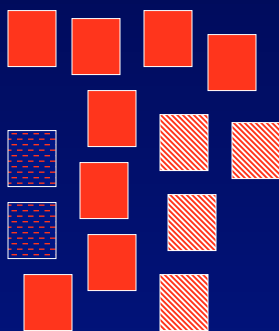
- Potential criteria to enrich genetic effect size
  - Low risk of disease rather than population-based samples
  - Same ancestry as cases
  - Matched to cases on age, sex, demographics

### Controls

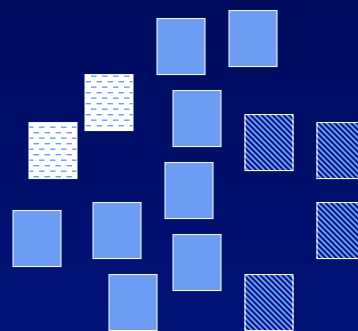


## Matched ancestry

### Cases

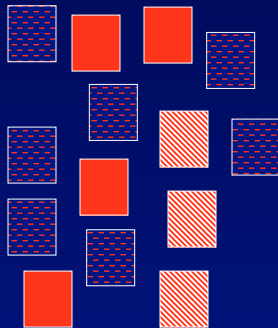


### Controls

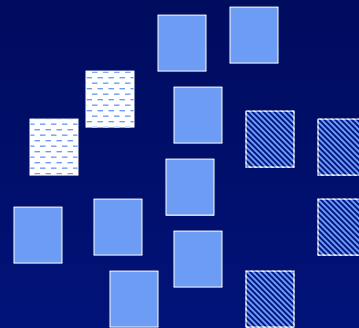


## **Unmatched ancestry**

### **Cases**



### **Controls**



May have inadequate ancestry information prior to genotyping

## **Population stratification and cryptic relatedness**

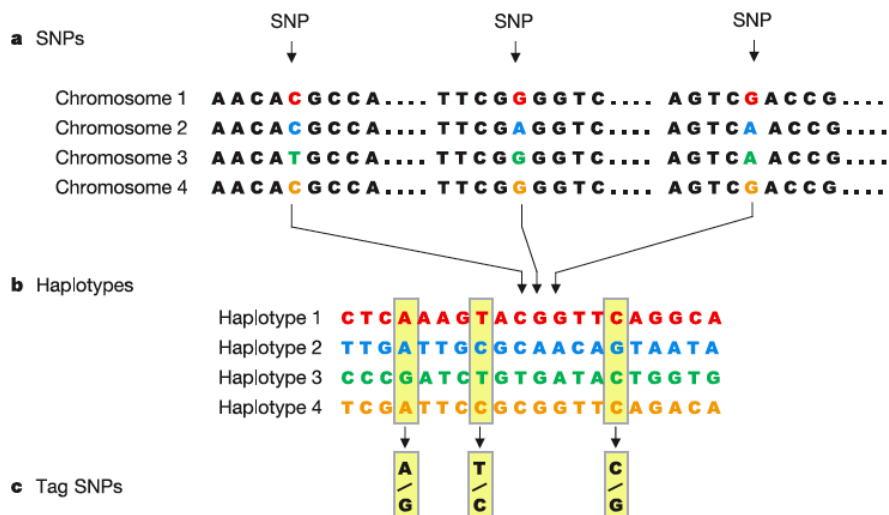
- Can produce spurious associations in case-control studies
- Account for or avoid
  - Genomic control
  - Principle components
  - Family-based study design



## Genome-wide SNP panels

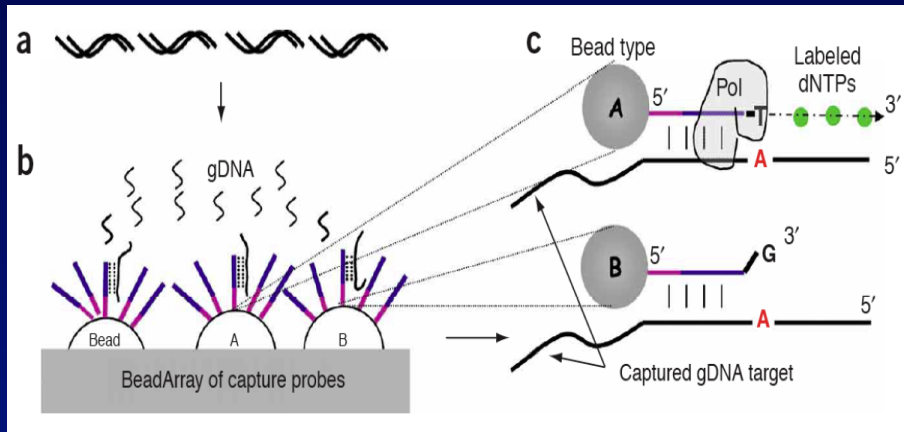
- 10,000 - 5 million SNPs
- Affymetrix, Illumina
  - Random SNPs
  - Selected haplotype tag SNPs
  - Copy number probes
  - Some arrays allow SNPs to be added

### Selecting 'haplotype tag' SNPs



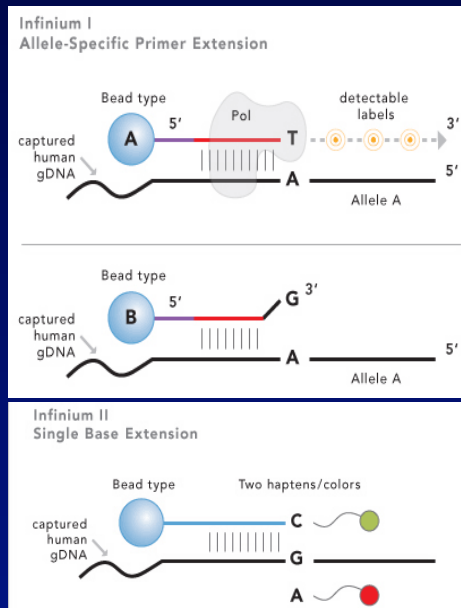
International HapMap Consortium (2003) Nature 426:789

## Illumina Infinium Assay



Gunderson et al. (2005) NatGen 37:549

## Illumina Infinium Assays



Infinium Assay BeadChips enable interrogation of >317,000 to over one million SNPs and offer comprehensive coverage of CNV regions. Shown above, from left to right, are the Human1M, Human150-Duo, and Human550-Duo BeadChips.

Illumina.com

## Affymetrix GeneChip Array

Figure 1: GeneChip® Mapping Assay Overview.

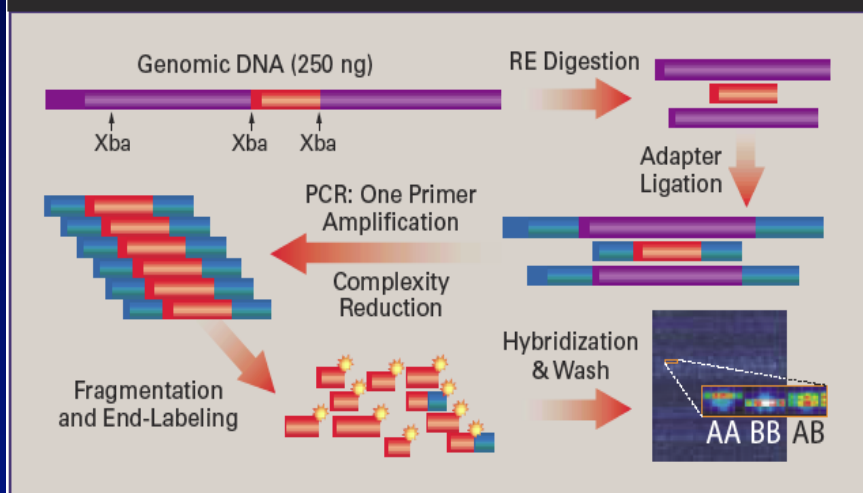
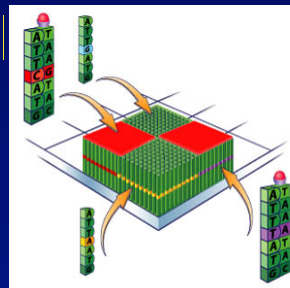
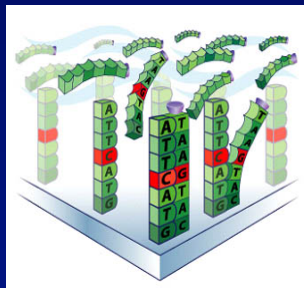
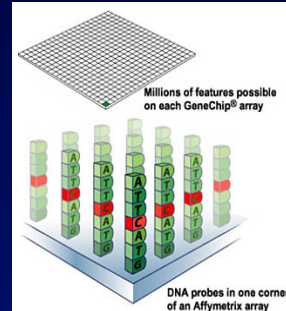
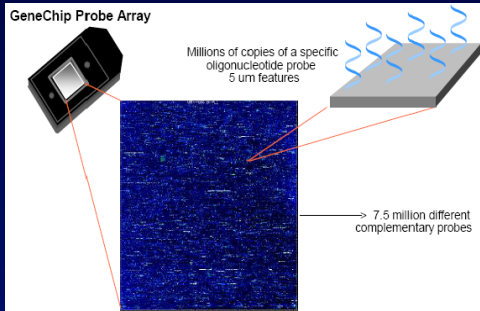
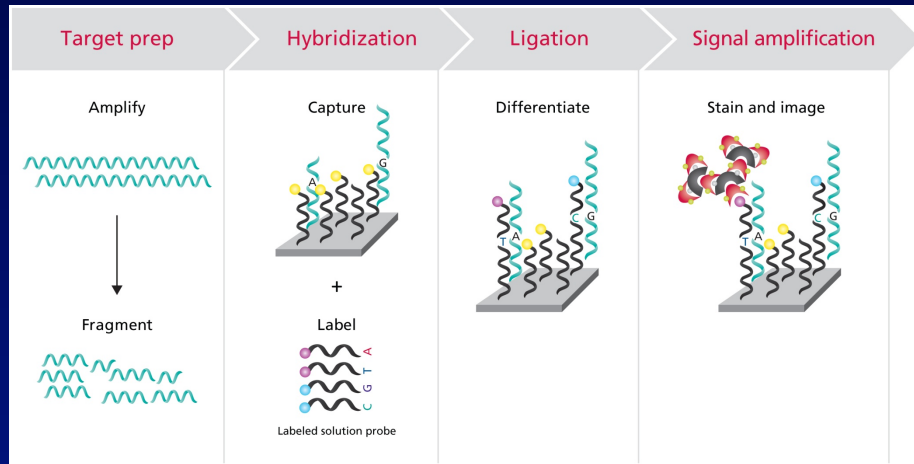


image from affymetrix.com



## Affymetrix Axiom Array



Affymetrix.com

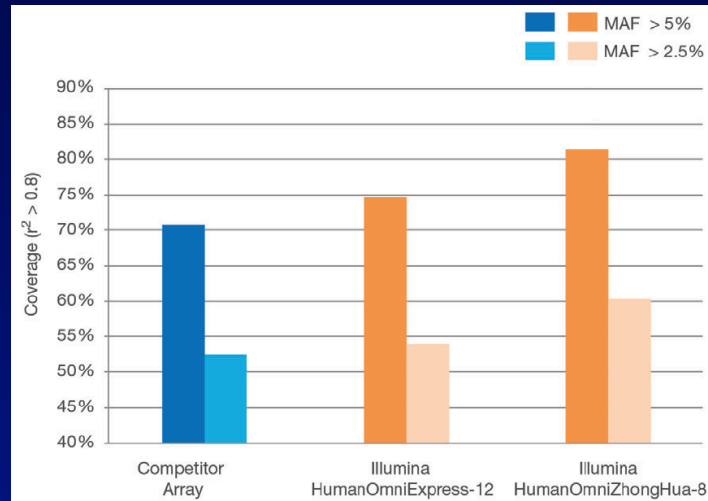
## Global genomic coverage

Global coverage (%) by SNP chips

SNP chip	CEU	CHB+JPT	YRI
SNP Array 5.0	64	66	41
SNP Array 6.0	83	84	62
HumanHap300	77	66	29
HumanHap550	87	83	50
HumanHap650Y	87	84	60
Human1M	93	92	68

Li (2008) EJHG 16:625

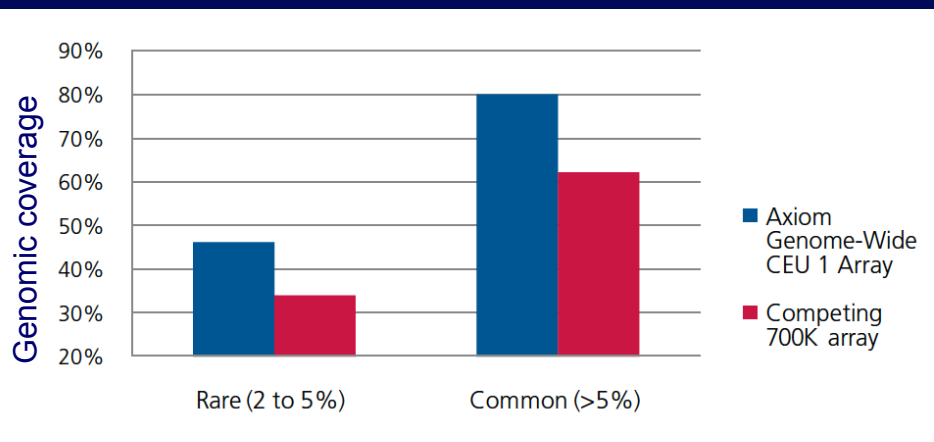
## Newer arrays improve coverage of less common variants



Coverage calculations based on known common and rare Chinese population variants identified in the HapMap and 1000 Genomes projects

Illumina.com

## Newer arrays improve coverage of less common variants

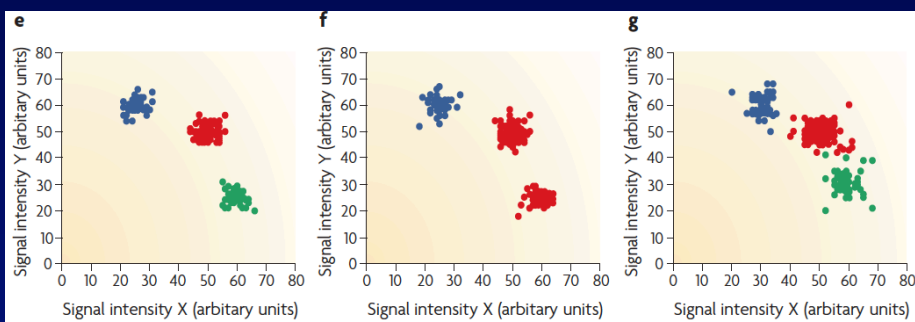


Affymetrix.com

## Quality control: Identify and remove bad samples

- **Poor quality samples**
  - Sample success rate < 95 %
  - Excess heterozygous genotypes
- **Sample switches**
  - Wrong sex
- **Unexpected related individuals**
  - Pair-wise comparisons of genotype similarity
  - Duplicates
- **Ancestry different from the rest of sample**

## Quality control: Identify and remove bad SNPs



Ideal genotyping plot

Clusters mis-called

Clusters overlap

McCarthy (2008) Nat Rev Gen 9:356

## **Quality control: Identify and remove bad SNPs**

- Genotyping success rate < 95%
- Different genotypes in duplicate samples
- Expected proportions of genotypes are not consistent with observed allele frequencies
- Non-Mendelian inheritance in trios
- Differential missingness in cases and controls

## **Test for association**

- Differences between cases & controls

	AA	AC	CC
Case			
Control			

- Ex. Cochran-Armitage test for trend
- Covariates (age, sex, ...)
- Other genetic models

## Odds ratio

- Surrogate measure of effect of allele on risk of developing disease

Allele	A	C	Total
Case	860	1140	2000
Control	1000	1000	2000
Total	1860	2140	4000

Odds of C allele given case status =  $\frac{\text{Case C}}{\text{Case A}}$

Odds of C allele given control status =  $\frac{\text{Control C}}{\text{Control A}}$

$$\text{Odds Ratio} = \frac{\text{Case C} / \text{Case A}}{\text{Control C} / \text{Control A}} = \frac{1140 / 860}{1000 / 1000} = 1.33$$

## Multiple testing

- Genotype and test > 300K – 5M SNPs
- Correct for the multiple tests

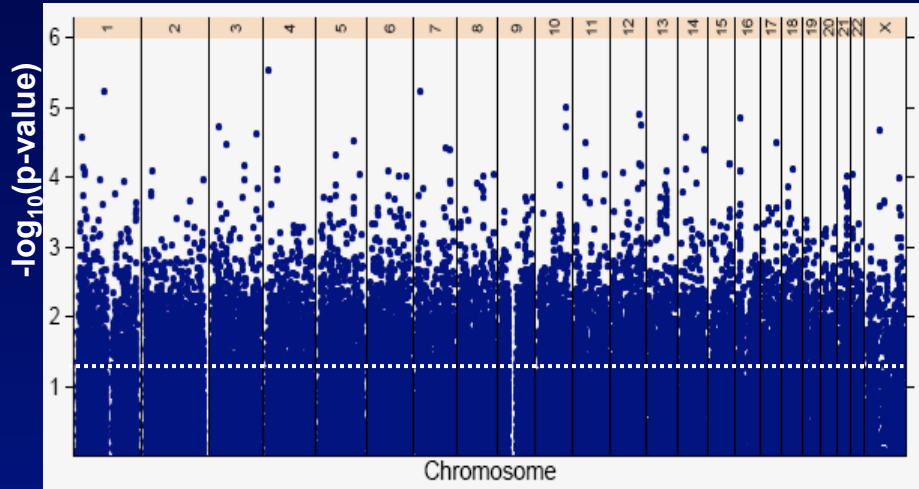
$$\frac{.05 \text{ P-value}}{1 \text{ million SNPs}} = 5 \times 10^{-8}$$

- Need large effect or large sample size



## Type 2 diabetes association results

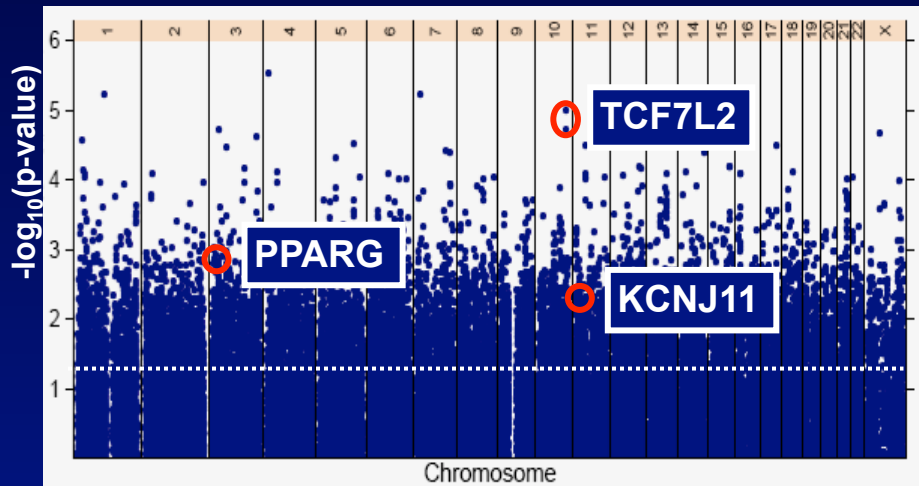
1161 Finnish T2D cases + 1174 Finnish normal glucose tolerant controls



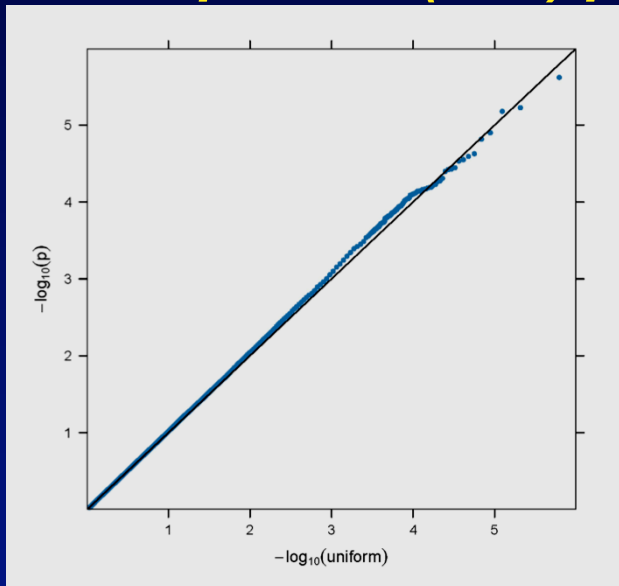
Logistic regression using additive model adjusted for age, gender, birth province

## Which results are true positives?

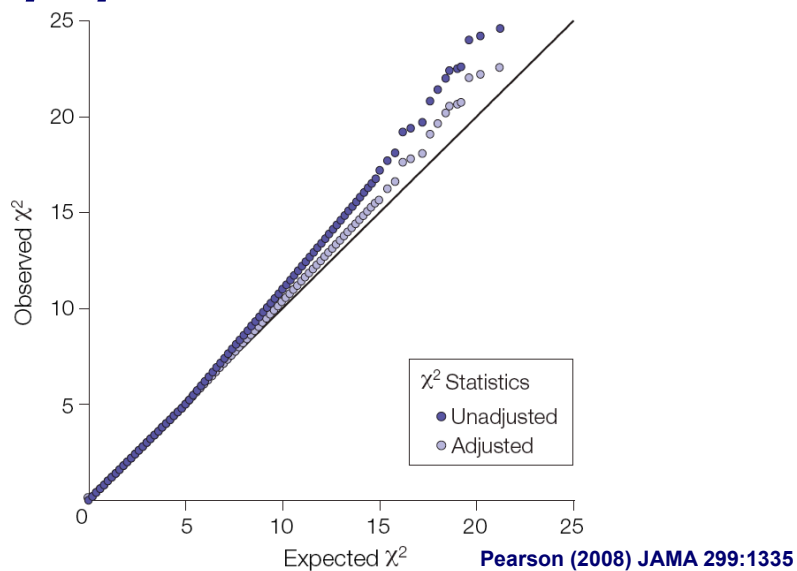
1161 Finnish T2D cases + 1174 Finnish normal glucose tolerant controls



## Quantile-quantile (Q-Q) plot



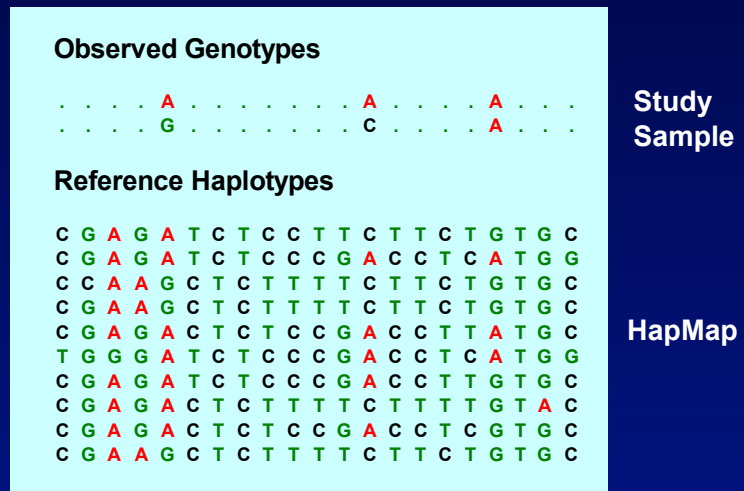
## Before and after adjustment of population stratification



## Gain power through collaboration

- Each study performs GWA
- Combine data from all studies by performing a meta-analysis
- Potential issues:
  - Different genotyping and analysis strategies
  - Case definitions are different

## Imputation: Observed genotypes



## Identify match among reference

### Observed Genotypes

. . . . . A . . . . . A . . . . . A . . . .  
 . . . . . G . . . . . C . . . . . A . . . .

### Reference Haplotypes

C G A G A T C T C C T T C T T C T G T G C  
 C G A G A T C T C C C G A C C T C A T G G  
 C C A A G C T C T T T T C T T C T G T G C  
 C G A A G C T C T T T T C T T C T G T G C  
 C G A G A C T C T C C G A C C T T A T G C  
 T G G G A T C T C C C G A C C T C A T G G  
 C G A G A T C T C C C G A C C T T G T G C  
 C G A G A C T C T T T T C T T T T G T A C  
 C G A G A C T C T C C G A C C T C G T G C  
 C G A A G C T C T T T T C T T C T G T G C

Li (2009) Ann Rev Gen Hum Genet 10:387

Gonçalo Abecasis

## Phase chromosomes, impute missing genotypes

### Observed Genotypes

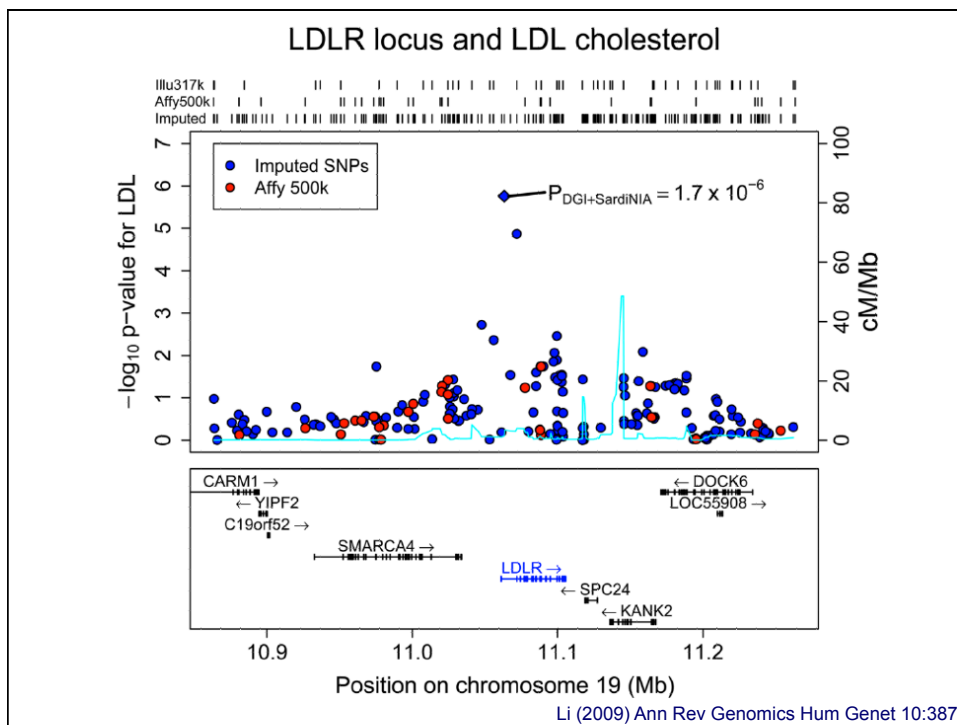
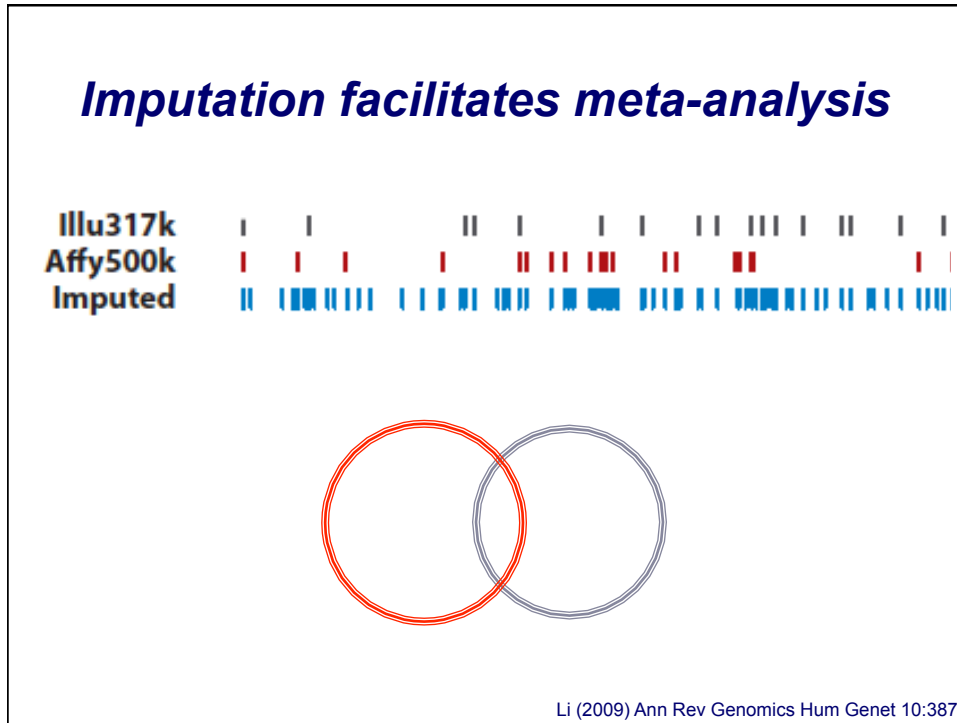
c g a g A t c t c c c g A c c t c A t g g  
 c g a a G c t c t t t t C t t t c A t g g

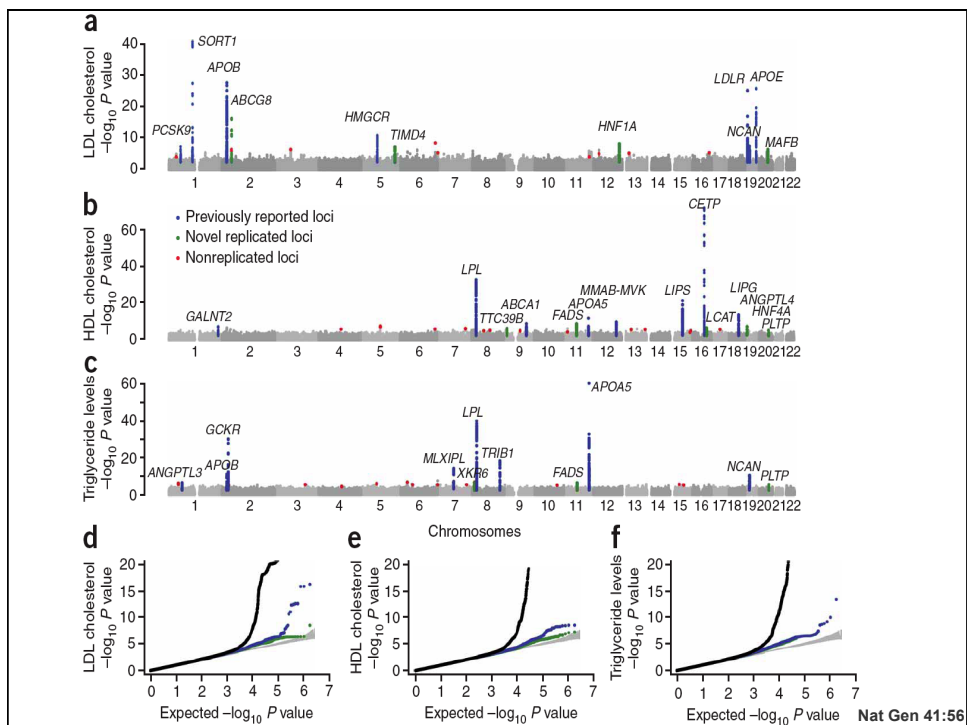
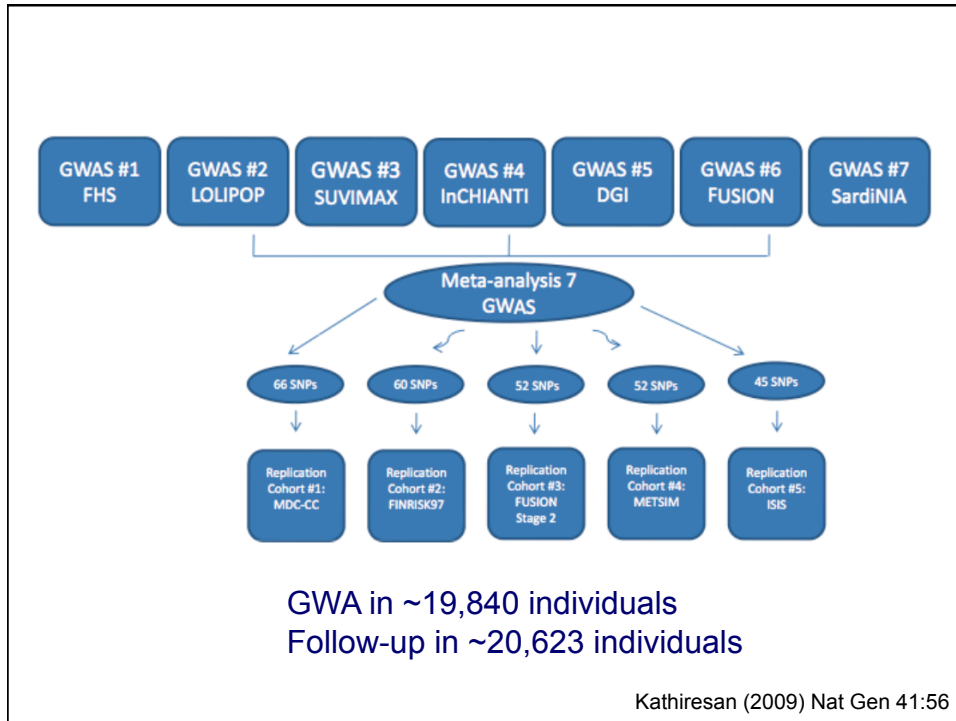
### Reference Haplotypes

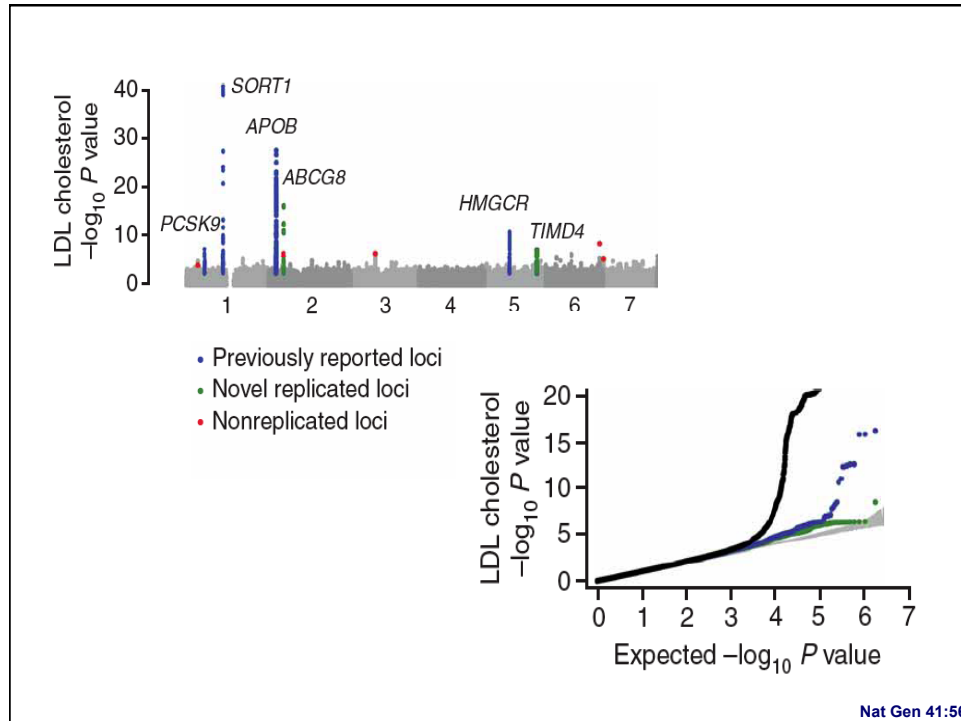
C G A G A T C T C C T T C T T C T G T G C  
 C G A G A T C T C C C G A C C T C A T G G  
 C C A A G C T C T T T T C T T C T G T G C  
 C G A A G C T C T T T T C T T C T G T G C  
 C G A G A C T C T C C G A C C T T A T G C  
 T G G G A T C T C C C G A C C T C A T G G  
 C G A G A T C T C C C G A C C T T G T G C  
 C G A G A C T C T T T T C T T T T G T A C  
 C G A G A C T C T C C G A C C T C G T G C  
 C G A A G C T C T T T T C T T C T G T G C

Li (2009) Ann Rev Gen Hum Genet 10:387

Gonçalo Abecasis





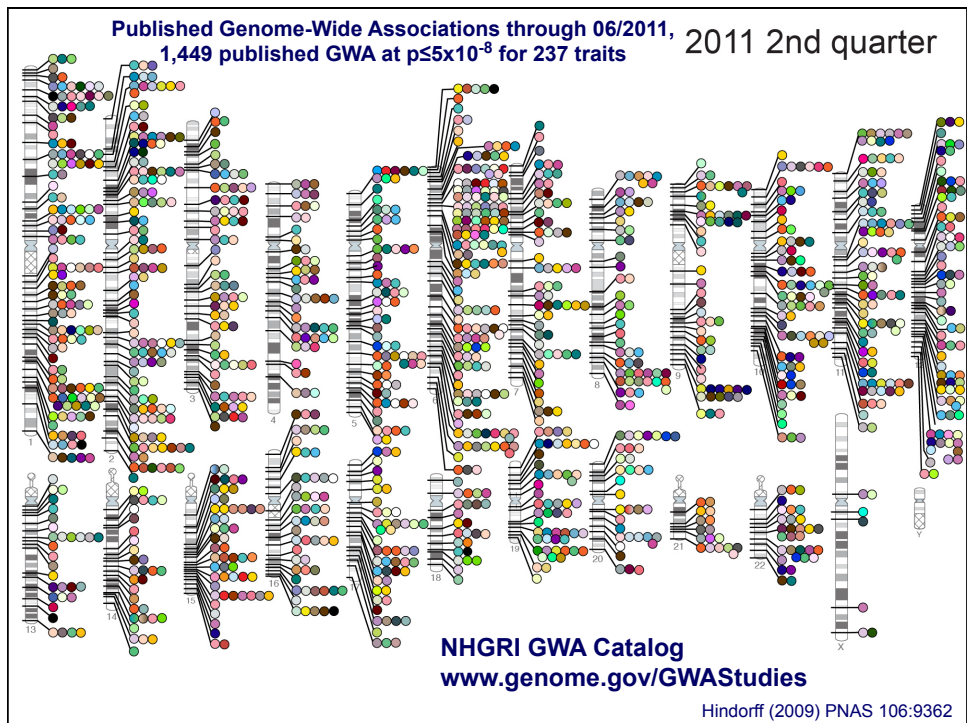
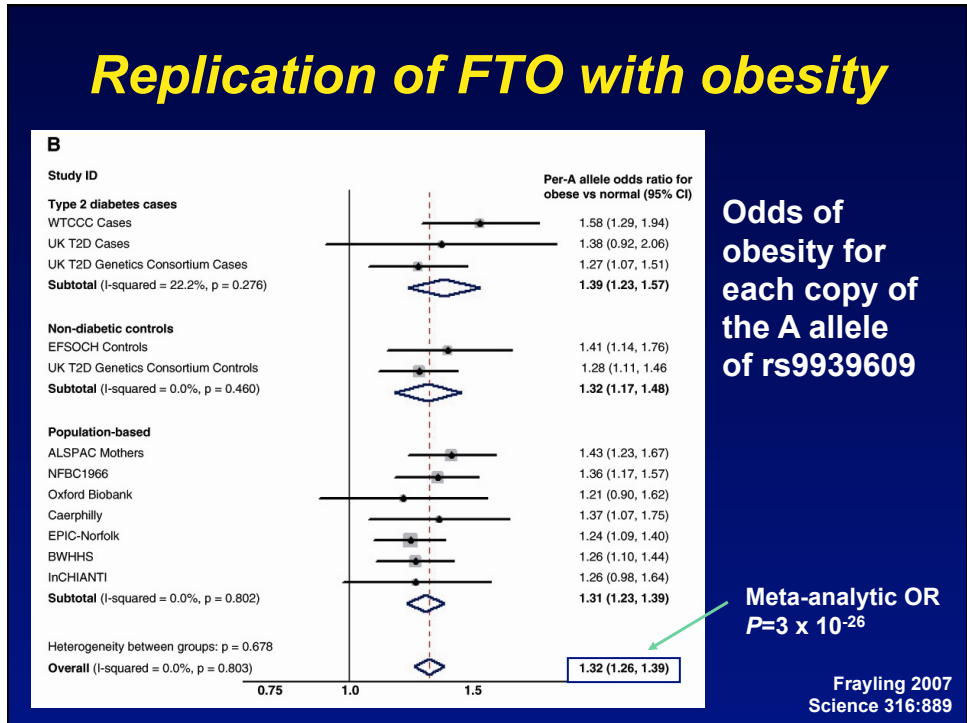


## Heterogeneity

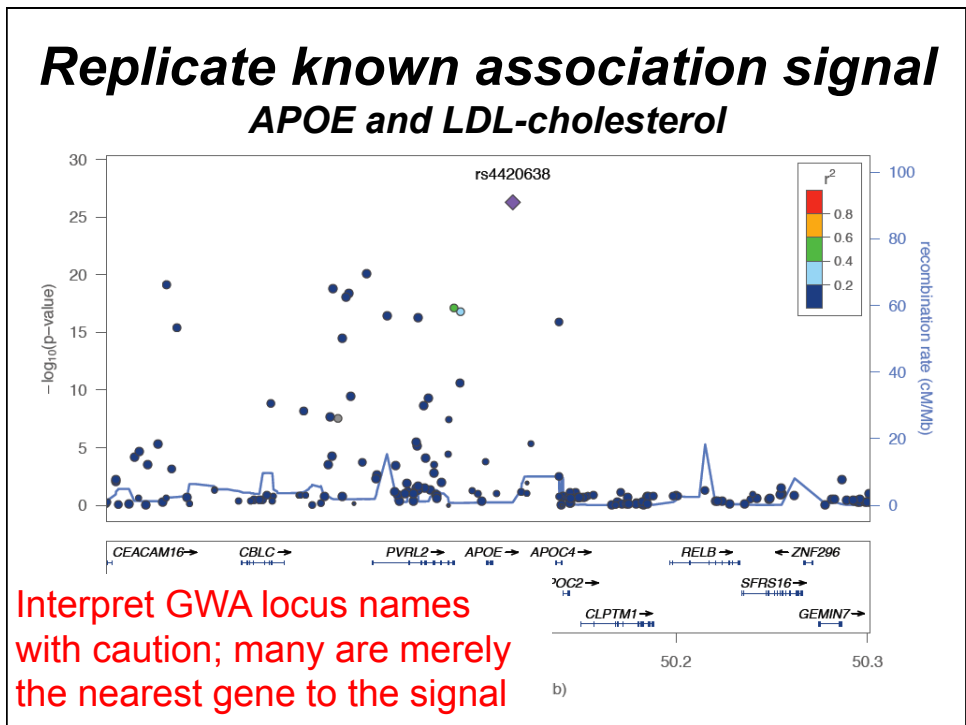
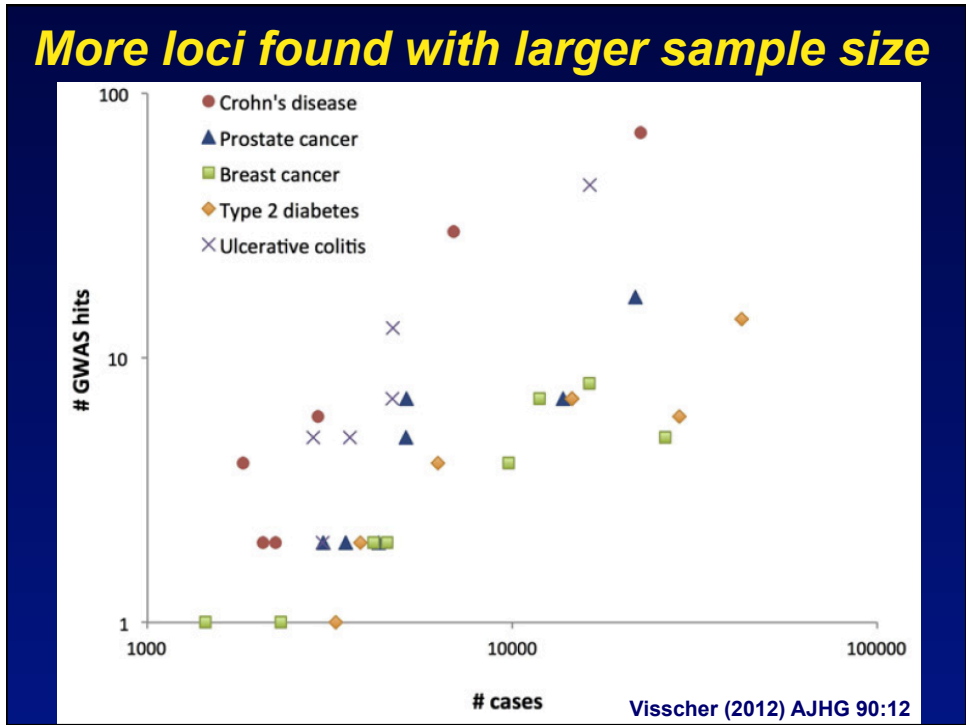
- *FTO* associated with type 2 diabetes in the Wellcome Trust Case-Control Consortium
- Mostly not observed in other diabetes studies
- WTCCC cases more obese than controls
- Diabetes signal abolished when adjust for BMI
- ID of heterogeneity source led to BMI gene

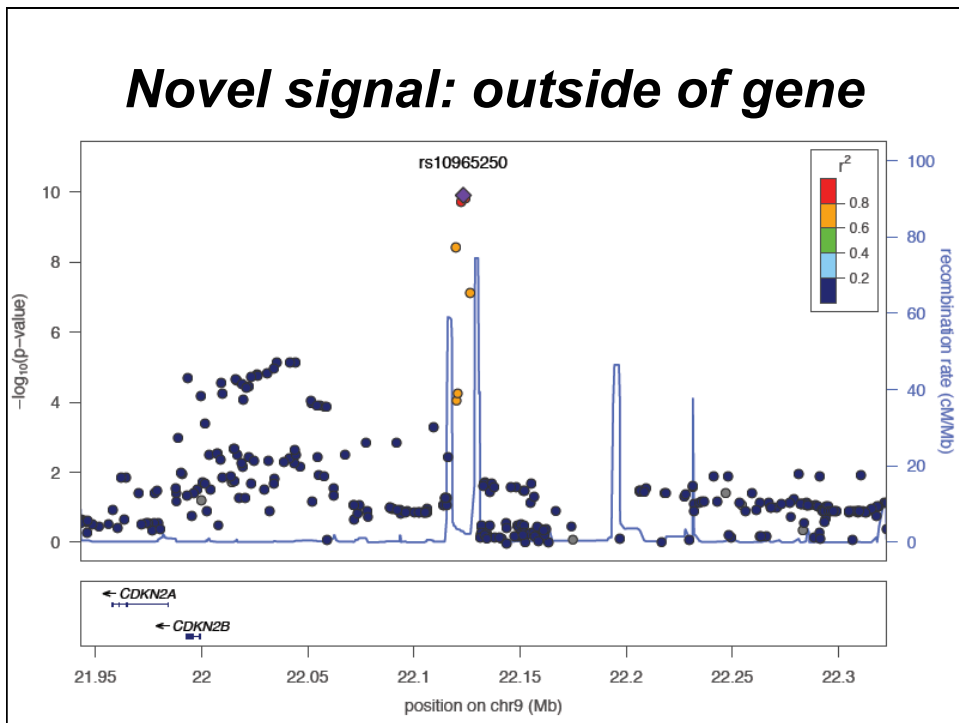
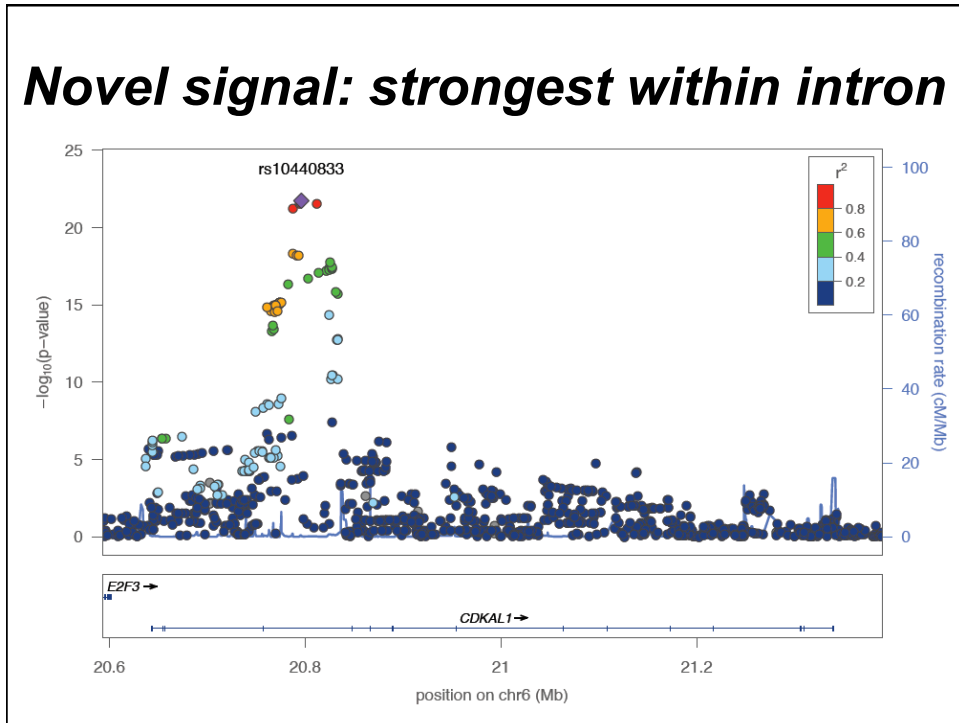
Frayling 2007 Science 316:889

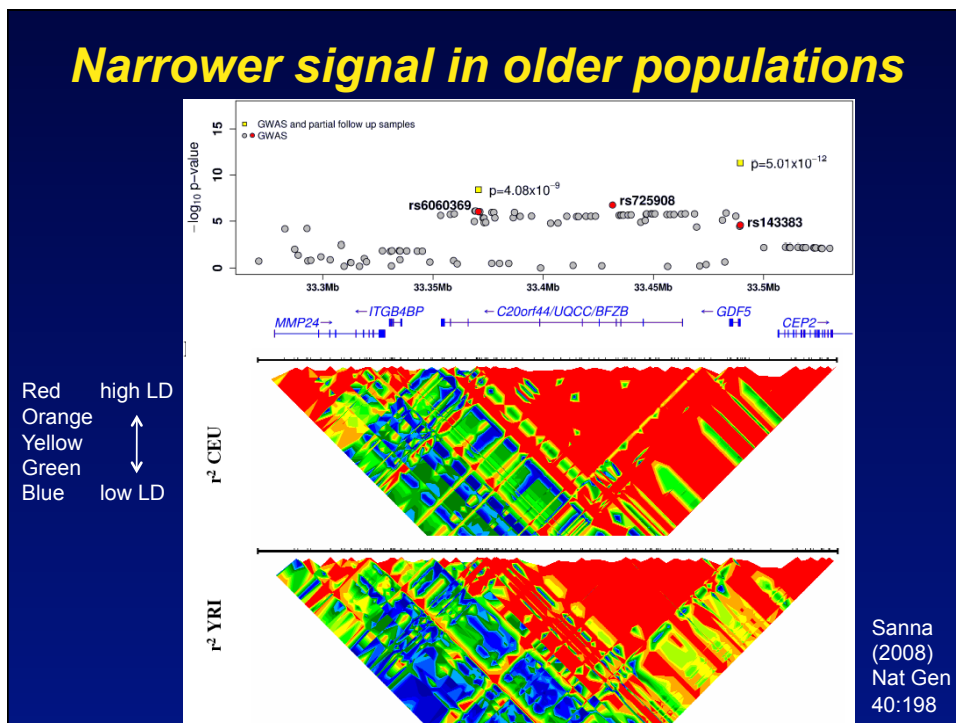
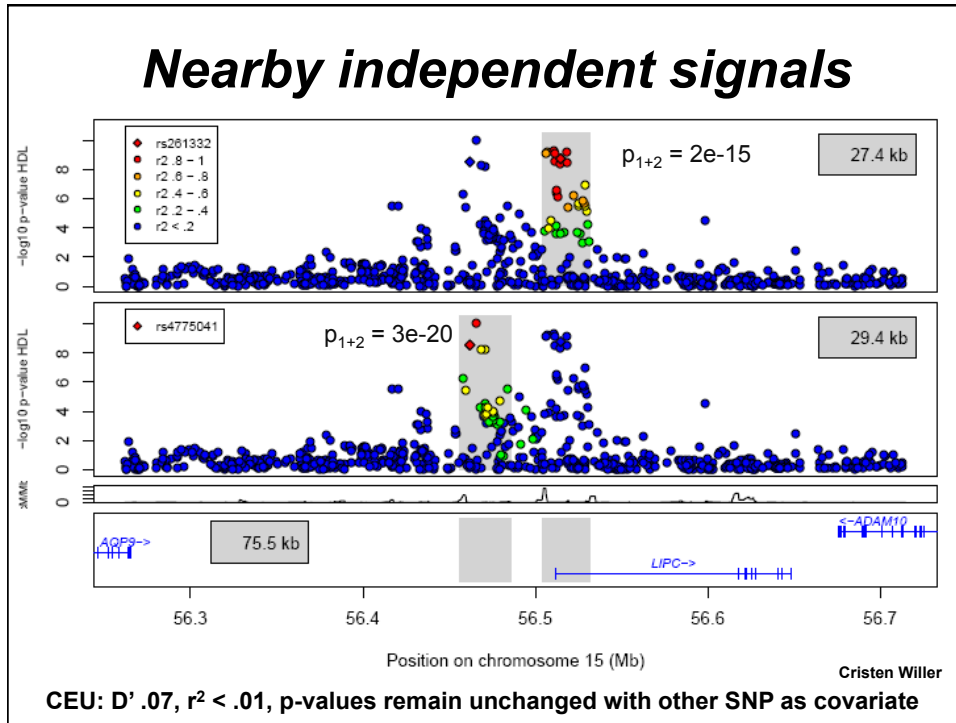
## Replication of *FTO* with obesity









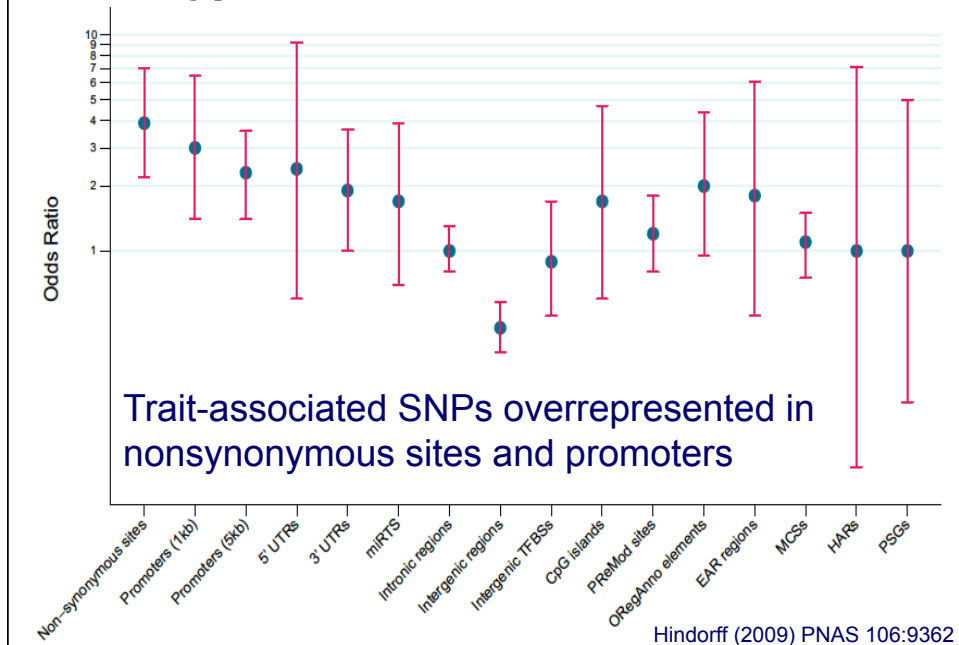


## Signals associated with $\geq 2$ traits

Attributed genes	Associated traits reported in catalog
<i>PTPN22</i>	Crohn's disease, type 1 diabetes, rheumatoid arthritis
<i>FCER1A</i>	Serum IgE levels, select biomarker traits (MCP1)
<i>BCL11A</i>	Fetal hemoglobin, F-cell distribution
<i>GCKR</i>	CRP, lipids, waist circumference
<i>HLA / MHC region</i>	Systemic lupus erythematosus, lung cancer, psoriasis, inflammatory bowel disease, ulcerative colitis, celiac disease, rheumatoid arthritis, juvenile idiopathic arthritis, multiple sclerosis, type 1 diabetes
<i>CDKAL1</i>	Crohn's disease, type 2 diabetes
<i>IRF4</i>	Freckles, hair color, chronic lymphocytic leukemia
<i>TNFAIP3</i>	Systemic lupus erythematosus, rheumatoid arthritis
<i>JAZF1</i>	Height, type 2 diabetes*
<i>Intergenic</i>	Prostate or colorectal cancer, breast cancer
<i>CDKN2A, CDKN2B</i>	Type 2 diabetes, intracranial aneurysm, myocardial infarction

Hindorff (2009) PNAS 106:9362

## Types of associated variants



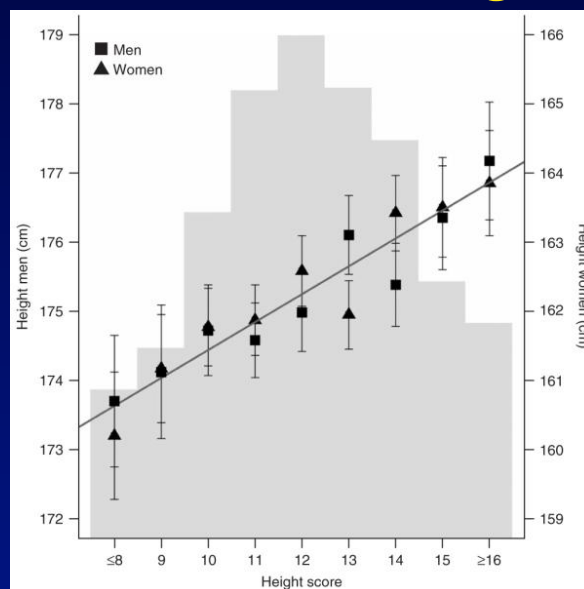
**Table 1. Population Variation Explained by GWAS for a Selected Number of Complex Traits**

Trait or Disease	$h^2$ Pedigree Studies	$h^2$ GWAS Hits <sup>a</sup>	$h^2$ All GWAS SNPs <sup>b</sup>
Type 1 diabetes	0.9 <sup>98</sup>	0.6 <sup>99, c</sup>	0.3 <sup>12</sup>
Type 2 diabetes	0.3–0.6 <sup>100</sup>	0.05–0.10 <sup>34</sup>	
Obesity (BMI)	0.4–0.6 <sup>101,102</sup>	0.01–0.02 <sup>36</sup>	0.2 <sup>14</sup>
Crohn's disease	0.6–0.8 <sup>103</sup>	0.1 <sup>11</sup>	0.4 <sup>12</sup>
Ulcerative colitis	0.5 <sup>103</sup>	0.05 <sup>12</sup>	
Multiple sclerosis	0.3–0.8 <sup>104</sup>	0.1 <sup>45</sup>	

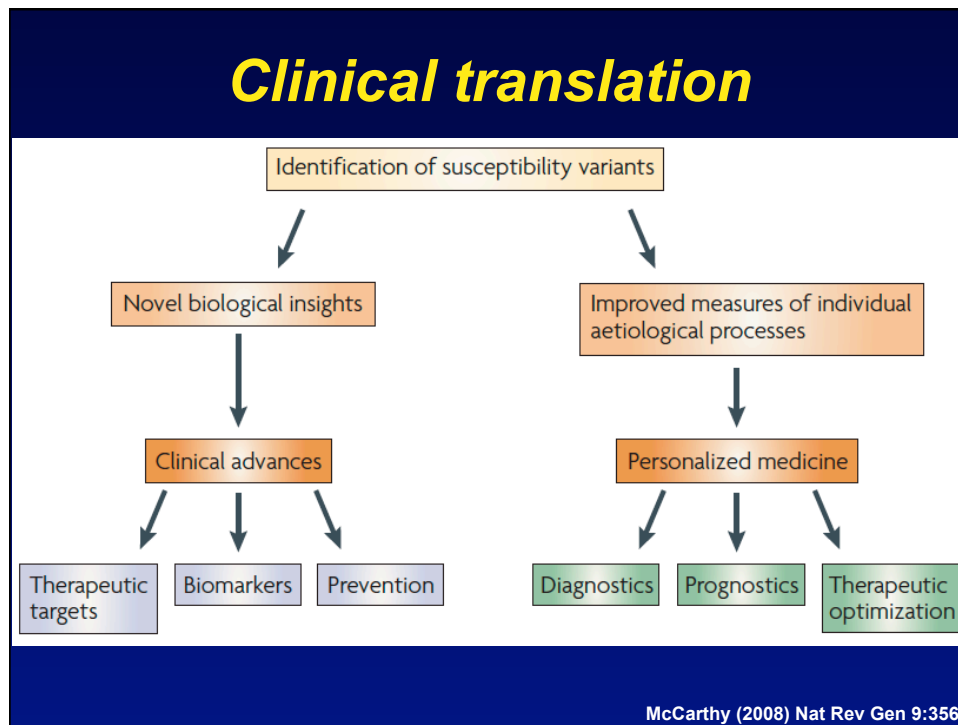
**Use of the current information in clinical practice will be disease dependent**

Partial table from Visscher (2012) AJHG 90:12

## Prediction of height



Lettre et al. (2008) Nat Gen 40:584–591



## Summary

- Need careful attention to design and QC
- Need large samples to find small signals
- 1,449 signals ( $P \leq 5 \times 10^{-8}$ ) and counting
- Finding an association signal does not immediately yield information on the underlying biology or clinical utility
- Time to changes in medical care based on GWA results may be many years

## ***Future of GWA***

- **More and more loci identified**
- **Larger meta-analyses**
- **Deeper follow-up of GWA signals**
- **Population-specific GWA panels**
- **More diverse populations**
- **Other sequence variants**
- **Multiple trait analysis**
- **Gene-gene and -environment interactions**
- **Molecular and biological mechanisms**