

# Technologic Issues in GWAS and Follow-up Studies

**Stephen Chanock, M.D.**

Senior Investigator, POB,CCR &

Director, Core Genotyping Facility, DCEG NCI

**May 22, 2007**

# Types of Polymorphisms

Single nucleotide polymorphisms (SNPs)

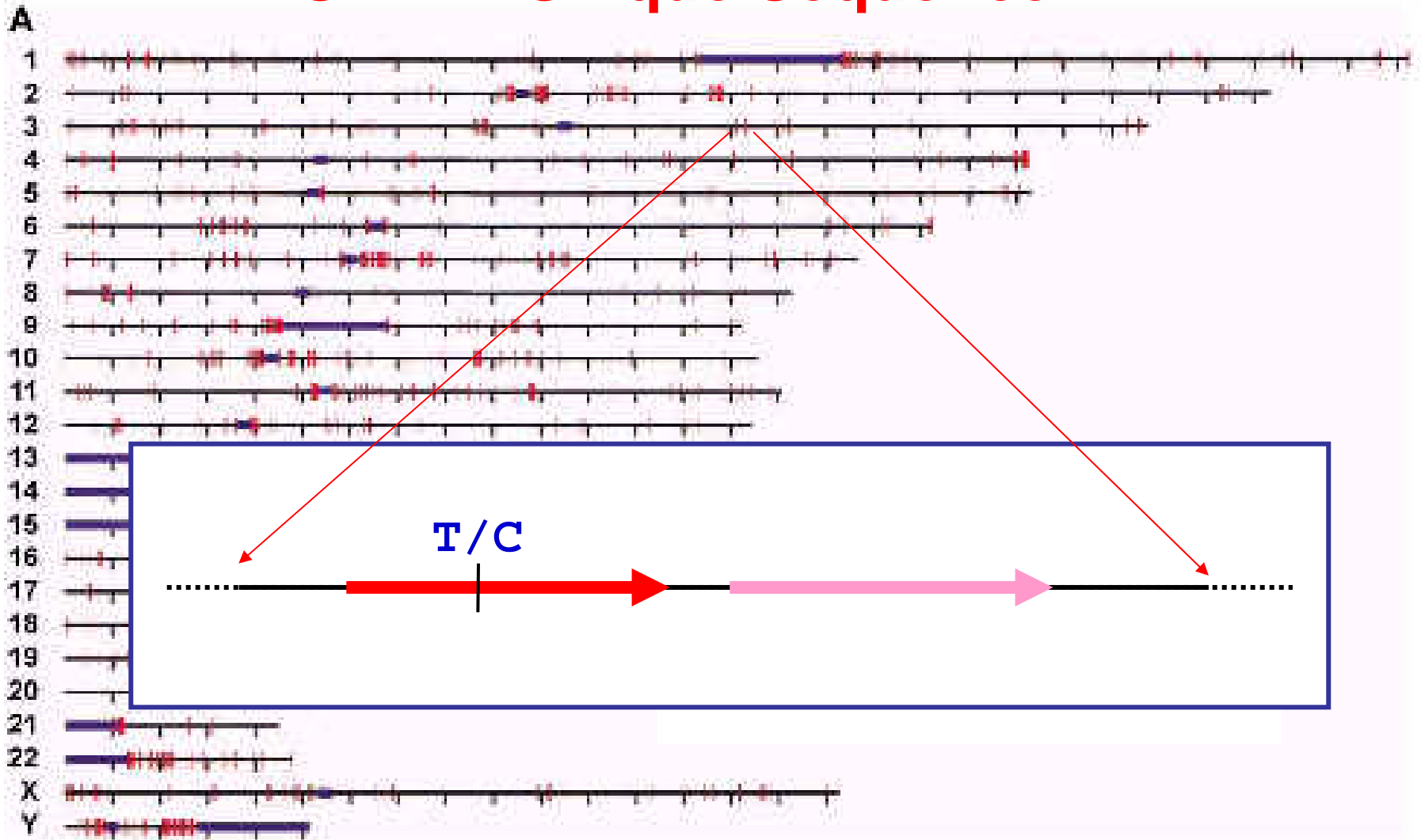
Common SNPs are defined as  $>1\%$  in at least one population

Rare SNPs are hard to identify and validate

*But*, it is estimated that there are a large number per individual

MAF= minor allele frequency

# SNP in Unique Sequence



# SNPs & Function: We know so little..

- Majority are “silent”
  - No known functional change
- Alter gene expression/regulation
  - Promoter/enhancer/silencer
  - mRNA stability
  - Small RNAs
- Alter function of gene product
  - Change sequence of protein

# SNPs in Genes: Take one

## Coding SNPs

**Synonymous-** no change in amino acid  
previously termed “silent” but.....

*Can alter mRNA stability*

*DRD2 (Duan et al 2002)*

**Nonsynonymous-** changes amino acid  
conservative and radical

**Nonsense-** insertion of stop codon

**Indel-** Disrupts codon sequence

Rare but disruptive

# SNPs Outside Genes: Take many....

- Majority distributed throughout genome are “silent” (excellent as markers)
- Alter transcription
  - Promoter, enhancer, silencer
- Regulate expression
  - Locus control region, mRNA stability
- Most are assumed to be ‘silent hitchhikers’
  - No function by predictive models or analysis

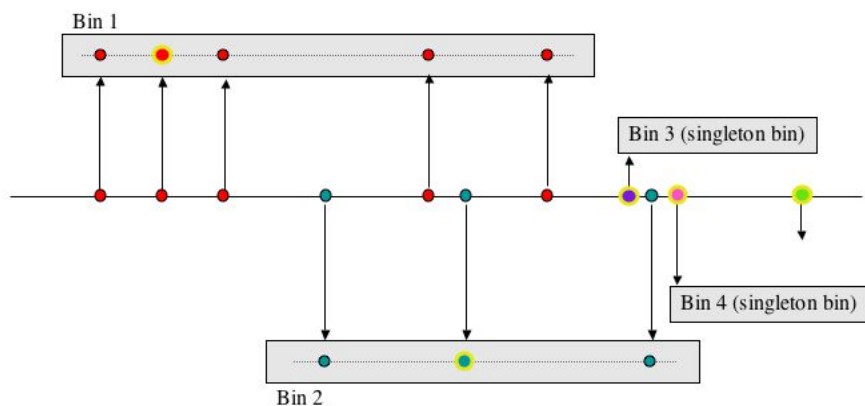
# Linkage disequilibrium (LD)

- The non-random association of alleles in the population
- Alleles at neighboring loci tend to cosegregate
- Linkage disequilibrium implies population allelic association

# Strategy for SNP Selection

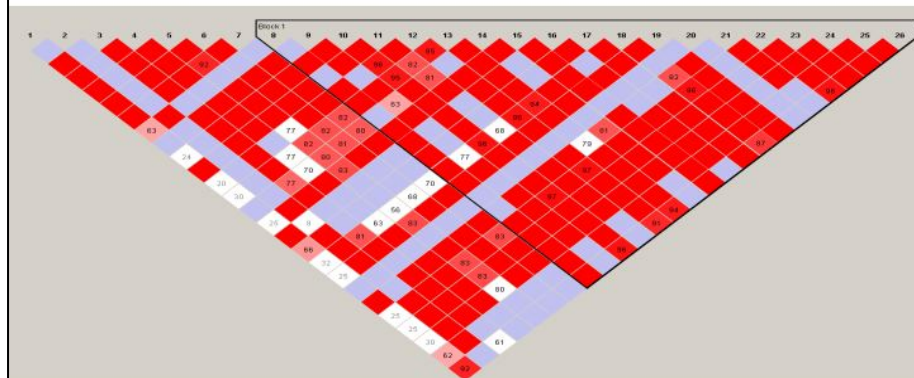
To test all SNPs is presently too costly  
Utilize a strategy that capitalizes on linkage disequilibrium between SNPs

Grouping of SNPs into bins based on pairwise  $r^2$ .



Carlson et al. *AJHG* 74:106 (2004)

Haplotype blocks defined by Gabriel et al  
Based on  $D'$  values for linkage disequilibrium

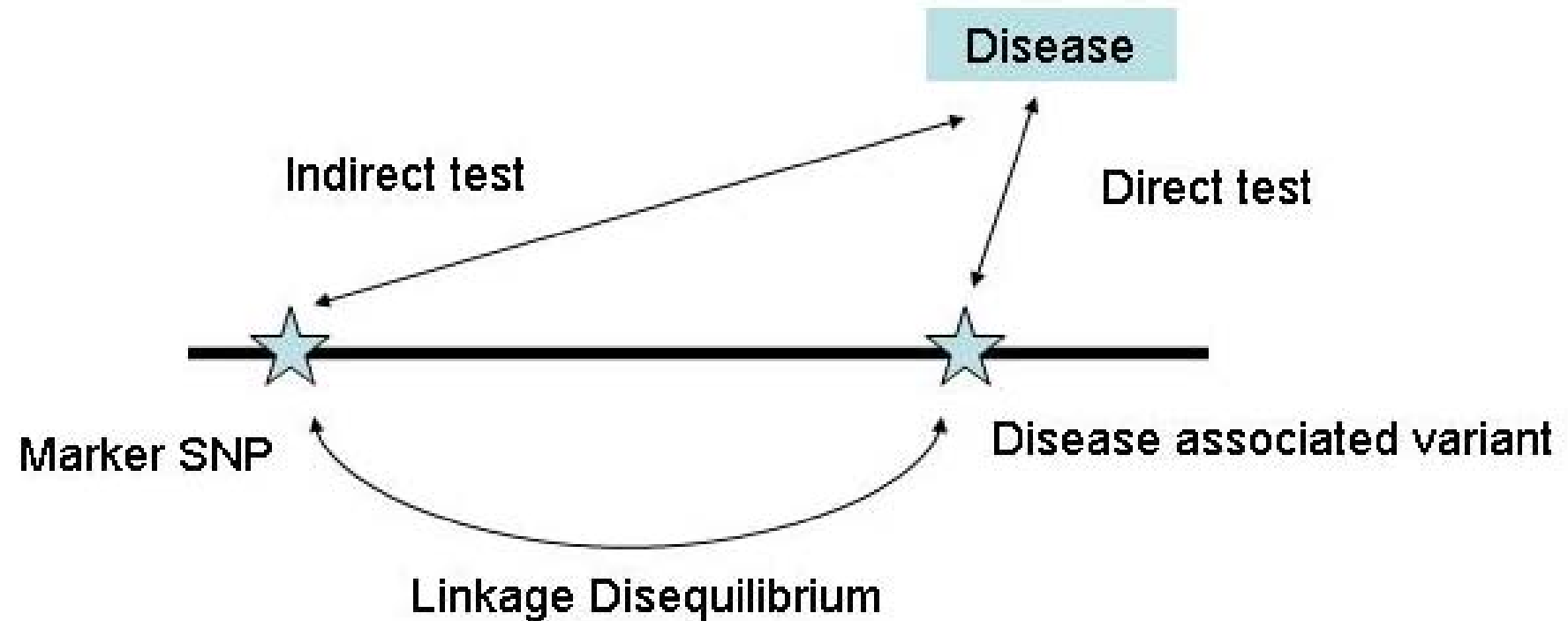




# What can LD do for me?

- Knowledge of patterns of linkage disequilibrium can be quite useful in the design and analysis of genetic data
- Design:
  - Estimation of theoretical power to detect associations
  - Evaluation of degree completeness of sampling of genetic variants
  - Choice of most informative genetic variants to genotype

# Genetic Association Testing: Finding Markers



# Large and Small Scale Polymorphisms

- **Copy Number of Polymorphisms**

- Regional “repeat” of sequence

- 10s to 100s kb of sequence

- Estimate of >10% of human genome

- Multi-copy in many individuals

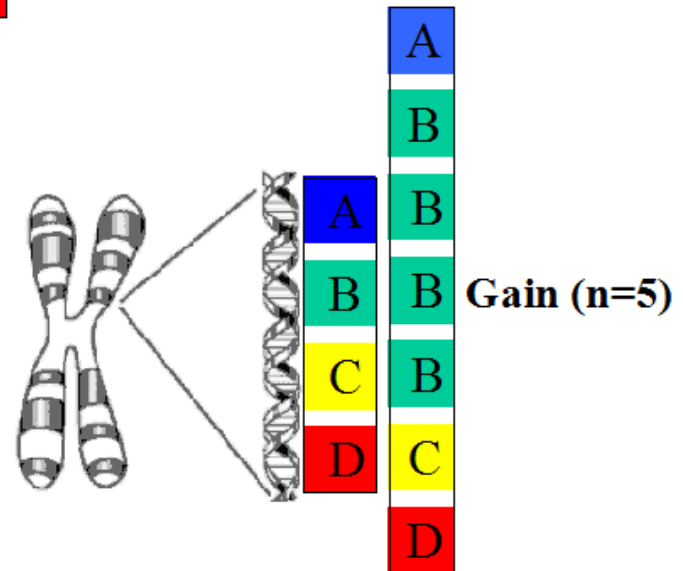
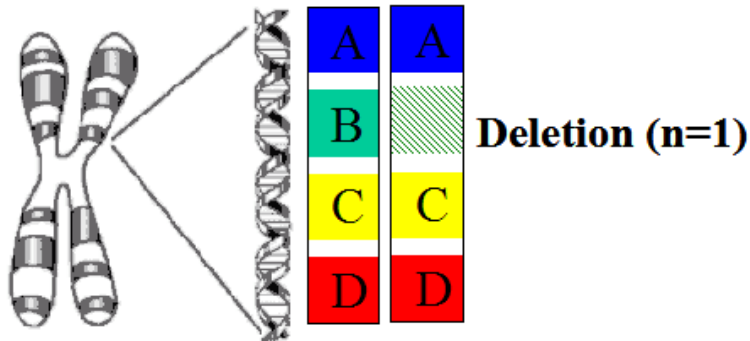
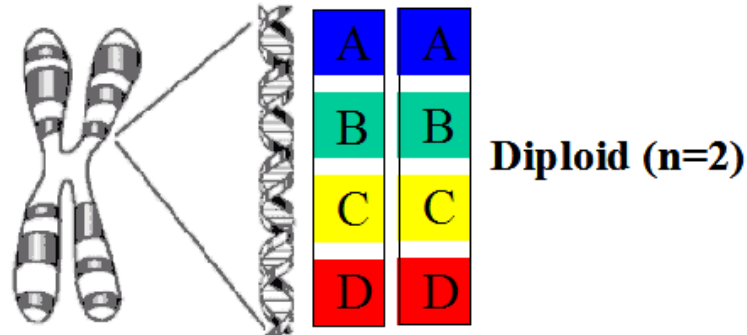
- **Duplicons**

- 90-100% similarity for >1 kb

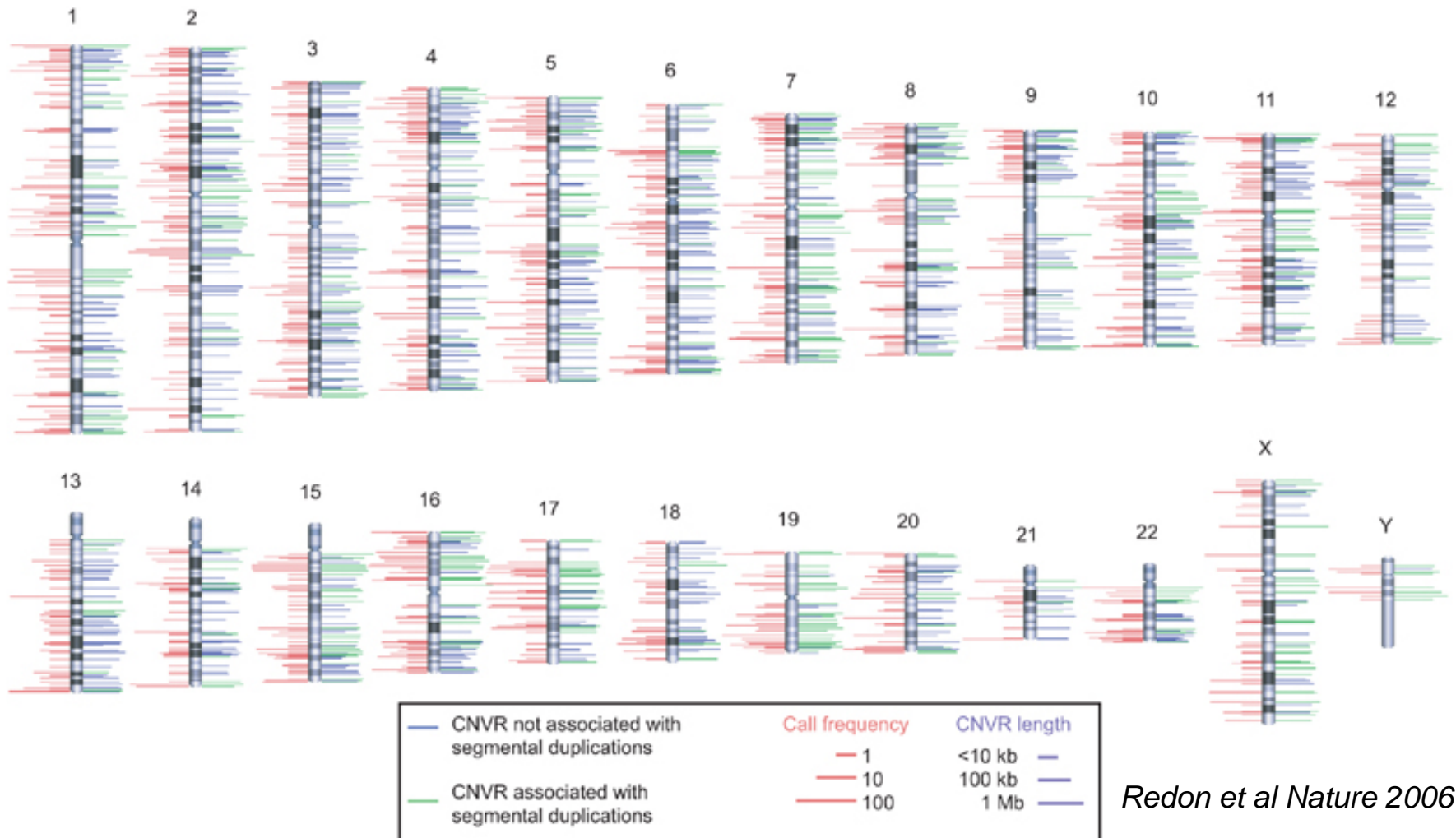
- 5-10% of genome (5% exons elsewhere)

- Multi-copy (high N) in all individuals

# Germ-line DNA Copy Number Variation(CNV)



# Copy Number Variation Across the Genome



# Copy Number Variation

Across the Genome

Frequency of CNVs

Most are uncommon (<5%)

Familial vs Unrelated Studies

Associated with Disease

*OPN1LW*

Red/green colorblind

*CCL3L1*

Reduced HIV Infection

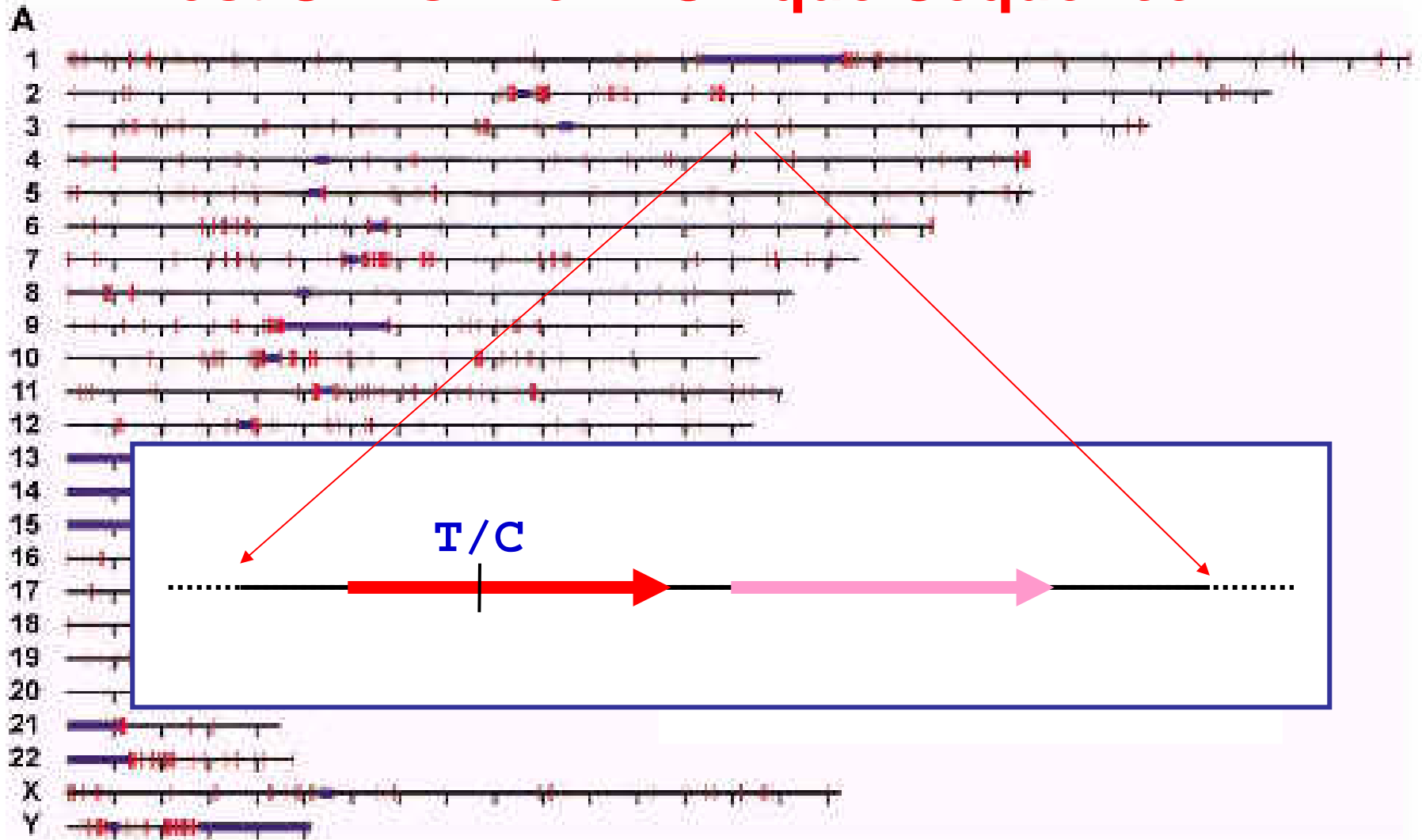
*CYP2A6*

Altered nicotine metabolism

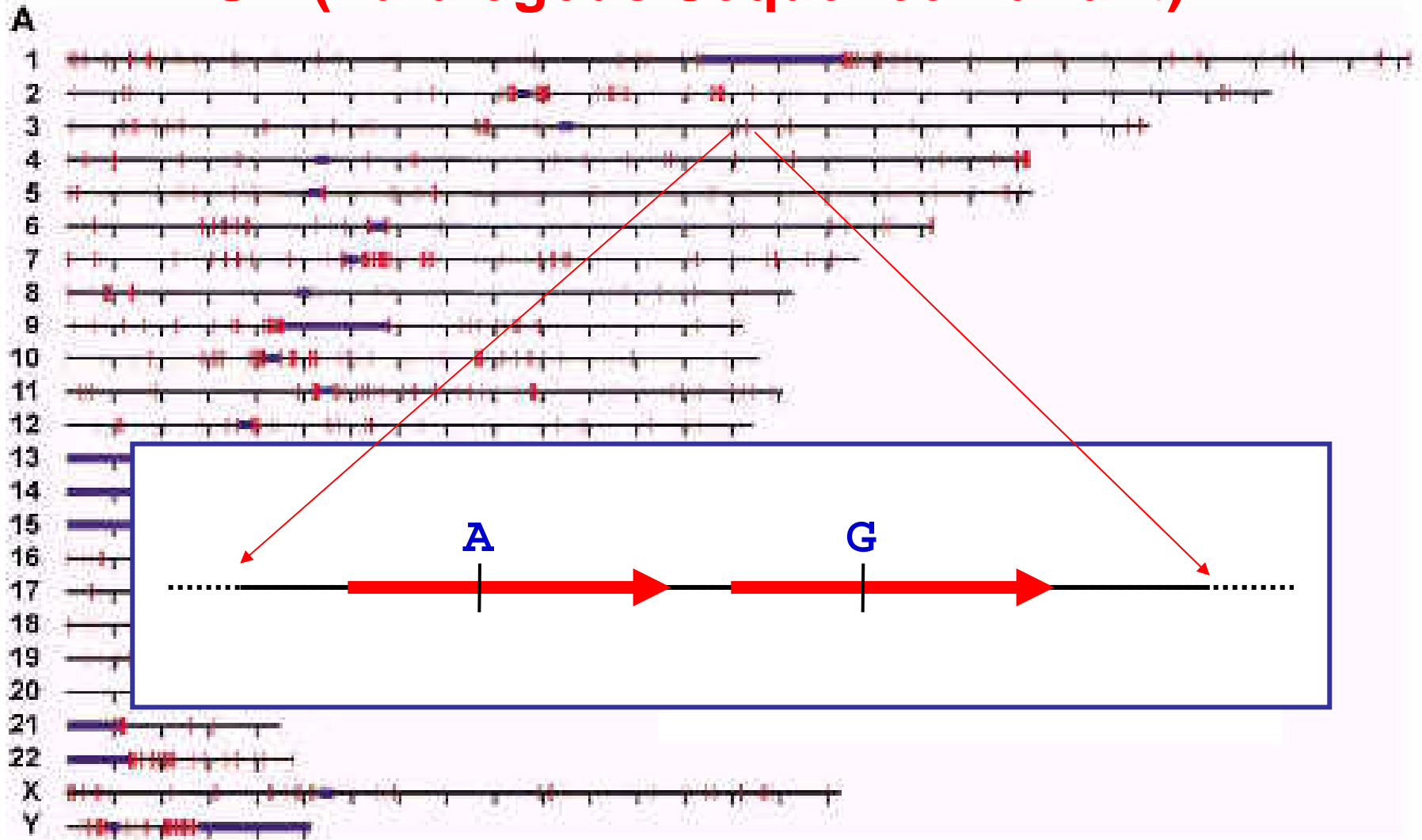
*VKORC1*

Warfarin metabolism

# Most SNPs Are In Unique Sequence !

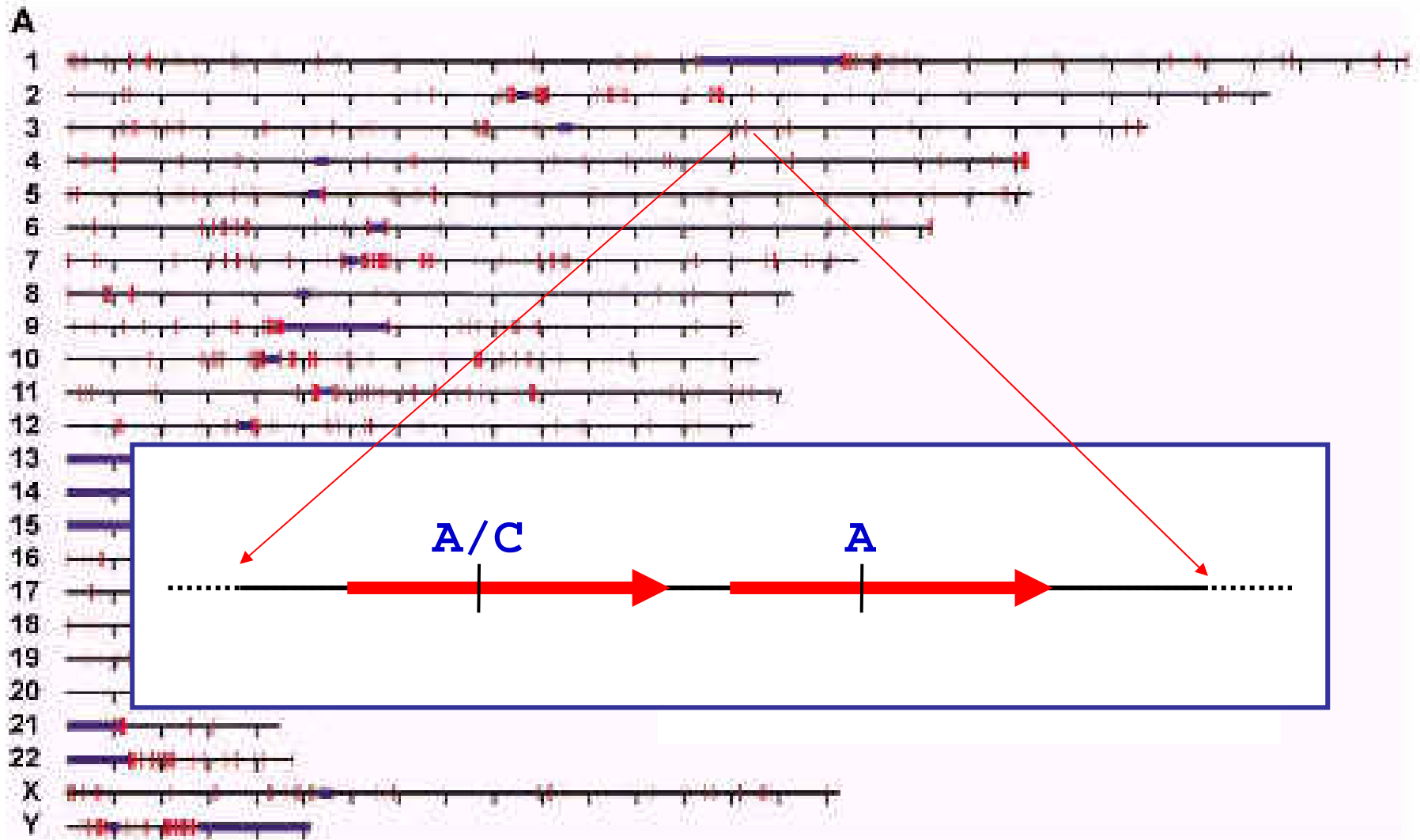


# PSV (Paralogous Sequence Variant)

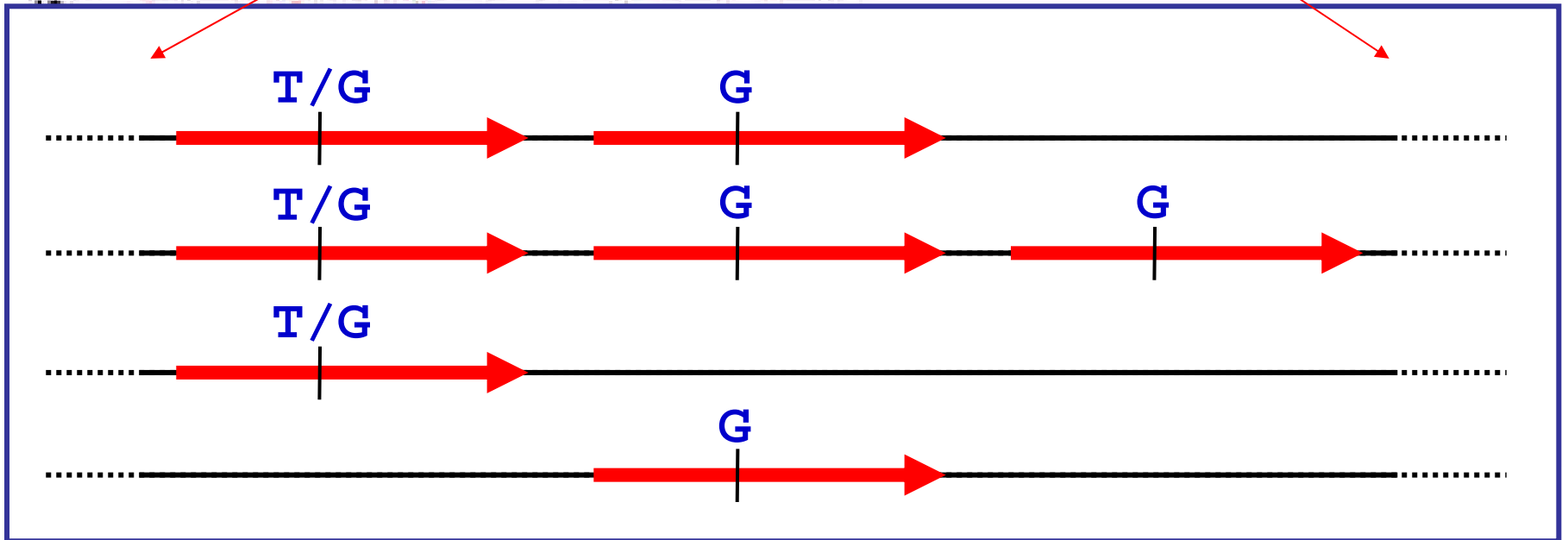
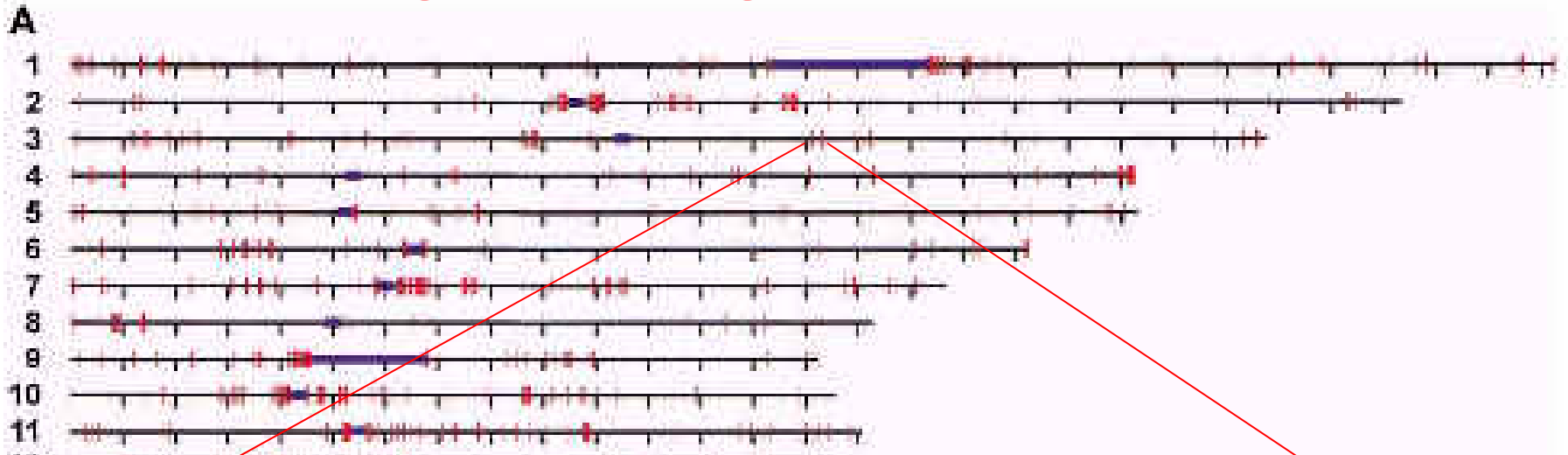




# SNP in Duplicon Sequence



# MSV: Multi-Site Variants



# Progress in Genotyping Technology

Cost per genotype  
Cents (USD)

$10^2$

**ABI  
TaqMan**  
**Sequenom**

10

**ABI  
SNPlex**

**Illumina  
Golden Gate**

**Affymetrix  
10K**     **Affymetrix  
MegAllele**

1

**Perlegen**

**Affymetrix  
100K/500K**     **Illumina  
Infinium/Sentrix**

1

10

$10^2$

$10^3$

$10^4$

$10^5$

$10^6$

# of  
SNPs

2001

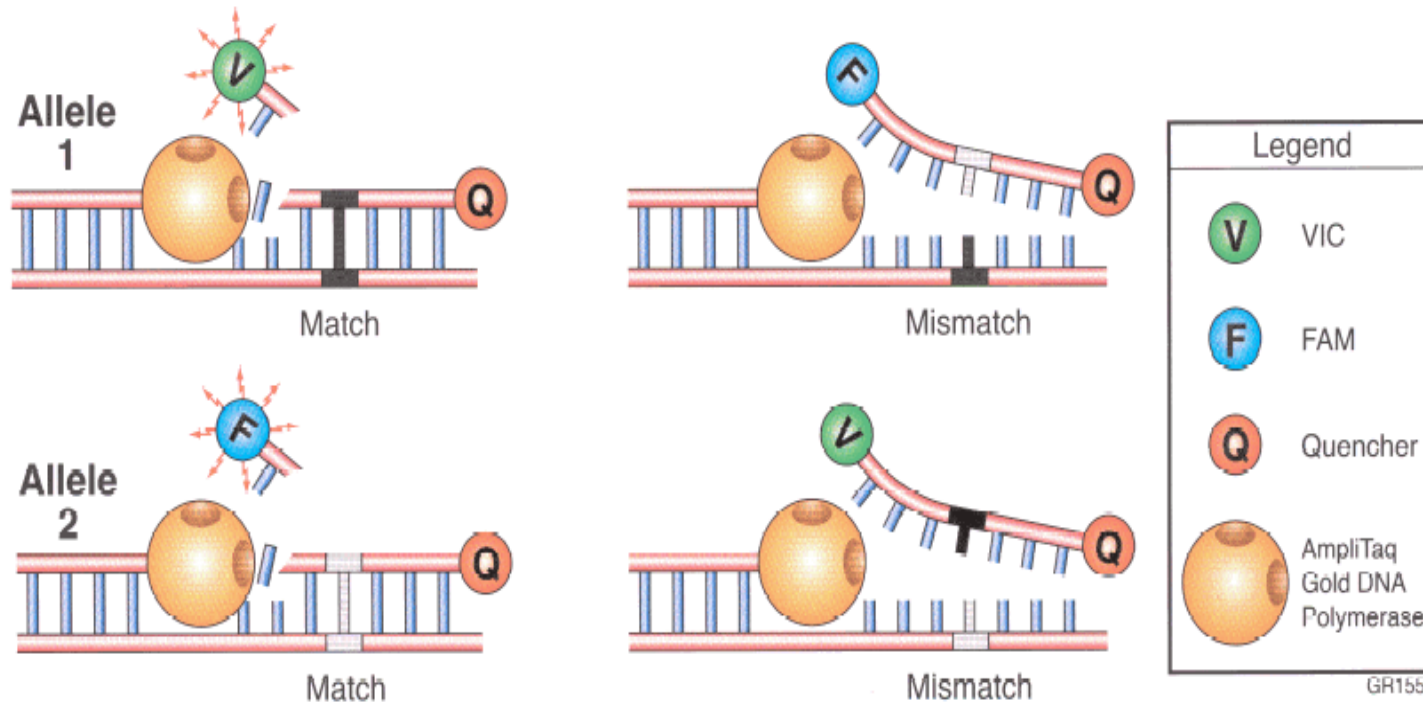


2007

# Genotype Technologies

- Dropping costs
- Smaller amounts of DNA
  - < 1ug for > 1 Million SNPs
- Economy of scale
  - Frequent Flier Paradigm
- Increased density but with fixed products
- Custom products bear high development cost
- Challenge of mid-range (50 to 500 SNPs)

# TaqMan™ (5' Exonuclease)



GR1556

A substantial increase in...

Indicates...

VIC™ fluorescence only

homozygosity for Allele 1

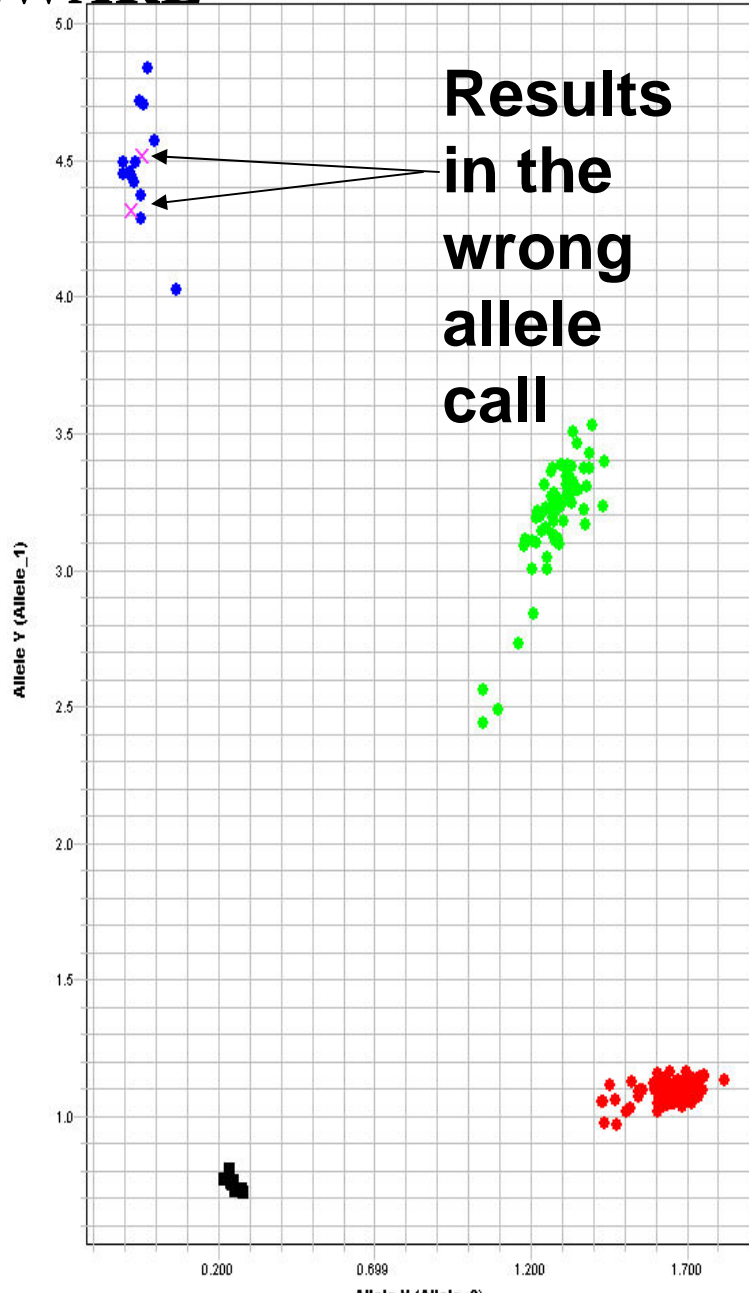
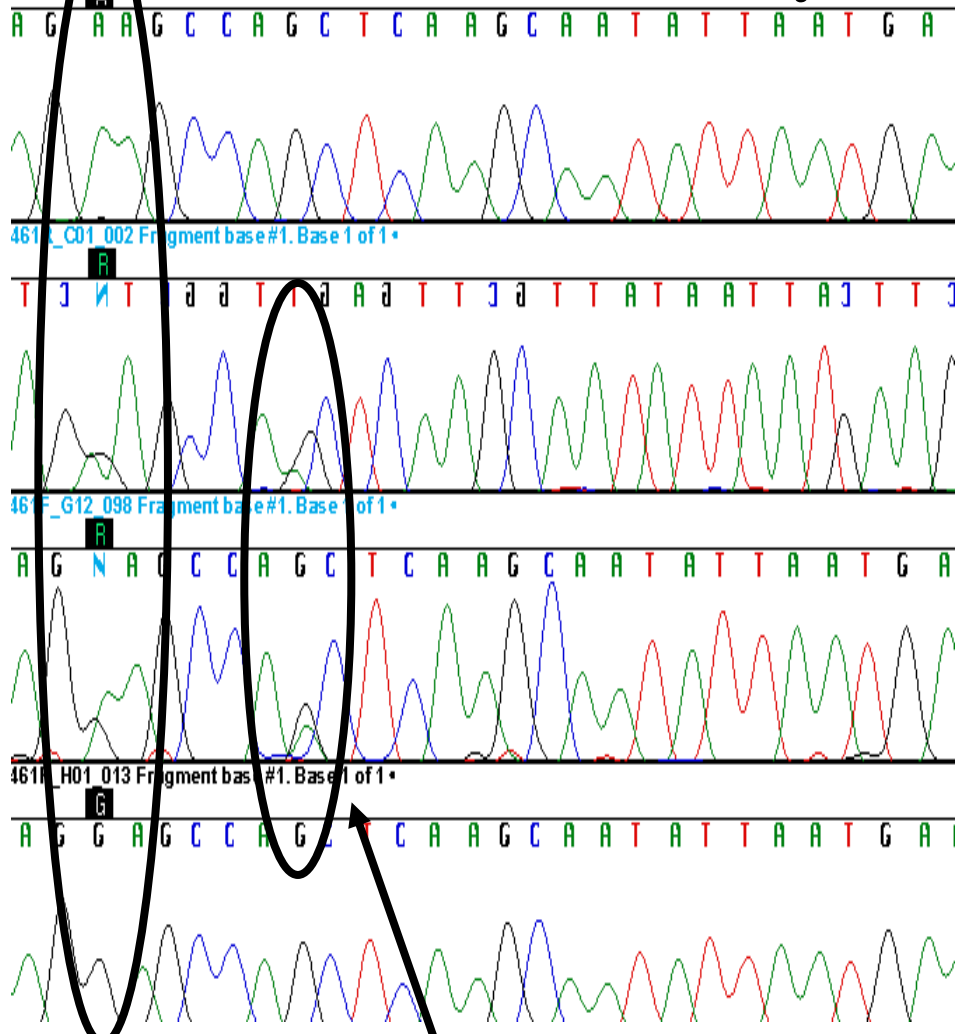
FAM™ fluorescence only

homozygosity for Allele 2

Both fluorescences

heterozygosity

# SNP Analysis- BEWARE



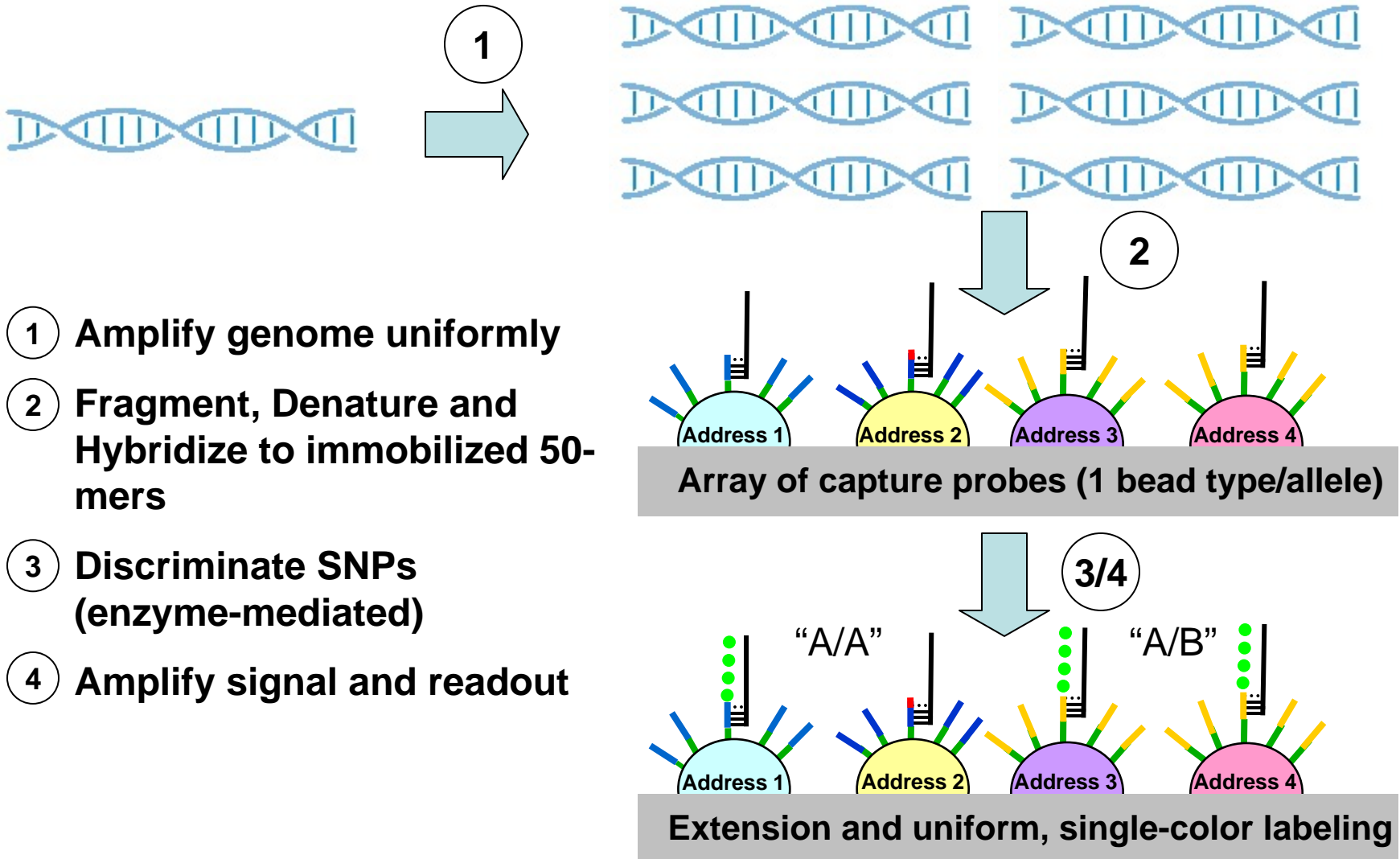
BRCA1-03

SNP under both TM probes

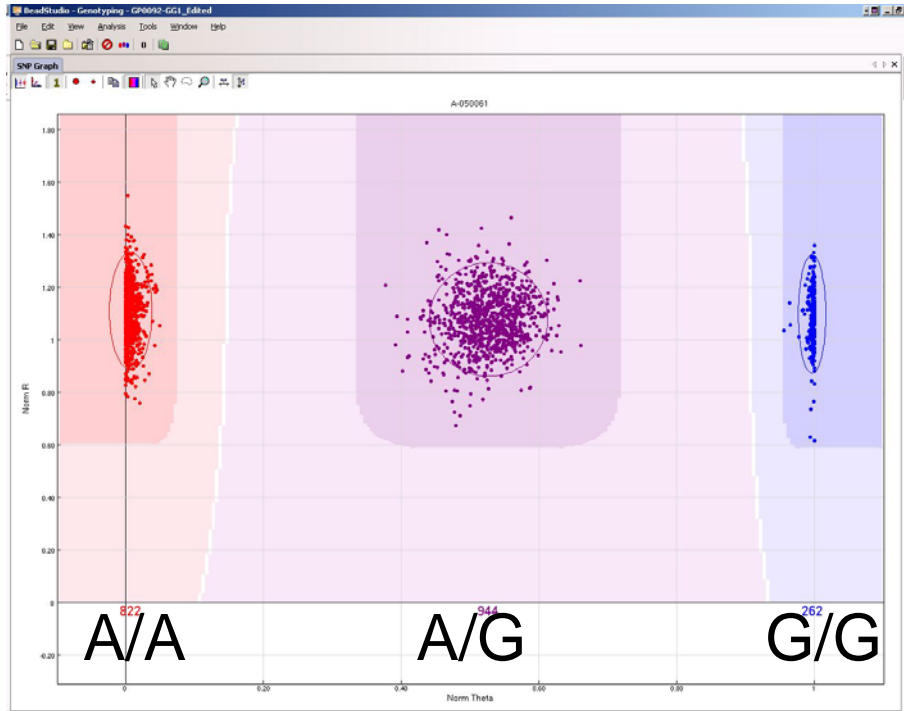
rs16941

<http://snp500cancer.nci.nih.gov>

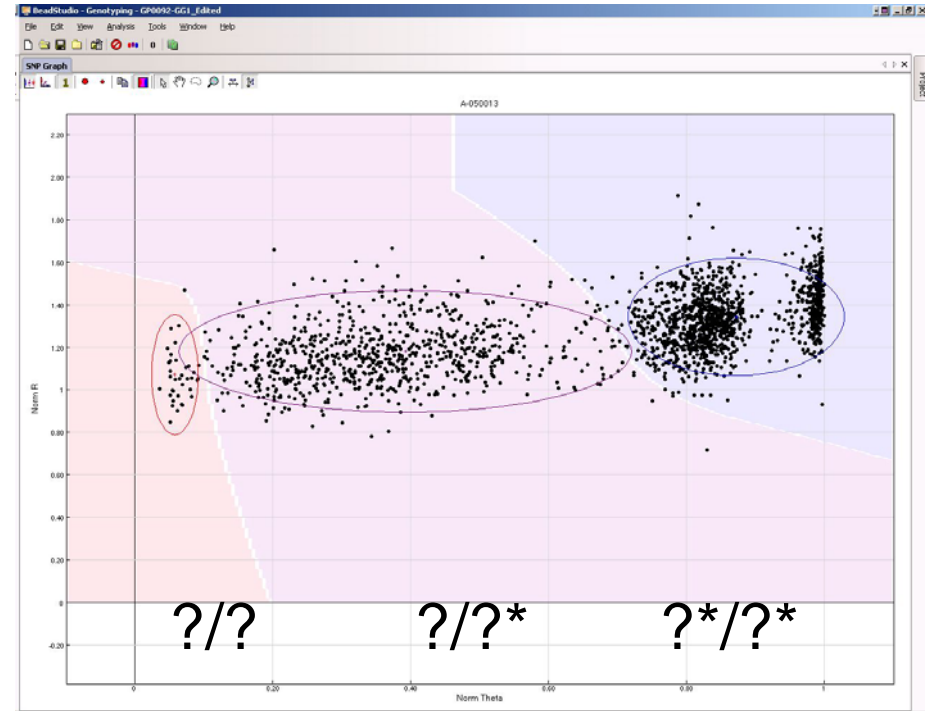
# Infinium™ Assay



# Illumina HumanHap500



Good Cluster

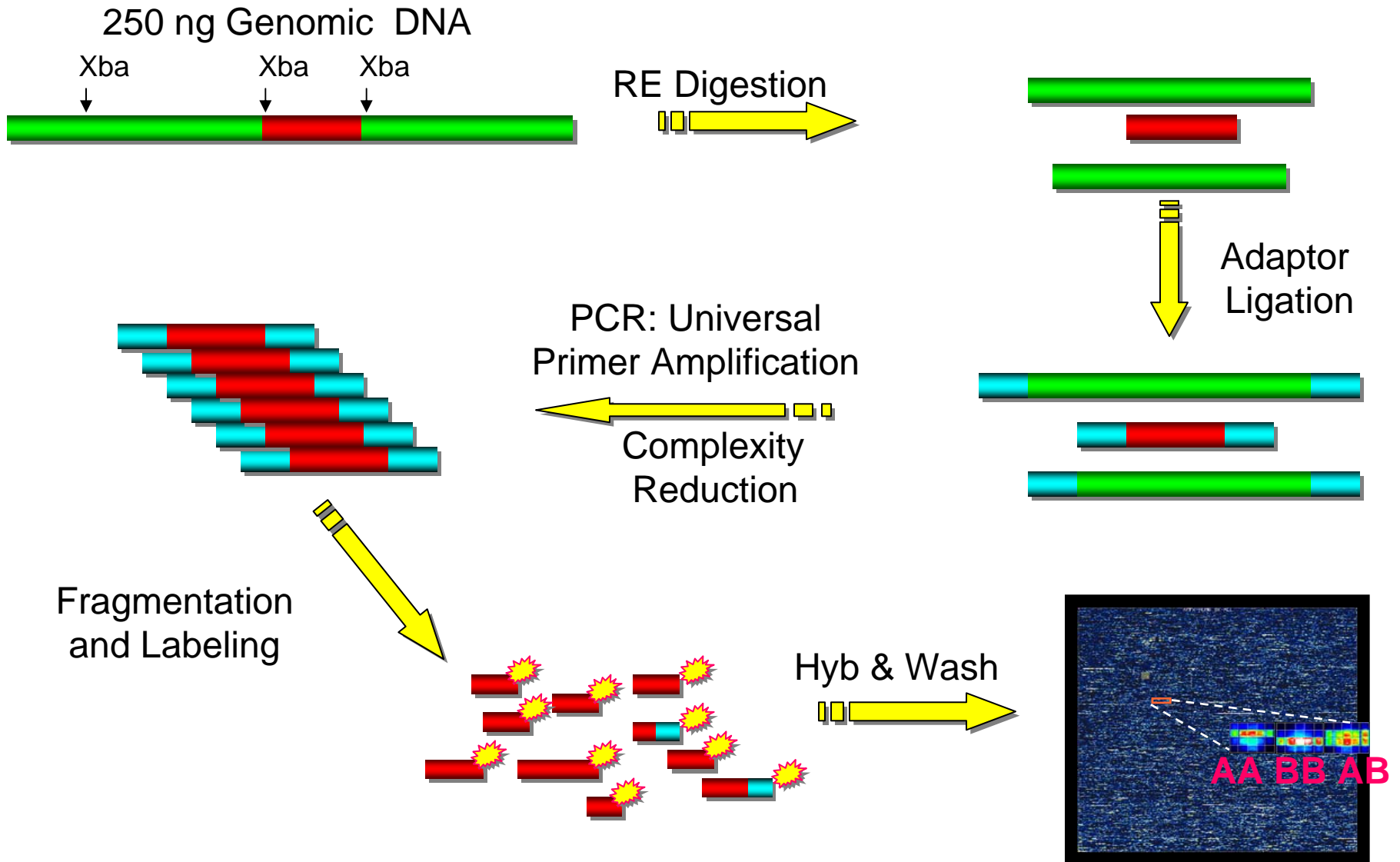


Poor Cluster

Read in BeadStudio™



# Affymetrix GeneChip<sup>®</sup> Mapping Assay



# Affymetrix 500K Chips

## Poor Quality Good Quality

GeneChip Operating Software - [Image Views]

File Edit View Run Tools Window Help

Open Save Print Tree Shorts Status Scans Expts Fluidics Start Add Resume Stop Analyze Tasks Image Help

Data Source: Local

- NA12750\_D8
- NA12751\_C7
- NA12753\_A1
- NA12762\_B4
- NA12763\_A6
- NA12801\_A3
- NA12801\_B9
- NA12801\_C3
- NA12801\_C9
- NA12812\_C10
- NA12813\_D3
- Image Data
  - NA06993\_C11.DAT
  - NA07022\_A4.DAT
  - NA07056\_A2.DAT
  - NA07345\_C4.DAT
  - NA07357\_A7.DAT
  - NA07357\_C1.DAT
  - NA10846\_A9.DAT
  - NA11831\_C8.DAT**
  - NA11832\_A10.DAT
  - NA11832\_C2.DAT
  - NA11832\_C2\_2.DAT
  - NA12144\_A8.DAT
  - NA12145\_A5.DAT
  - NA12155\_A11.DAT
  - NA12249\_C5.DAT
  - NA12740\_C6.DAT
  - NA12751\_C7.DAT
  - NA12753\_A1.DAT

GeneChip Software

Publish

Sample History

Instrument Control

Settings

NA11831\_C8 - NA11831\_C8

NA11832\_C2 - NA11832\_C2\_2

Auto Adj. Grid Mask Avg In Out Full

Position	Experiment Name	Probe Array Type	Barcode ID	User	Date & Time	Scan Status
9	NA11832_C2	GenomeWideSNP_5	36976	cgfbio	May 07 2007 04:37PM	Failed to align one or more subgrids.
10	NA07345_C4	GenomeWideSNP_5	40232	cgfbio	May 07 2007 04:59PM	Scan complete

2 Cartridges Loaded      Autoloader Door: Unlocked

Pixel X = 9116, Y = 7627, Intensity = 167      Standby      Filters applied

05/07/07 14:59:47 - Initializing the fluidics station  
 05/07/07 14:59:47 - Initializing the scanner  
 05/07/07 14:59:51 - GeneChip Operating Software initialization complete  
 05/07/07 17:00:06 - Data from experiment NA07345\_C4 added to the analysis queue  
 05/07/07 17:00:09 - Running analysis for D:\Program Files\Affymetrix\GeneChip\Affy  
 05/07/07 17:00:56 - Analysis completed for D:\Program Files\Affymetrix\GeneChip\A

start      GeneChip Operating ...      HapMapTraining\_WK...      2:23 PM

# Important Points

- Too many data points to review individually
- Iterative algorithm for analysis
  - Still undergoing improvements
- Validation of notable SNPs with second technology
  - “Neighboring SNP-land mines”
- Do not do this at home- Only for highly trained personnel

# Choice of Dense SNP Platforms

## Affymetrix

### Basic Points

Based on 'Spacing'  
100k, 500k, 1 Million  
CNV Analysis  
WGA compatible

### Issues

Lower price  
2 Enzyme Problem  
Calling Algorithms  
Redundancy (useful)

## Illumina

### Basic Points

Based on 'tagging'  
317k, 550k, 1 Million  
CNV Analysis  
WGA not yet rec'd

### Issues

Higher price  
HapMap II Based

# 2007 What is Available for Whole Genome Scans

- Coverage analysis based on HapMap II Data

- Build 20 MAF  $\geq 5\%$ ,  $r^2 \geq 0.8$  (pair-wise)

		CEU	YRI	JPT/CHB
• Illumina	HumanHap300	80%	35%	40%
• Illumina	HumanHap500	91%	58%	88%
• Affymetrix*	500k Mapping	63*%	41%	63%

\*77% (with 50k MegA)

**Quality control of  
genotype calls  
&  
DNA handling**

# Quality Control for Called Genotypes

## PURPOSE:

Identification of unreliable SNPs and DNAs to be *entirely* removed from the analysis .

Evaluation of completion rate (DNAs)

Evaluation of call rate (SNPs)

Evaluation of discordance rate (error rate)

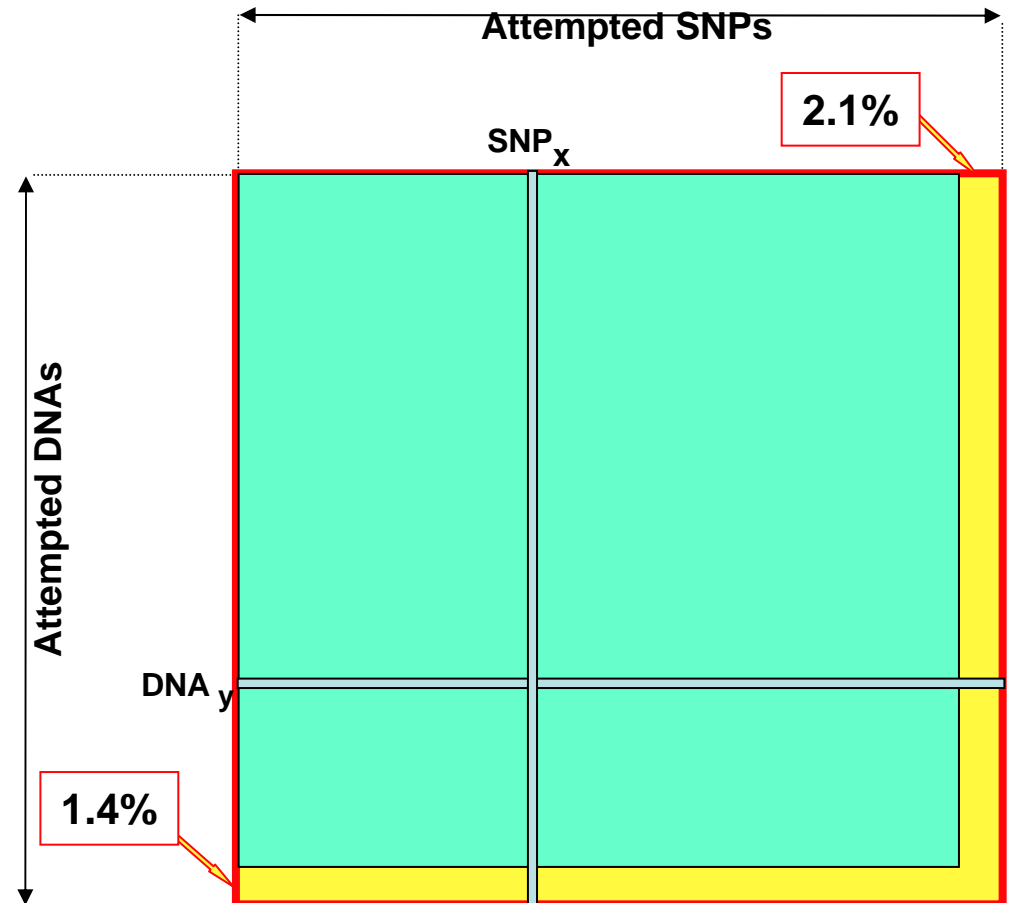
# CGEMS SNP DNA Success Rates

**Criterion**  
SNP call rate > 90%

	Number of individuals attempted	failed	success rate
PLCO	561,494	1,490	0.973
NHS	555,352	8,706	0.984

**Threshold**  
DNA completion rate > 94%

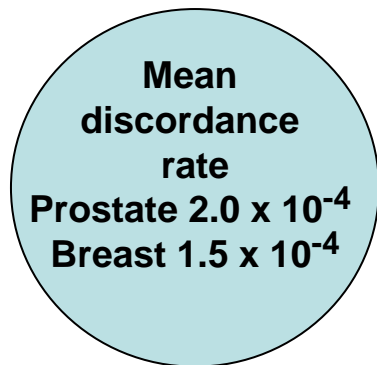
Number of individuals attempted	failed	Success rate
4696	66	0.986



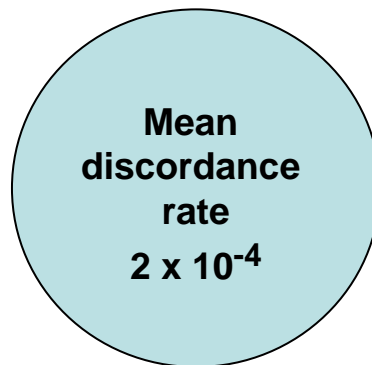


# Discordance rate for CGEMS: HumanHap500 (Illumina)

**Participants**  
142 duplicate pairs



**CEPH-CGEMS**  
74 duplicate pairs

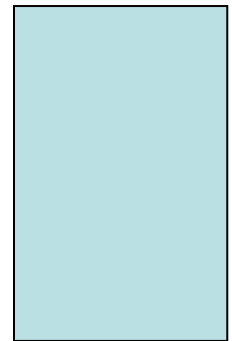


**28 individuals**  
(with 24 duplicates)



Mean discordance rate  
 $1.4 \times 10^{-3}$

**CEPH-HapMap**



Concordance rates >99.5%  
Subtle Differences in Quality of DNA

# **Quality Control for Recruitment DNA handling**

## **Checking for:**

**Chromosome X Ploidy**

**Identification of Familial Relationships**

**Evaluation of Continental Admixture**

**Population Stratification**

**Principal Component Analysis**

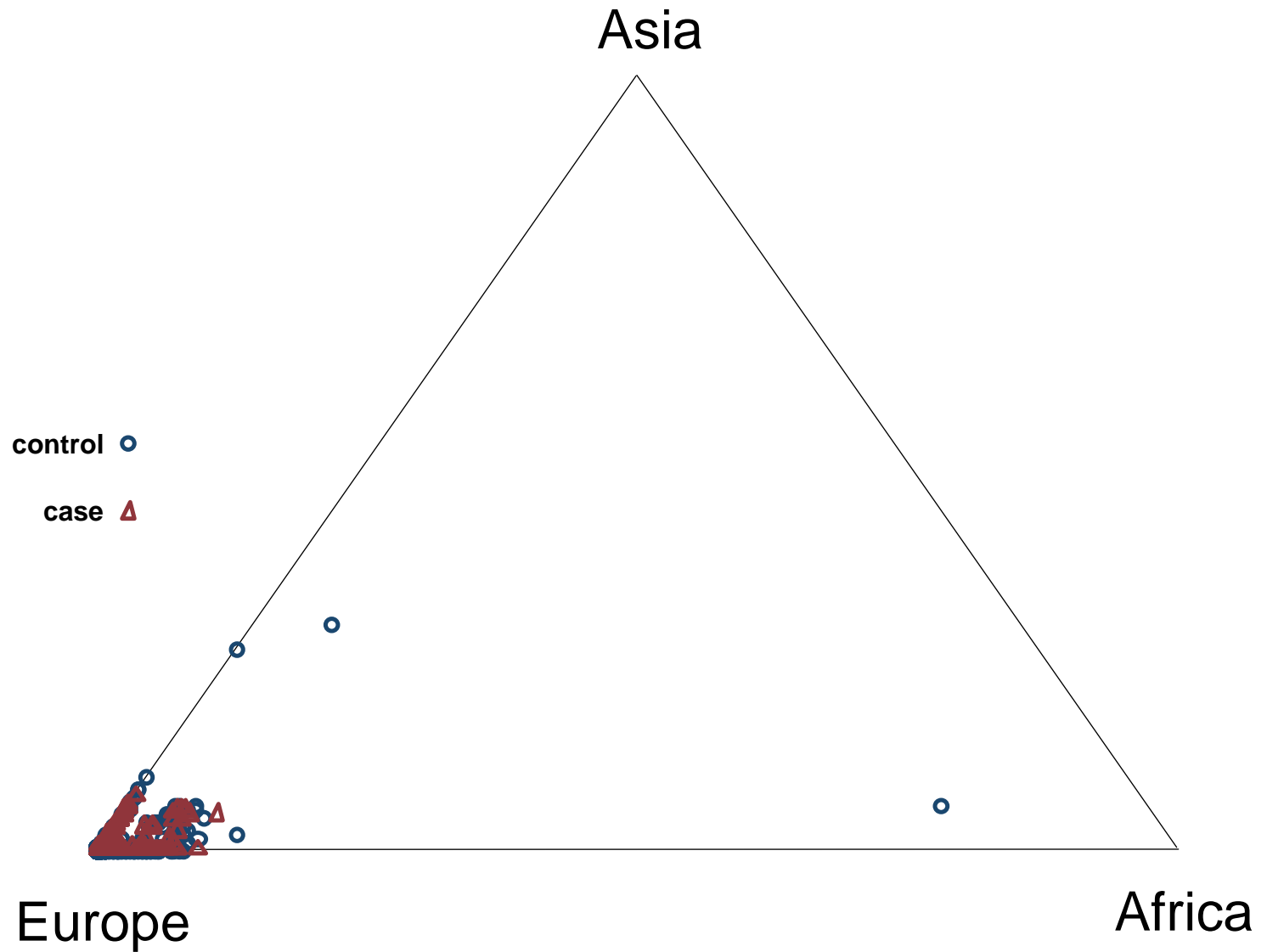
**(Hardy Weinberg Statistics)**

# Analysis of CGEMS Data Sets

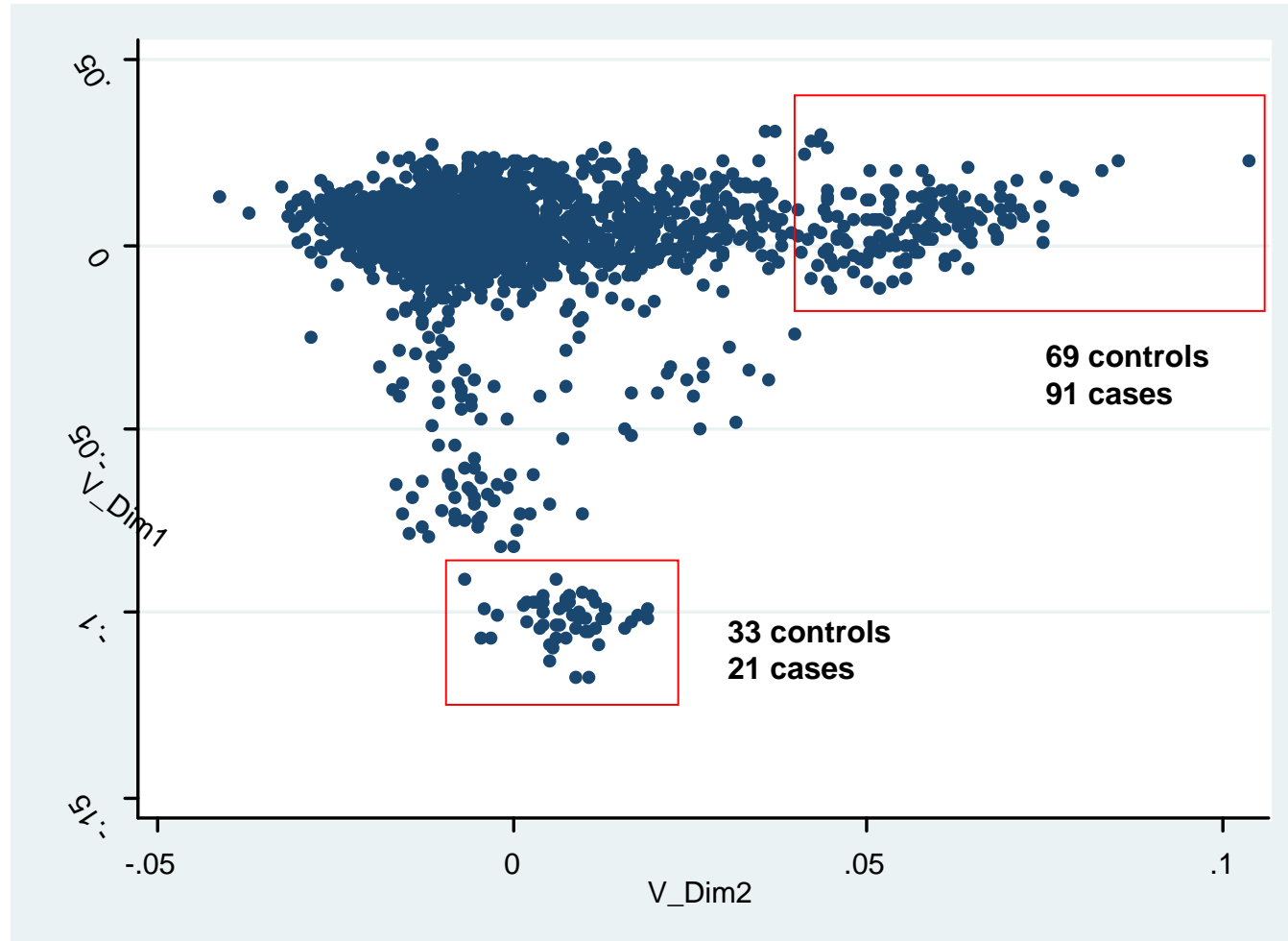
	Chromosome X		Unexpected duplicates	1st & 2nd degree relatives	Other 3rd to 5th degree relatives
	1 copy	2 copies			
Prostate	2279	3	3 pairs	5	20*
Breast	0	2299	3 pairs	1	to be done

\*It was noted subsequently that both members of each pair had been recruited in the same center.

# Admixture coefficient in PLCO samples



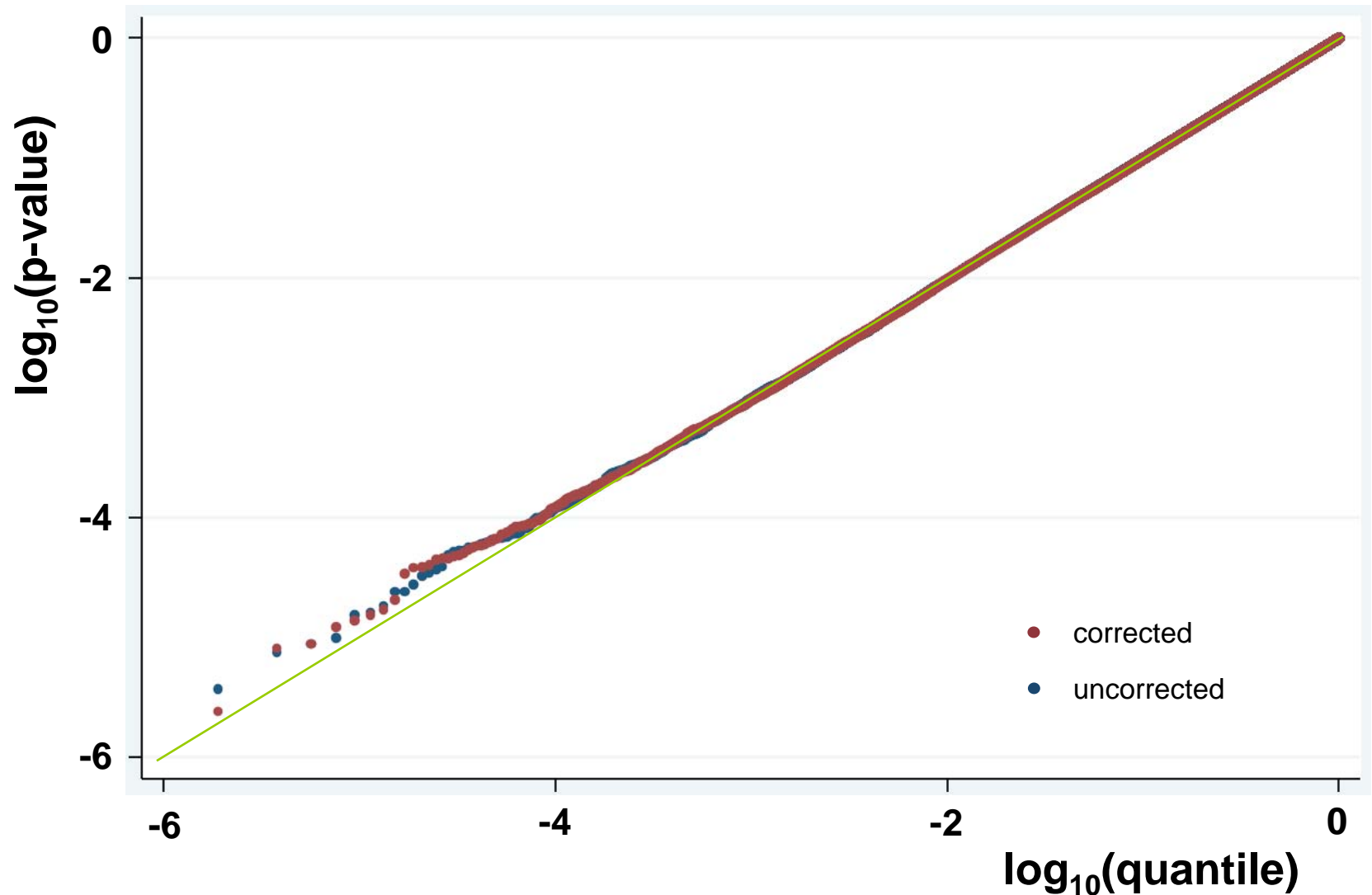
# CGEMS Prostate Cases & Controls Principal Component Analysis



*Based on Price et al Nat Genet 2006*

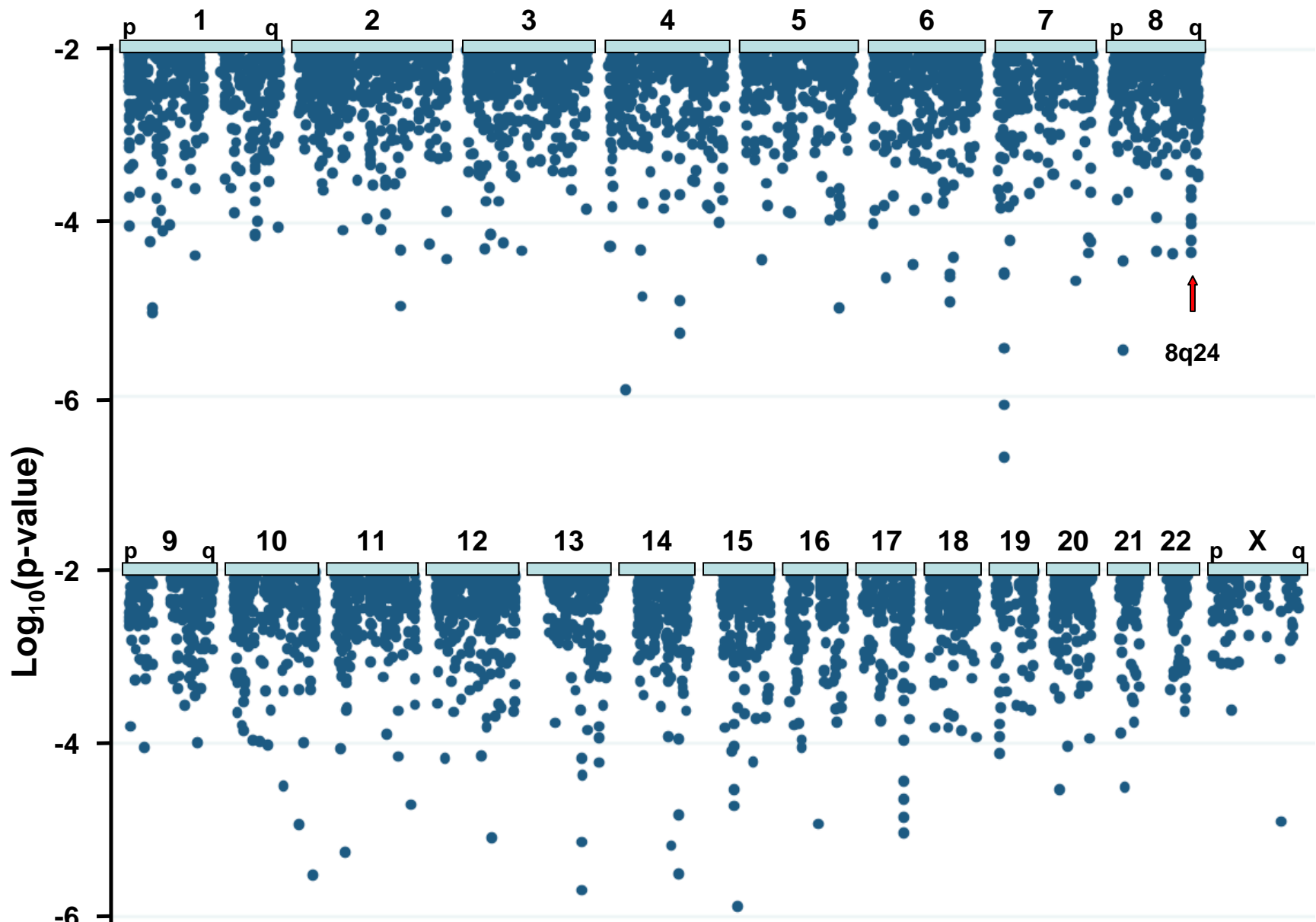
# CGEMS Breast Cancer Scan

log quantile plot of p-values for the Entire Set of Markers



# CGEMS Prostate Cancer GWAS

Chromosomes



8q24

*P values < 0.01*

# Replication Studies in CGEMS Prostate Cancer GWAS

	Subjects		rs6983267			rs1447295		
			Predisposing allele frequency		P-value	Predisposing allele frequency		P-value
			Cases	Cont.		Cases	Cont.	
PLCO	1157	1172	0.55	0.49	$2.4 \times 10^{-05}$	0.14	0.10	$9.8 \times 10^{-05}$
ACS	1151	1150	0.55	0.50	$3.2 \times 10^{-03}$	0.12	0.08	$2.7 \times 10^{-05}$
ATBC	896	894	0.57	0.51	$1.9 \times 10^{-03}$	0.21	0.17	$2.9 \times 10^{-02}$
FPCC	459	455	0.56	0.51	$1.2 \times 10^{-01}$	0.12	0.07	$4.4 \times 10^{-03}$
HPFS	636	625	0.57	0.51	$1.0 \times 10^{-02}$	0.13	0.09	$2.7 \times 10^{-03}$
ALL	4299	4296	0.56	0.50	$9.4 \times 10^{-13}$	0.15	0.11	$1.5 \times 10^{-14}$

Estimated Odds Ratios Overall

Heterozygotes

1.26

1.43

Homozygotes

1.58

2.23



# Meta-Analysis of 8q24 papers in *Nature Genetics*: J Witte

