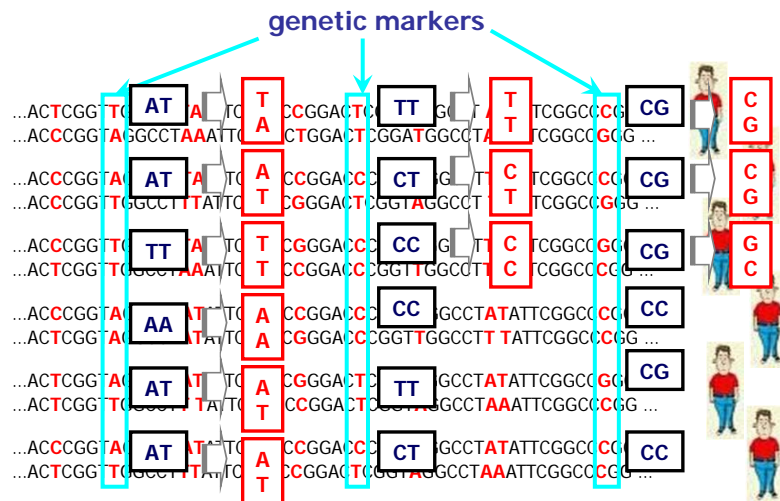# Disease Association Studies

Lectures 7 – Oct 19, 2011
CSE 527 Computational Biology, Fall 2011

Instructor: Su-In Lee
TA: Christopher Miles

Monday & Wednesday  12:00-1:20
Johnson Hall (JHN) 022

1

---

# Last Class …

- Haplotype reconstruction



**genetic markers**

**Single nucleotide polymorphism (SNP) [snip] = a variation at a single site in DNA**

2

*1*

# Outline

- **Application to disease association analysis**
  - Single marker based association tests
  - Haplotype-based approach
  - Indirect association – predicting unobserved SNPs
  - Selection of tag SNPs

- **Genetic linkage analysis**
  - Pedigree-based gene mapping
  - Elston-Stewart algorithm
  - Association vs linkage

3

# A single marker association test

- **Data**
  - Genotype data from case/control individuals
    - e.g. case: patients, control: healthy individuals

- **Goals**
  - Compare frequencies of particular alleles, or genotypes, in set of cases and controls
  - Typically, relies on standard contingency table tests
    - Chi-square goodness-of-fit test
    - Likelihood ratio test
    - Fisher's exact test

4

# Construct contingency table

- Organize genotype counts in a simple table
    - Rows: one row for cases, another for controls
    - Columns: one of each genotype (or allele)
    - Individual cells: count of observations

| i: case, control j: 0/0, 0/1, 1/1 | j=1 0/0 | j=2 0/1 | j=3 1/1 | |
|---|---|---|---|---|
| i=1 Case (affected) | $O_{1,1}$ | $O_{1,2}$ | $O_{1,3}$ | $O_{1,\cdot}=O_{1,1}+O_{1,2}+O_{1,3}$ |
| i=2 Control (unaffected) | $O_{2,1}$ | $O_{2,2}$ | $O_{2,3}$ | $O_{2,\cdot}=O_{2,1}+O_{2,2}+O_{2,3}$ |
| | $O_{\cdot,1}=O_{1,1}+O_{2,1}$ | $O_{\cdot,2}=O_{1,2}+O_{2,2}$ | $O_{\cdot,3}=O_{1,3}+O_{2,3}$ | |

- Notation
    - Let $O_{ij}$ denote the observed counts in each cell
    - Let $E_{ij}$ denote the expected counts in each cell
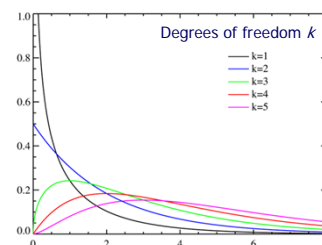        - $E_{ij} = O_{i,\cdot} \cdot O_{\cdot,j} / O_{\cdot,\cdot}$

5

---

# Goodness of fit tests (1/2)

- Null hypothesis
    - There is no statistical dependency between the genotypes and the phenotype (case/control)

- P-value
    - Probability of obtaining a test statistic at least as extreme as the one that was actually observed

- Chi-square test

$$\chi^2 = \sum_{i,j} \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}$$



- If counts are large, compare statistic to chi-squared distribution
    - p = 0.05 threshold is 5.99 for 2 df (degrees of freedom, e.g. genotype test)
    - p = 0.05 threshold is 3.84 for 1 df (e.g. allele test)
- If counts are small, exact or permutation tests are better

6

# Goodness of fit tests (2/2)

- Likelihood ratio test
  - The test statistics (usually denoted D) is twice the difference in the log-likelihoods:

$$D = -2\ln\left(\frac{\text{likelihood for null model}}{\text{likelihood for alternative model}}\right)$$

$$= -2\ln\frac{\prod_{i,j}\left(E_{i,j}/O\right)^{O_{i,j}}}{\prod_{i,j}\left(O_{i,j}/O\right)^{O_{i,j}}} = 2\sum_{i,j}O_{i,j}\ln\frac{O_{i,j}}{E_{i,j}}$$

- How about we do this for haplotypes?
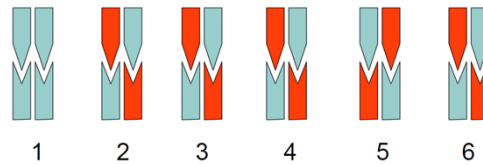  - When does it out-perform the single marker association test?

# Haplotype association tests

- Calculate haplotype frequencies in each group

- Find most likely haplotype for each group

- Fill in contingency table to compare haplotypes in the two groups (case, control)
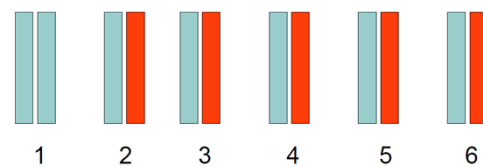
- Not recommended!

# Case genotypes & haplotypes

- Observed case genotypes



1   2   3   4   5   6

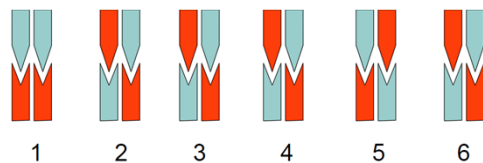  - The phase reconstruction in the five ambiguous individuals will be driven by the haplotypes observed in individual 1 ...

- Inferred case haplotypes



1   2   3   4   5   6

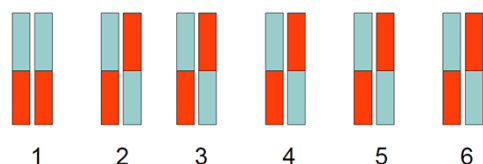  - This kind of phenomenon will occur with nearly all population based haplotyping methods!

9

# Control genotypes & haplotypes

- Observed control genotypes



1   2   3   4   5   6

  - Note these are identical, except for the single homozygous individual ...

- Inferred case haplotypes



1   2   3   4   5   6

  - Oops... The difference in a single genotype in the original data has been greatly amplified by estimating haplotypes...

10

# Haplotype association tests

- Never impute haplotypes in two groups separately

- Alternatively,
  - Consider both samples jointly
    - Schaid et al (2002) *Am J Hum Genet* **70**:425-34
    - Zaytkin et al (2002) *Hum Hered.* **53**:79-91
  - Use maximum likelihood

$$L = \prod_i \sum_{H \sim G_i} P(H)$$

individuals

Haplotype pair frequency

Possible haplotype pairs, conditional on genotype

11

# Likelihood-based test

- Calculate 3 likelihoods
  - Maximum likelihood for combined samples, $L_A$
  - Maximum likelihood for control sample, $L_B$
  - Maximum likelihood for case sample, $L_C$

$$D = 2\ln\left(\frac{L_B L_C}{L_A}\right) \sim \chi^2_{df}$$

  - df (degrees of freedom) corresponds to number of non-zero haplotype frequencies in large samples

12
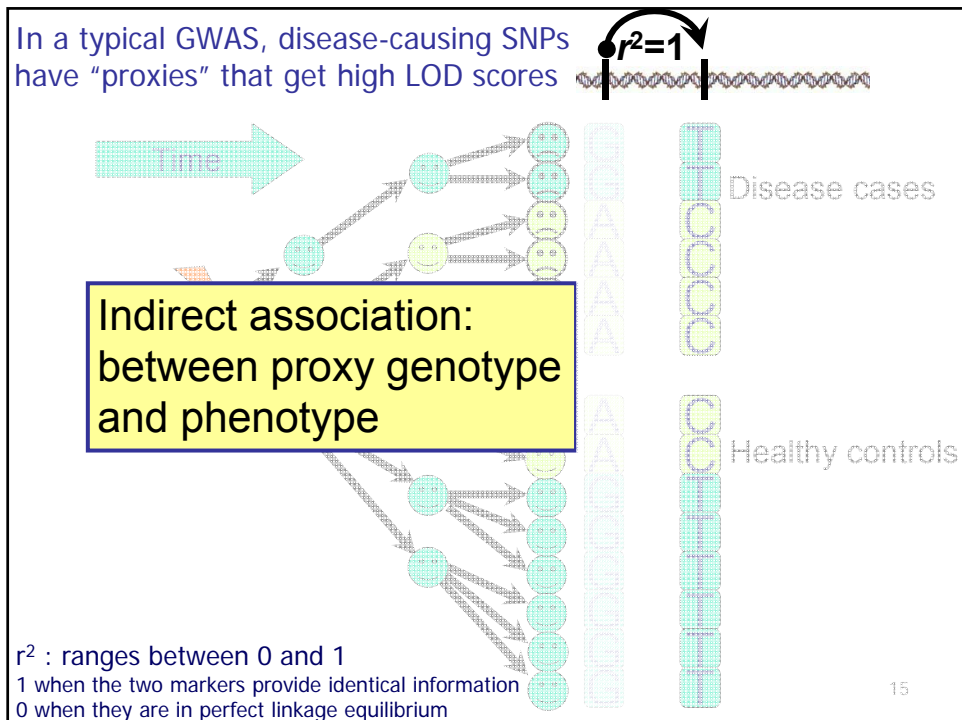
*6*

# Significance in small samples

- In reality sample sizes, it is hard to estimate the number of *df* accurately

- Instead, use a permutation approach to calculate empirical significance levels
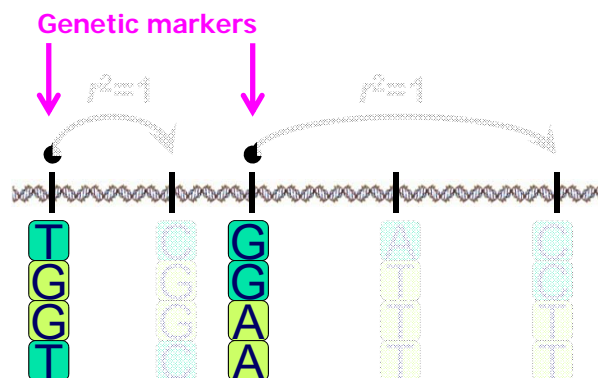
- How?

13

# Outline

- Application to disease association analysis
  - Single marker based association tests
  - Haplotype-based approach
  - Indirect association – predicting unobserved SNPs
  - Selection of tag SNPs

- Genetic linkage analysis
  - Pedigree-based gene mapping
  - Elston-Stewart algorithm
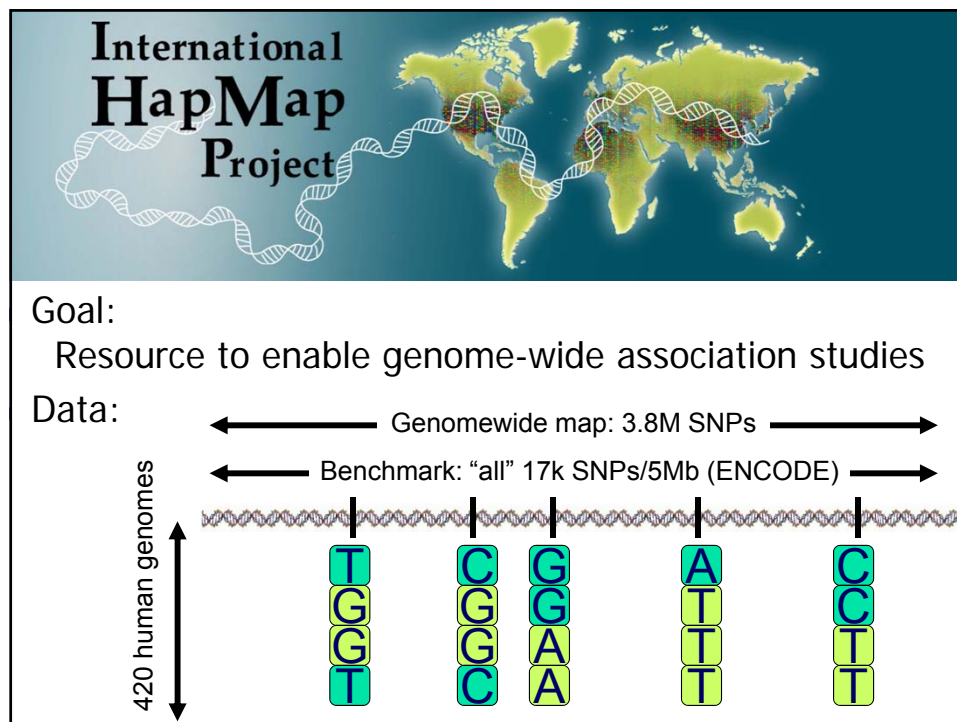  - Association vs linkage

14

In a typical GWAS, disease-causing SNPs have "proxies" that get high LOD scores

$r^2=1$

Time

Disease cases

Indirect association: between proxy genotype and phenotype

Healthy controls

$r^2$ : ranges between 0 and 1
1 when the two markers provide identical information
0 when they are in perfect linkage equilibrium

15



# Pre-requisite for association studies

Genetic markers

$r^2=1$        $r^2=1$

- How can we know which SNP pairs?
  - Very dense genotype data
  - Learn correlation between SNPs – haplotype structures
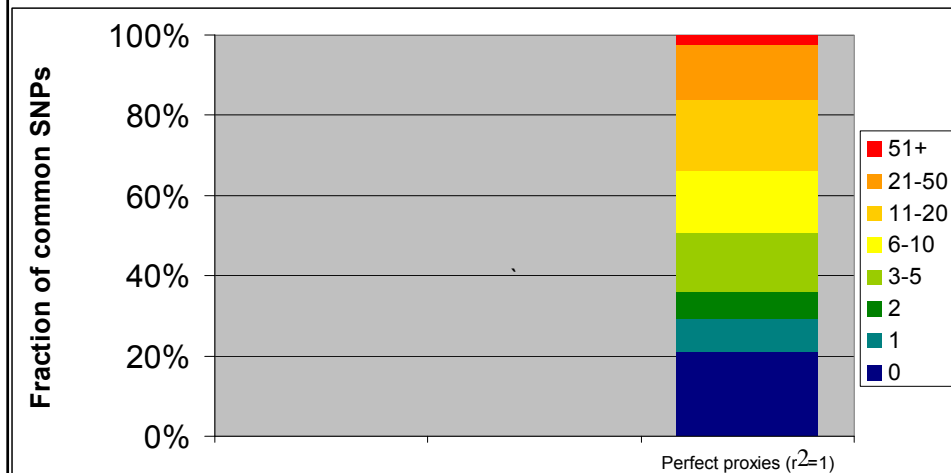- Goal: dense genome-wide association scan

International **HapMap** Project

Goal:

Resource to enable genome-wide association studies

Data:

Genomewide map: 3.8M SNPs

Benchmark: "all" 17k SNPs/5Mb (ENCODE)

420 human genomes

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

---

# Main question for HapMap:



- Are genomewide association studies doable?

or

- Do SNPs have enough proxies?
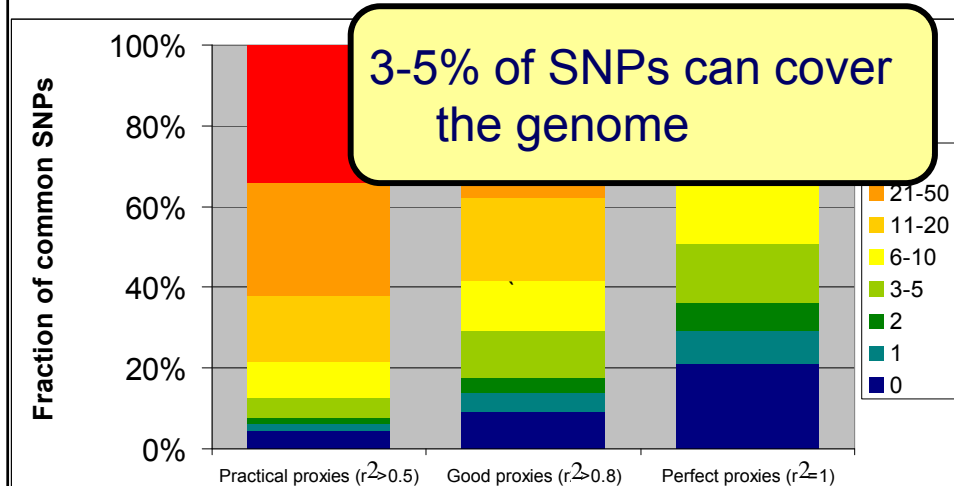
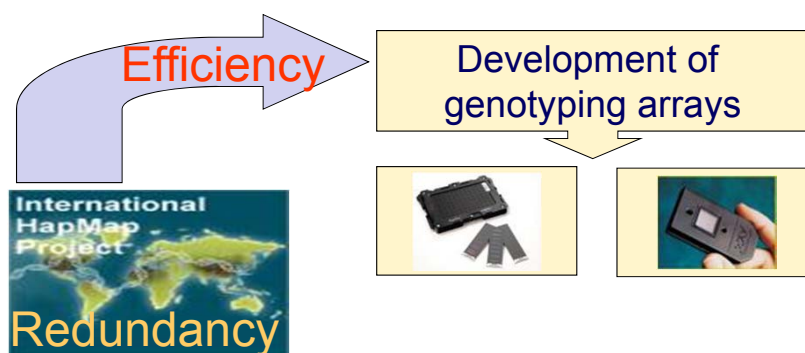18

*9*

# How many proxies will my causal SNP have?

Fraction of common SNPs

- 51+
- 21-50
- 11-20
- 6-10
- 3-5
- 2
- 1
- 0

Perfect proxies ($r^2=1$)

19

# Imperfect proxies

$r^2=1$   $r^2=0.75$

Disease cases

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

Healthy controls

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

20

# How many proxies will my causal SNP have?



3-5% of SNPs can cover the genome

Fraction of common SNPs

100% / 80% / 60% / 40% / 20% / 0%

Legend: 21-50 / 11-20 / 6-10 / 3-5 / 2 / 1 / 0

Practical proxies ($r^2>0.5$) — Good proxies ($r^2>0.8$) — Perfect proxies ($r^2=1$)

21

# Computational challenges



Efficiency

Development of genotyping arrays

International HapMap Project

Redundancy

22

# Optimizing SNP-set efficiency

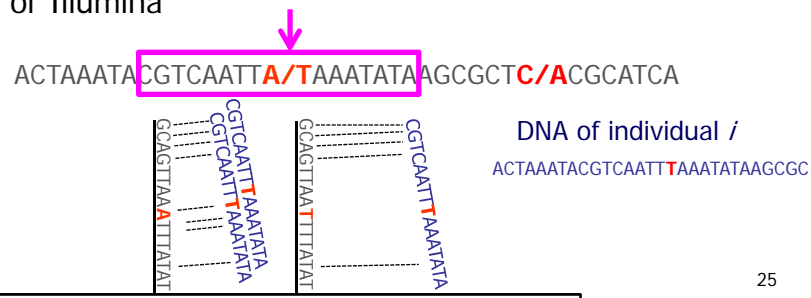- Select "tag" SNPs that maximize the number of other SNPs whose alleles are revealed by them

Markers tested:

high $r^2$    high $r^2$

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

Markers captured:  ☑    ☑  ☑    ☑    ☑

- **How?**

23

---

# Computational challenges

Efficiency
(tag SNP selection)

Development of genotyping arrays

International HapMap Project

Redundancy

Genotyping study cohort

Power
(predicting unobserved SNPs)

Analysis

24

# Analysis questions

- Can we quantify the coverage of common sequence variations measured by genome-wide SNP genotyping arrays?

- SNP genotyping arrays
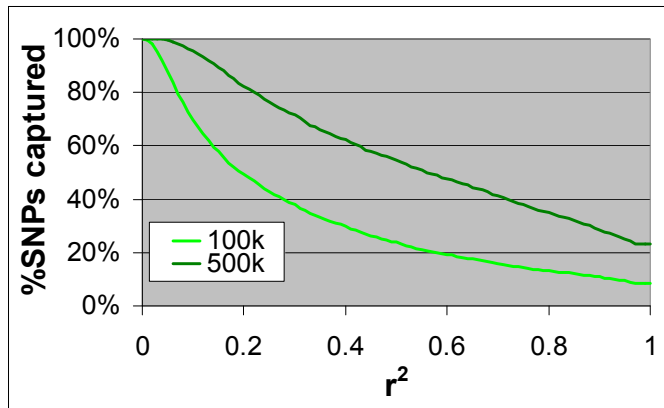  - Arrays covering 100K/500K/1M SNPs from Affymetrix or Illumina

ACTAAATACGTCAATT**A/T**AAATATAAGCGCT**C/A**CGCATCA

DNA of individual *i*

ACTAAATACGTCAATT**T**AAATATAAGCGC

25

---

# Association tests with fixed markers

Tests of association:

high *r²*     high *r²*

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

SNPs captured:   ☑     ☑     ☑           ☑

26

---

*13*

# Arrays cover many common alleles



**Panel:**

African (most diverse)

27

# Arrays cover many common alleles
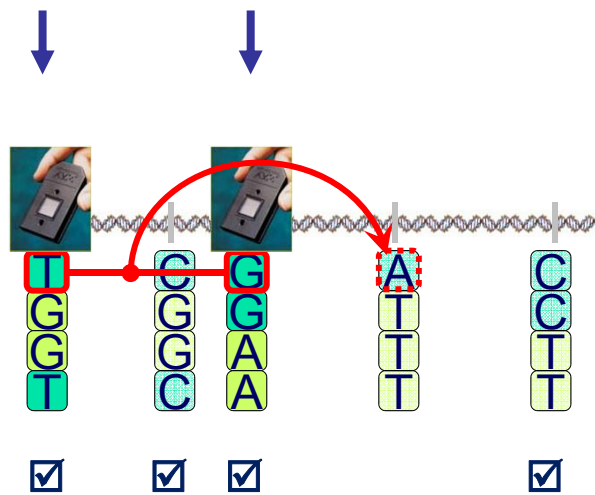


**Panel:**

European

28

# Analysis questions

- Can we quantify the coverage of common sequence variations measured by genome-wide SNP genotyping arrays?
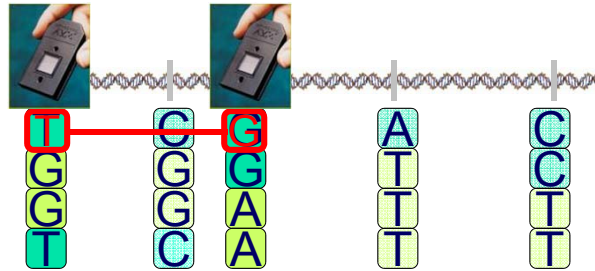
- Can we do better?

# Association with haplotypes
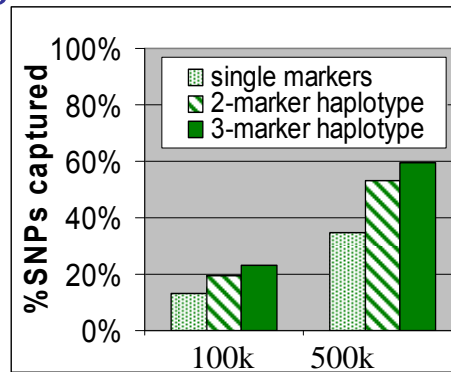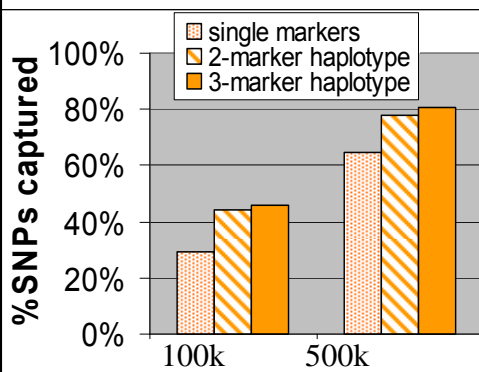
Tests of association:

SNPs captured:

Association with haplotypes



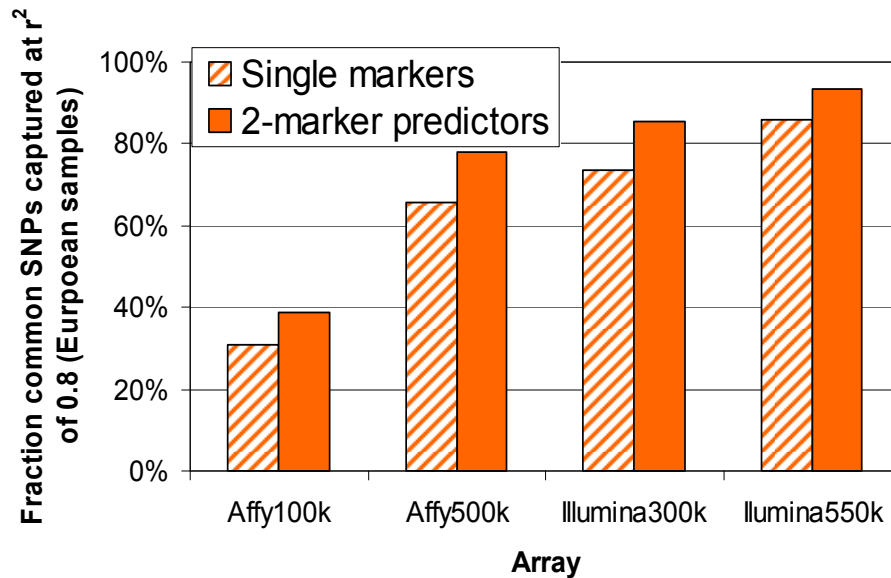Increasing coverage ($r^2=0.8$) by specified haplotypes

Panel: European

Panel: African (most diverse)

# Which platform to use?



# Summary

- Association analysis is a powerful strategy for common disease research

- HapMap and genomewide technologies enable whole-genome association scans

34