

Linkage Disequilibrium

Biostatistics 666

Logistics: Office Hours

- Office hours on Mondays at 4 pm.
- Room 4614
School of Public Health Tower

Previously ...

- Basic properties of a locus
 - Allele Frequencies
 - Genotype Frequencies
- Hardy-Weinberg Equilibrium
 - Relationship between allele and genotype frequencies that holds for most genetic markers
- Exact Tests for Hardy-Weinberg Equilibrium

Today ...

- We'll consider properties of pairs of alleles
- Haplotype frequencies
- Linkage equilibrium
- Linkage disequilibrium

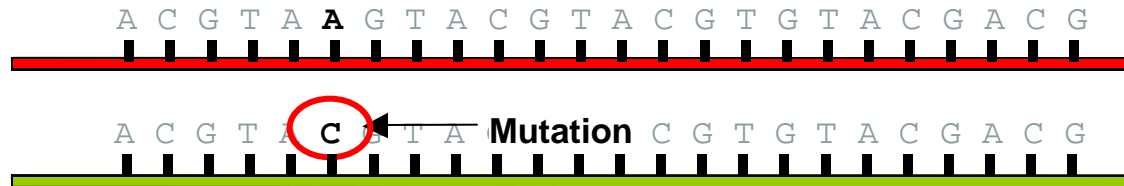
*Let's consider the history of
two neighboring alleles...*

Alleles that exist today arose through ancient mutation events...

Before Mutation



After Mutation

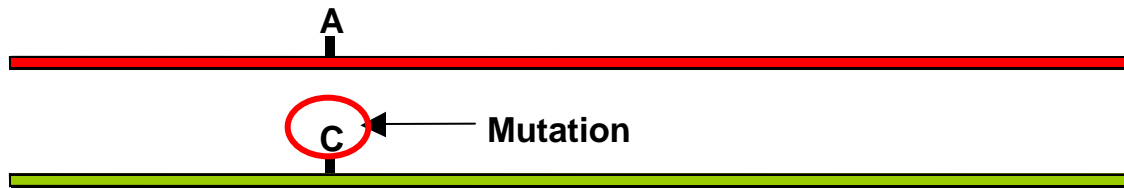


Alleles that exist today arose through ancient mutation events...

Before Mutation

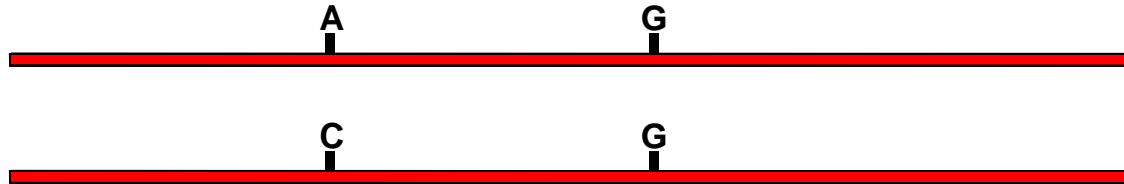


After Mutation

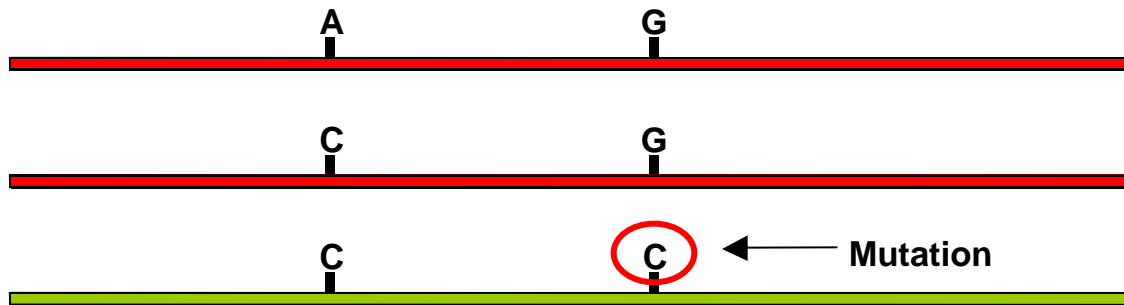


One allele arose first,
and then the other...

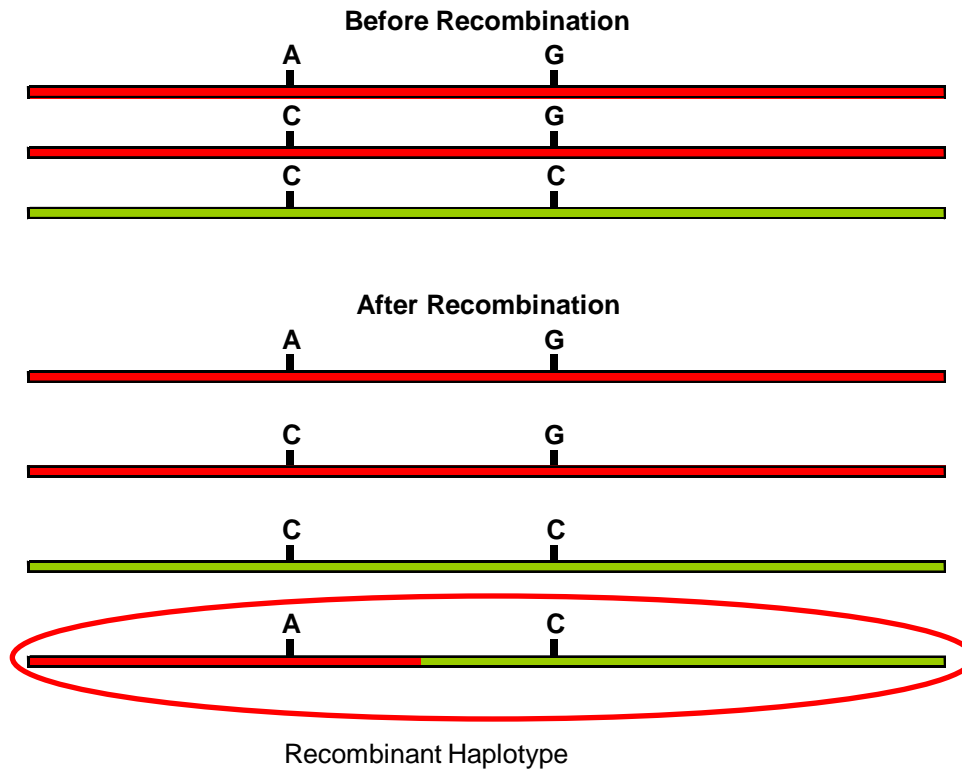
Before Mutation



After Mutation



Recombination generates new arrangements for ancestral alleles



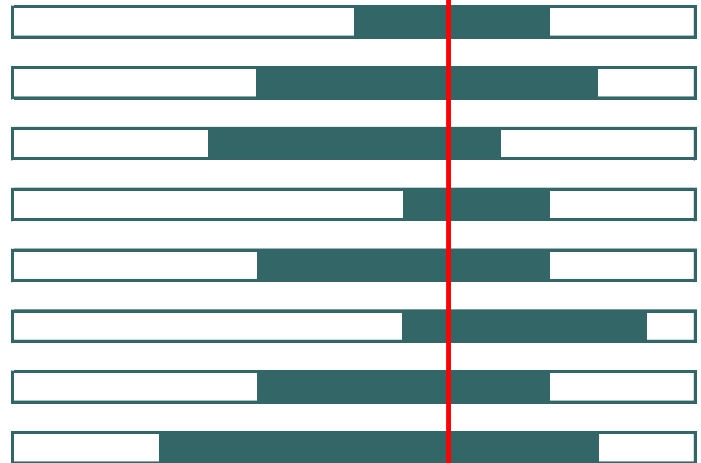
Linkage Disequilibrium

- Chromosomes are mosaics
- Extent and conservation of mosaic pieces depends on
 - Recombination rate
 - Mutation rate
 - Population size
 - Natural selection
- Combinations of alleles at very close markers reflect ancestral haplotypes

Ancestor



Present-day



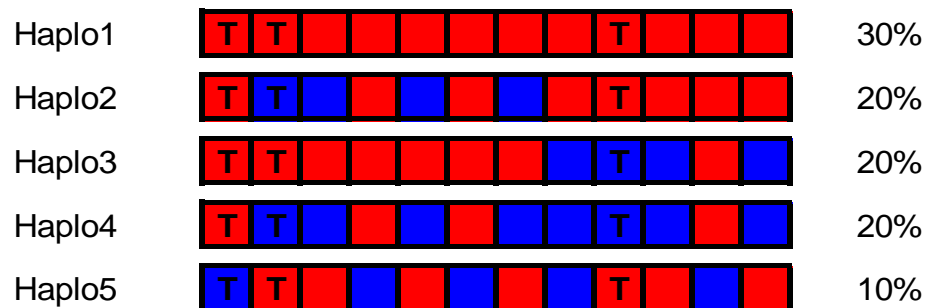
Why is linkage disequilibrium important for gene mapping?

Association Studies and Linkage Disequilibrium

- If all polymorphisms were independent at the population level, association studies would have to examine every one of them...
- Linkage disequilibrium makes tightly linked variants strongly correlated producing cost savings for association studies

Tagging SNPs

- In a typical short chromosome segment, there are only a few distinct haplotypes
- Carefully selected SNPs can determine status of other SNPs



*Basic Descriptors of
Linkage Disequilibrium*

Commonly Used Descriptors

- Haplotype Frequencies
 - The frequency of each type of chromosome
 - Contain all the information provided by other summary measures

- Commonly used summaries
 - D
 - D'
 - r^2 or Δ^2

Haplotype Frequencies

		<u>Locus B</u>		Totals
		<i>B</i>	<i>b</i>	
<u>Locus A</u>	<i>A</i>	p_{AB}	p_{Ab}	p_A
	<i>a</i>	p_{aB}	p_{ab}	p_a
Totals		p_B	p_b	1.0

Linkage Equilibrium Expected for Distant Loci

$$P_{AB} = P_A P_B$$

$$P_{Ab} = P_A P_b = P_A (1 - P_B)$$

$$P_{aB} = P_a P_B = (1 - P_A) P_B$$

$$P_{ab} = P_a P_b = (1 - P_A)(1 - P_B)$$

Linkage Disequilibrium Expected for Nearby Loci

$$p_{AB} \neq p_A p_B$$

$$p_{Ab} \neq p_A p_b = p_A(1 - p_B)$$

$$p_{aB} \neq p_a p_B = (1 - p_A)p_B$$

$$p_{ab} \neq p_a p_b = (1 - p_A)(1 - p_B)$$

Dis-equilibrium Coefficient D_{AB}

$$D_{AB} = p_{AB} - p_A p_B$$

$$p_{AB} = p_A p_B + D_{AB}$$

$$p_{Ab} = p_A p_b - D_{AB}$$

$$p_{aB} = p_a p_B - D_{AB}$$

$$p_{ab} = p_a p_b + D_{AB}$$

D_{AB} is hard to interpret

- Sign is arbitrary ...
 - A common convention is to set A, B to be the common allele and a, b to be the rare allele
- Range depends on allele frequencies
 - Hard to compare between markers

What is the range of D_{AB} ?

- What are the maximum and minimum possible values of D_{AB} when
 - $p_A = 0.3$ and $p_B = 0.3$
 - $p_A = 0.2$ and $p_B = 0.1$
- Can you derive a general formula for this range?

D' – A scaled version of D

$$D'_{AB} = \begin{cases} \frac{D_{AB}}{\min(p_A p_B, p_a p_b)} & D_{AB} < 0 \\ \frac{D_{AB}}{\min(p_A p_b, p_a p_B)} & D_{AB} > 0 \end{cases}$$

- Ranges between -1 and $+1$
 - More likely to take extreme values when allele frequencies are small
 - ± 1 implies at least one of the observed haplotypes was not observed

More on D'

- **Pluses:**

- $D' = 1$ or $D' = -1$ means no evidence for recombination between the markers
- If allele frequencies are similar, high D' means the markers are good surrogates for each other

- **Minuses:**

- D' estimates inflated in small samples
- D' estimates inflated when one allele is rare

Δ^2 (also called r^2)

$$\Delta^2 = \frac{D_{AB}^2}{p_A(1-p_A)p_B(1-p_B)}$$
$$= \frac{\chi^2}{2n}$$

- Ranges between 0 and 1
 - 1 when the two markers provide identical information
 - 0 when they are in perfect equilibrium
- Expected value is $1/2n$

More on r^2

- $r^2 = 1$ implies the markers provide exactly the same information
- The measure preferred by population geneticists
- Measures loss in efficiency when marker A is replaced with marker B in an association study
 - With some simplifying assumptions (e.g. see Pritchard and Przeworski, 2001)

*When does
linkage equilibrium hold?*

Equilibrium or Disequilibrium?

- We will present simple argument for why linkage equilibrium holds for most loci
- Balance of factors
 - Genetic drift (a function of population size)
 - Random mating
 - Distance between markers
 - ...

Why Equilibrium is Reached...

- Eventually, random mating and recombination should ensure that mutations spread from original haplotype to all haplotypes in the population...
- Simple argument:
 - Assume fixed allele frequencies over time

Generation t, Initial Configuration

	B	b	
A	$p_A p_B + D_{AB}$	$p_A p_b - D_{AB}$	p_A
a	$p_a p_B - D_{AB}$	$p_a p_b + D_{AB}$	p_a
	p_B	p_b	

Assume arbitrary values for the allele frequencies
and disequilibrium coefficient

Generation $t+1$, Without Recombination

	B	b	
A	$p_A p_B + D_{AB}$	$p_A p_b - D_{AB}$	p_A
a	$p_a p_B - D_{AB}$	$p_a p_b + D_{AB}$	p_a
	p_B	p_b	

Haplotype Frequencies Remain Stable Over Time
Outcome has probability $1 - R$

Generation t+1, With Recombination

	B	b	
A	$p_A p_B$	$p_A p_b$	p_A
a	$p_a p_b$	$p_a p_B$	p_a
	p_B	p_b	

Haplotype Frequencies Are Function of Allele Frequencies
Outcome has probability R

Generation t+1, Overall

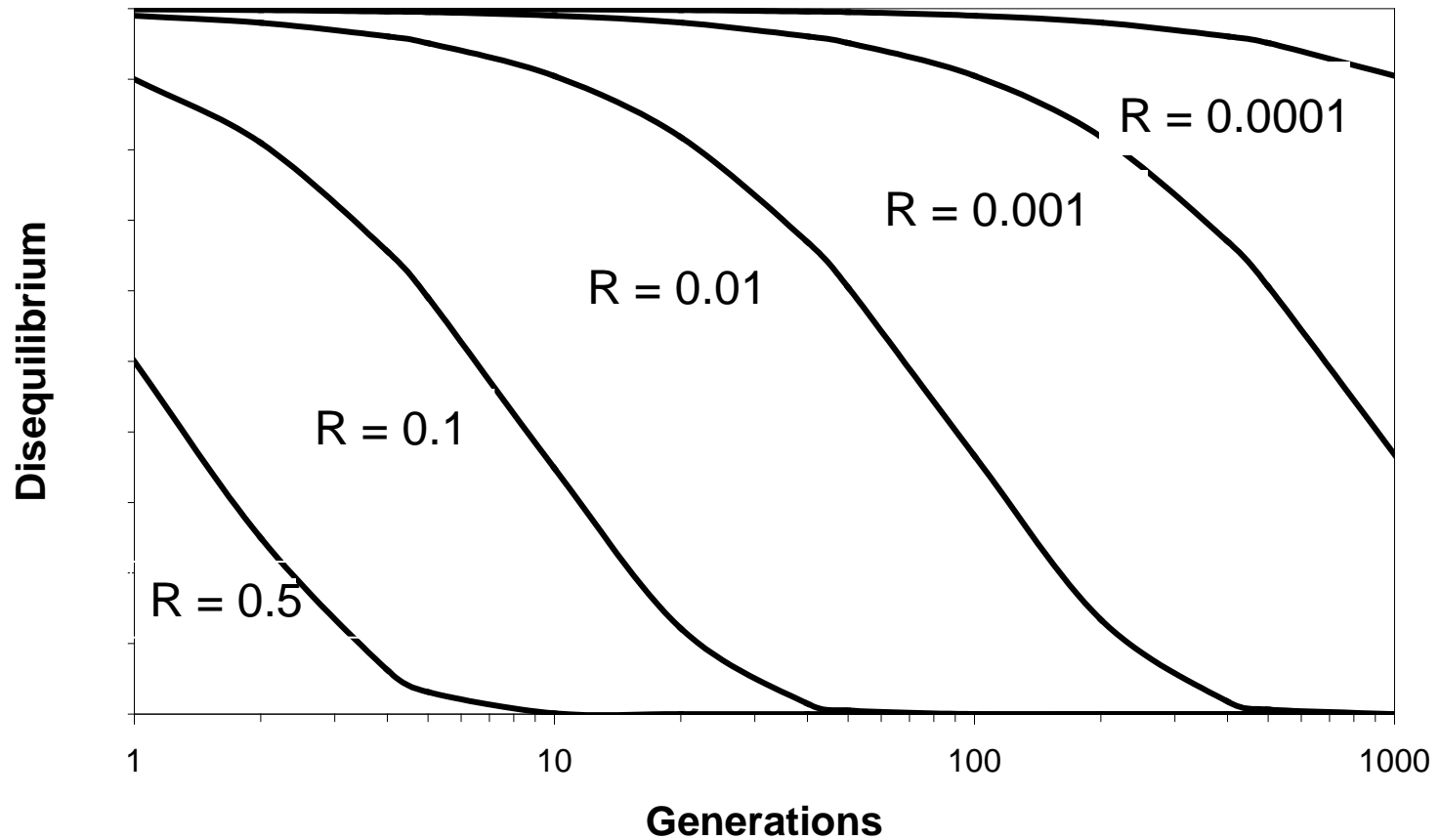
	B	b	
A	$p_A p_B + (1-r)D_{AB}$	$p_A p_b - (1-r)D_{AB}$	p_A
a	$p_A p_b - (1-r)D_{AB}$	$p_a p_b + (1-r)D_{AB}$	p_a
	p_B	p_b	

Disequilibrium Decreases...

Recombination Rate (R)

- Probability of an odd number of crossovers between two loci
- Proportion of time alleles from two different grand-parents occur in the same gamete
- Increases with physical (base-pair) distance, but rate of increase varies across genome

Decay of D with Time



Predictions

- Disequilibrium will decay each generation
 - In a large population
- After t generations...
 - $D_{AB}^t = (1-R)^t D_{AB}^0$
- A better model should allow for changes in allele frequencies over time...

Linkage Equilibrium

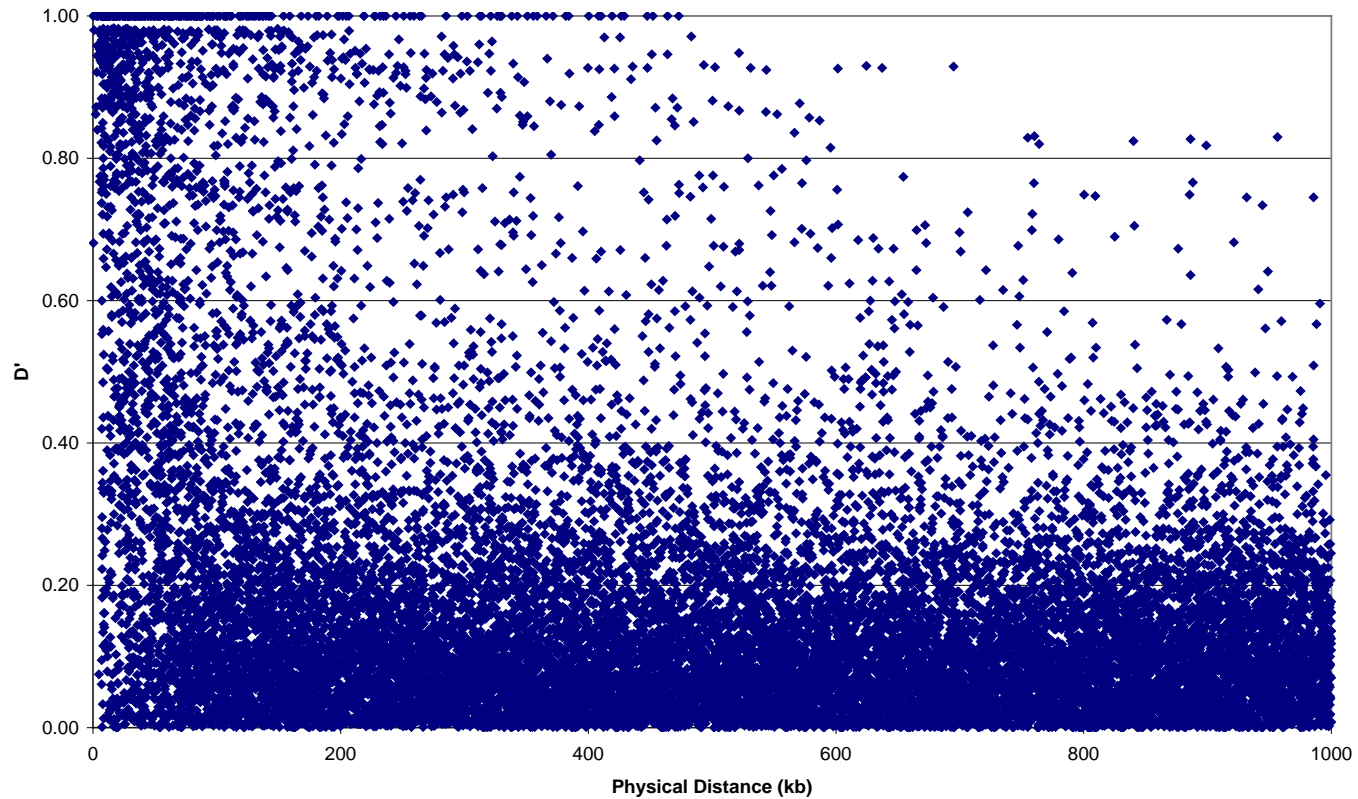
- In a large random mating population haplotype frequencies converge to a simple function of allele frequencies

*Some Examples of Linkage
Disequilibrium Data*

Summary of Disequilibrium in the Genome

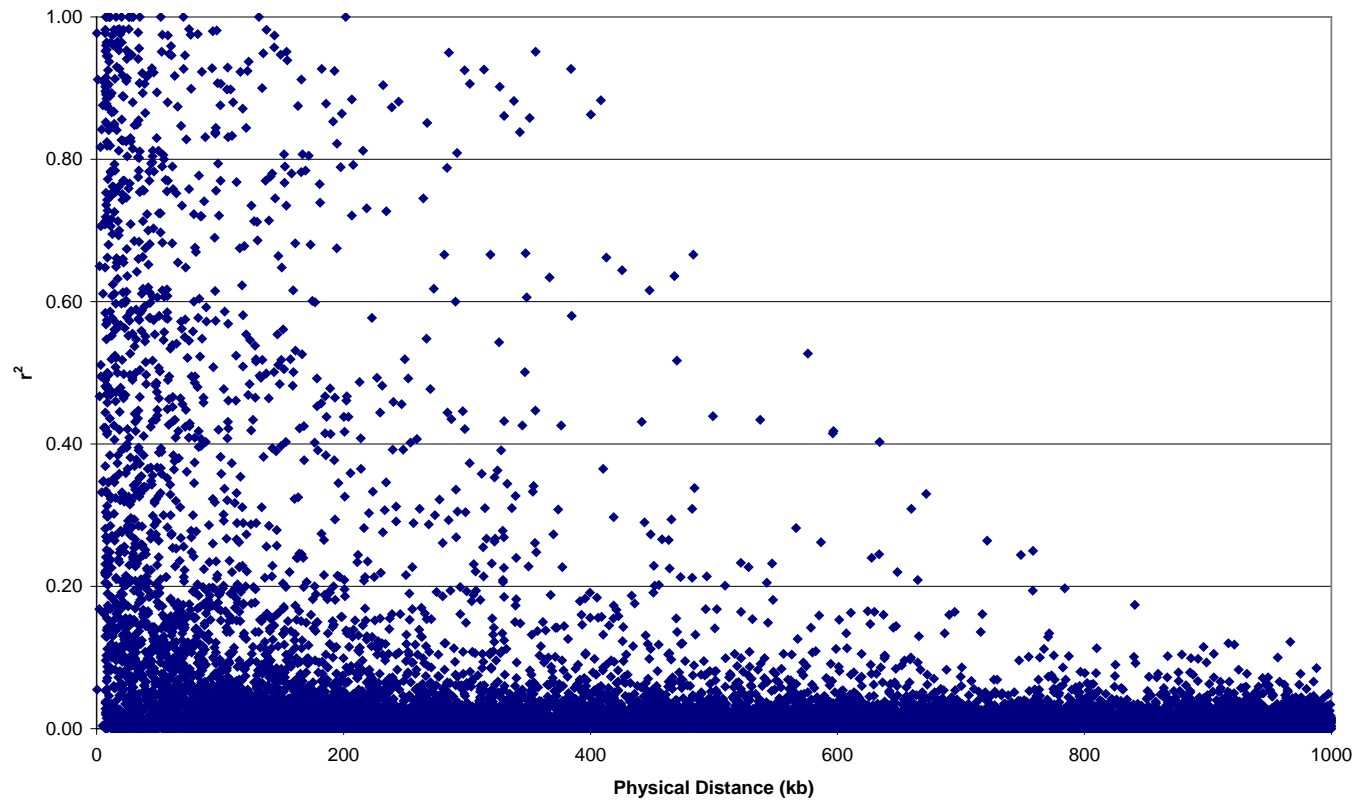
- How much disequilibrium is there?
- What are good predictors of disequilibrium?

Raw $|D'|$ data from Chr22



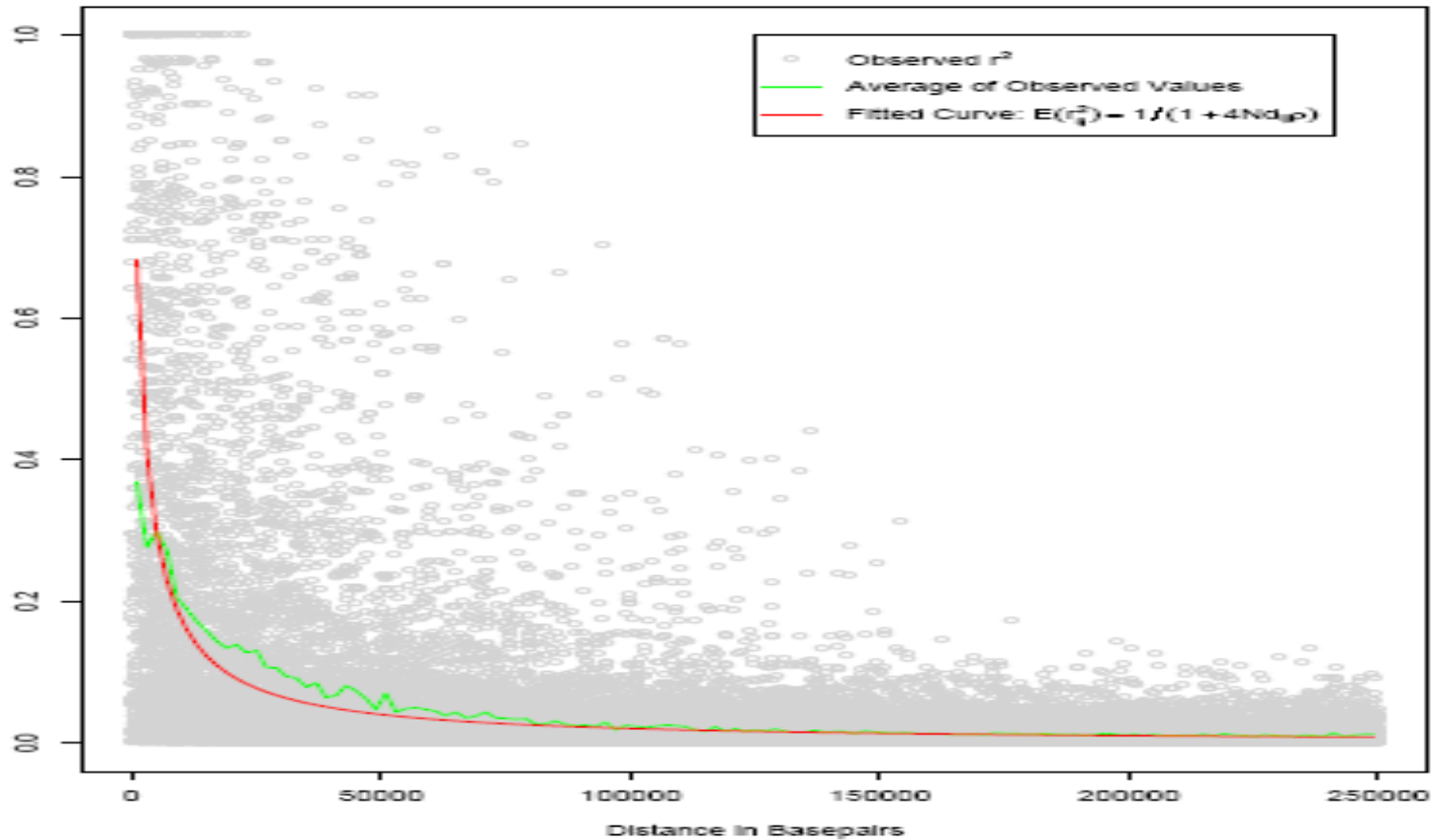
Dawson et al, *Nature*, 2002

Raw Δ^2 data from Chr22

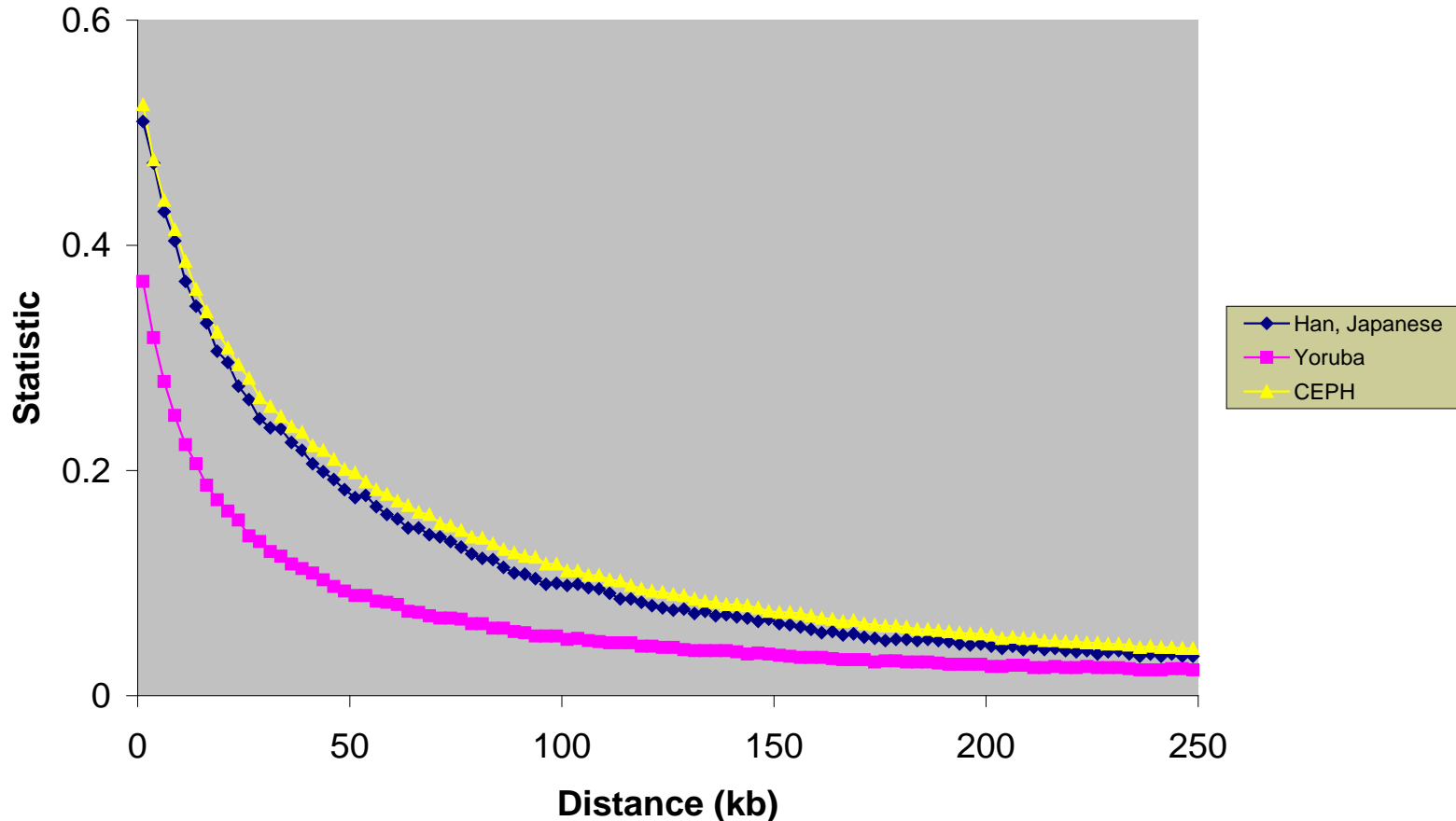


Dawson et al, *Nature*, 2002

Summarizing Disequilibrium



Comparing Populations ...



LD extends further in CEPH and the Han/Japanese than in the Yoruba

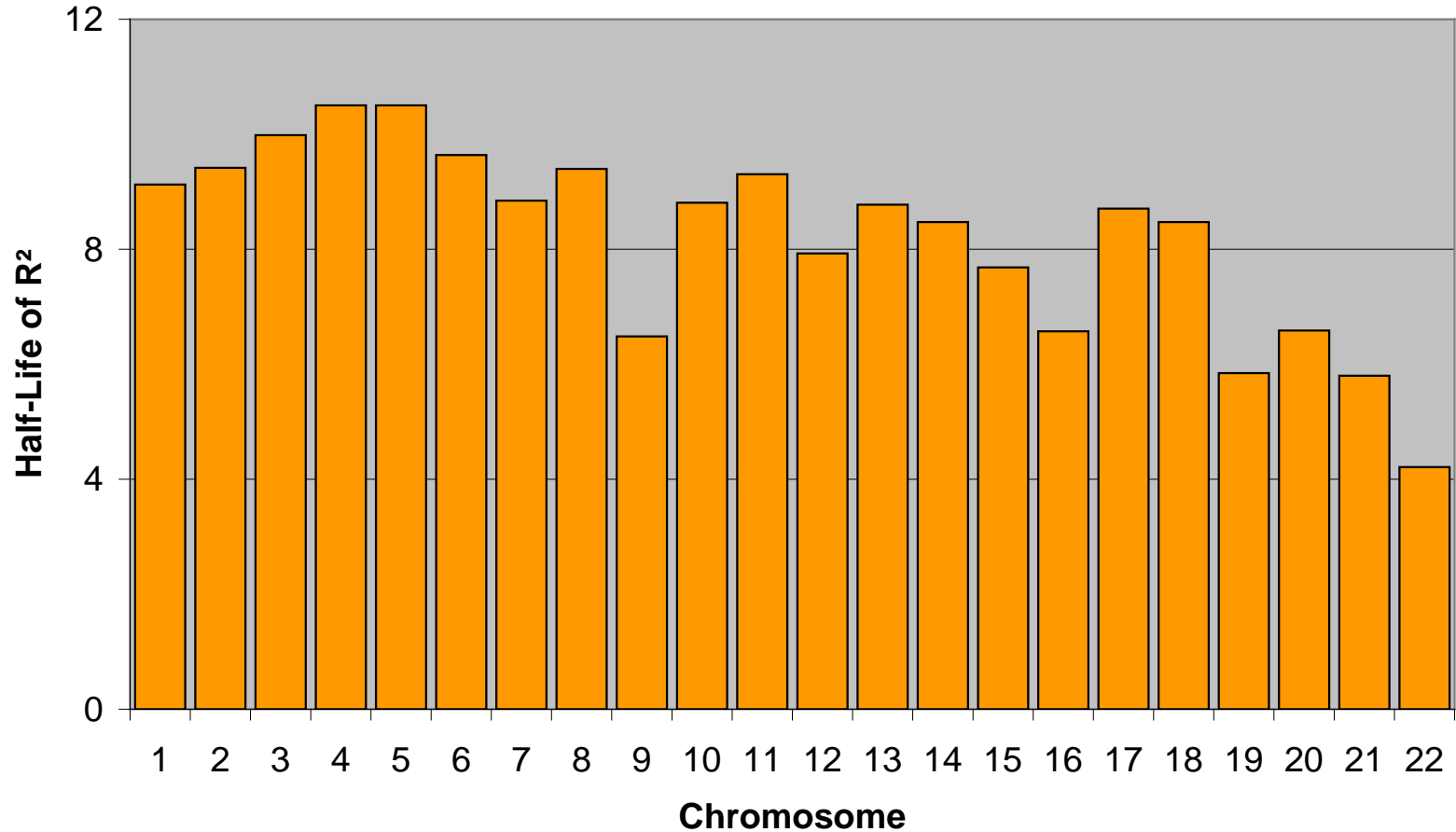
International HapMap Consortium, *Nature*, 2005

Variation in Linkage Disequilibrium Along The Genome

Comparing Genomic Regions ...

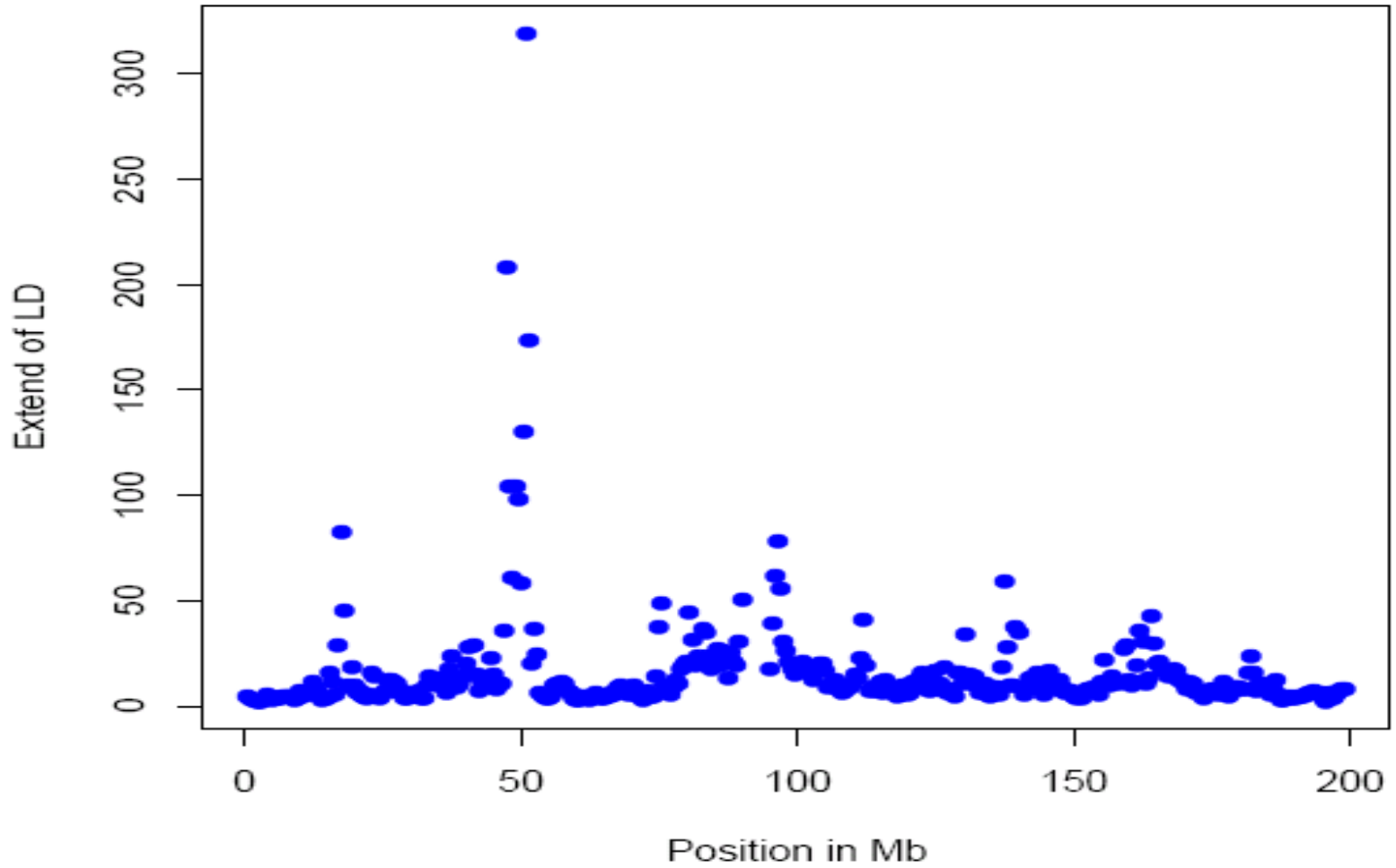
- Rather than compare curves directly, it is convenient to pick a summary for the decay curves
- One common summary is the distance at which the curve crosses a threshold of interest (say 0.50)

Extent of Linkage Disequilibrium



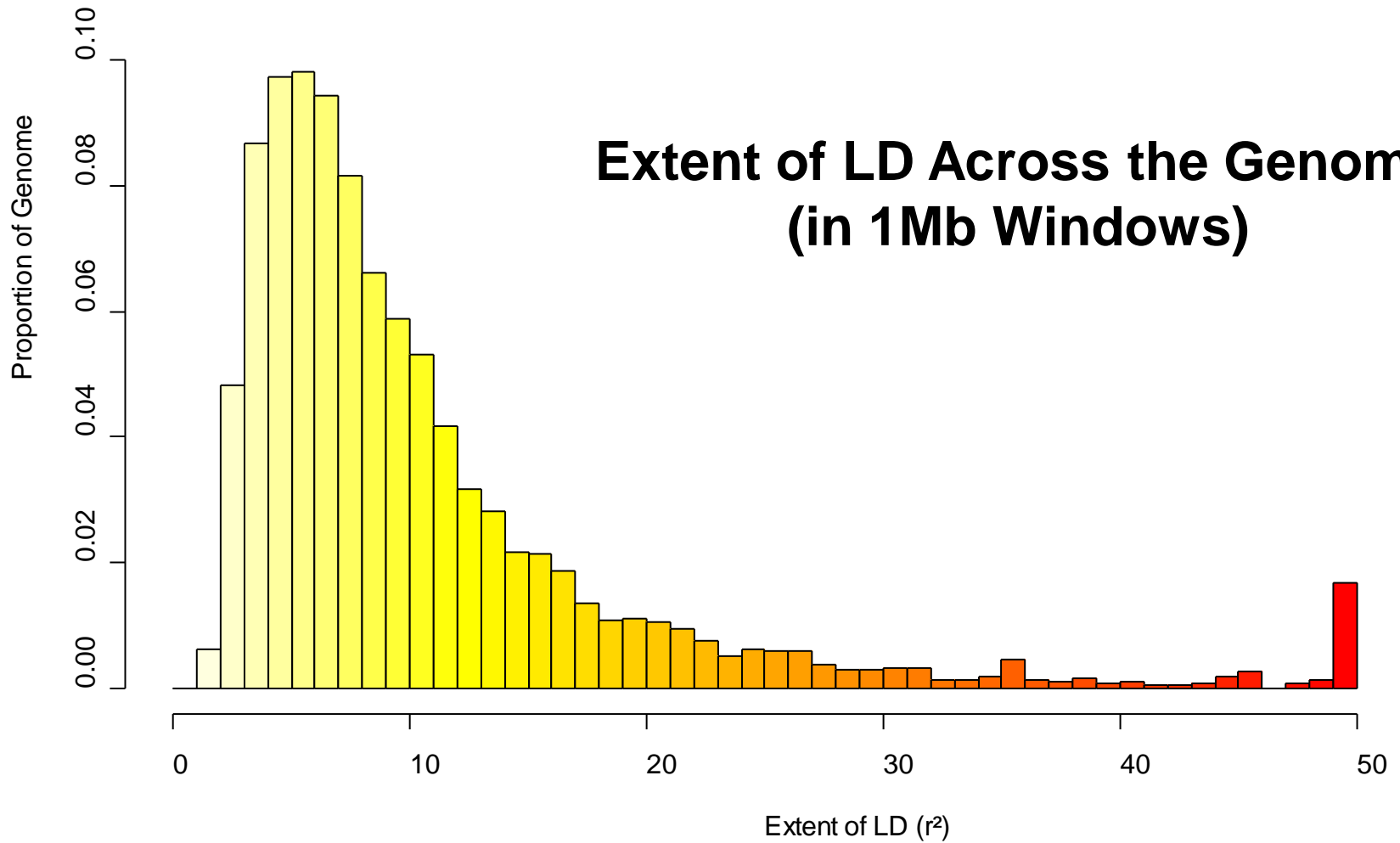
LD extends further in the larger chromosomes, which have lower recombination rates

Chromosome 3



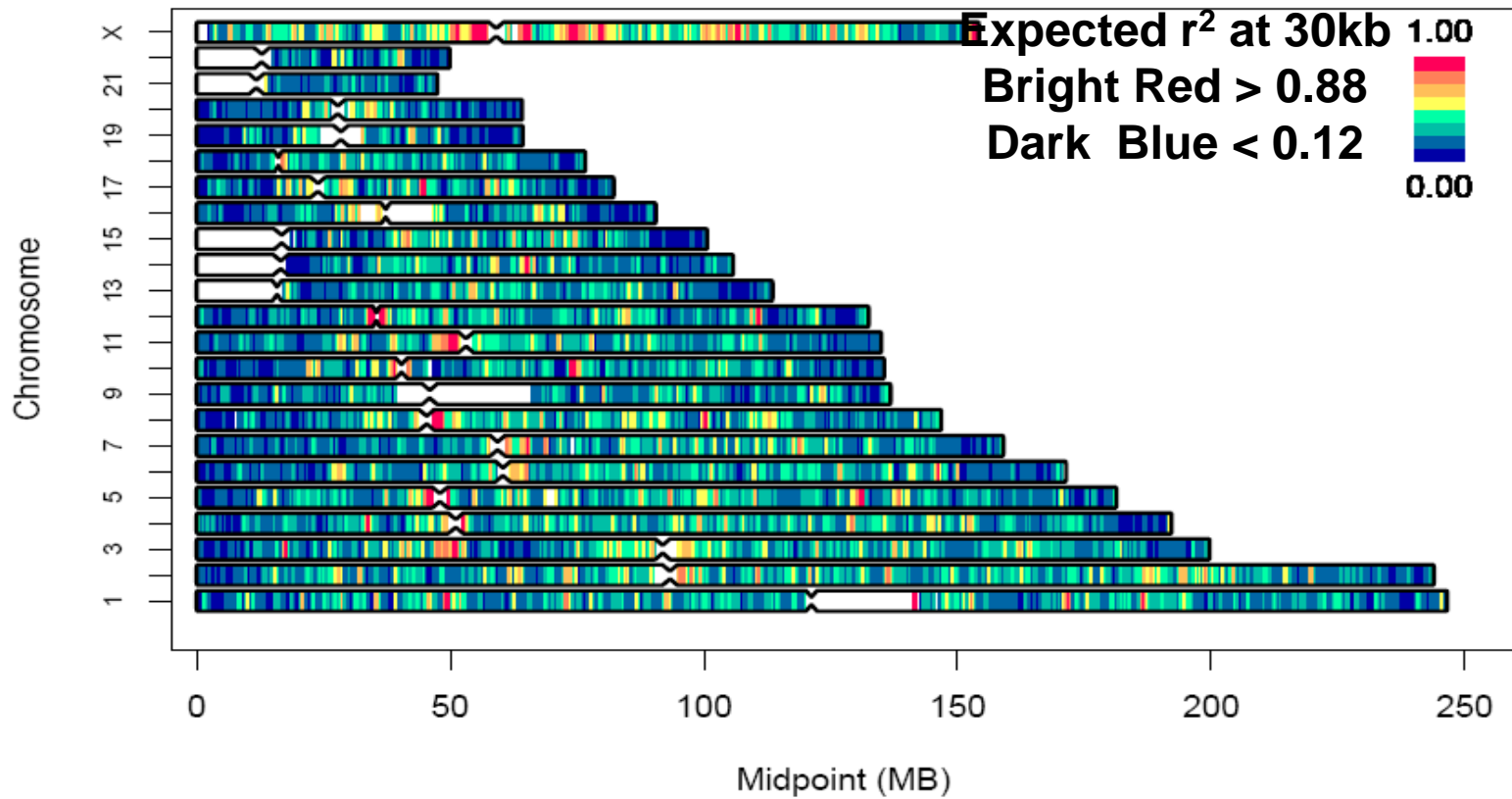
But within each chromosome, there is still huge variability!

Extent of LD



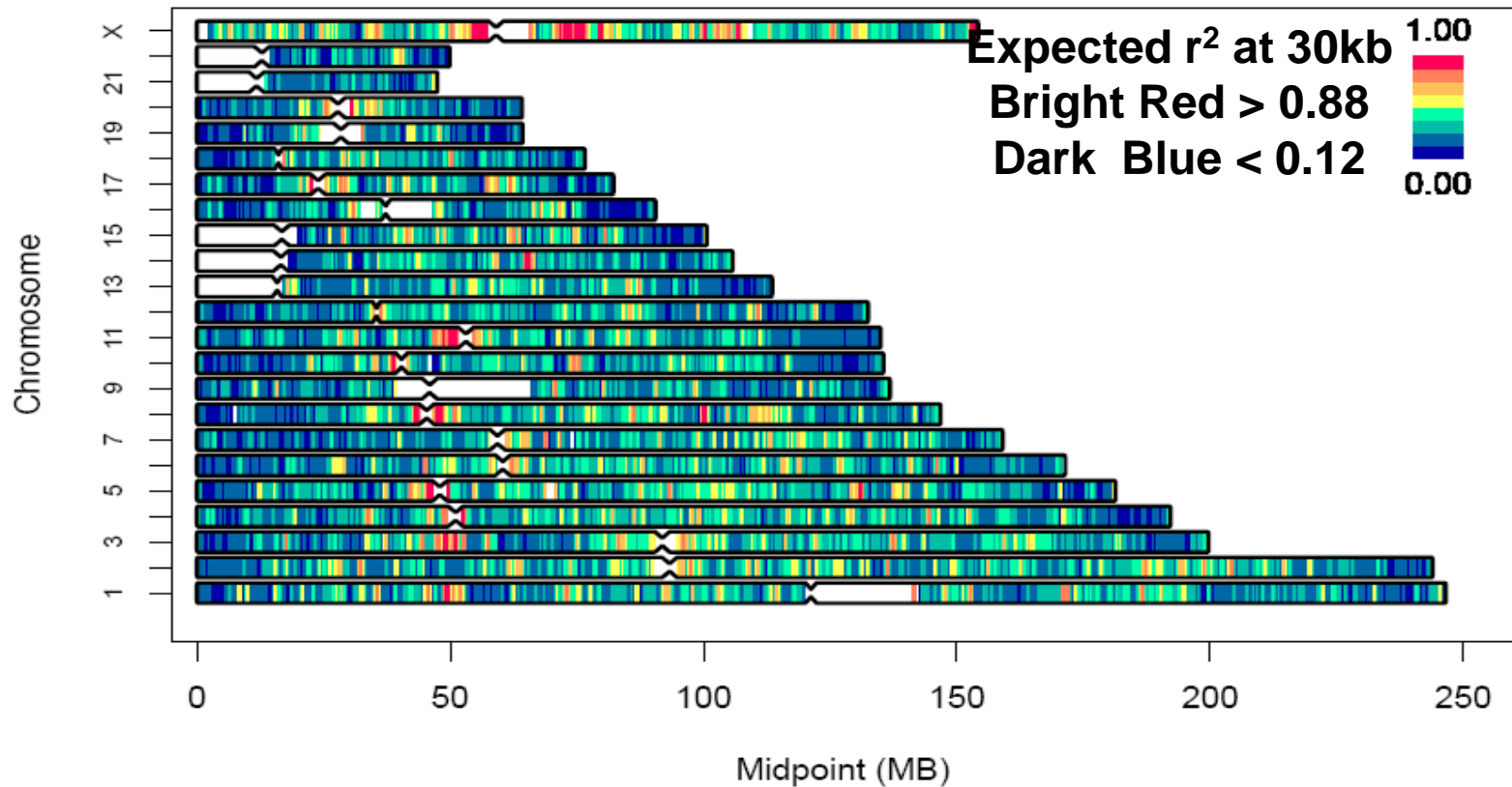
Average Extent: 11.9 kb
Median Extent: 7.8 kb
10th percentile: 3.5 kb
90th percentile: 20.9 kb

Genomic Variation in LD (CEPH)



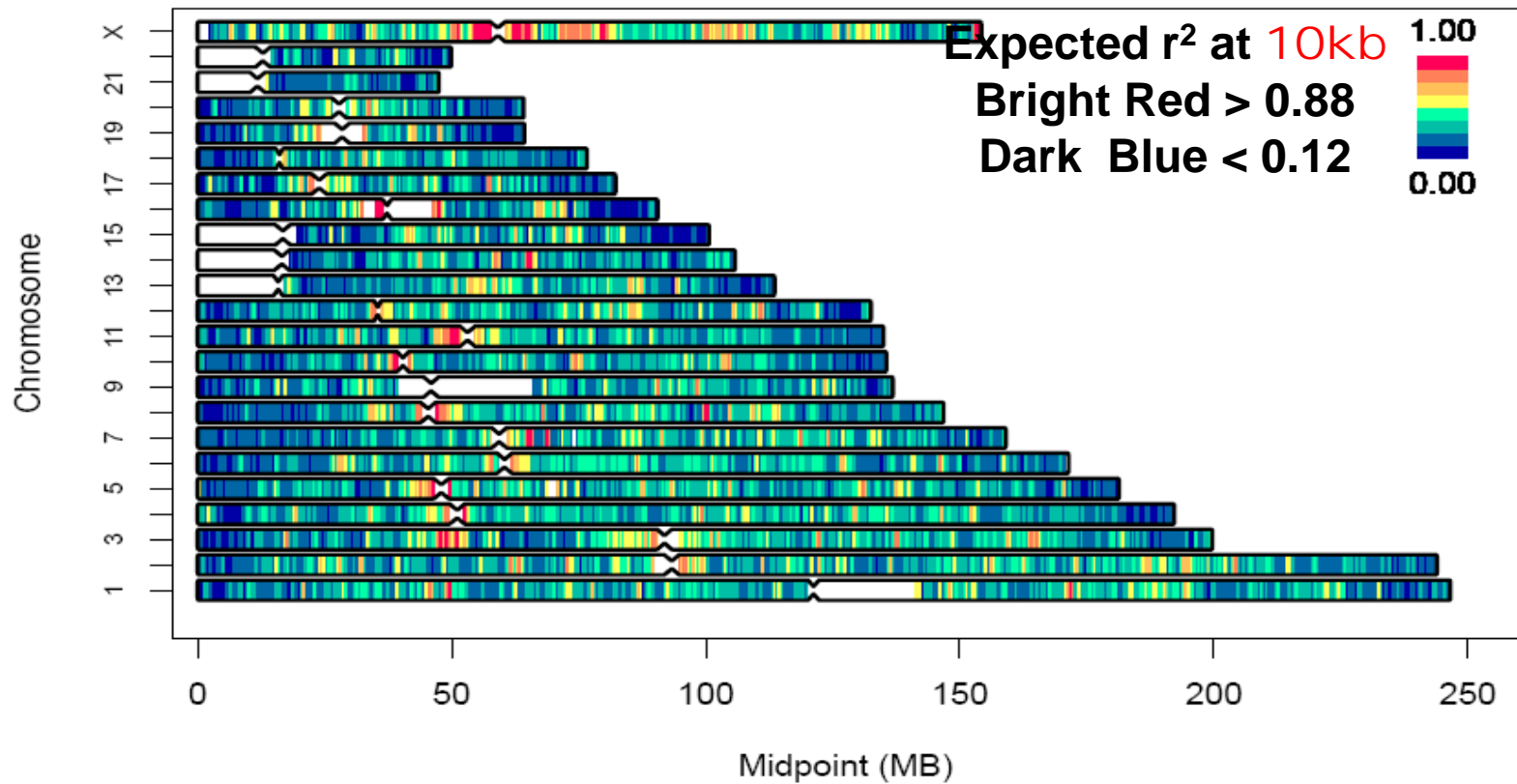
Smith et al, *Genome Research*, 2005

Genomic Variation in LD (JPT + CHB)



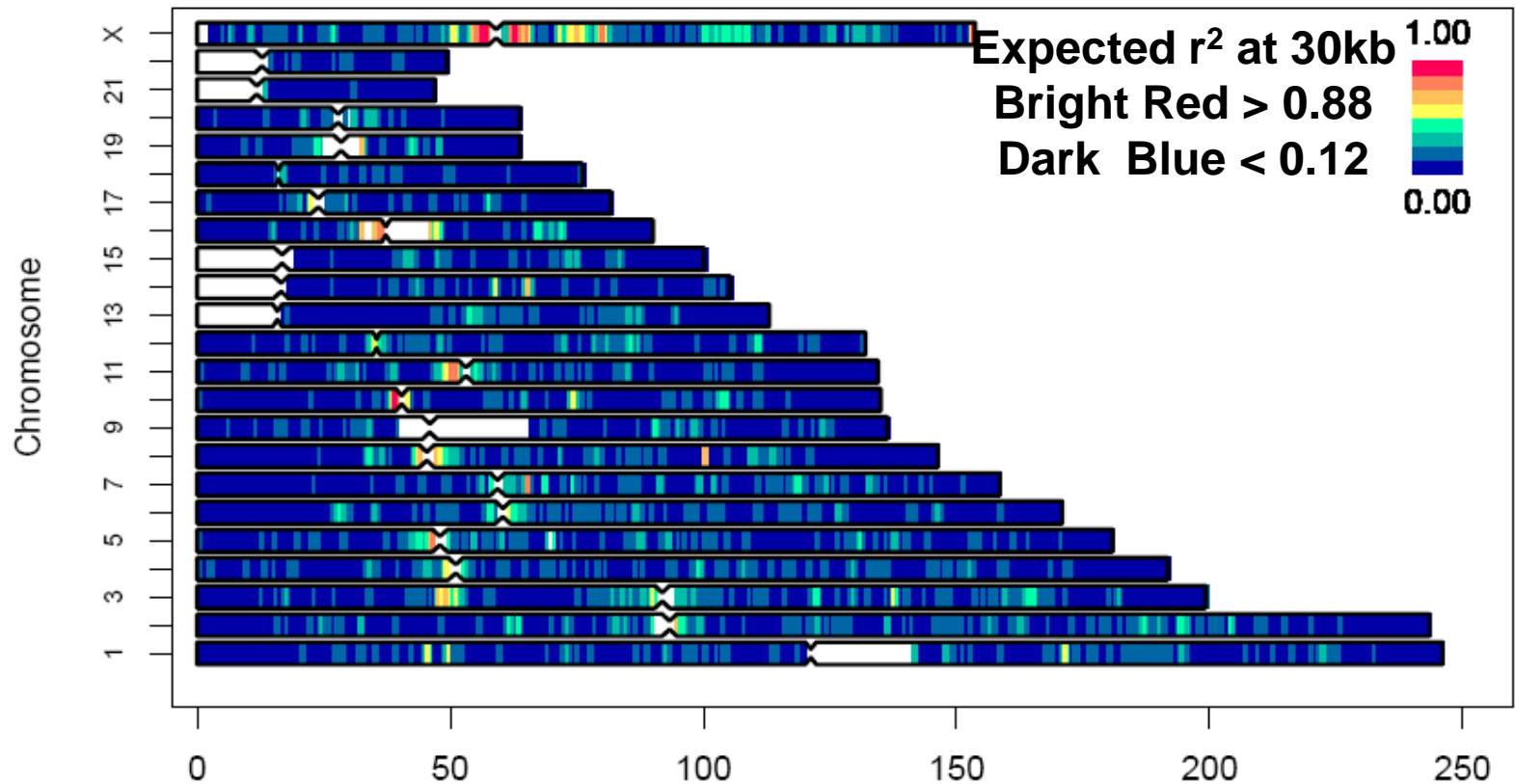
Smith et al, *Genome Research*, 2005

Genomic Variation in LD (YRI)



Smith et al, *Genome Research*, 2005

Genomic Distribution of LD (YRI)



Smith et al, *Genome Research*, 2005

Sequence Composition vs. LD

(some selected comparisons)

	Genome Quartiles, Defined Using LD					Trend
	Genome	(Low LD) Q1	Q2	Q3	(High LD) Q4	
Basic Sequence Features						
GC Bases (%)	40.8	43.5	41.0	39.6	39.0	Decreases With LD
Bases in CpG Islands (%)	0.7	0.9	0.7	0.6	0.6	Decreases With LD
Polymorphism (Π * 10,000)	10.1	11.9	10.6	9.6	8.3	Decreases With LD
Genes and Related Features						
Known Genes (per 1000 kb)	6.4	6.6	6.1	6.2	6.7	U shaped
Genic Bases (Exon, Intron, UTR, %)	38.5	37.6	34.6	36.0	45.9	U shaped
Coding Bases (%)	1.2	1.1	1.0	1.1	1.4	U shaped
Conserved Non-Coding Sequence (%)	1.4	1.6	1.5	1.3	1.1	Decreases with LD
Repeat Content						
Total Bases in Repeats (%)	47.9	44.2	46.4	48.6	52.3	Increases with LD
Bases in LINE repeats (%)	20.9	16.5	19.9	22.4	24.9	Increases with LD
Bases in SINE repeats (%)	13.6	14.7	13.1	12.6	14.0	U shaped

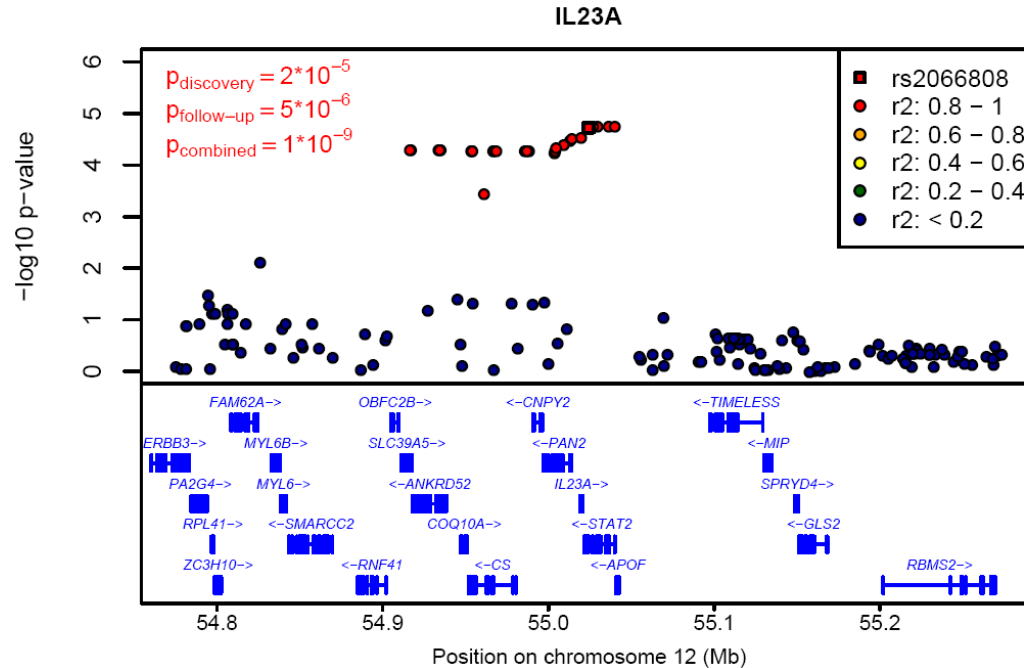
Gene Function in Regions of High and Low Disequilibrium

Gene Function (GO Term)	Annotated			χ^2	P-value
	Genes	Low LD	High LD		
All Swissprot Entries Examined	7520	2305	2045	-	-
DNA metabolism	366	74	139	35.37	<.0001
Immune Response	622	232	94	34.36	<.0001
Cell cycle	493	119	177	26.24	<.0001
Protein Metabolism	1193	318	375	23.33	<.0001
Organelle Organization and Biogenesis	444	107	152	19.65	<.0001
Intracellular Transport	263	56	95	19.61	<.0001
Organogenesis	805	294	162	16.49	0.00005
Cell Organization and Metabolism	545	138	178	16.43	0.00005
RNA Metabolism	208	41	71	15.33	0.00009

Results from a comparison of the distribution of 40 most common gene classifications in the GENE Ontology Database

Implications for Association Studies

Linkage Disequilibrium in Association Studies: Psoriasis and IL23A Example



Multiple nearby SNPs show evidence for association with psoriasis.
The SNPs are all in linkage disequilibrium, so it is hard to pinpoint causal SNP.

Nair et al, *Nature Genetics*, 2009

Linkage Disequilibrium in Association Studies: Tag SNP Picking

- Many nearby SNPs will typically provide similar evidence for association
- To decrease genotyping costs, most association studies will examine selected “tag SNPs” in each region
- The most common tagging strategy focuses on pairwise r^2 between SNPs

Linkage Disequilibrium in Association Studies: Pairwise Tagging Algorithm

- Select an r^2 threshold, typically 0.5 or 0.8
 - SNPs with r^2 above threshold can serve as proxies for each other
- For each marker being considered, count the number of SNPs with r^2 above threshold
- Genotype SNP with the largest number of pairwise “proxies”
- Remove SNP and all the SNPs it tags from consideration
- Repeat the previous three steps until there are no more SNPs to pick or genotyping budget is exhausted

Potential Number of tag SNPs

Table 3 | Number of tag SNPs required to capture common (MAF \geq 0.05) Phase II SNPs

Threshold	YRI	CEU	CHB+JPT
$r^2 \geq 0.5$	627,458	290,969	277,831
$r^2 \geq 0.8$	1,093,422	552,853	520,111
$r^2 = 1.0$	1,616,739	1,024,665	1,078,959

Current tag SNP panels typically examine 500,000 – 1,000,000 SNPs for a cost of \$200 - \$500 per sample.

It is anticipated that future panels will allow for as many as 5 million SNPs.

The International HapMap Consortium, *Nature*, 2007

Today ...

- Basic descriptors of linkage disequilibrium
- Learn when linkage disequilibrium is expected to hold (or not!)

Additional Reading I

- Dawson E et al (2002). A first-generation linkage disequilibrium map of human chromosome 22. *Nature* **418**:544-548
- The International HapMap Consortium. (2005). A haplotype map of the human genome. *Nature* 437:1299-320
- Carlson CS et al (2004). Selecting a maximally informative set of single-nucleotide polymorphisms for association analyses using linkage disequilibrium. *Am J Hum Genet* 74:106-120

Additional Reading II

- Cardon and Bell (2001) Association study designs for complex diseases. *Nature Reviews Genetics* 2:91-99
- Surveys important issues in analyzing population data.
- Illustrates shift from focus on linkage to association mapping for complex traits.