

SNPs and Haplotypes

Bioinformatics
NSF



SNP Concepts

- SNPs – what are they?
- Why are SNPs important?
- SNPs and linkage disequilibrium
- Common SNPs / haplotype blocks
- SNPs / Haplotype block – navigation
- Building complex traits and ‘the myth of race?’, the role of SNPs and haplotypes

2

SNPs - Introduction

- Single Nucleotide Polymorphism
- Occurs **once** in human evolution
- Estimate of 1 bp in 600 – 1000 bp
- Occur mostly in introns (2/3)
 - Regulatory regions, leading to cancer
- Often lethal when in exons (1/3)
 - Leading to a fatal amino acid substitution

3

What are Single Nucleotide Polymorphisms (SNPs)?

```
ATGGTAAGCCTGAGCTGACTTAGCGT-AT
ATGGTAAACCTGAGTTGACTTAGCGTCAT
      ↑       ↑           ↑
      SNP    SNP         indel
```

SNPs result from replication errors and DNA damage
They are a ‘polymorphic’ bit state at a nucleoside address

4

What (exactly) is a SNP ?

- A **SNP** is defined as a **single base change** in a DNA sequence that occurs in a significant proportion (more than 1 percent) of a large population.
- Occurs exactly **once** in human evolution
- Degenerate ‘**bit state**’ at a genomic address (A, T, C, or G) usually n=2
- Degenerate ‘bit state’ = **polymorphism**

5

SNPs / Polymorphisms

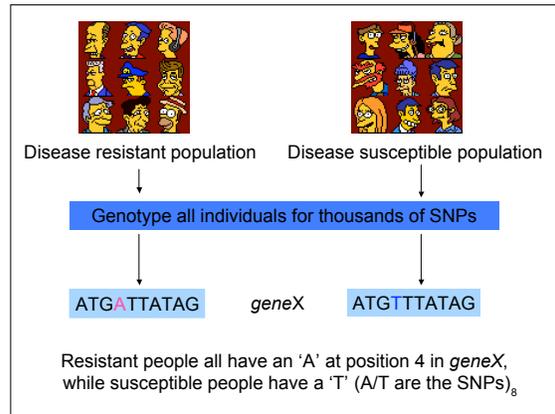
- A Single Nucleotide Polymorphism is a **source of variance** in a genome. A SNP (“snip”) is a single base mutation in DNA.
- SNPs are ‘**conserved**’ across the genome, often in patterns called ‘**haplotype blocks**’
- SNPs are the most simple form and most common source of **genetic polymorphism** in the human genome (90% of all human DNA **polymorphisms** are associated with SNPs).

6

Why are 'SNP' Polymorphisms Useful?

- It's sometimes possible to **correlate a SNP** with a particular trait or disease
 - This is known as **association genetics**
- **Susceptibility to disease** may also be described as an '**unfortunate trait**'
- **Traits** are also '**larger**' than genes
- SNPs in (regulatory) intronic regions may be as important as (coding) exons

7



SNP Applications in Medicine

- Gene discovery and allele mapping
- Association-based (drug) candidate
 - polymorphism testing of a trait pool
- Diagnostics / risk profiling
- Drug response prediction
- Homogeneity testing / study design
- Gene function identification

9

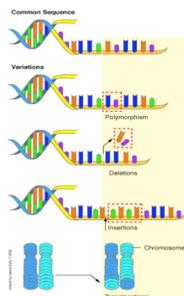
Genetic Polymorphism

- **Genetic Polymorphism:** A difference in DNA sequence among individuals, groups, or populations. **One or more SNPs.**
- **Genetic Mutation:** A change in the nucleotide sequence of a DNA molecule. Genetic mutations are one type of genetic polymorphism (but less than 1%).
- **Polymorphism is common, mutation rare**
 - Polymorphism is the 'stuff of variation'

10

Variations in Genomes

- SNP variations
- Polymorphisms
- Deletions
- Insertions
- Translocations



11

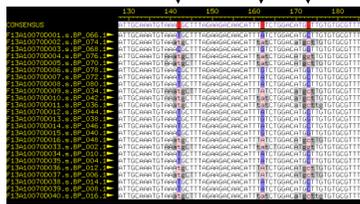
Single Nucleotide Polymorphisms (SNPs), Haplotypes, Linkage Disequilibrium, and the Human Genome

Coding Region SNPs

- Types of coding region SNPs:
 - **Synonymous:** the substitution causes no amino acid change to the protein it produces. This is also called a **silent** mutation.
 - **Non-Synonymous:** the substitution results in an alteration of the encoded amino acid. A **missense** mutation changes the protein by causing a change of codon. A **nonsense** mutation results in a misplaced termination.
 - One half of all **coding sequence SNPs** result in **non-synonymous codon changes.** (but half do)

12

Sequence-Based Detection and Genotyping of SNPs



Jim Sloan, Tushar Bhangle (PolyPhred)
 Matthew Stephens, Paul Scheet (Quality Scores for SNPs)
 Phil Green, Brent Ewing, David Gordon (Phred, Phrap, Consed)

SNP Discovery and Genotyping Workshop

19

SNPs and Variation

- In human beings, 99.9 percent of bases are same
- Remaining 0.1 percent makes a person unique.
 - Different attributes / characteristics / traits
 - how a person looks,
 - diseases he or she develops.
- These variations can be:
 - Harmless (change in phenotype)
 - Harmful (diabetes, cancer, heart disease, Huntington's disease, and hemophilia)
 - Latent (variations found in coding and regulatory regions, are not harmful on their own, and the change in each gene only becomes apparent under certain conditions e.g. susceptibility to lung cancer)

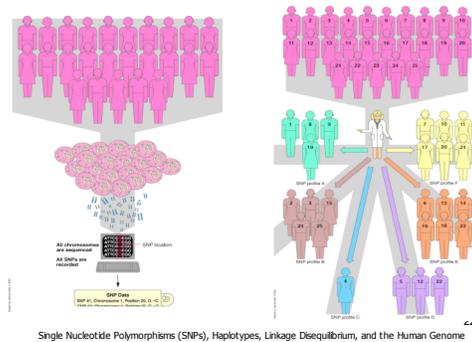
20

Why Create SNP Profiles?

- Genome of each individual contains a *distinct SNP pattern* (haplotype block).
- People can be grouped based on their SNP profiles (*association studies*).
- SNP profiles may be important for identifying response to drug therapy.
- Correlations might emerge between certain SNP profiles and specific responses to treatment (good and bad).

21

Populations Based on SNPs

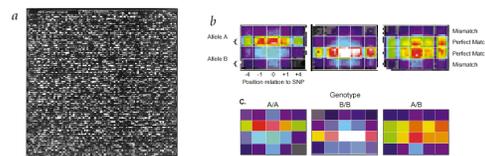


Techniques to Detect Known Polymorphisms

- Hybridization Techniques
 - Microarrays
 - Real time PCR
 - HTS SNP arrays
- Enzyme based Techniques
 - Nucleotide extension
 - Cleavage
 - Ligation
 - Reaction product detection and display
- [Comparison of SNP Assay Techniques](#)

23

SNP Genotyping



SNPs can be measured using High Throughput Screening with custom (HTS) microarray technology (Applied Biosystems)

SNP Discovery and Genotyping Workshop

24

Techniques to Detect Unknown Polymorphisms

- Direct Sequencing
 - HTS SNP sequencing
- SNP Microarray
 - Rapid SNP genotyping
- Cleavage / Ligation
- Electrophoretic mobility assays
- [Comparison of Techniques](#) used to detect unknown (SNP) polymorphisms

25

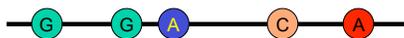
More SNP Terminology

- Polymorphism / Haplotype
 - Correlation of characters states among polymorphic sites (across the genome). **SNP patterns = 'blocks'**
 - Insufficient passage of time to randomize character states by meiotic recombination – patterns conserved
- Haplotypes are **'blocks'** of associated SNPs
- Haplotypes may be *'too recent'*...
 - Not enough time for recombination to merge SNPs
 - Or SNP recombination may be a more difficult process

28

SNPs and Haplotypes

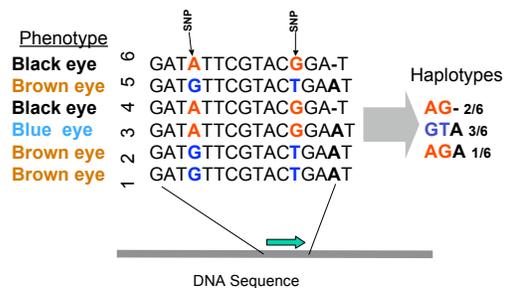
- **SNP**: Single Nucleotide Polymorphism
- **Haplotype**: A set of closely linked genetic markers present on one chromosome which tend to be inherited together (not easily separable by recombination).



Set of SNP polymorphisms: a SNP haplotype

Seoul National University School of Public Health³⁰

From SNP to Haplotype



Seoul National University School of Public Health³⁰

SNPs and Linkage Disequilibrium (LD)

- Recent Resurgence of LD Study Motivated by SNP Markers / haplotypes (HapMap project)
- High density: ~ one SNP in every 600 bp in the human genome make them easy to map.
- Simple SNPs: Biallelic (occur on both alleles)
- Common: ~93% are found globally (among human populations); ~7% are restricted to local populations. (NHGRI, 2001)

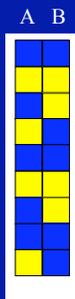
31

SNPs / Linkage Disequilibrium

- Generally speaking...
- SNPs should be inherited *independently*
 - Following Mendelian inheritance
- Many SNPs appear to be **co-inherited**
 - Creating 'hot spots' in the human genome
- Blocks of SNPs – 'SNP haplotypes'
 - We don't why or how, **but they exist**

32

Linkage Disequilibrium



Haplotype is the pattern of alleles on a single chromosome
– 4 possible haplotypes

Linkage Disequilibrium (LD) describes the allelic association between two SNPs

Two popular LD statistics:

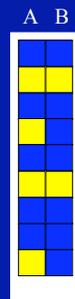
$$D'$$

$$r^2$$

33

Identifying Haplotypes for Genotype-Phenotype Analysis Dana C. Crawford dcrawfo@gs.washington.edu

Complete LD



Unequal allele frequency
Allelic association is as strong as possible

- 3 haplotypes observed
- No detected recombination between SNPs
- Genotype is not perfectly correlated

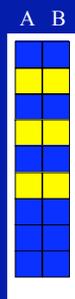
$$D' = 1$$

$$r^2 < 1$$

34

Identifying Haplotypes for Genotype-Phenotype Analysis Dana C. Crawford dcrawfo@gs.washington.edu

Perfect LD



Equal allele frequency
Allelic association is as strong as possible

- 2 haplotypes observed
- No detected recombination between SNPs
- Genotype is perfectly correlated

$$D' = 1$$

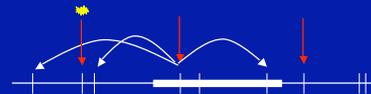
$$r^2 = 1$$

35

Identifying Haplotypes for Genotype-Phenotype Analysis Dana C. Crawford dcrawfo@gs.washington.edu

Rational SNP Selection

- Select SNPs to genotype on the basis of LD
- Some SNPs are in LD with many other SNPs
- Some SNPs are in LD with no other SNPs
- SNPs between a pair of associated SNPs are not necessarily associated with the flanking SNPs



36

Identifying Haplotypes for Genotype-Phenotype Analysis Dana C. Crawford dcrawfo@gs.washington.edu

SNP HapMap Project

<http://www.HapMap.org/>

- Sequence genomes of a large number of people ($n > 10,000$ ethnically diverse)
- Compare the base sequences to discover SNPs, their location, and their frequency.
- Generate a single map of the human genome containing all possible SNPs => SNP maps

43



Perlegen Haplotype Study

"Global patterns of human DNA sequence variation (haplotypes) defined by common single nucleotide polymorphisms (SNPs) have important implications for identifying disease associations and human traits. We have used high-density oligonucleotide arrays, in combination with somatic cell genetics, to identify a large fraction of all common human chromosome 21 SNPs and to directly observe the haplotype structure defined by these SNPs. This structure reveals blocks of limited haplotype diversity in which more than 80% of a global human sample can typically be characterized by only three common haplotypes."

Abstract from seminal article by Cox *Science* 294: 1719-1723 (2001)

45

Haplotypes in the Genome

- Defined as a pattern of SNPs that appears as an 'associated block' on one or more chromosomes
- There are (estimated*) to be roughly 200K to 300K haplotypes in genome
- Of these, most can be defined by the identity of 3 SNPs, and always < 10

* Cox et.al. Perlegen Haplotyping of chromosome 21

46

Golden Path / Human Variation

- 3 billion base pairs
- 6 to 10 million SNPs
- 200K – 300K interrelated haplotype 'blocks' in the human population
- Each block is about 7.8 K bp
– (this is an average)
- Each block contains roughly 10 SNPs
– (of which 1 to 3 define the haplotype)

47

Haplotype Questions

- David Cox – traits vs. genes?
- Disequilibrium – patterns of inheritance?
- 'Hot spots' – defined by haplotypes?
- Are SNPs more definable than 'race'?
 - One race, 200,000 'haplotype / traits'?
 - One race, common haplotype patterns?
- How and why do haplotypes occur?
 - Is there a benefit to SNP patterns?

48

SNPs per Haplotype Block

- In common haplotypes (>80%) 3 SNPs can determine the identity of the block, and often just 1 SNP determines the haplotype identity (3 SNPs / haplotype)
- At most, only 10 SNPs (10% of roughly 100 SNPs) will define the 'identity' for the majority of observed haplotypes
- These SNPs are often called 'markers'

49

Identification

- In half of all haplotype blocks, 3 or fewer SNPs are needed for an identification
- In haplotype blocks with 3 SNPs, only one SNP is needed to identify a block
- In other cases, 10 SNPs define a block
- In all cases, 10% SNP identification is enough to create a haplotype definition

50

Haplotyping Tools (Now/Later)

- High density nucleotide arrays
 - Rapid SNP genotyping
- High resolution maps of chromosome
 - Data mining tools for pattern recognition
- Key word / ontologies for disease, traits, other definitions of human variability
- Eventually haplotype expression data

51

Haplotype Utility

- Are SNP patterns as important as SNPs in characterizing disease susceptibility?
- Are haplotype patterns a better way to define an individual's (UID) genome?
- Is this a better tool to understand the 'evolution of race', or 'myth of race'?
- Is this a method to identify proteomes?
 - Proteome = (UID) genome * expression?

52

Haplotype Exercise Overview

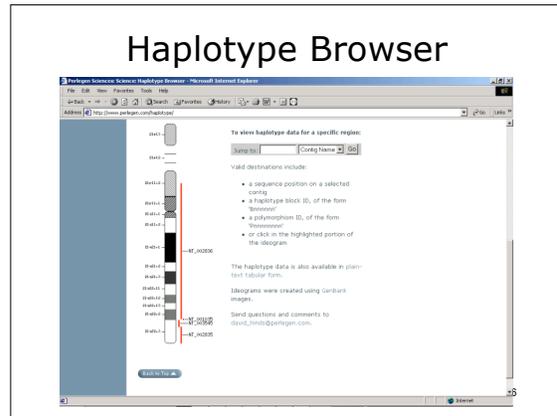
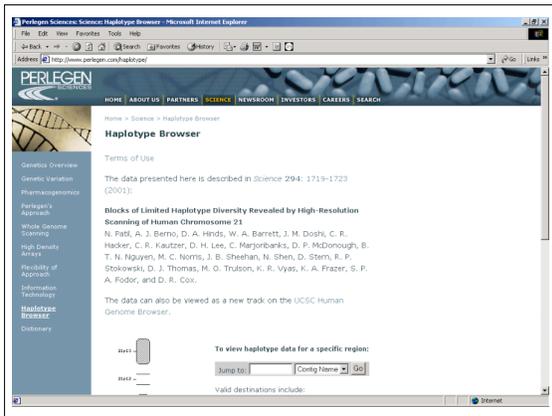
- Perlegen haplotype browser
- Database of haplotypes
- Navigate along a chromosome
- Pull up a haplotype block
- Examine blocks and SNPs
- Follow the SNP links into NCBI
- Perlegen genotype browser

53

Haplotype Exercise (I)

- Go to Perlegen haplotype browser
- Download the paper by David R Cox
- [Science 294: 1719-1723 \(2001\)](#)
- Scan the paper, which is a tough to read (technically), and use the next slides as a review of the high points
- The image of haplotype blocks is key
 - (This was shown earlier in the presentation)

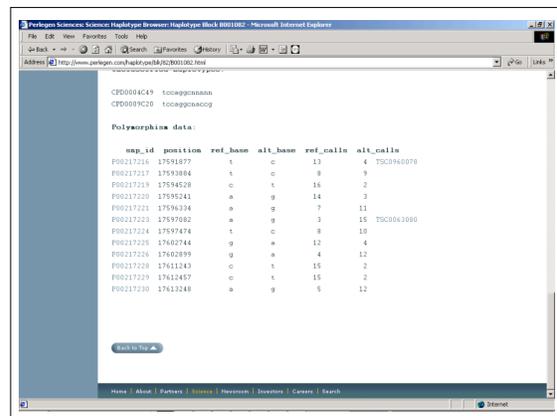
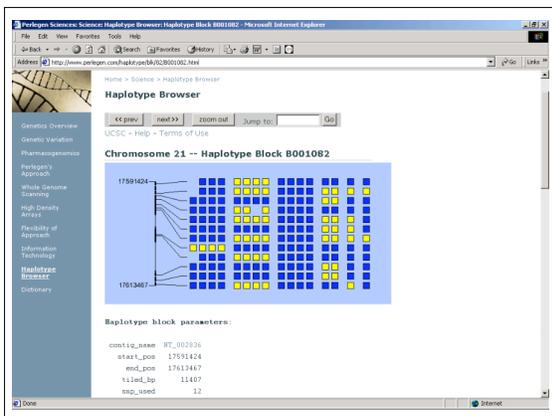
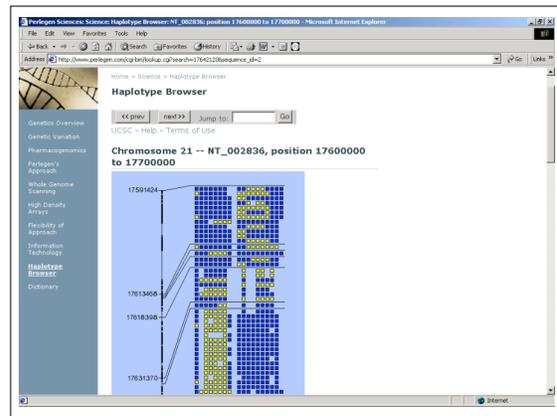
54

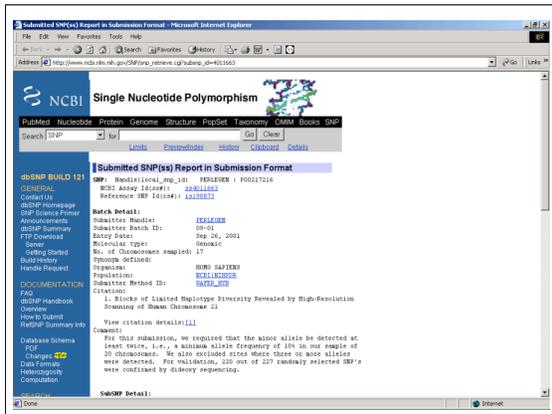


Haplotype Exercise (II)

- Click in the middle of 21q21+2 (or anywhere along the chromosome)
- That will bring up a haplotype block
- Scroll down that page, and look at the entries for each of the rows in the block
- At the bottom of the page, look at the base calls for each (block row) entry
- Follow the bottom links to NCBI dbSNP

57





Perlegen Cox Paper (I)

- 24 separate ethnically diverse individuals
- 35,989 SNPs were identified on chromosome 21 (32,397,439 bp)
- 47% of 53,000 common SNPs occur with an allele frequency of 10%
 - 24,000 SNPs are highly reproducible from human to human - polymorphism

62

Perlegen Cox Paper (II)

- Used the set of 24,000 SNPs (with a minor allele) to define a haplotype
- Blocks define haplotype patterns
- The four most common haplotypes account for 80% of all occurrences
 - These 'universal haplotypes' defined by an address in the genome with SNP variation
 - Their distribution follows a 'normal curve'

63

Perlegen Cox Paper (III)

- Alleles making up blocks of such SNPs are often correlated – resulting in reduced genetic variability and defining a limited number of 'SNP haplotypes'.
- Leads to 'linkage disequilibrium'
- 80% of the haplotype structure is defined by < 10% of the SNPs in a haplotype block. This is universal.

64

Perlegen Cox Paper (IV)

- Most common haplotype pattern (of the four universal haplotypes) is found in:
 - 1st – 50% of all individuals
 - 2nd – 25% of all individuals
 - 3rd – 12.5% of all individuals
 - 4th – ~ 6% of all individuals (estimate)
- *These total 93% of all human (genomes)*

65

Perlegen Cox Paper (V)

- Average size of a block is 7.8 Kbytes
- Some blocks are up to 114 SNPs and 115 Kbytes (1 SNP per 1000 bp)
- 10% of 114 SNPs or only 11 SNPs are needed to define the largest haplotypes
- Average distance among all SNPs was 900 bp, and among 24,000 common SNPs, was 1,300 bp

66

Perlegen Cox Paper (VI)

- On average, there are 2.7 common haplotypes per block, defined as haplotypes observed on multiple chromosomes
- 94% of blocks contained fewer than 3 SNPs
- Exonic bases:
 - > 10 SNPs
 - 3 to 10 SNPs
 - < 3 SNPs

67

Perlegen Cox Paper (VII)

- Haplotype blocks are defined by their genetic information, and not on knowledge of how this information originated, or why it exists.
- Perlegen (David Cox) suggests that perhaps 'human traits' associated with haplotypes are a better way to approach genetic variation in the human genome

68

The Myth of Race?

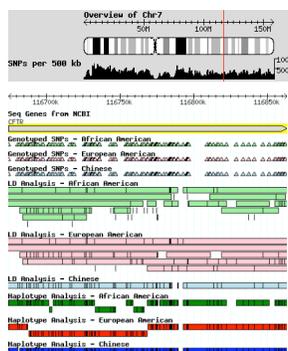
- 93% of SNPs are 'universal' - 'public SNPs'
- 80% of SNP variation defined by four or less haplotypes within a haplotype block
- Haplotype patterns and 'multigenic', meaning that blocks can span genes (CDS and introns)
- Traits may be multigenic and multi-block.
- We are more complex as individuals than we are differentiated by race, but 'race' can be studied
- <http://genome.perlegen.com/browser/index.html>

69

Perlegen Genome Browser

- Database of haplotypes / SNPs and LD
- Data separated by 'race'
 - European
 - African American
 - Han Chinese
- Navigate along a chromosome
- Examine **SNP frequencies** by 'race'
 - Study **alleles** and **population distributions**
- <http://genome.perlegen.com/browser/>

70



<http://genome.perlegen.com/browser/>

Whole-Genome Patterns of Common DNA Variation in Three Human Populations

David A. Hinds,¹ Laura L. Stuve,¹ Geoffrey B. Nilsen,¹ Eran Halperin,² Eleazar Eskin,³ Dennis G. Ballinger,¹ Kelly A. Frazer,¹ David R. Cox^{1*}

Individual differences in DNA sequence are the genetic basis of human variability. We have characterized whole-genome patterns of common human DNA variation by genotyping 1,586,383 single-nucleotide polymorphisms (SNPs) in 71 Americans of European, African, and Asian ancestry. Our results indicate that these SNPs capture most common genetic variation as a result of linkage disequilibrium, the correlation among common SNP alleles. We observe a strong correlation between extended regions of linkage disequilibrium and functional genomic elements. Our data provide a tool for exploring many questions that remain regarding the causal role of common human DNA variation in complex human traits and for investigating the nature of genetic variation within and between human populations.

www.sciencemag.org SCIENCE VOL 307 18 FEBRUARY 2005

73

Are We Really That Different?

Table 1. SNPs segregating in the three genotyped populations. Percentages are of 1,586,383 genotyped SNPs or of 291,012 private SNPs.

Population	Segregating		MAF > 0.05		MAF > 0.10	
	SNPs	%	SNPs	%	SNPs	%
All SNPs						
African-American	1,483,594	93.5	1,267,594	79.9	1,083,652	68.3
European-American	1,286,277	81.1	1,123,765	70.8	991,046	62.5
Han Chinese	1,168,029	73.6	1,027,109	64.7	910,451	57.4
Private SNPs						
African-American	218,500	75.1	139,536	47.9	88,525	30.4
European-American	44,555	15.3	18,284	6.3	8,062	2.8
Han Chinese	27,957	9.6	15,946	5.5	9,817	3.4

www.sciencemag.org SCIENCE VOL 307 18 FEBRUARY 2005

74

Table 3. Distribution of genic, synonymous, and nonsynonymous coding SNPs spanned by bins of extended LD in any of the three population samples. Genic SNPs are defined as within 10 kb of a protein-coding gene annotation.

Longest spanning LD bin (kb)	SNPs	Genic		Synonymous		Nonsynonymous	
		SNPs	%	SNPs	%	SNPs	%
<500	1,336,094	707,950	46.1	10,330	0.67	8,898	0.58
500 to 1000	42,492	22,189	52.3	347	0.82	302	0.71
>1000	7857	4,955	63.1	120	1.52	171	2.17

Table 4. LD statistics for common SNPs genotyped in this study, with common variants identified by complete resequencing in 152 genes.

Subset*	Yield (%)	r ² ₁	r ² > 0.5 (%)	r ² > 0.8 (%)	r ² = 1.0 (%)
African-American					
All	23.3	0.715	70.9	53.7	41.5
Tag	12.3	0.698	70.1	51.9	33.2
European-American					
All	25.0	0.841	86.5	72.6	62.4
Tag	8.1	0.810	85.6	69.7	44.9

*SNPs from the current study within all common SNPs or minimal tagging subset. †Percentage of all Seattle SNP PGAs variants that were in the selected set. ‡Across all PGA variants, the mean maximum r² with a selected SNP in the same locus. §Percentages of PGA variants having an r² greater than the specified threshold with any selected SNP in the same locus.

www.sciencemag.org SCIENCE VOL 307 18 FEBRUARY 2005

75

Perlegen Next Steps

- Identify SNP patterns in thousands
- Map SNPs / haplotypes to disease
- Map SNPs / haplotypes to human 'traits'
- Association maps for haplotypes
 - How do blocks interrelate?
- Create a 'Hap Map' for human genome
- Ability to survey the entire genome will dramatically increase the power of genetic association analysis.

79

Conclusions

- SNPs (single nucleotide polymorphisms) are abundant and useful genetic markers
 - *Disease, drug resistance or adverse effects, etc*
- Linkage disequilibrium describes the tendency of many SNPs to be inherited in patterns or blocks
 - *Typical blocks include 7 to 8 common SNPs / 7.8 Kb*
- These haplotype patterns, or blocks, appear in >90% of the population, in over 250,000 identified blocks
 - *Often 3 or fewer SNPs can define a haplotype block*
- Association studies of individuals and populations that show disease or drug response behaviour in progress
 - *Where we'll see the benefit of haplotype research!*

80

Further Reading / Resources

- Applied Biosystems <http://www.appliedbiosystems.com/>
- Perlegen - <http://www.perlegen.com/>
- SNP Consortium - <http://snp.cshl.org/>
- [International SNP Working Group Data \(Nature\)](#)
- HapMap Project – <http://www.hapmap.org/>
- Google 'haplotypes, SNPs, haplotype maps, SNP genotyping, linkage disequilibrium', HapMap, and all (word / keyword) combinations of the above

81

Presentation References

- Introduction to SNPs: Discovery of Markers of Disease <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1152833/>
- SNP mapping using Genome-wide Unique Sequences <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC1152833/>
- The Structure of Haplotypes Blocks in Human Genome
- Using Haplotype blocks to map human complex trait loci
- High Resolution haplotype structure in human genome
- Detection of regulatory variation in mouse genes <http://linkage.rockefeller.edu/will01.html>
- <http://statewww.epfl.ch/revision/teaching/Microarrays/snp.pdf>
- <http://www.genome.gov/10001885>
- Resolution of Haplotypes and Haplotype Frequencies from SNP Genotype of Pooled Samples http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed&list_uids=11752031
- http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed&list_uids=11752031
- <http://www.genome.gov/10001885>
- <http://www.genome.gov/10001885>
- Visualizing Haplotype Blocks <http://www.media.mit.edu/people/fr/haplotype/>

82

Presentation Acknowledgements

1. Single Nucleotide Polymorphisms (SNPs), Haplotypes, Linkage Disequilibrium, and the Human Genome Manish Anand, Nihar Sheth Jim Costello
2. High-resolution haplotype structure in the human genome Mark J. Daly, John D. Rioux, Stephen F. Schaffner, Thomas J. Hudson & Eric S. Lander
3. Blocks of Limited Haplotype Diversity Revealed by High-Resolution Scanning of Human Chromosome 21 David R. Cox* et. Al.
4. Selecting SNPs for Genotype-Phenotype Analysis Using Allelic Association Linkage Disequilibrium Christopher Carlson csc47@u.washington.edu

83

Copyright Status

- Slides used in this presentation include material from several presentations downloaded from **Programs for Genomic Applications Educational Activities**. This material is used consistent with the educational license.
- **Web citations / links are used whenever possible (see the presentation acknowledgements slide).**
- Screen grabs from the Perlegen [website](#) are included in the tutorial section on the Perlegen haplotype browser.
- Inquiries for **educational** use of this material should be forwarded to me at rdcormia@earthlink.net Thank you!