

## Lecture 29 - Polymorphisms in Human DNA Sequences

- SNPs
- SSRs

### Eukaryotic Genes and Genomes

**genome** = DNA content of a complete haploid set of chromosomes  
 = DNA content of a gamete (sperm or egg)

Species	Chromosomes	cM	DNA content/ haploid(Mb)	year sequence completed	genes/ haploid
<i>E. coli</i>	1	N/A	5	1997	4,200
<i>S. cerevisiae</i>	16	4000	12	1997	5,800
<i>C. elegans</i>	6	300	100	1998	19,000
<i>D. melanogaster</i>	4	280	180	2000	14,000
<i>M. musculus</i>	20	1700	3000	2002 draft 2005 finished?	30,000?
<i>H. sapiens</i>	23	3300	3000	2001 draft 2003 finished	30,000?

Note: cM = centi Morgan = 1% recombination  
 Mb = megabase = 1 million base-pairs of DNA  
 Kb = kilobase = 1 thousand base-pairs of DNA

Species	cM	DNA content/ haploid (Mb)	generation time	design crosses?	true breeding strains?
<i>E. coli</i>	N/A	5	30 min	yes	yes
<i>S. cerevisiae</i>	4000	12	90 min	yes	yes
<i>C. elegans</i>	300	100	4 d	yes	yes
<i>D. melanogaster</i>	280	180	2 wk	yes	yes
<i>M. musculus</i>	1700	3000	3 mo	yes	yes
<i>H. sapiens</i>	3300	3000	20 yr	no	no

- Human genetics is retrospective (vs prospective). Human geneticists cannot test hypotheses prospectively. The mouse provides a prospective surrogate.

- Can't do selections

- Meager amounts of data Human geneticists typically rely upon statistical arguments as opposed to overwhelming amounts of data in drawing connections between genotype and phenotype.

- Highly dependent on DNA-based maps and DNA-based analysis

The unique advantages of human genetics:

- A large population which is self-screening to a considerable degree
- Phenotypic subtlety is not lost on the observer
- The self interest of our species

A locus is said to be polymorphic if two or more alleles are each present at a frequency of at least 1% in a population of animals.

1) SNPs = single nucleotide polymorphisms = single nucleotide substitutions

$H_{nuc}$  = average heterozygosity per nucleotide site = 0.001

In human populations:

		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys STOP Trp	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G
						3rd base in codon

### SYNONOMOUS CHANGES

TTT GCT GGC CAC

Phe Ala Gly His

TTT GCT GGA CAC

Phe Ala Gly His

### NON-SYNONOMOUS CHANGES

TTT GCT GGC CAC

Phe Ala Gly His

TTT GCT TGC CAC

Phe Ala Cys His

The great majority (probably 99%) of SNPs are selectively “neutral” changes of little or no functional consequence:

- outside coding or gene regulatory regions (>97% of human genome)
- silent substitutions in coding sequences
- some amino acid substitutions do not affect protein stability or function
- disadvantageous SNPs selected against --> further underrepresentation

A small minority of SNPs are of functional consequence and are selectively advantageous or disadvantageous.

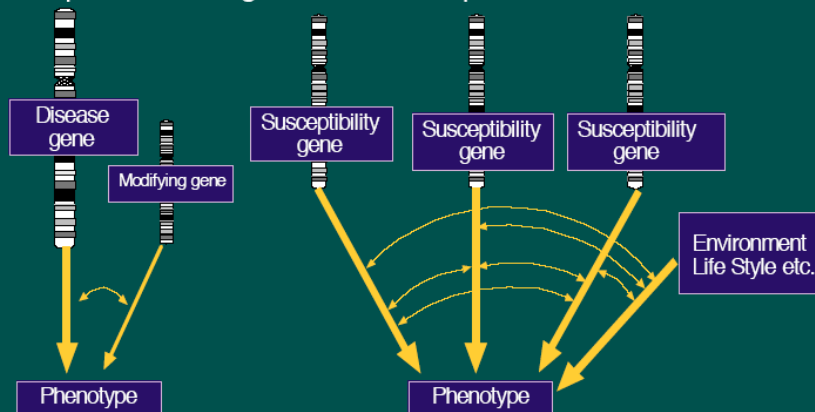
## Affymetrix chip

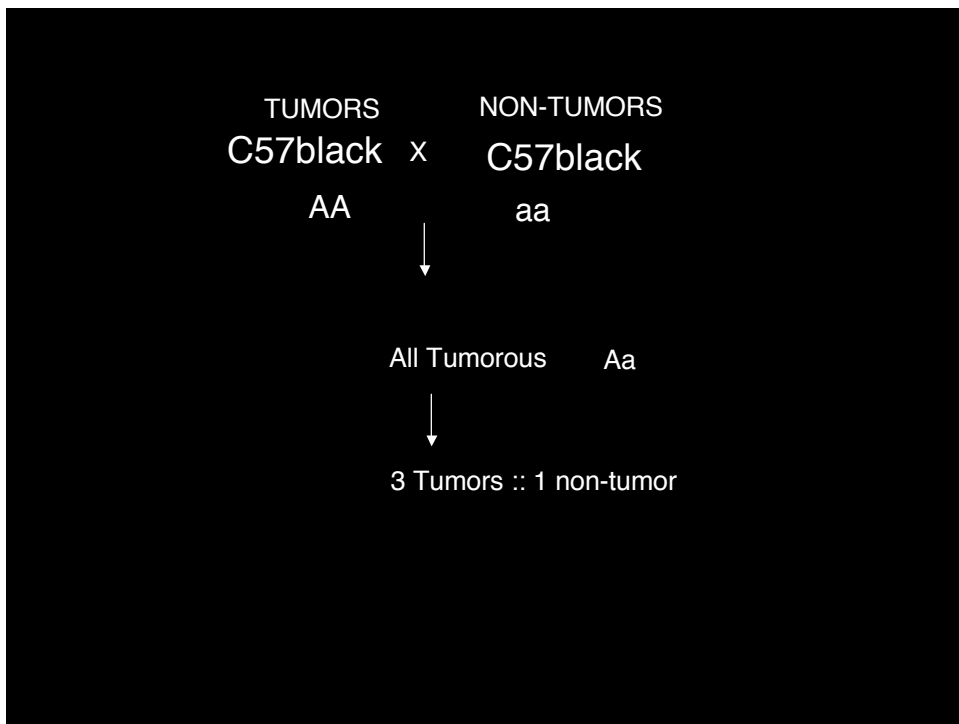
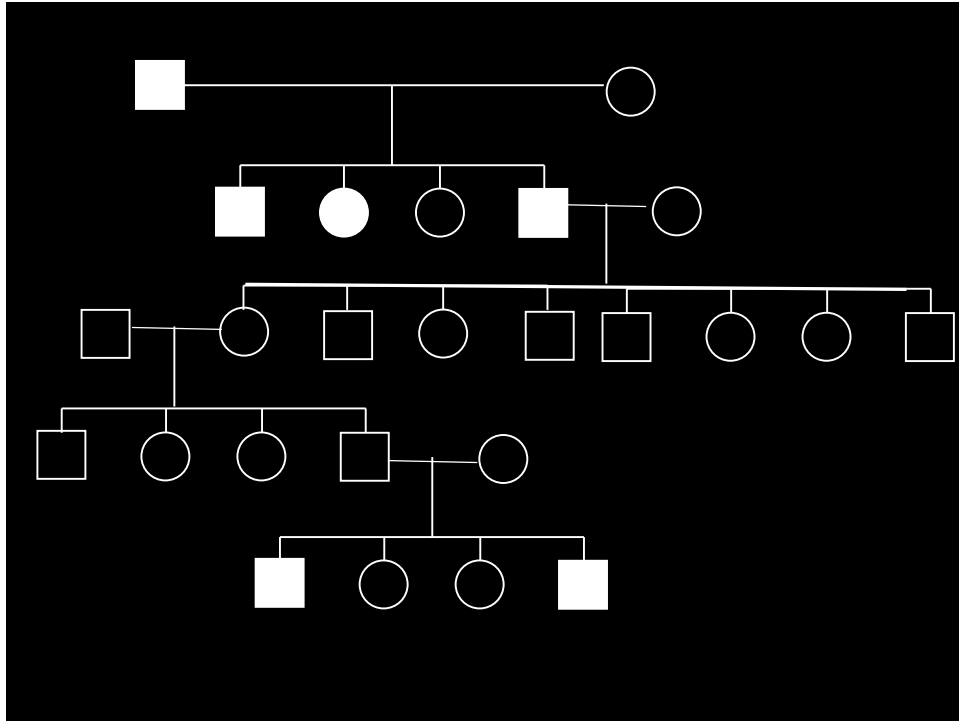


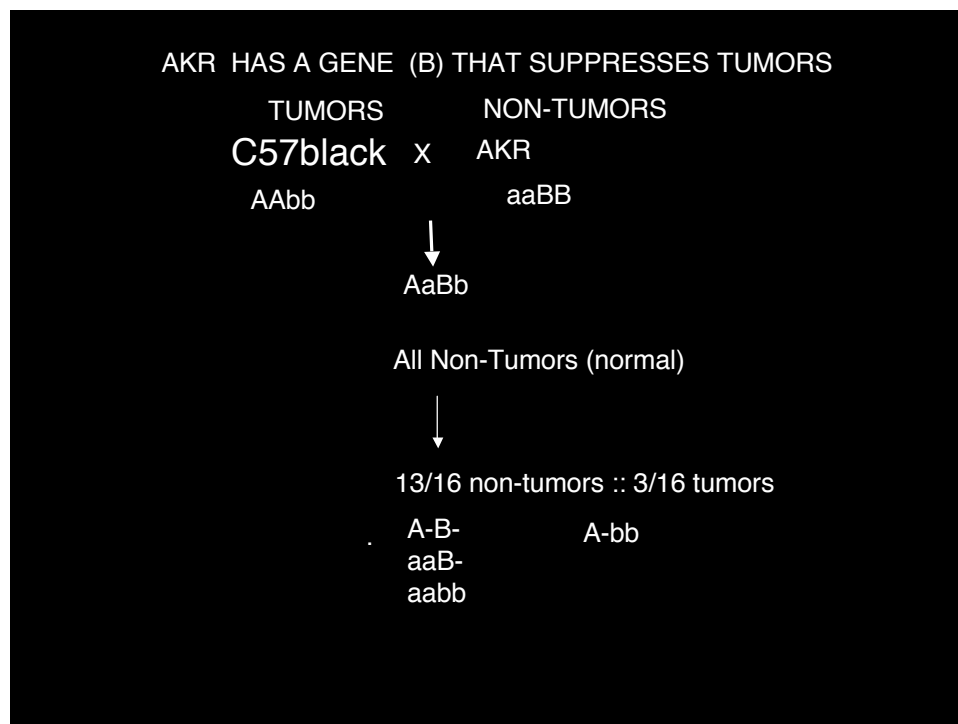
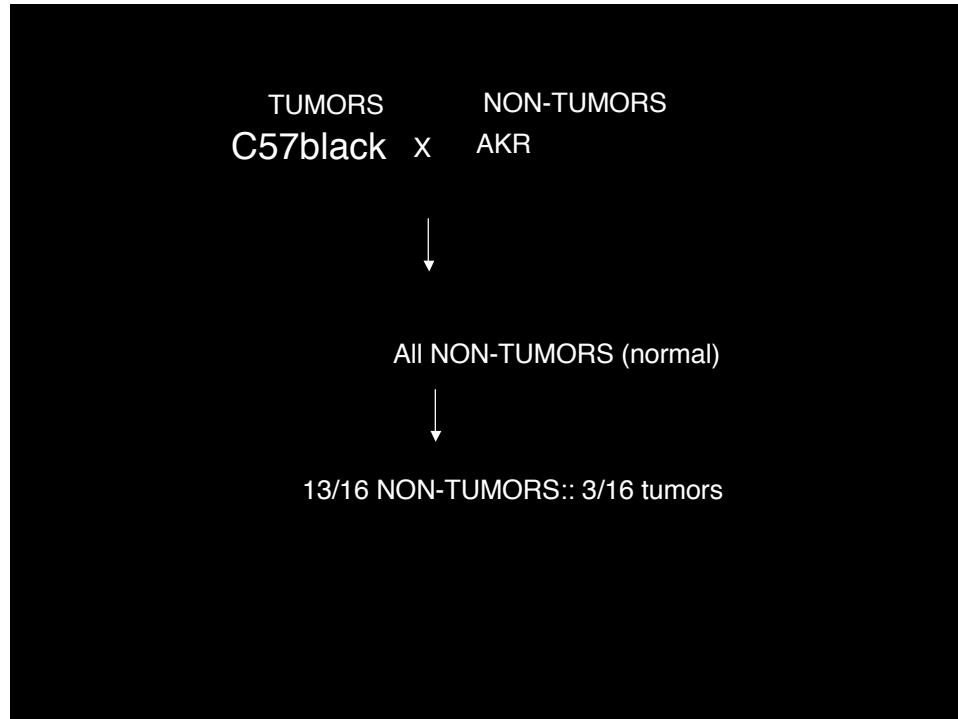
## Genetic Traits

Simplex or monogenic

Complex or multifactorial







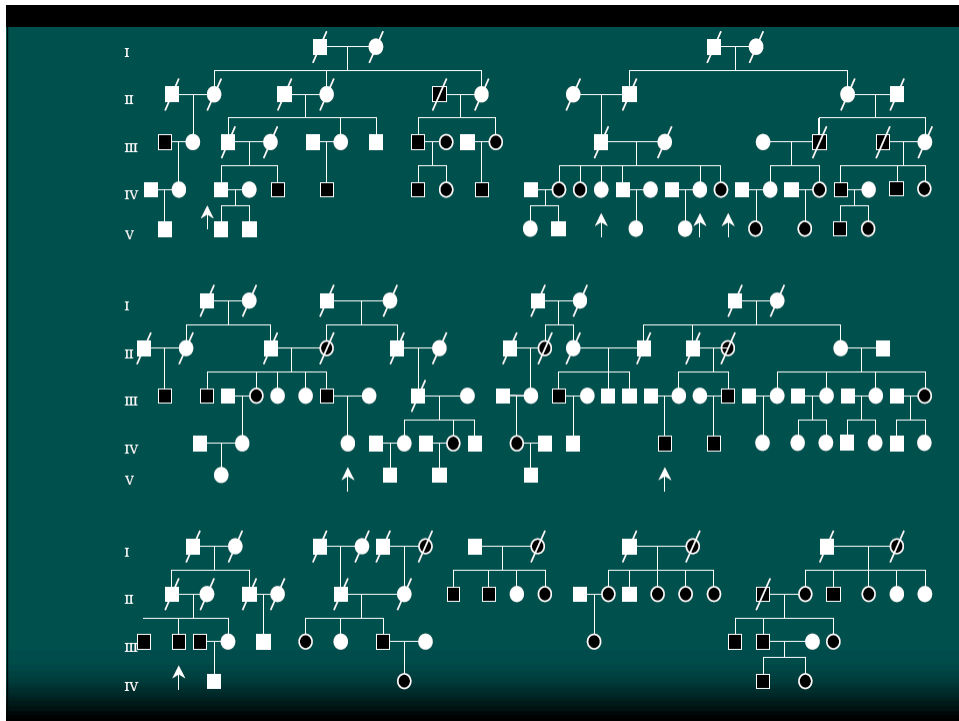
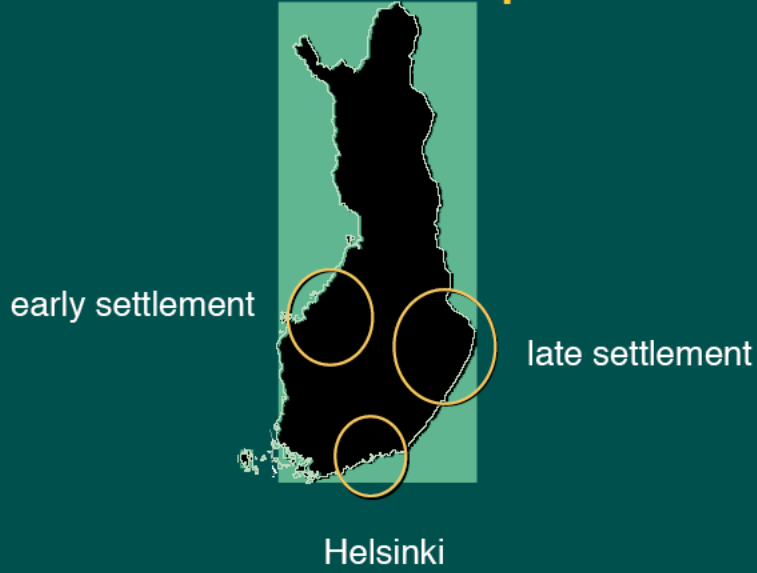
## Efforts to Simplify the Complex Genetic Background of Common Diseases

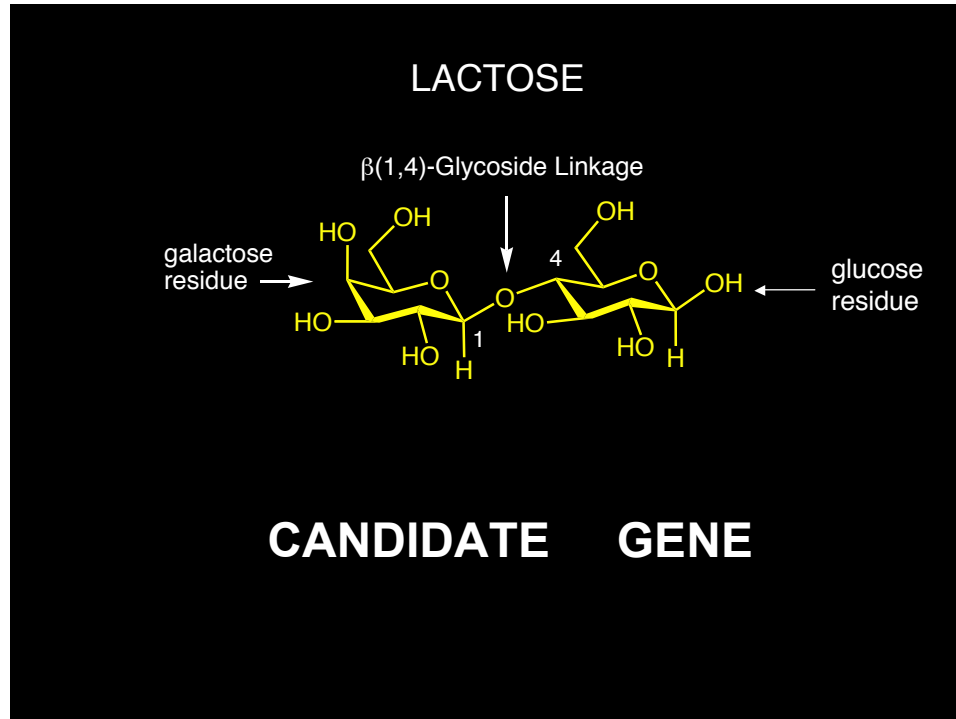
- Familial cases
- Population isolates
- Defined clinical phenotype
- Animal models

## The Effect of Population Bottlenecks to Disease Alleles

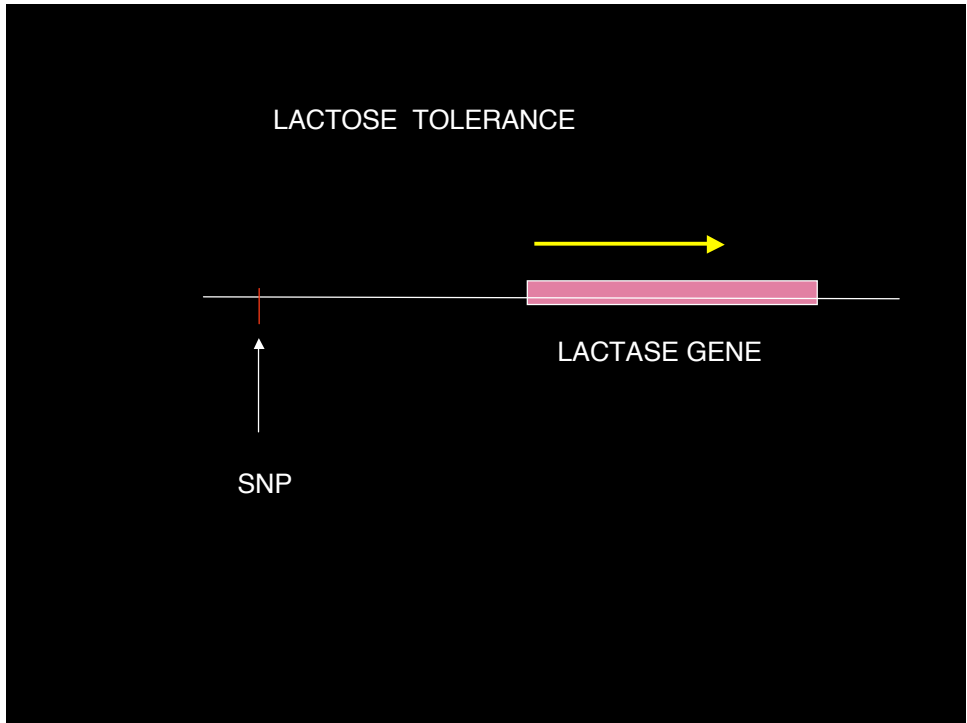


# Carrier Frequencies





The enzyme lactase that is located in the villus enterocytes of the small intestine is responsible for digestion of lactose in milk. Lactase activity is high and vital during infancy, but in most mammals, including most humans, lactase activity declines after the weaning phase. In other healthy humans, lactase activity persists at a high level throughout adult life, enabling them to digest lactose as adults. This dominantly inherited genetic trait is known as lactase persistence. The distribution of these different lactase phenotypes in human populations is highly variable and is controlled by a polymorphic element *cis*-acting to the lactase gene. A putative causal nucleotide change has been identified and occurs on the background of a very extended haplotype that is frequent in Northern Europeans, where lactase persistence is frequent. This single nucleotide polymorphism is located 14 kb upstream from the start of transcription of lactase in an intron of the adjacent gene *MCM6*. This change does not, however, explain all the variation in lactase expression.



2) **SSRs** = simple sequence repeat polymorphisms = "microsatellites"

Most common type in mammalian genomes is **CA repeat**:

The diagram shows a DNA strand with a (CA)<sub>n</sub> repeat region. Primer #1 is on the left, pointing right, and primer #2 is on the right, pointing left. Below the repeat, the text 'PCR gel electrophoresis' is shown with a downward arrow.

alleles	n
A	11
B	12
C	13
D	14
E	15
F	16

n	16	15	14	13	12	11
F	—	—	—	—	—	—
E	—	—	—	—	—	—
D	—	—	—	—	—	—
C	—	—	—	—	—	—
B	—	—	—	—	—	—
A	—	—	—	—	—	—

Genotype

AB CD EF AD CF

SSRs are extremely useful as genetic markers in human studies because:

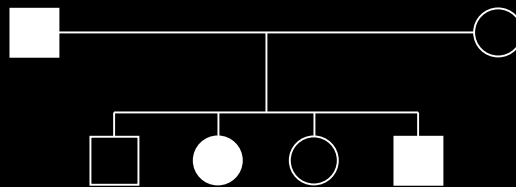
- they are easily scored (by PCR)
- they are codominant
- many SSRs exhibit very high average heterozygosities:  $H_{SSR} = 0.7$  to  $0.9$

A randomly selected person is likely to be heterozygous.

- SSRs are abundant

SSRs occur, on average, about once every 30 kb in the human (or mouse) genomes. > 20,000 SSRs have been identified and mapped within the human genome.

### Huntington's disease (HD)

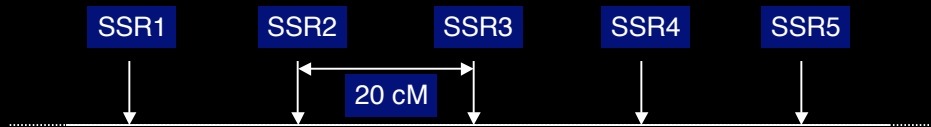


HD: autosomal dominant affecting 1/20,000 individuals

Phenotype: Loss of neurons → personality change, memory loss, motor problem

# genetic linkage mapping

We genotype the six members of the family for SSRs scattered throughout the genome (which spans 3300 cM)—perhaps 165 different SSRs distributed at 20 cM intervals so that one SSR must be within 10 cM of the Huntington's gene:



We obtain potentially exciting results with SSR37, on chromosome 4:

SSR37	A	—	—	—	—	—	—
	B	—	—	—	—	—	—
	C	—	—	—	—	—	—
	D	—	—	—	—	—	—
Genotypes:	HD SSR37	HD/+ AB	+/+ BD	HD/+ AC	+/+ BC	HD/+ AD	+/+ CD
Paternal alleles:	HD SSR37	+	B	HD A	+	B	HD A