



# **“Inferring parental genomes from offspring's DNA”**

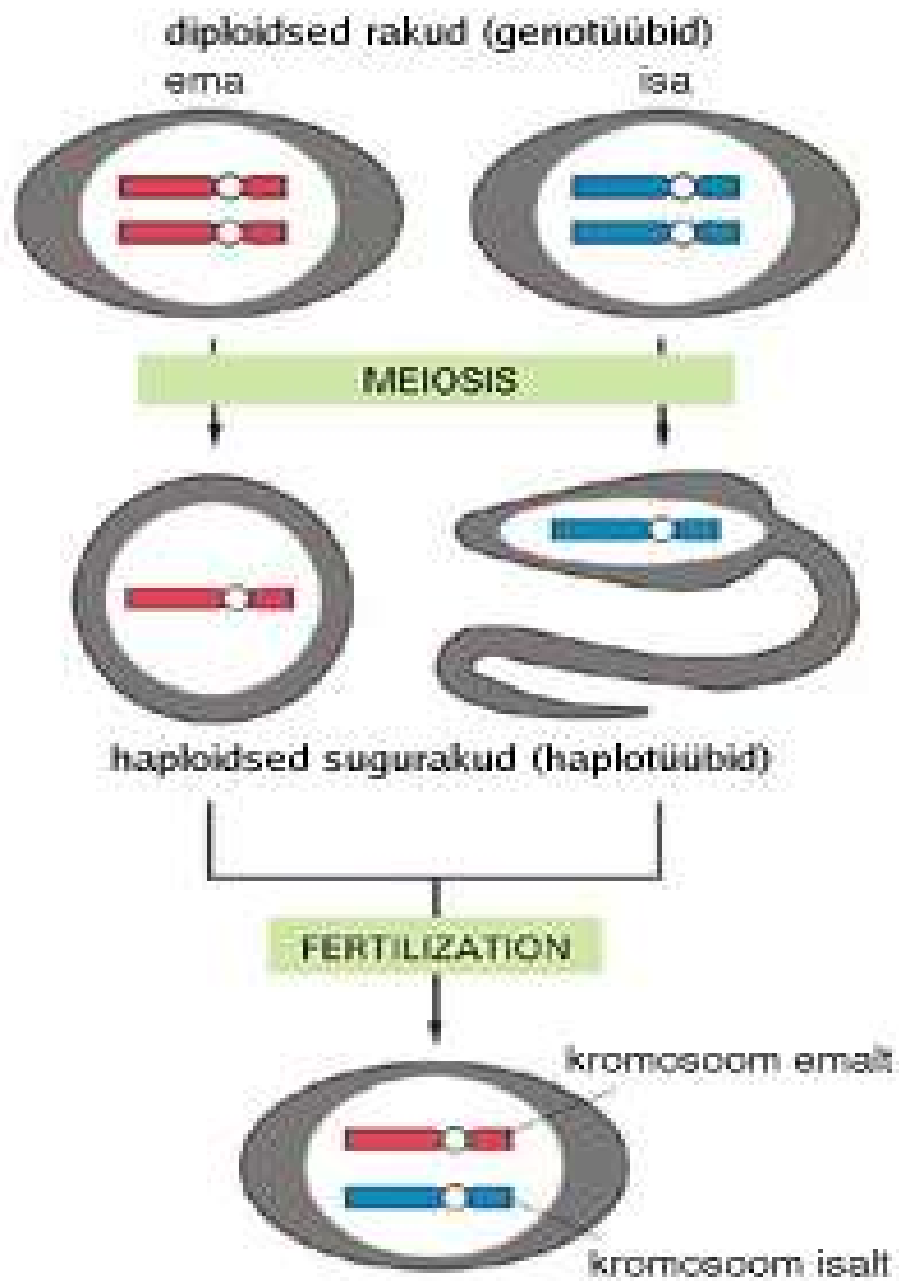
*Kristo Käärmann,  
t-days @ Veskisilla '04  
joint work with Sven Laur*

# Talk outline

- ♦ What is that `haplotype` and why do we need to know about it?
- ♦ *Task 1*: Identify haploid genomes for person X
  - Family method
  - Computational methods
- ♦ *Task 2*: Reconstruct block structure of human genome
- ♦ Further ideas to look into ...

# Human genome

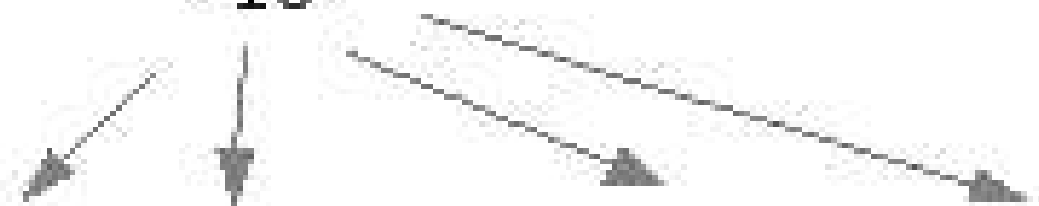
- ♦ DNA ~ 3.4 billion base pairs
- ♦ 99.6% same for all
  - Genetic variation in 10 million SNP-s
    - SNP - Single Nucleotide Polymorphism
    - More or less equally distributed
- ♦ Genome comes in  $2 \times 23$  chromosomes
  - Every cell possesses two versions of each chromosome – one from both parents
  - Recombination merges two genomes into one (takes chunks from each)



# Block structure ?

- ♦ Recombination in large chunks (10 Mb)
  - Size and position of chunks varies
  - Rubik's cube effect
  - *hotspots*, conserved regions
- ♦ Haplotype block: conserved region in DNA (~20-100Kb)
  - *tagSNPs*: set of polymorphisms that uniquely specify all block alleles
- ♦ HapMap project (250 × 400 000 SNPs)
  - Map the block structure of human genome

SNP  
~107

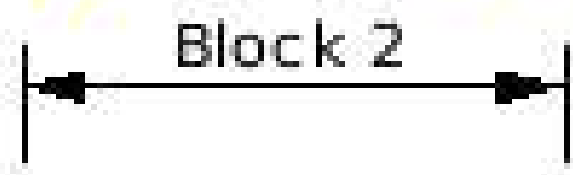


population

AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC  
AC**T**CTAG**A**TTAGACTG**C**ATGAGGAG**A**GC  
AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC  
AC**T**CTAG**A**TTAGACTG**G**ATGAGGAG**T**GC  
AC**T**CTAG**A**TTAGACTG**G**ATGAGGAG**T**GC  
AC**A**CTAG**C**TTAGACTG**G**ATGAGGAG**T**GC  
AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC

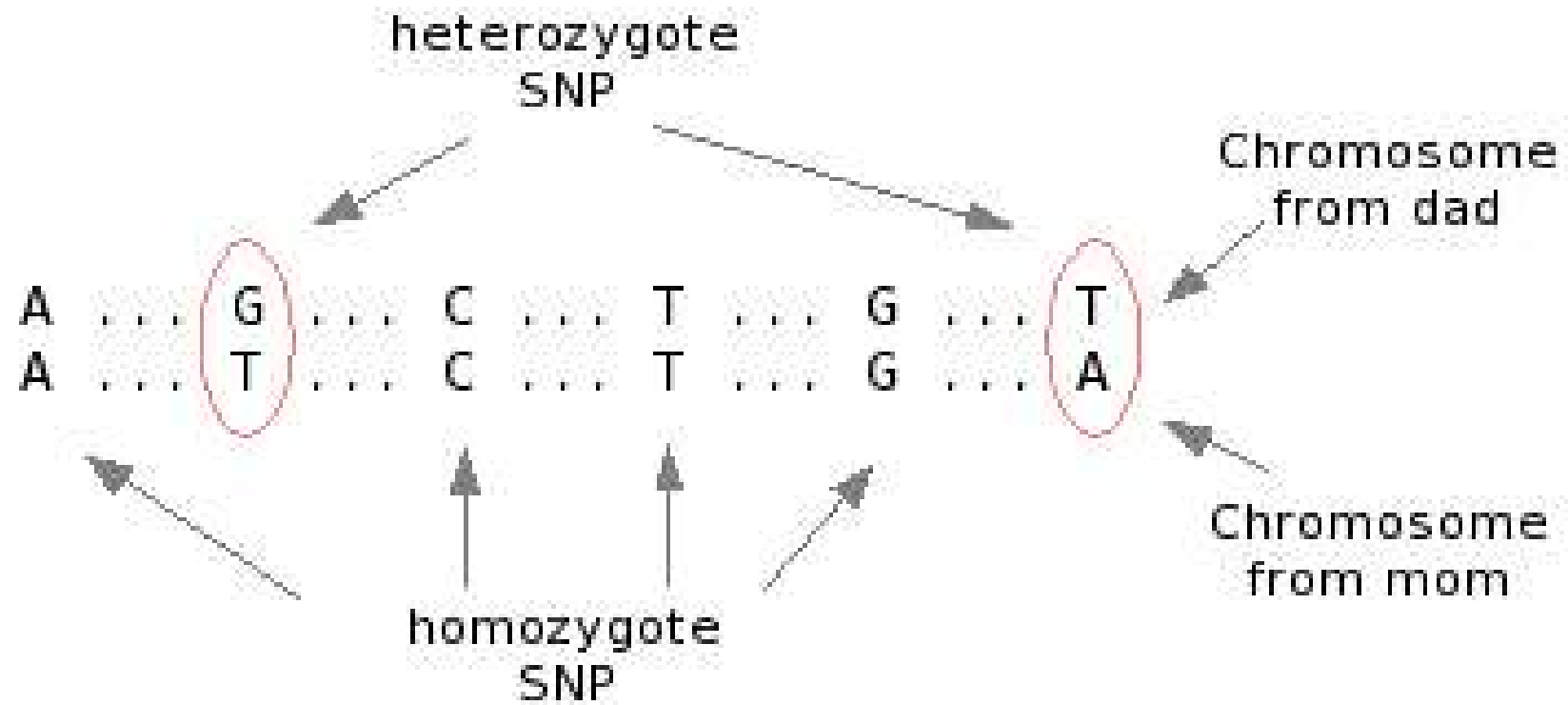


Haploid DNA  
(haplotype)  
 $3.4 \times 10^9$



# Haplotype & genotype

- ♦ Genotyping – identifying a base pair (\$0.1 per SNP)
  - In position  $x$  is a pair  $\{A,G\}$
- ♦ No reasonable lab method for reading haplotypes – that is two separate sequences  $\{ATG \dots TGA, TAG \dots GCC\}$
- ♦ Positions that differ in chromosomes are heterozygote SNP-s
- ♦ Block structure in haplotypes!

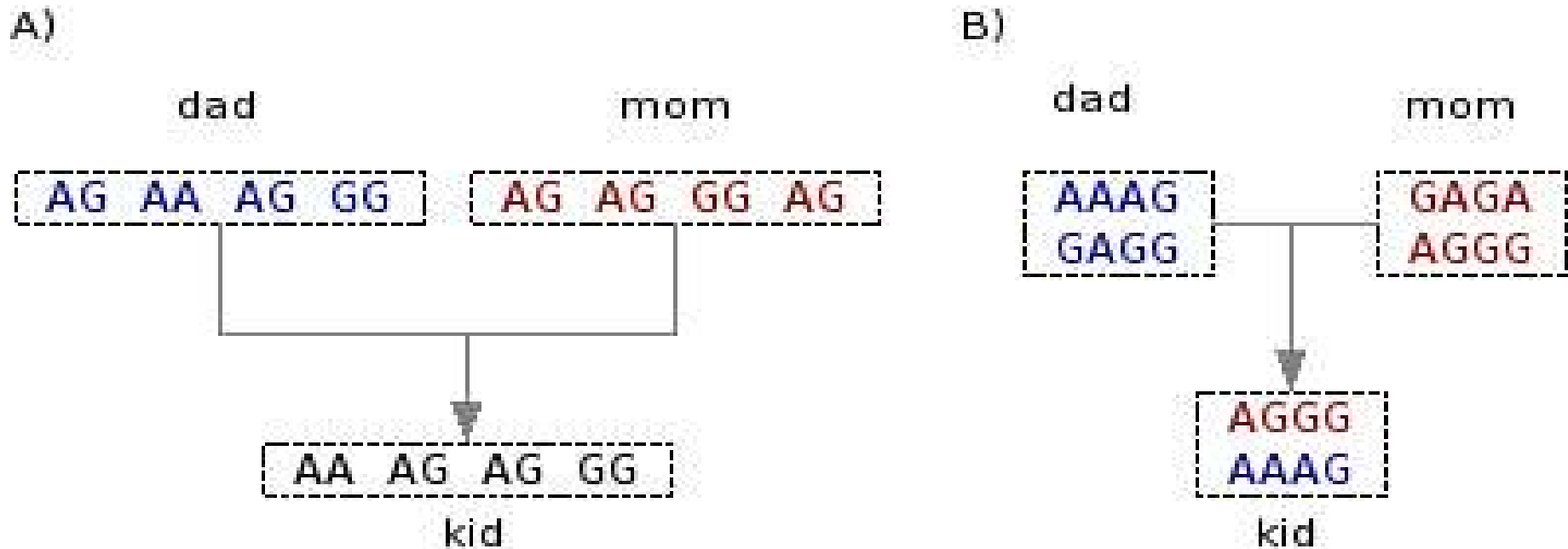


	emalt	isalt
variant 1	GA	TT
variant 2	TA	GT
variant 3	GT	TA
variant 4	TT	GA



# Haplotypes by family method

- Requires genotypes from mom-dad-child trios
- 12.5% SNP-s unsolvable by this method
- Expensive, not very accurate



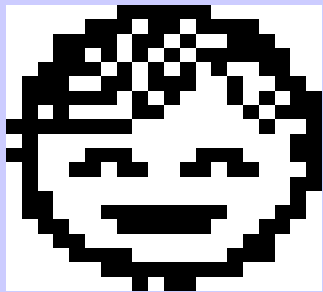
# Haplotypes by computation

- No family data required
- Methods require a population to run on
  - Greedy algorithm (*parsimony*)
  - EM algorithm ('95)
    - Haplotype frequencies in population
  - Markov chains ('01)
    - Gibbs sampling (PHASE)
    - Haplotype inference for every single individual
    - Accurate, slow
  - Phylogeny tree
    - Reconstructing evolution

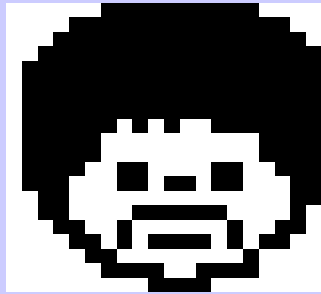
# Testing results

- ♦ How to test produced haplotypes?
  - No lab method to test against!
  - Simulating populations over generations
    - 10 random genotypes x10 random mating
    - Thinning
  - Not a trivial task
- ♦ PHASE
  - Accuracy on HapMap data: ~98.8%
  - 100 SNP: 3h, 500 SNP: 48h

# Sample population



# Illustrative experiment



20 x 20 monochrome bitmap

```
1111111100001111111111110000110
00011111111100111111110111111100
101000010100111110011011111101100
111001111000000111101110011111111
111111001000111001001001110010000
110011110011000100000111111111100
001000001111111111000010000001010
110000000100000000000000000001000
000000000000000001000000000000000
00010000000000000000000001100000000
000000000011100000000000000001111
100000000000000011111111000000000
1111
```

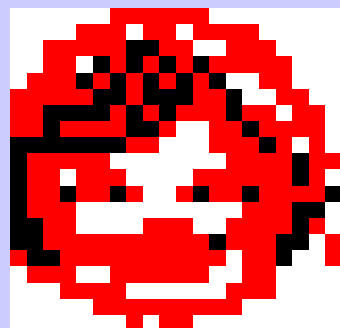
400 bits

```
ATAGCTAGACGATAGAATGCTCCCTAG
ATACTTGTTTCGCTAGCTAGCTTAGATCG
ATTCGAGGATCGTAGATGCCCATGATC
GATCGAATGCATGCAGGGAGGAAATCG
ACTGACTGATGCATGCATCTTACGTACG
ATCATCACTAGTCGAGTCAGCAGCATC
GACTGACATGCTTGACTCGATGTTGACT
ACGTTGCTAGCATCCTCTAGCATCGATC
TAGCTGATCGTAGTACGACTGACTGAT
GCATTAGCATGC.....
```

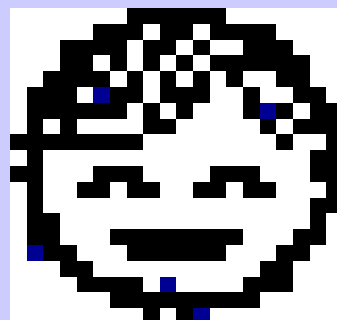
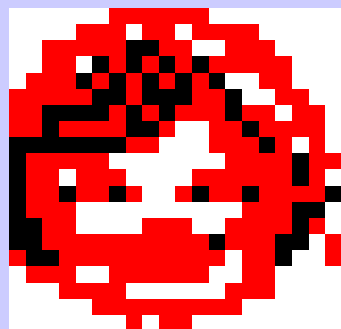
400 SNPs = 100 kB region



# Offsprings from random xxx



# Haplotypes by computation



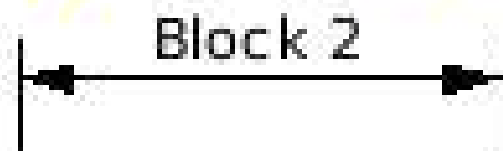
Noise rate: 1.2% = 5 pixels

SNP  
~10<sup>7</sup>

AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC  
AC**T**CTAG**A**TTAGACTG**C**ATGAGGAG**A**GC  
AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC  
AC**T**CTAG**A**TTAGACTG**G**ATGAGGAG**T**GC  
AC**T**CTAG**A**TTAGACTG**G**ATGAGGAG**T**GC  
AC**A**CTAG**C**TTAGACTG**G**ATGAGGAG**T**GC  
AC**A**CTAG**C**TTAGACTG**C**ATGAGGAG**A**GC

Haploid DNA  
(haplotype)  
3.4 x 10<sup>9</sup>

population





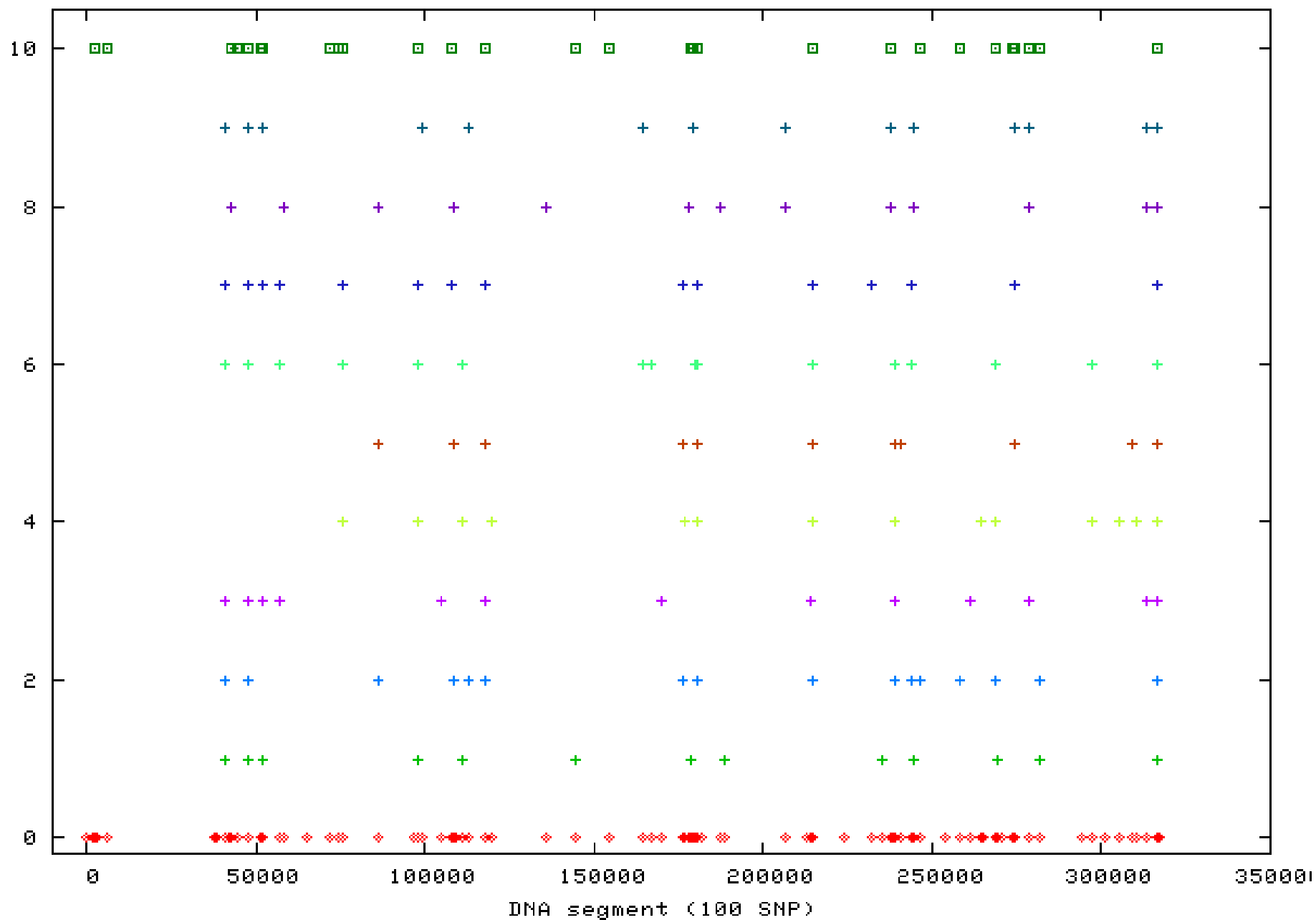
# Mapping block structure

- ♦ Find conserved regions in DNA, claim tagSNP-s
- ♦ Methods
  - Correlation based (four-gamete test)
  - Dynamic programming
  - Minimum description length (MDL)
- ♦ How to test the results?
  - Unsupervised learning ... (generalisation!)

# Dynamic progr. (Zhang)

- ♦  $S_j = \min_{i < j} [S_i + f(i, j)]$ 
  - Minimising additive cost function
  - $f(i, j)$  = number of *tag*SNPs, which uniquely identify x% of block [ i, j ] alleles
  - Block is defined by *tag*SNPs
- ♦ Experimental input:
  - 190 haplotypes (95 individuals)
  - 100 SNPs over 350Kb (~0.01% of DNA)

190 haplotyybi blokistruktuur (10 fold katse)



# What is to be done ...

- ♦ Haplotype inference
  - Scalability! Processing time exp to segment length
  - Partition-Ligation
    - Hard to merge, overlapping?
  - Parallel processing
- ♦ Mapping into blocks
  - Model-based approach
  - How well can we generalise?
- ♦ Blocks from genotypes?!