

**Analysis of Oligonucleotide
Single Nucleotide Polymorphism (SNP)
Array Data**

Cheng Li

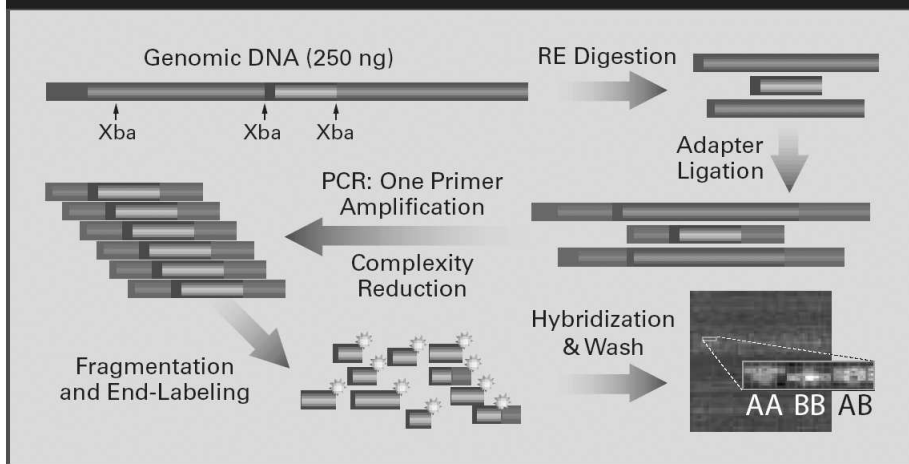
2/12/04

Department of Biostatistical Science, Dana-Farber Cancer Institute
Department of Biostatistics, Harvard School of Public Health

SNP array technology

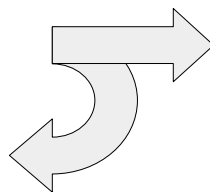
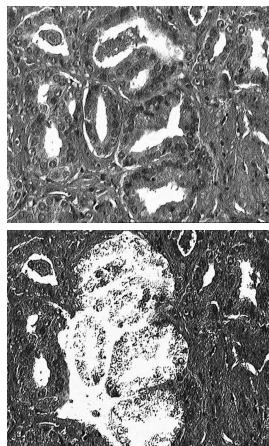
GeneChip® Human Mapping 10K Array

Figure 1: GeneChip® Mapping Assay Overview.



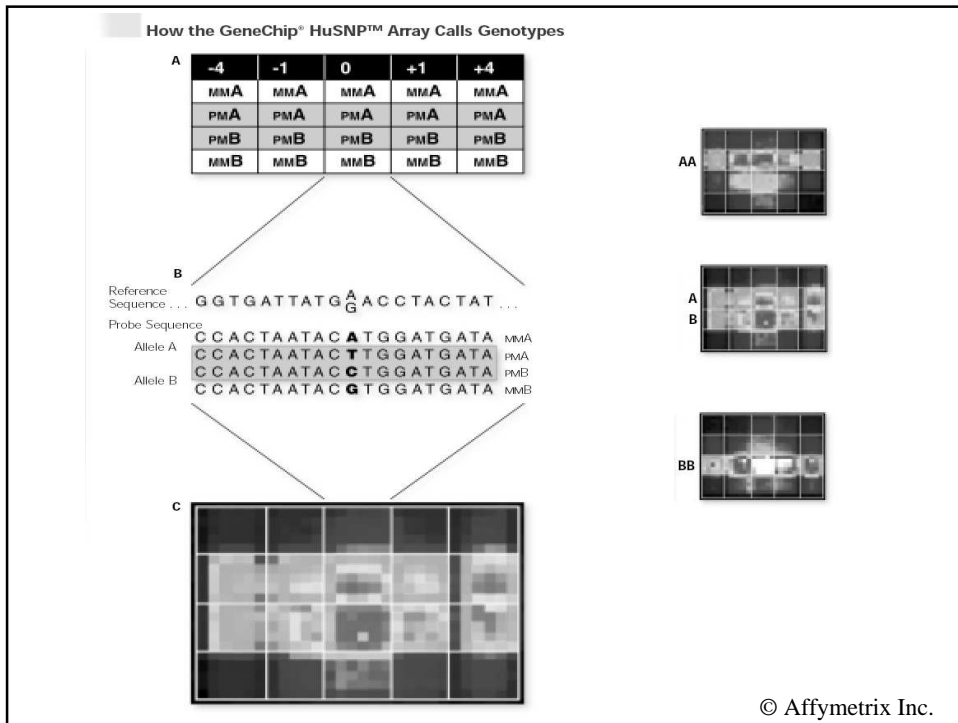
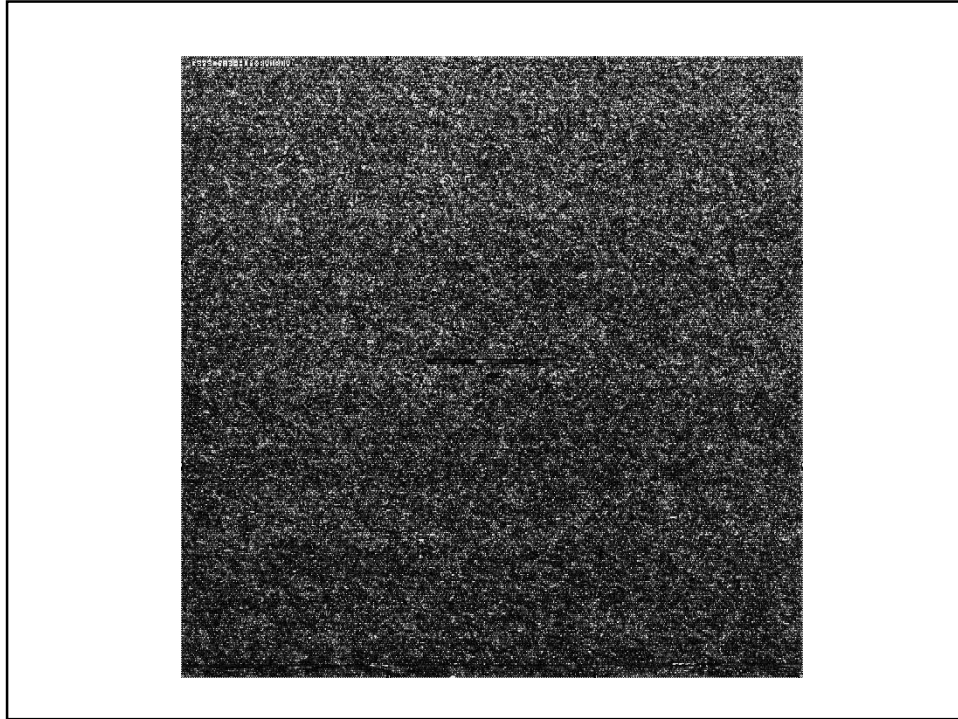
© Affymetrix Inc.

Laser Capture Microdissection

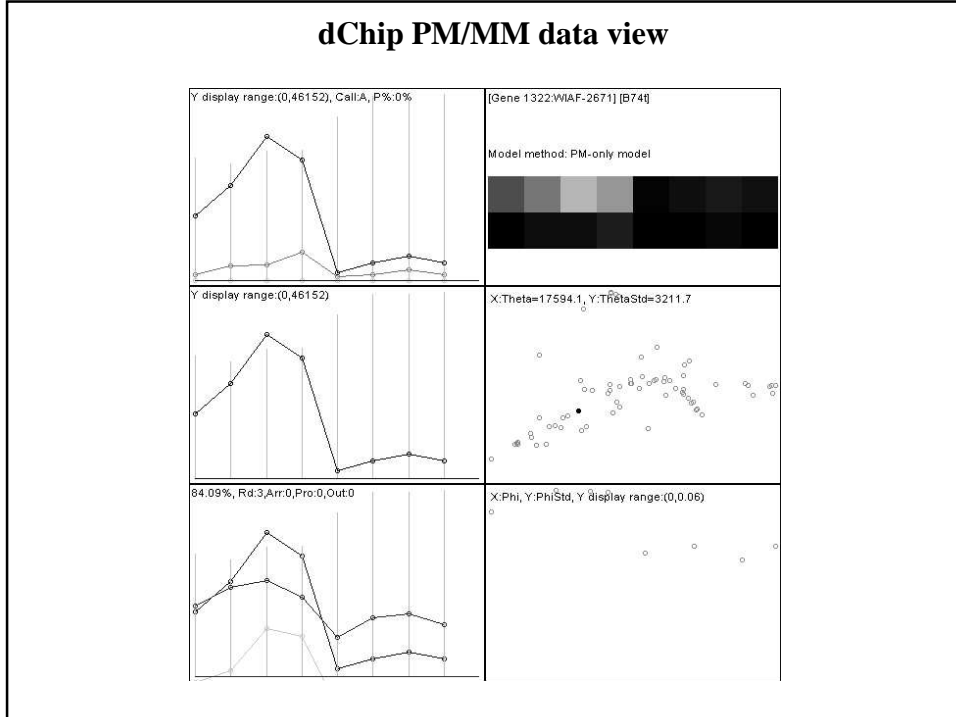


Prostate carcinoma: Formaldehyde fixed, Embedded in paraffin, H&E stained

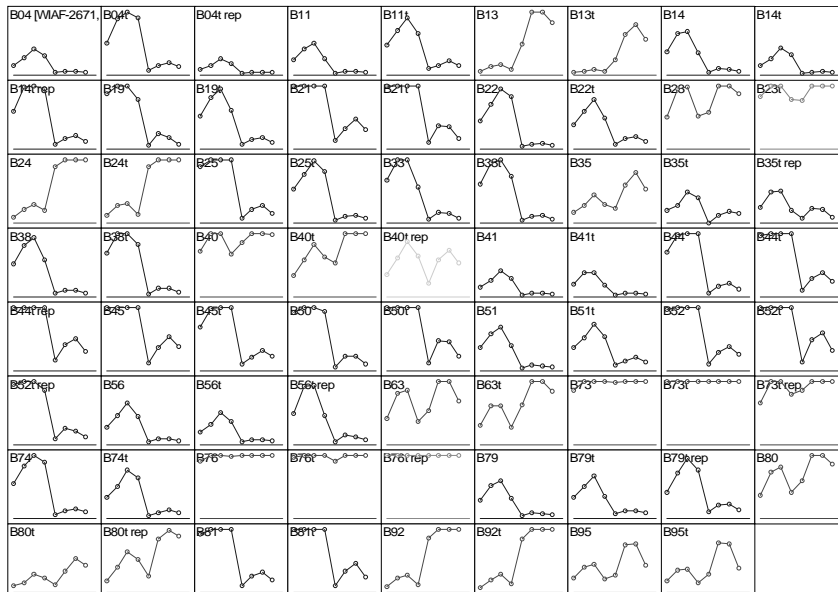
© M.E. Lieberfarb



dChip PM/MM data view



dChip Variation View: Unsupervised clustering may recognize different PM patterns or SNP genotypes



WIAF-2671, Breast cancer data

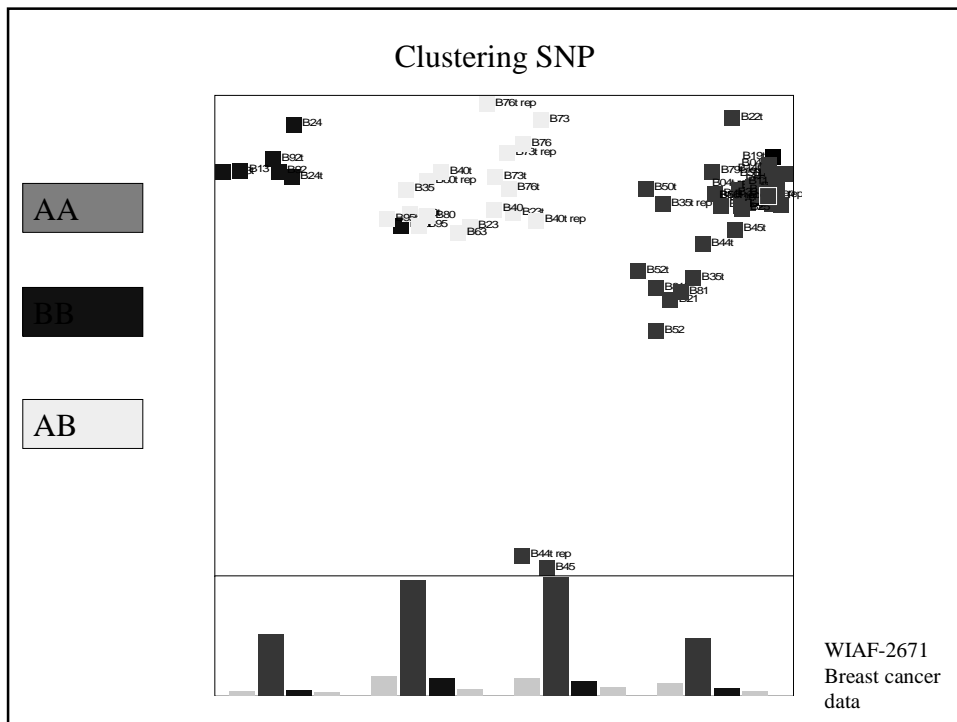
dChip SNP view: project probe data of a SNP to 2D

- For each MiniBlock $i = 1 \dots M$, compute

$$\text{Diff_A} = \max(\text{pmA} - \text{mmA}, 1)$$

$$R_i = \text{Diff_A} / (\text{Diff_A} + \text{Diff_B})$$

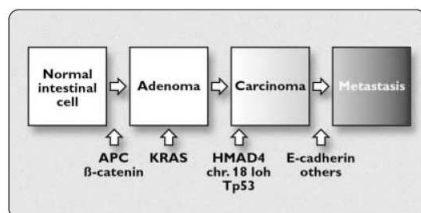
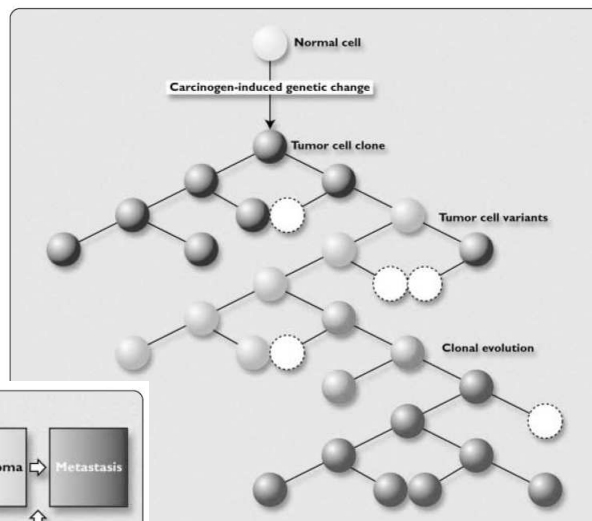
- The data of one SNP in one sample is (R_1, R_2, \dots, R_M)
- Use principle component analysis (PCA) to project S data points (for S samples) into two dimension to visualize



Loss of heterozygosity (LOH) by SNP array

Genetic evolution model

Cancer is caused by deregulation of growth controlling molecular pathways.



© C. J. Cornelisse

Motivation

- Despite apparent locally confined disease, up to 30% of prostate cancer patients undergoing radical prostatectomy will develop recurrence.
- Initial Hypothesis: The differing clinical outcomes of prostate cancer arise from differentially expressed genes
- New hypothesis: Correlating gene expression data with specific genetic alterations may identify biologically relevant gene expression patterns.

© M.Meyerson

Is there evidence that a genetic lesion can produce a global alteration in gene expression?

- BRCA

-Mixed Lineage Leukemia

- ALL, MLL, AML

- Upregulation of genes in MLL (FLT3)

- Expression profile classifies independent tumor set

© M.Meyerson

Chromosome Alterations

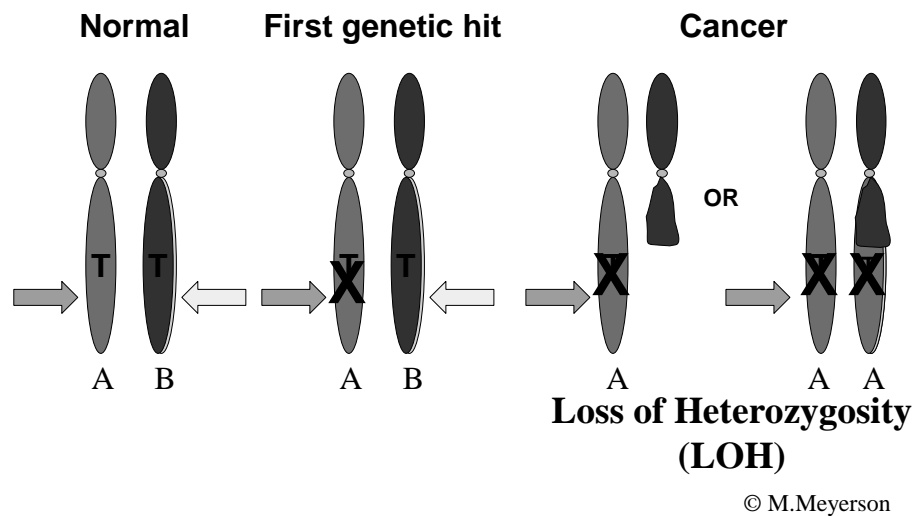
- Non-reciprocal translocations
- Aneuploidy
- Chromosomal amplifications
- Chromosomal deletions
 - loss of tumor suppressor via a point mutation followed by a deletion

© M.E. Lieberfarb

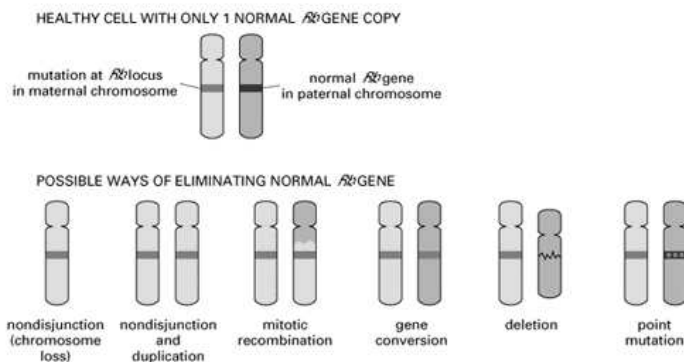
Loss-of-Heterozygosity (LOH)

- If a marker (SNP, micro-satellite) has heterozygous genotype in the normal sample but has homozygous genotype in the tumor sample from the same patient.
- Indicates chromosomal alteration; often related to tumor-suppressor genes

Paradigm for Tumor Suppressor Gene Inactivation by Allelic Loss in Cancers



Six ways of losing the remaining good copy of a tumor suppressor gene (Rb)



W.K. Cavenee et al., *Nature* 305:779-784, 1983.

Alberts et al. 1994 *Molecular Biology of the Cell*, 3rd ed.

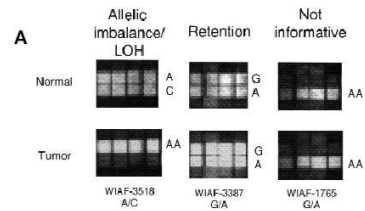
Loss-of-heterozygosity analysis of small-cell lung carcinomas using single-nucleotide polymorphism arrays

Kerstin Lindblad-Toh^{1*}, David M. Tanenbaum^{2-4*}, Mark J. Daly¹, Ellen Winchester¹, Weng-Onn Lui⁵, Anuradha Villapakkam¹, Sasha E. Stanton², Catharina Larsson⁶, Thomas J. Hudson^{1,6}, Bruce E. Johnson^{2,3}, Eric S. Lander^{1,7} and Matthew Meyerson^{2,4}

¹Whitehead Institute/MIT Center for Genome Research, Whitehead Institute for Biomedical Research, Cambridge, MA 02139. ²Department of Adult Oncology, Dana-Farber Cancer Institute, Boston, MA. ³Departments of Medicine and ⁴Pathology, Harvard Medical School, Boston, MA. ⁵Department of Molecular Medicine, CMM, Karolinska Hospital, Stockholm, S-171 76 Sweden. ⁶Montreal Genome Centre, McGill University Health Centre, Montréal, Québec. ⁷Department of Biology, Massachusetts Institute of Technology, Cambridge, MA 02139.

the SNP array, stained with streptavidin-phycoerythrin, and assayed by fluorescence detection. The principles underlying genotyping with SNP arrays were described in an earlier study⁴. Briefly, the detector for each SNP locus contains four rows of 25-mer oligonucleotides, two of which contain oligonucleotides that perfectly match either SNP allele A or SNP allele B, whereas the other two contain single-base mismatches at various positions. The allele-type at a locus is determined by fluorescence intensity ratios in an automated fashion. The approach dramatically decreases the work involved in assaying 1,500 loci, as well as the amount of DNA required (to a total of only 120 ng DNA, corresponding to ~20,000 diploid human genomes) in comparison to both SSLPs and CGH.

The call rate (the proportion of loci to which genotypes could be assigned) was 80.7% ± 3.0% over all samples, yielding ~1,205 SNPs scored per sample (Table 1). The rate did not differ between normal and tumor samples. Many SNPs performed in a robust fashion,



Lindblad-Toh et al. *Nature Biotechnology* 2000

Compare normal and tumor samples

Uninvolved SV (snap frozen)

Normal



Genomic DNA
Purification



paraffin-embedded tumor



Laser-Capture Microdissection



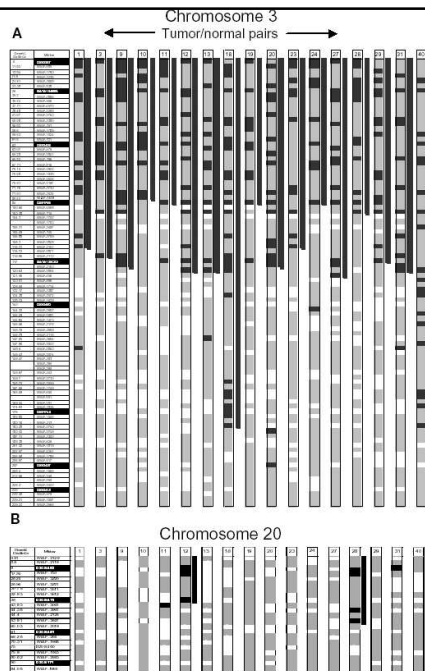
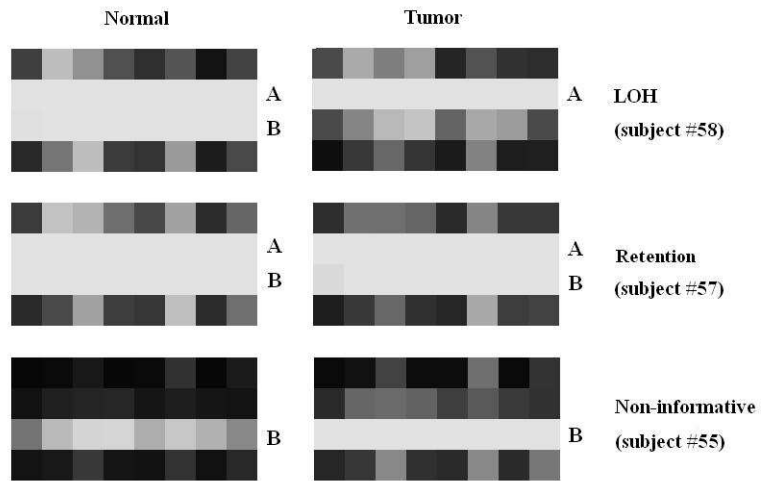
Genomic DNA Purification



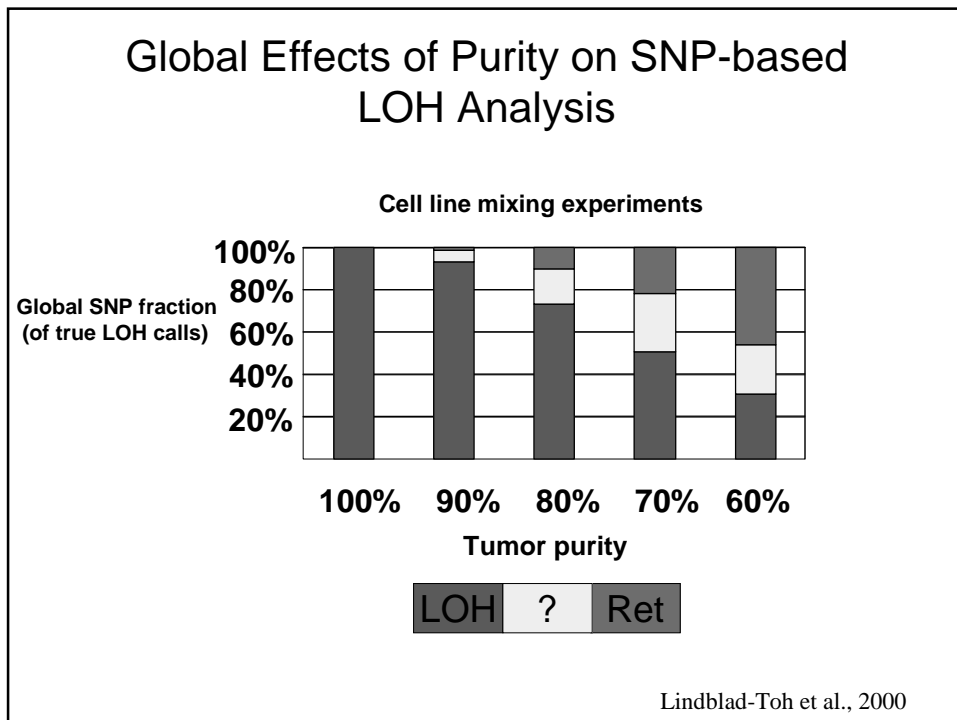
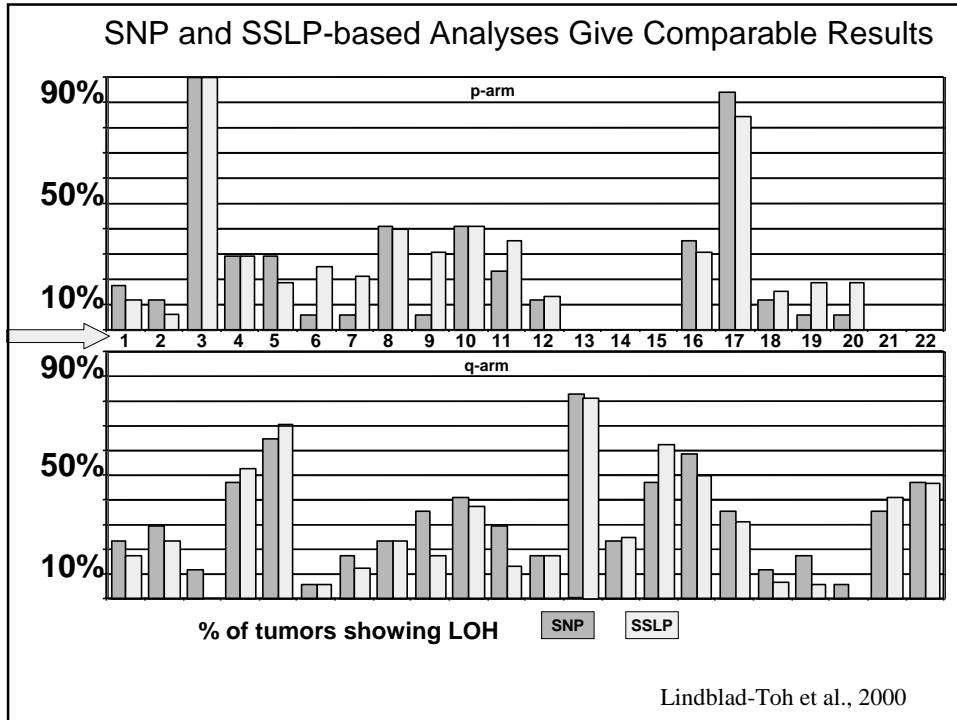
SNP array analysis

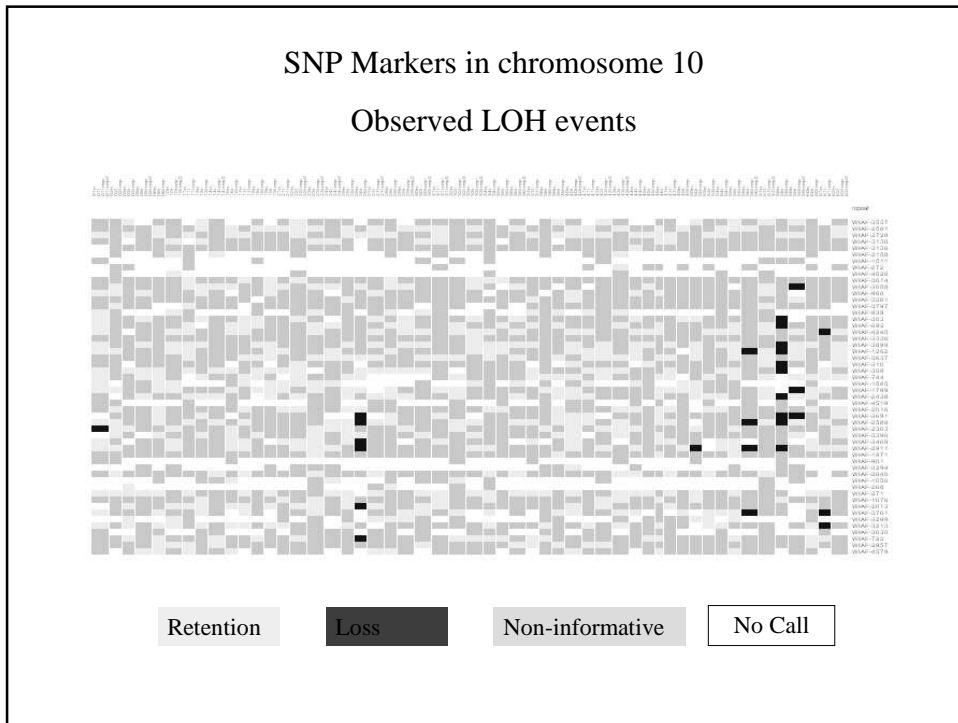
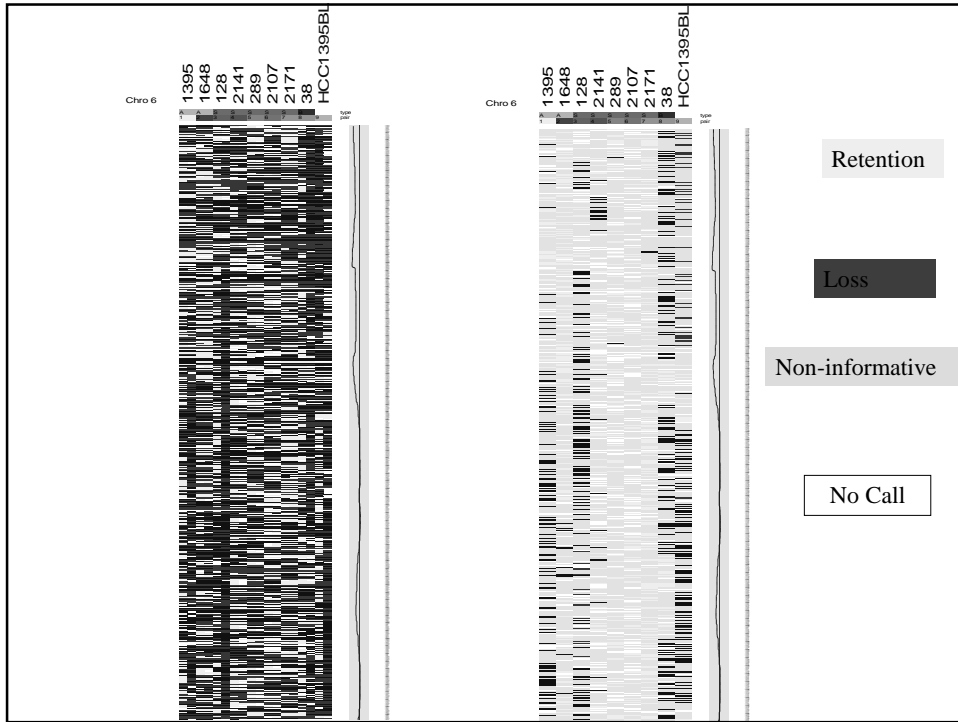
© M.E. Lieberfarb

Making LOH calls

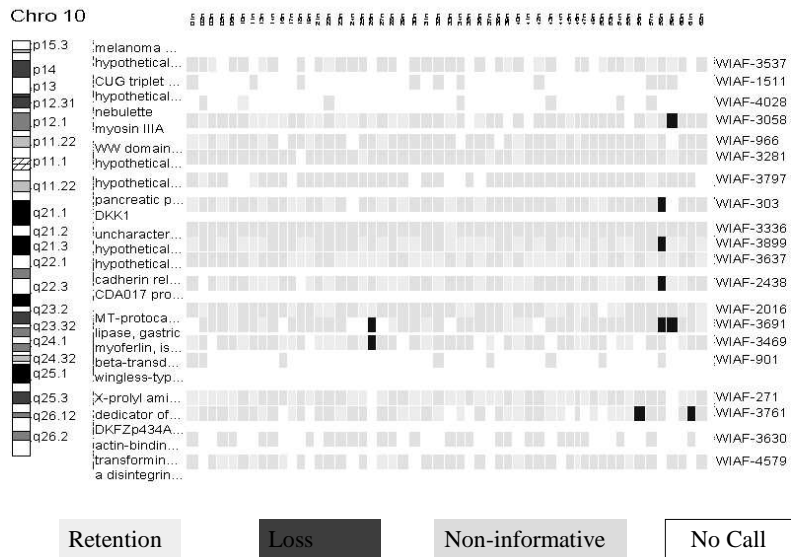


Lindblad-Toh et al., 2000

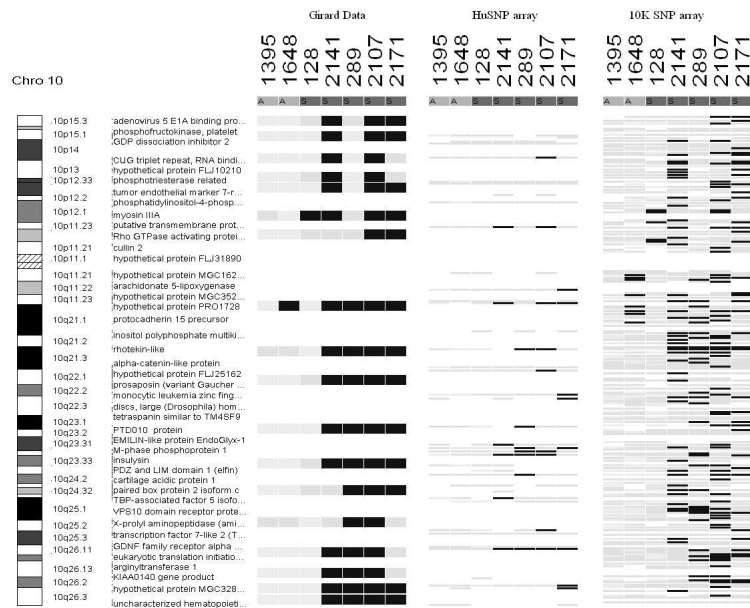


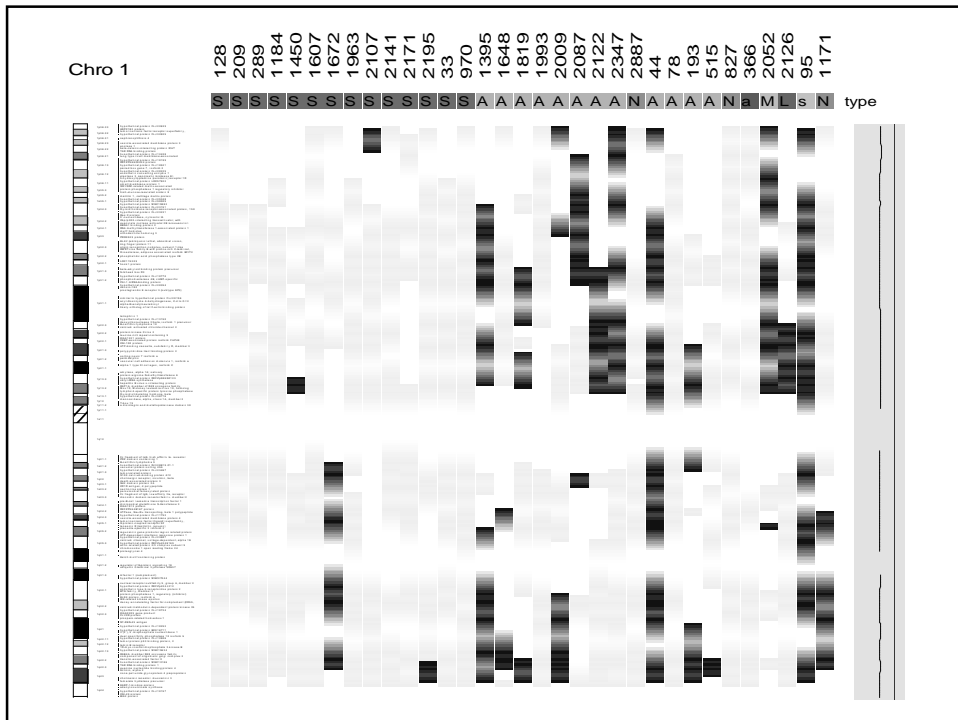
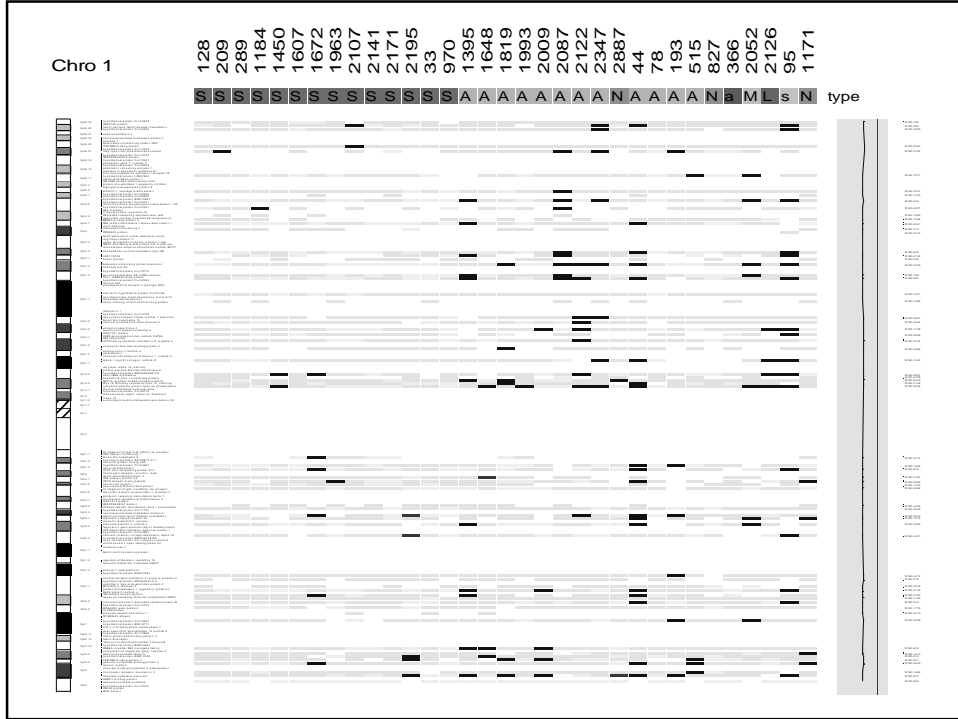


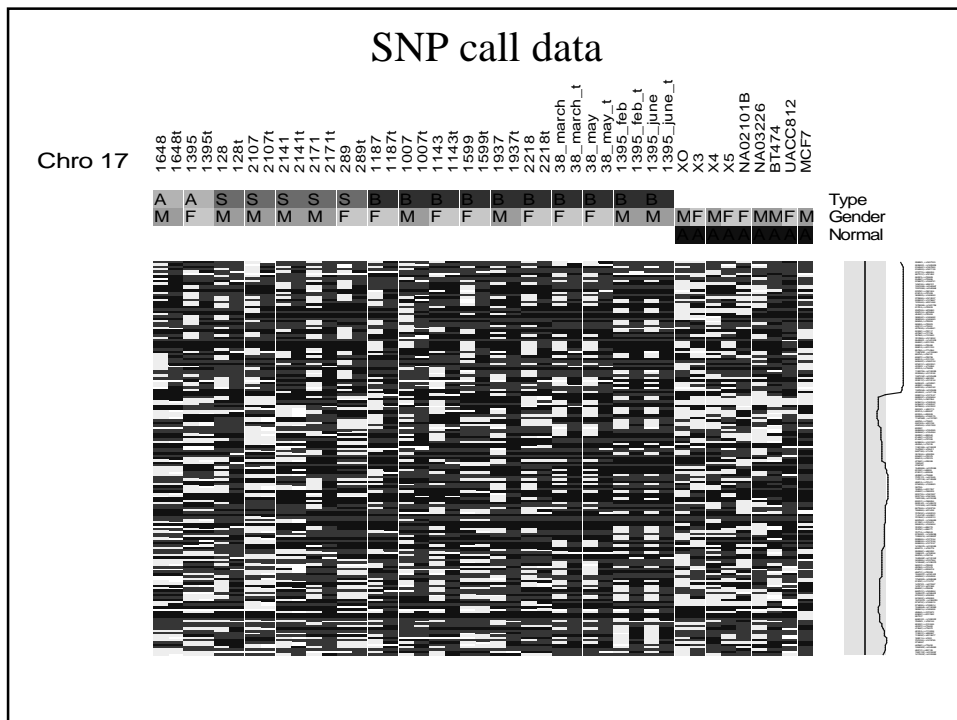
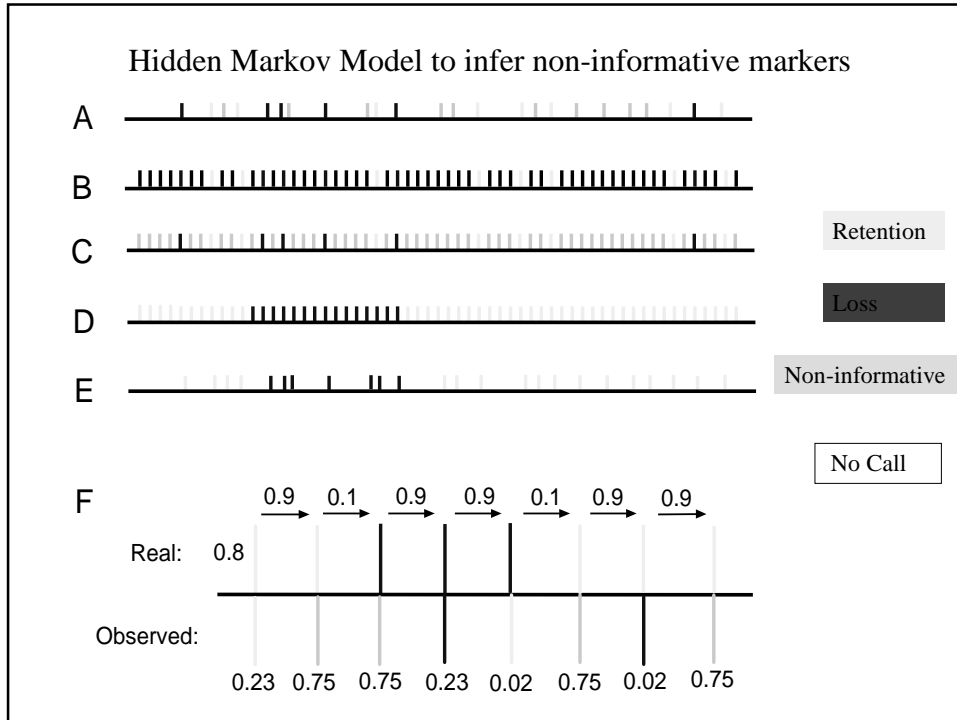
dChip Chromosome View with proportional distance



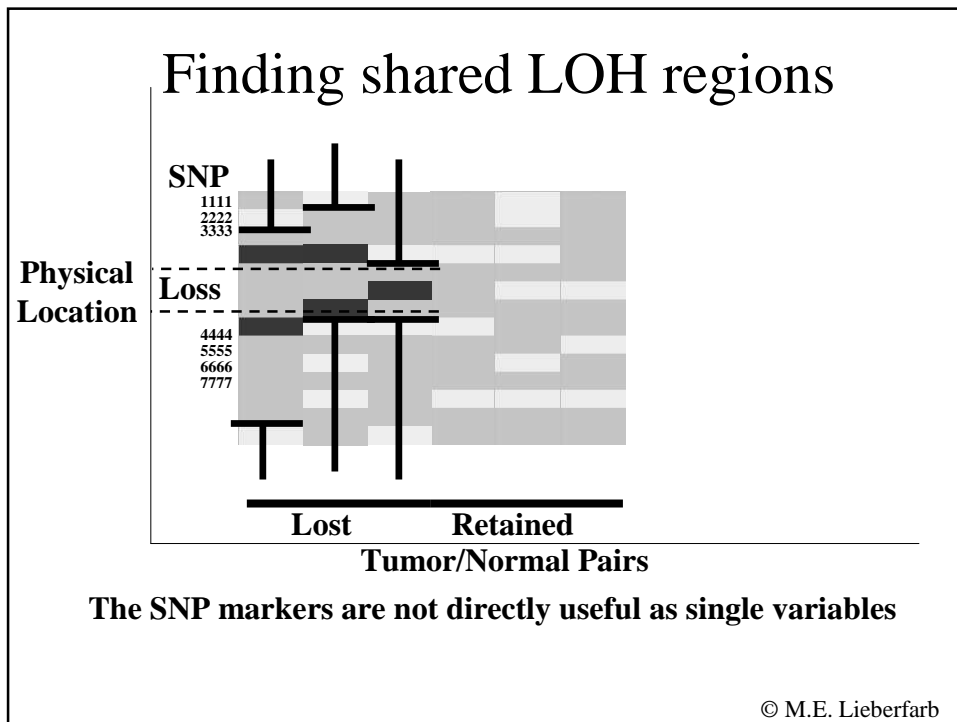
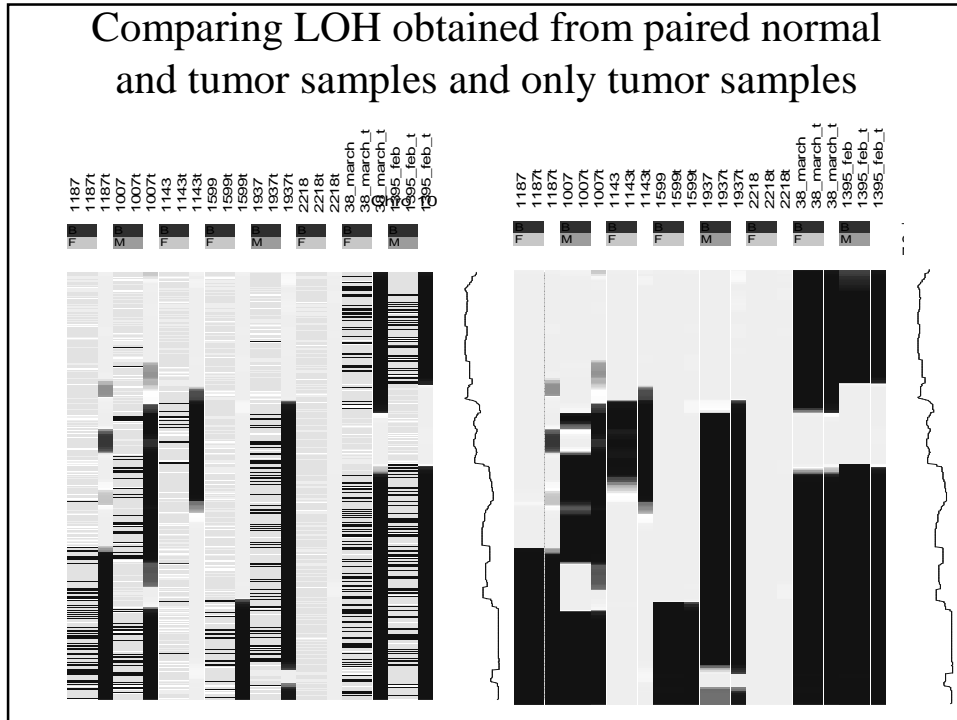
Compare LOH based on SSLP, HuSNP and 10K SNP







Chro 10

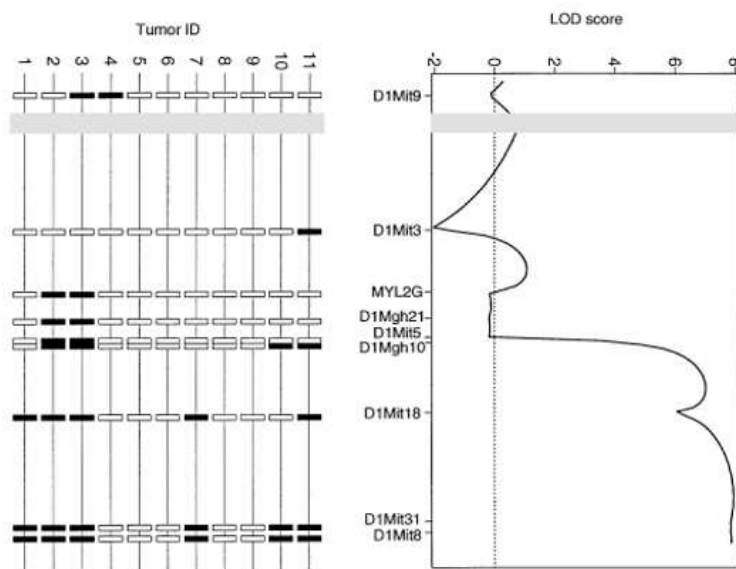


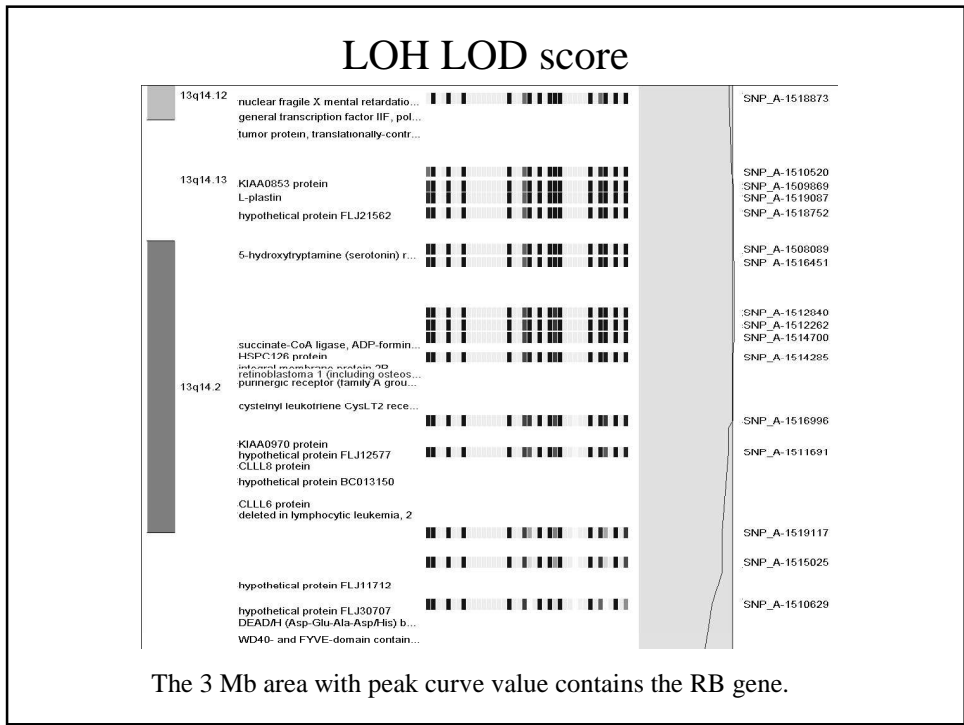
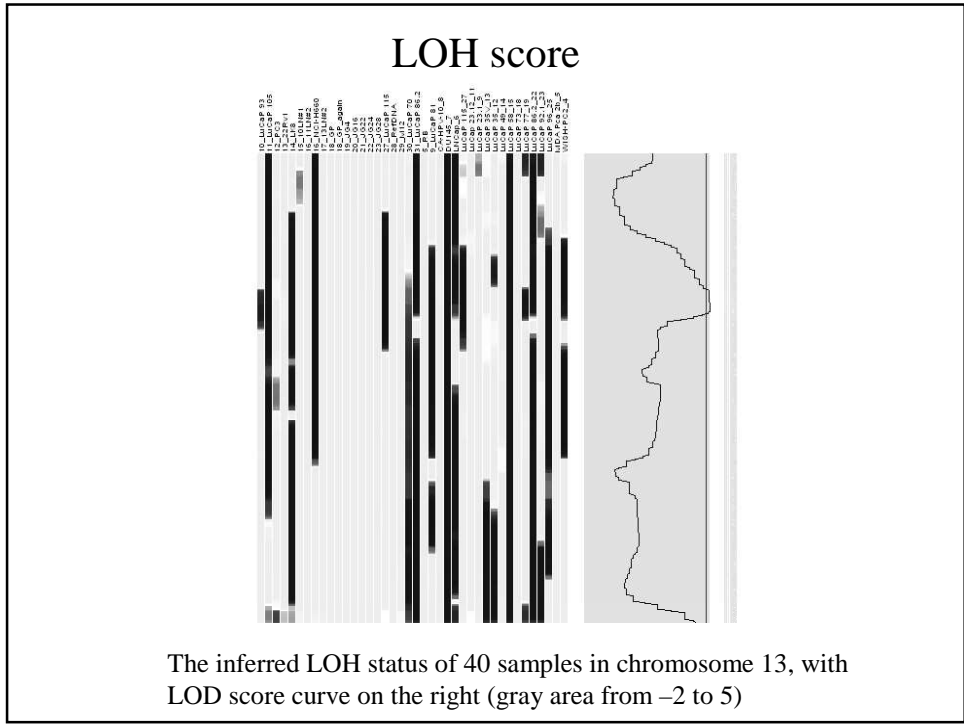
Finding shared LOH regions

- Shared LOH regions might contain tumor suppressor genes
- Complication:
 1. Markers are 300Kb apart, and many of them are non-informative
 2. Call errors, mapping errors
 3. Observed LOH events may be due to genetic instability of the tumor and are not cancer-related.

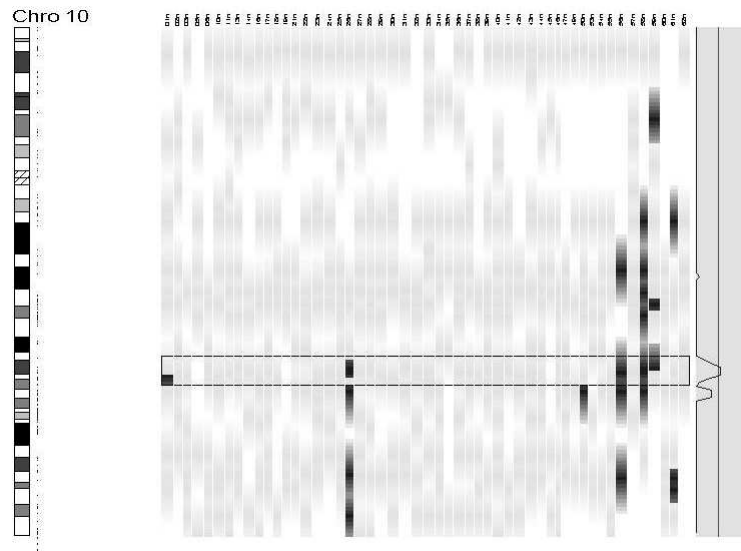
LOH LOD score

M. Newton et al. 1998 *Statist. Med*



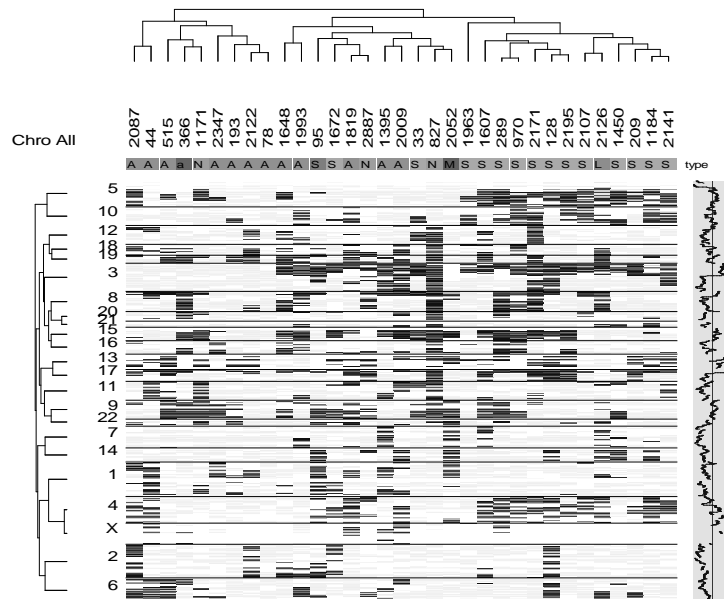


Permutation to compute p-value of cancer-relatedness

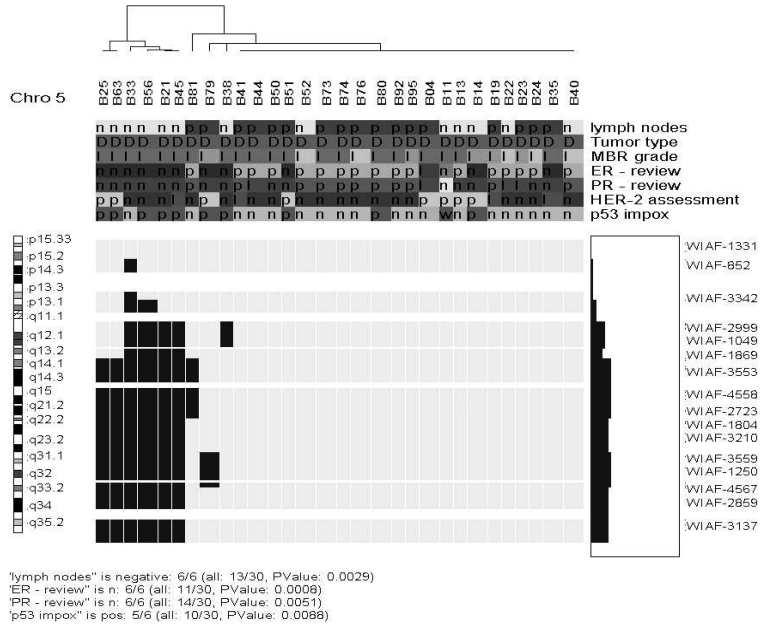


M. Lin et al. *Bioinformatics*, in press

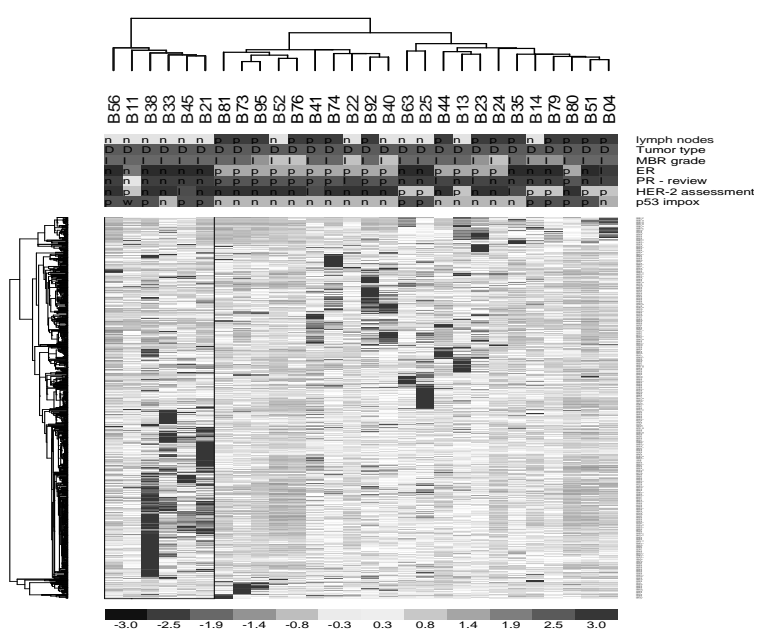
Clustering samples and chromosomes



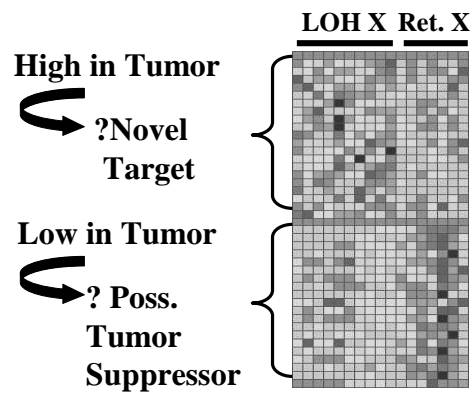
Correlate LOH data with expression data



Correlate LOH data with expression data



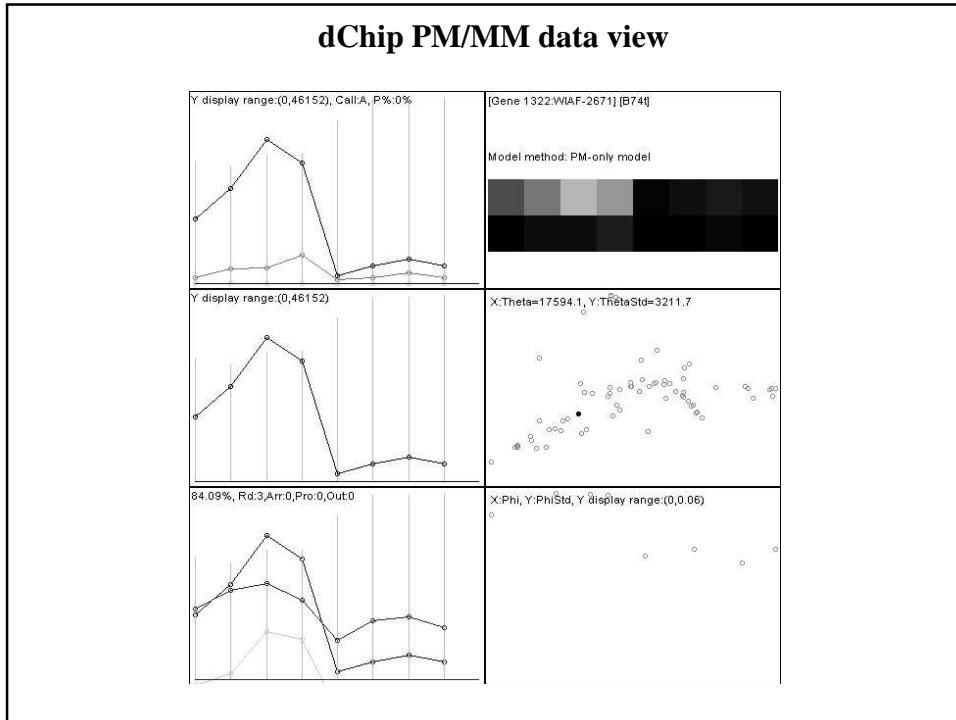
Applying supervised clustering using an LOH
as a class distinction



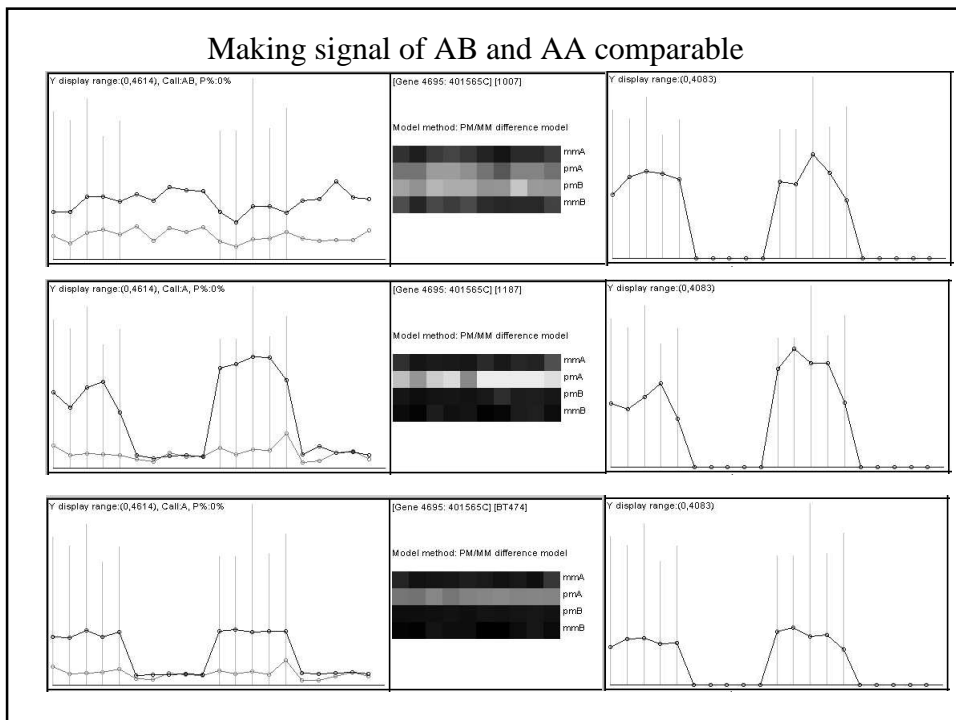
© M.E. Lieberfarb

Array CGH by SNP array

dChip PM/MM data view



Making signal of AB and AA comparable

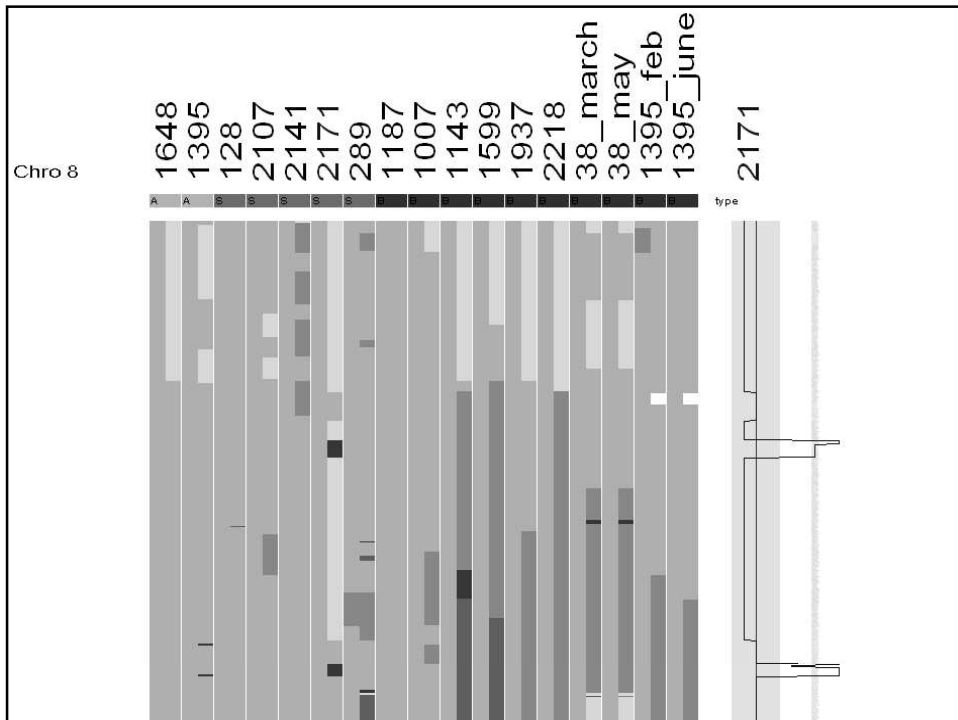
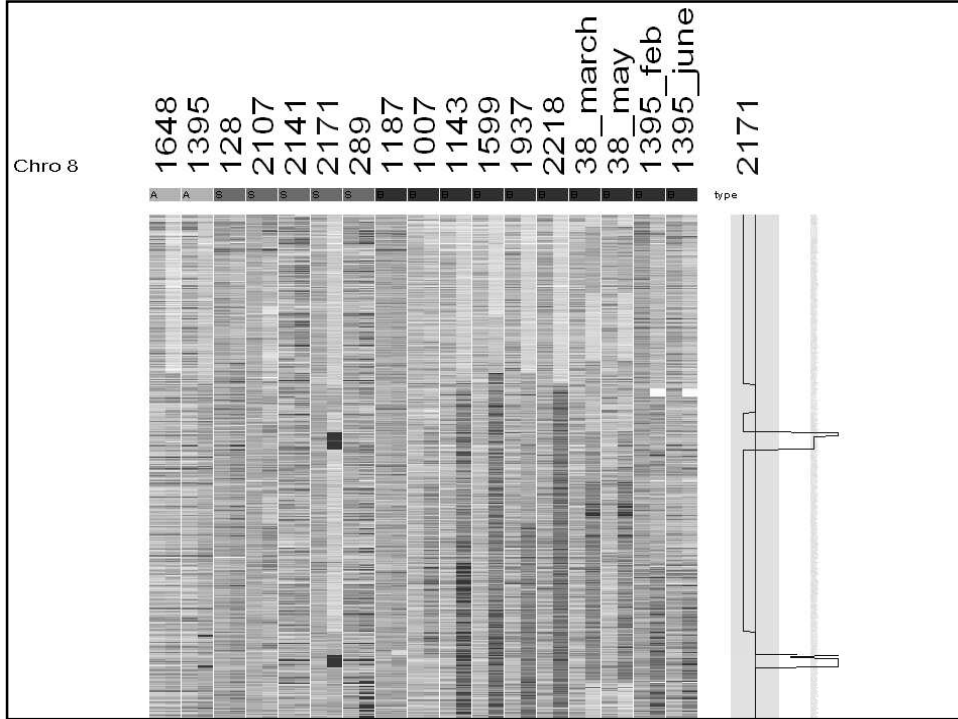


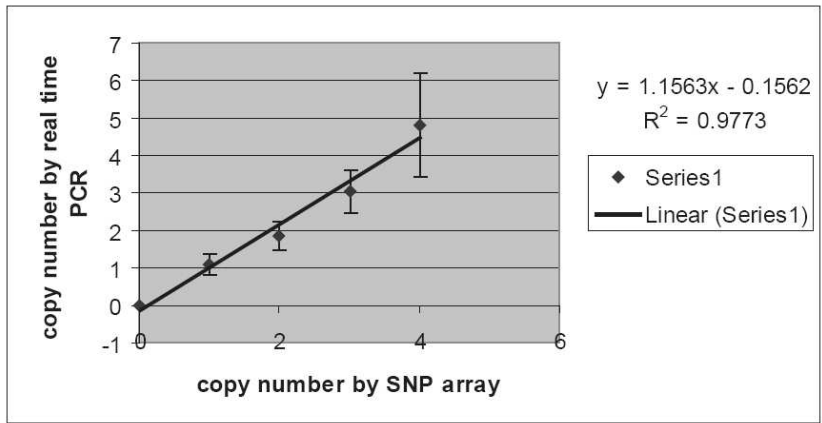
Copy number analysis using SNP array

- Normalization and model-based signal for each array and SNP
- For a SNP, the signal values of all normal cell lines were averaged to obtain the mean signal of 2 copy; observed copy number = (observed signal / mean signal of two copy) * 2
- HMM to infer real copy number by best path

Copy number analysis using SNP array

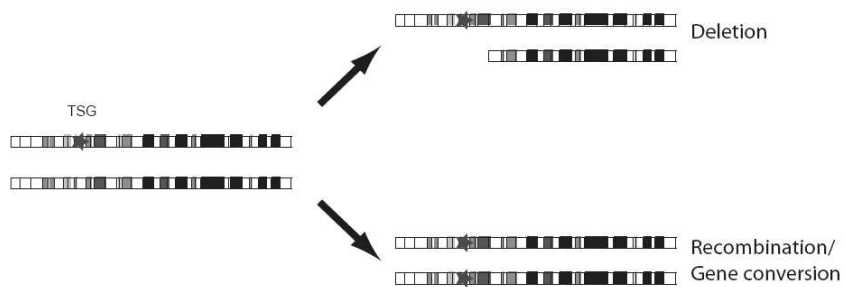
- Emission probabilities: for a SNP the observed signal values are random values drawn from $\left(\frac{\text{Signal} - \text{Mean} \cdot \text{Fold}}{\text{Std} \cdot \text{Fold}}\right) \sim t(40)$
- Transition probabilities: The Haldane's map function $\theta = \frac{1}{2}(1 - e^{-2d})$ convert the genetic distance d between two SNP markers to the probability ($2 \cdot \theta$) that the copy number of the 2nd marker will return to the background distribution of copy numbers in this sample and thus independent from the copy number of the 1st marker
- Initial probabilities: The proportion of chromosome regions that have a particular copy number is set to fixed values in the first round (0.9 for 2 copy, $0.1/(N-1)$ for copy 0 to MaxCopy except 2)



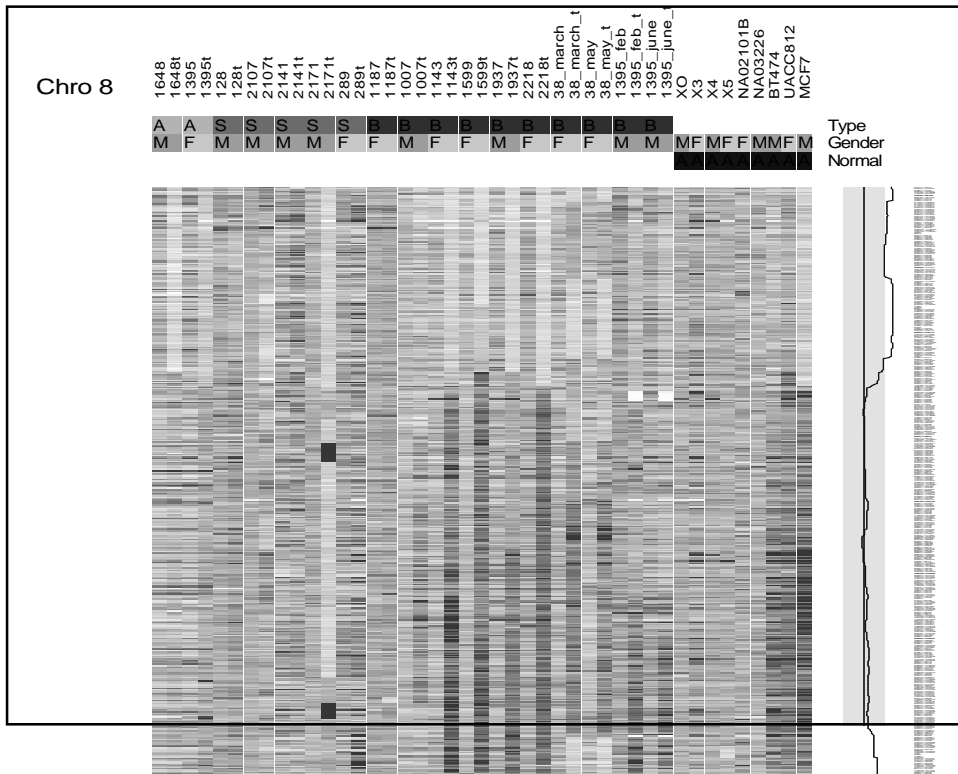
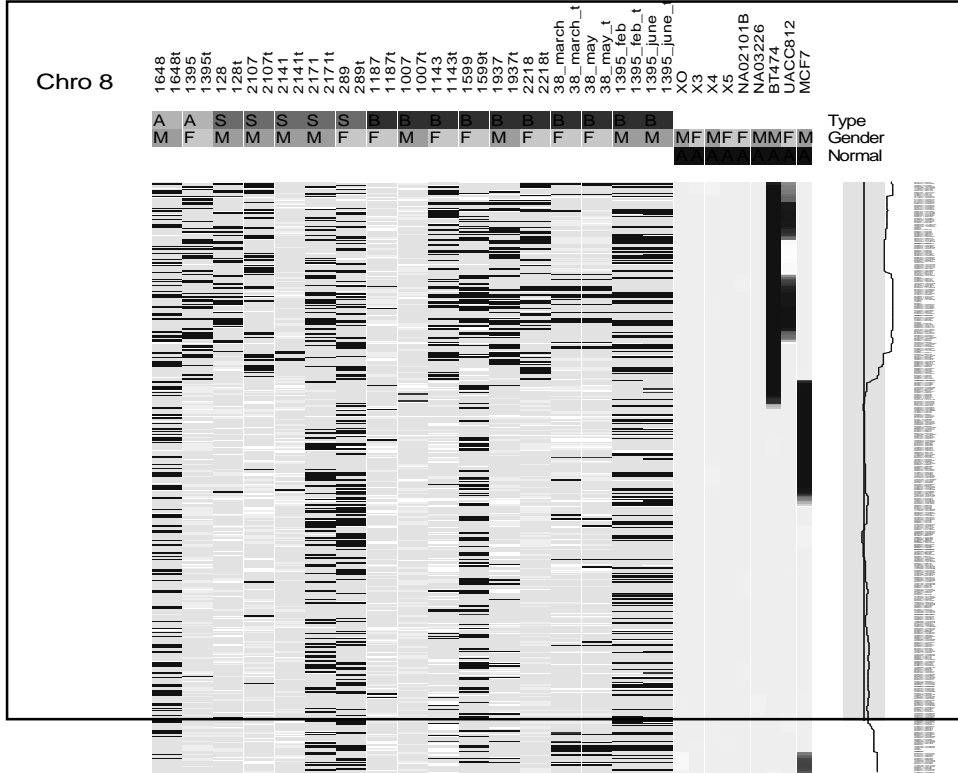


© X. Zhao

Two mechanisms of LOH



© X. Zhao



Acknowledgment

Harvard School of Public Health

Ming Lin

Ke Hao

Lee-Jen Wei

Wing Hung Wong

Carsten Rosenow (Affymetrix)

Richard Smith (U. of Iowa)

Dana-Farber Cancer Institute

Matthew Meyerson

William R. Sellers

Pasi Antero Janne

Marshall Lieberfarb

Andrea Richardson

Zhigang C. Wang

Edward Fox

Xiaojun Zhao



Cheng Li: cli@hsph.harvard.edu