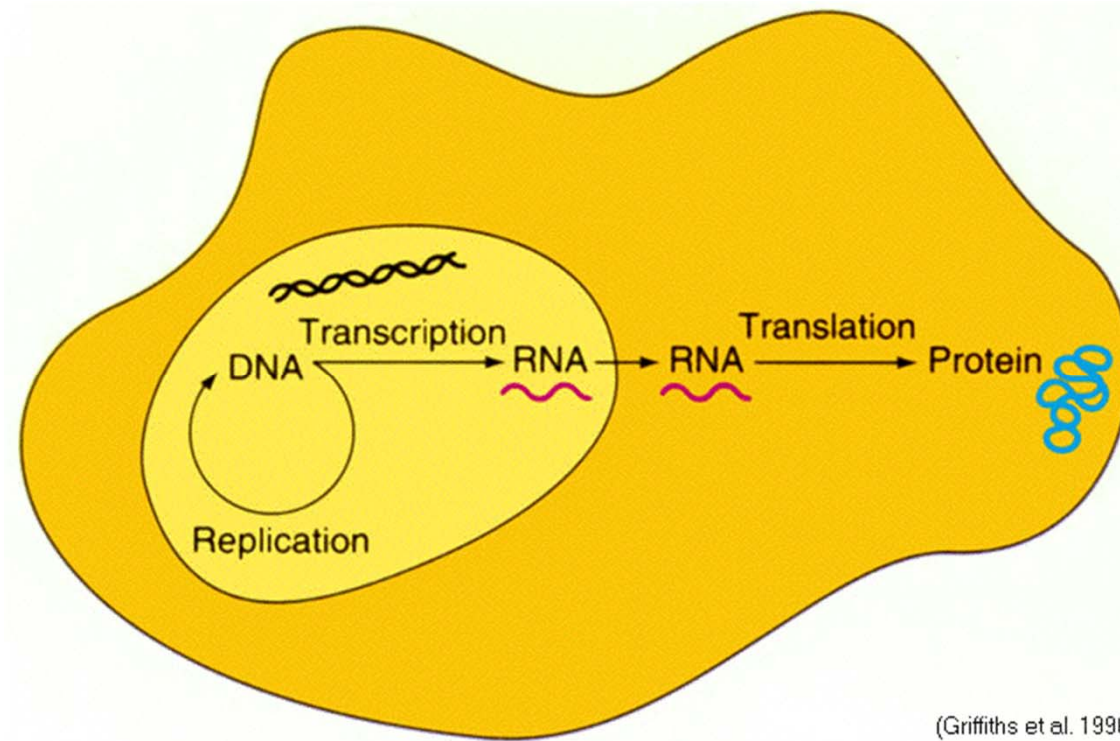


第6-1章: Microarray and Expression Measurements

- Introduction to Microarray
- Expression measurements

Transcriptome



Gene Expression

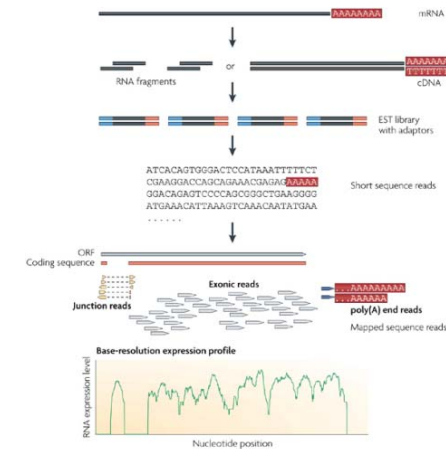
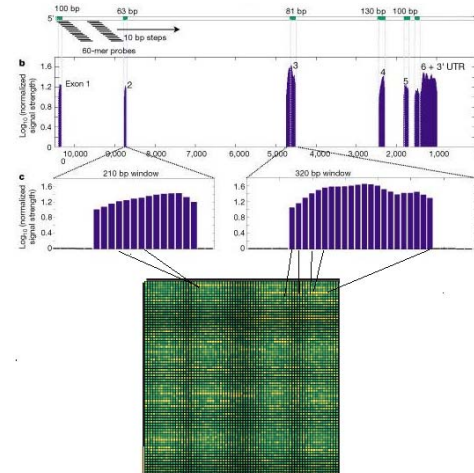
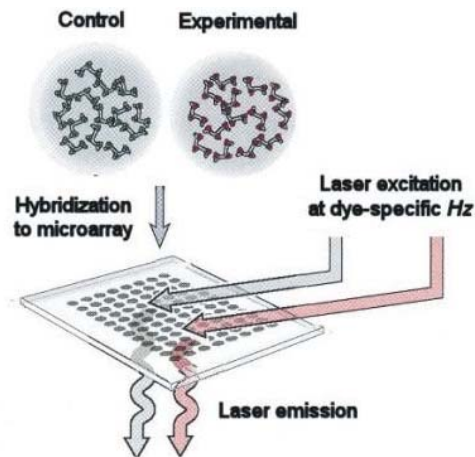
- Each cell contains a complete set of DNA.
- Only a fraction of these are used (or “expressed”) in any particular cell at any given time. For example, genes specific for erythroid cells, such as the hemoglobin genes, are not expressed in brain cells.

What is a DNA Microarray?

- Also known as DNA Chip
- Allows simultaneous measurement of the level of transcription for every gene in a genome (gene expression)
- Transcription?
 - Process of copying of DNA into messenger RNA (mRNA)
 - Environment dependent!
- Microarray detects mRNA, or rather the more stable cDNA

The Evolution of Transcriptomics

Hybridization-based



Nature Reviews | Genetics

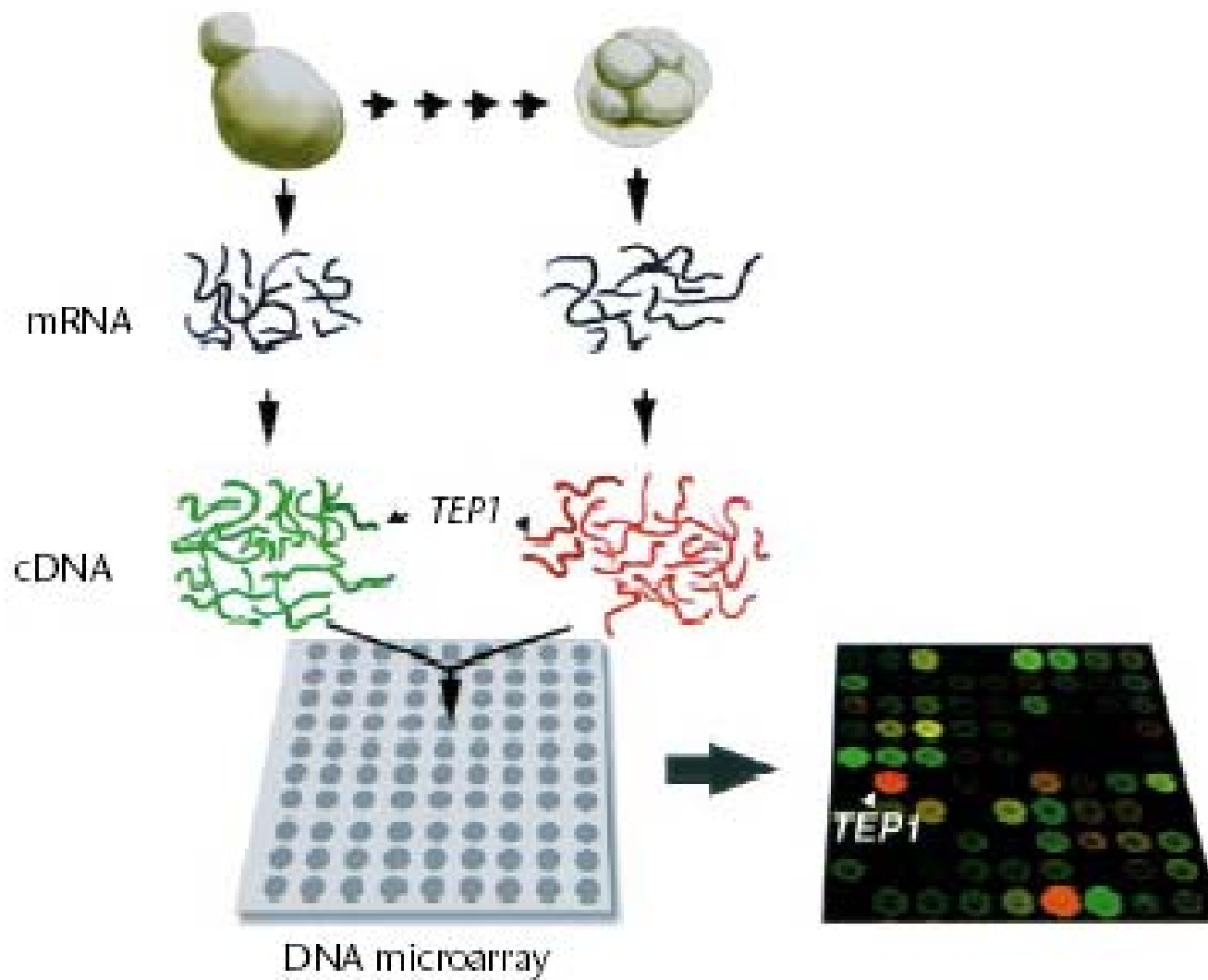
RNA-seq is still a technology under active development

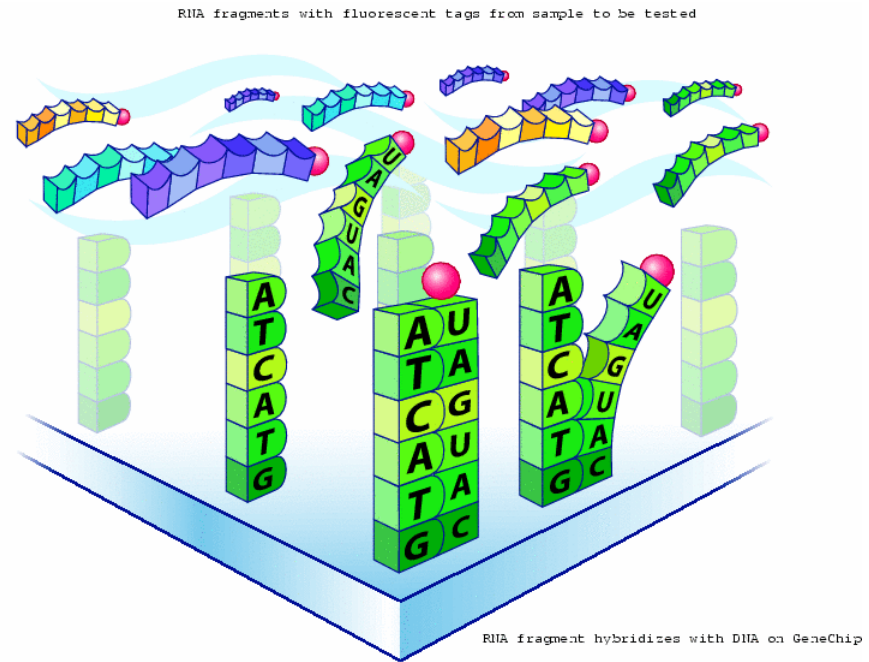
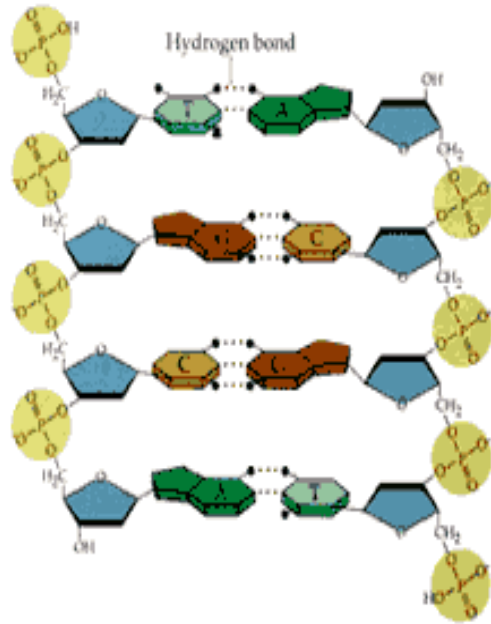
1995 P. Brown, et. al.
Gene expression profiling
using spotted cDNA
microarray: expression levels
of known genes

2002 Affymetrix, whole
genome expression profiling
using tiling array: identifying
and profiling novel genes and
splicing variants

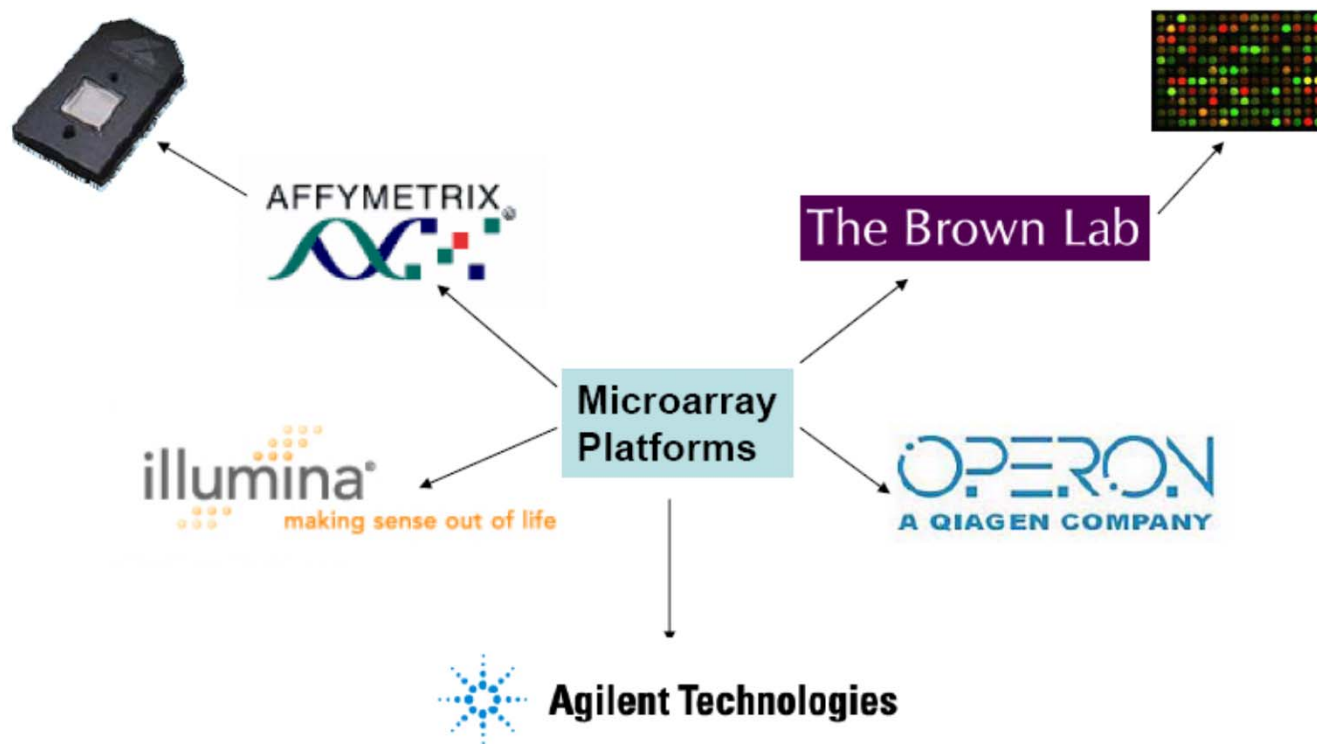
2008 many groups, mRNA-seq:
direct sequencing of mRNAs
using next generation
sequencing techniques (NGS)

cDNA Microarray

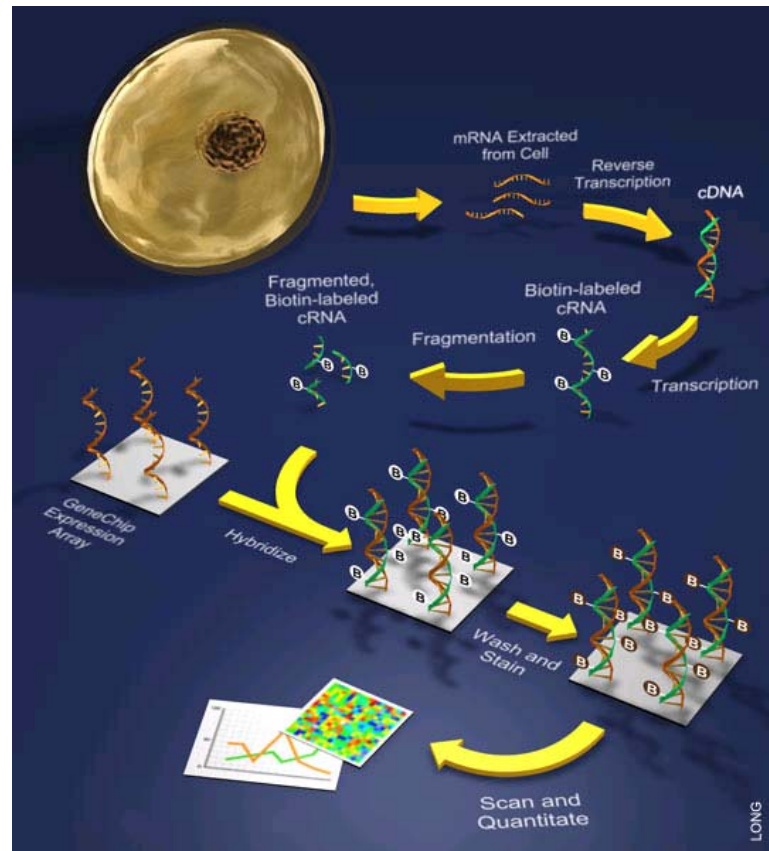




生物学芯片



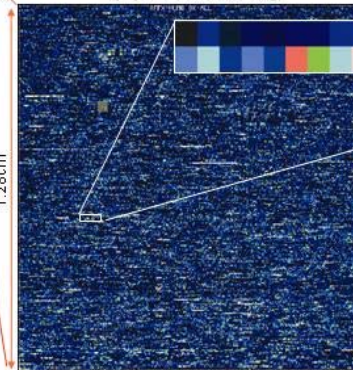
Affymetrix 表达芯片



Human Genome U133A GeneChip® Array

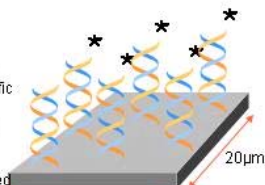


(1) Probe Array



(4) Probe Cell

Each Probe Cell contains $\sim 40 \times 10^4$ copies of a specific probe complementary to genetic information of interest
probe: single stranded, sense, fluorescently labeled oligonucleotide (25 mers)



(2) Probe Set

Each Probe Set contains 11 Probe Pairs (PM:MM) of different probes

(3) Probe Pair

Each Perfect Match (PM) and Mismatch (MM) Probe Cells are associated by pairs

The Human Genome U133 A GeneChip® array represents more than 22,000 full-length genes and EST clusters.

Glossary

- Probe: gene specific oligonucleotides
- 20 - 25 bases in each probe, i.e., 20-25 mer
- Perfect match (PM)
- Miss match (MM)
- Probe pair: a (PM, MM) pair
- Probe-pair set: a collection of probe-pairs (usually 20) represent each gene
- “20 digit barcode”

GeneChip® Expression Array Design

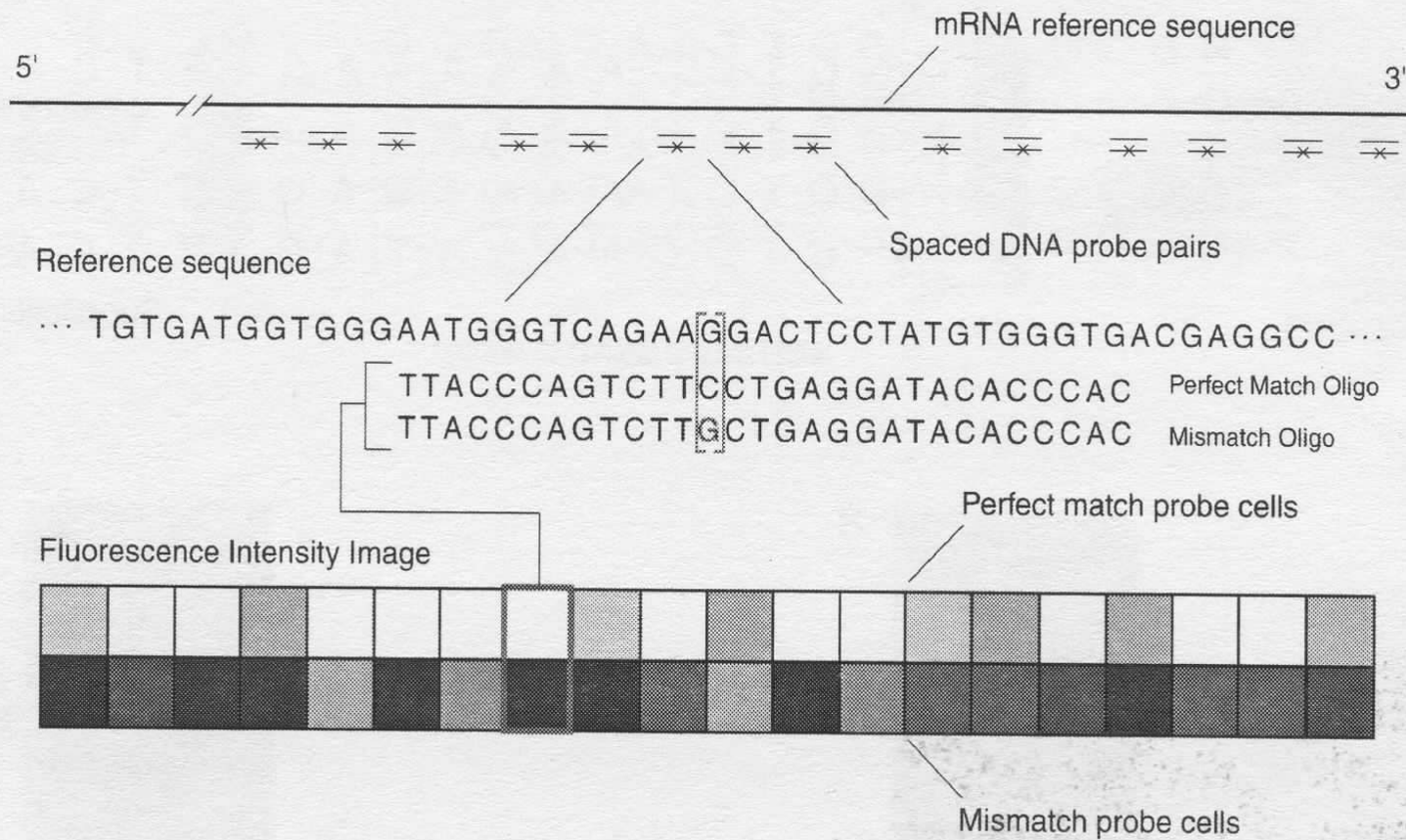
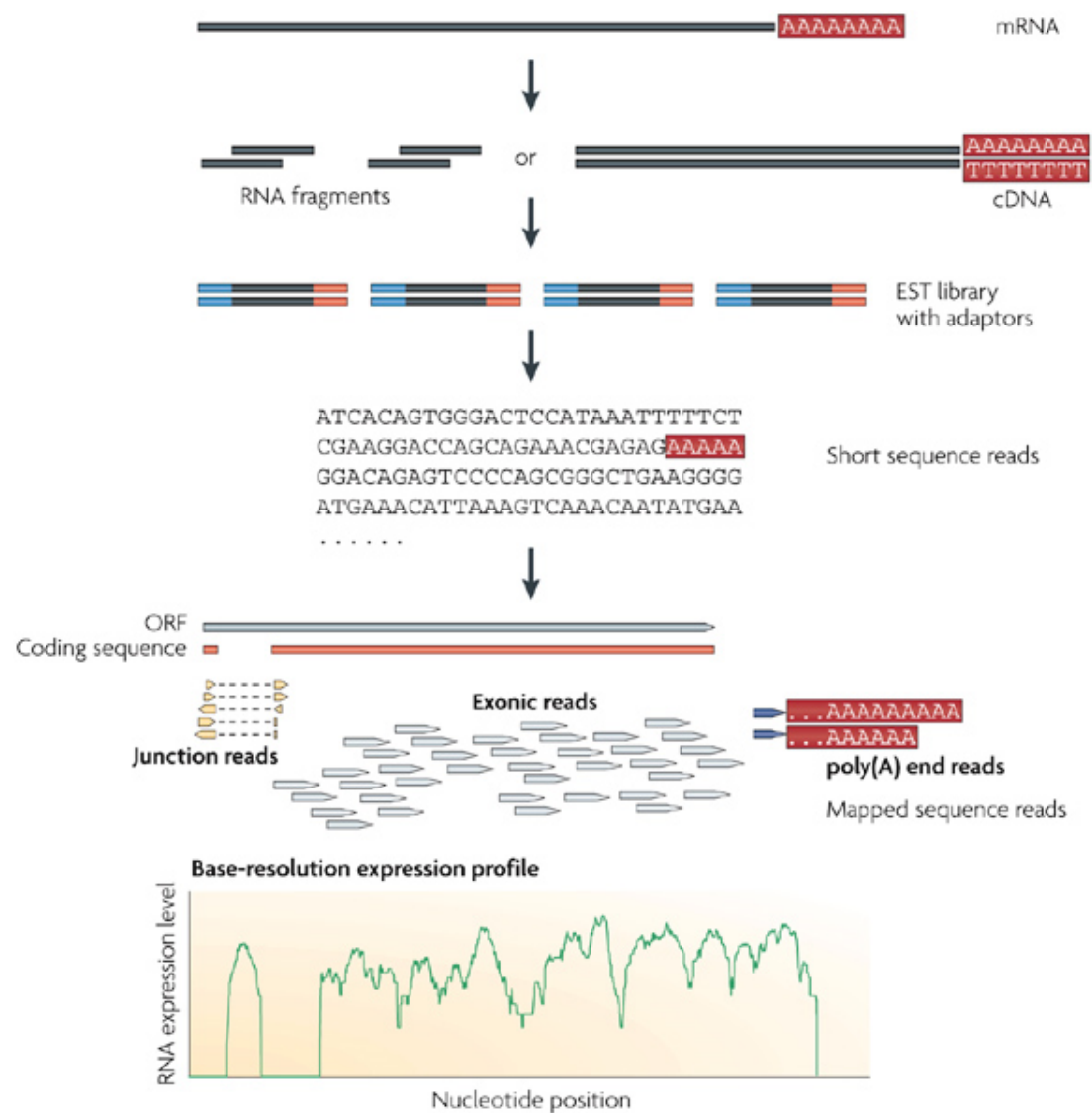


Figure 1-3 Expression tiling strategy



How RNA-seq works

Sample preparation

Next generation sequencing (NGS)

Data analysis:

- ✓ Mapping reads
- ✓ Visualization (Gbrowser)
- ✓ De novo assembly
- ✓ Quantification

FPKM (RPKM): Expression Values

- ▶ Fragments ~~Reads~~ Per Kilobase of exon model per Million mapped fragments
- ▶ Nat Methods. 2008, Mapping and quantifying mammalian transcriptomes by RNA-Seq.
Mortazavi A et al.

$$FPKM = 10^9 \times \frac{C}{NL}$$

C= the number of reads mapped onto the gene's exons

N= total number of reads in the experiment

L= the sum of the exons in base pairs.

Part I: Expression Measurement

1. MEBI
2. PDNN

References

- Cheng Li, Wing H. Wong. Model-based analysis of oligonucleotide arrays Expression index computation and outlier detection. PNAS 98:31-36, 2001.
- Zhang L, Miles MF, Aldape KD. A model of molecular interactions on short oligonucleotide microarrays. Nat Biotechnology 21(7):818-821, 2003.

Expression Measurement

- Affymetrix average approach
- Model Based Expression Index approach (Li & Wong)
- Robust Multi-Array approach (Irizarry & Speed)
- Position dependent nearest neighbor model (Zhang et al)

Data and Notation

- Probe intensity in chip i , probe j , and gene n

$$PM_{ijn}, MM_{ijn}$$

- $i=1,\dots,I$ (ranging from 1 to hundreds)
- $j=1,\dots,J$ (usually 16 or 20)
- $n=1,\dots,N$ (between 8,000 to 12,000)

Affymetrix Average Approach

- Affymetrix's Genechip@ software use Avg.diff

$$Avg.diff = \frac{1}{|A|} \sum_{j \in A} (PM_j - MM_j)$$

where A is a set of suitable pairs chosen by the software

Affymetrix's MAS5.0

- Affymetrix's new analysis approach

$$signal = TukeyBiweight\{\log(PM_j - MM_j^*)\}$$

where MM^* a version of MM that is never bigger than PM

dChip---MBEI

Model Base Expression Index

- Li-Wong full

$$\begin{cases} PM_{ij} = v_j + \alpha_j \theta_i + \phi_j \theta_i + \epsilon \\ MM_{ij} = v_j + \alpha_j \theta_i + \phi_j \theta_i + \epsilon \\ \epsilon \sim N(0, \xi^2) \end{cases}$$

with identifiability constraint

$$\sum_j \phi_j^2 = J$$

dChip---MBEI

Model-based Expression Index

- Li-Wong reduced

$$\begin{cases} y_{ij} = PM_{ij} - MM_{ij} = \phi_j \theta_i + \epsilon \\ \epsilon \sim N(0, \sigma^2), \quad \sigma^2 = 2\xi^2 \end{cases}$$

with identifiability constraint

$$\sum_j \phi_j^2 = J$$

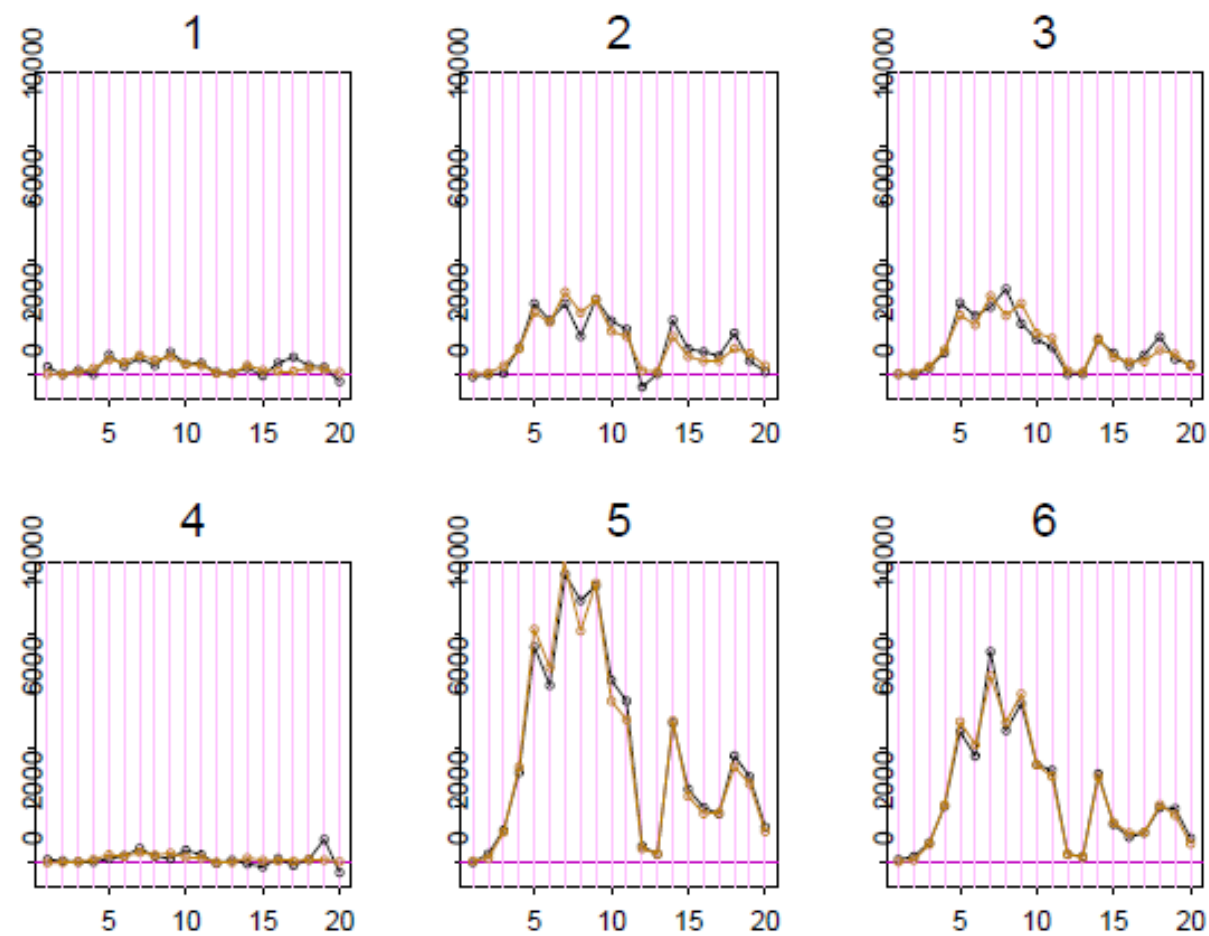


Figure 1.2. Black curves are the PM-MM difference data of gene A in the first 6 arrays. Light curves are the fitted values to model (2).

PDNN Model

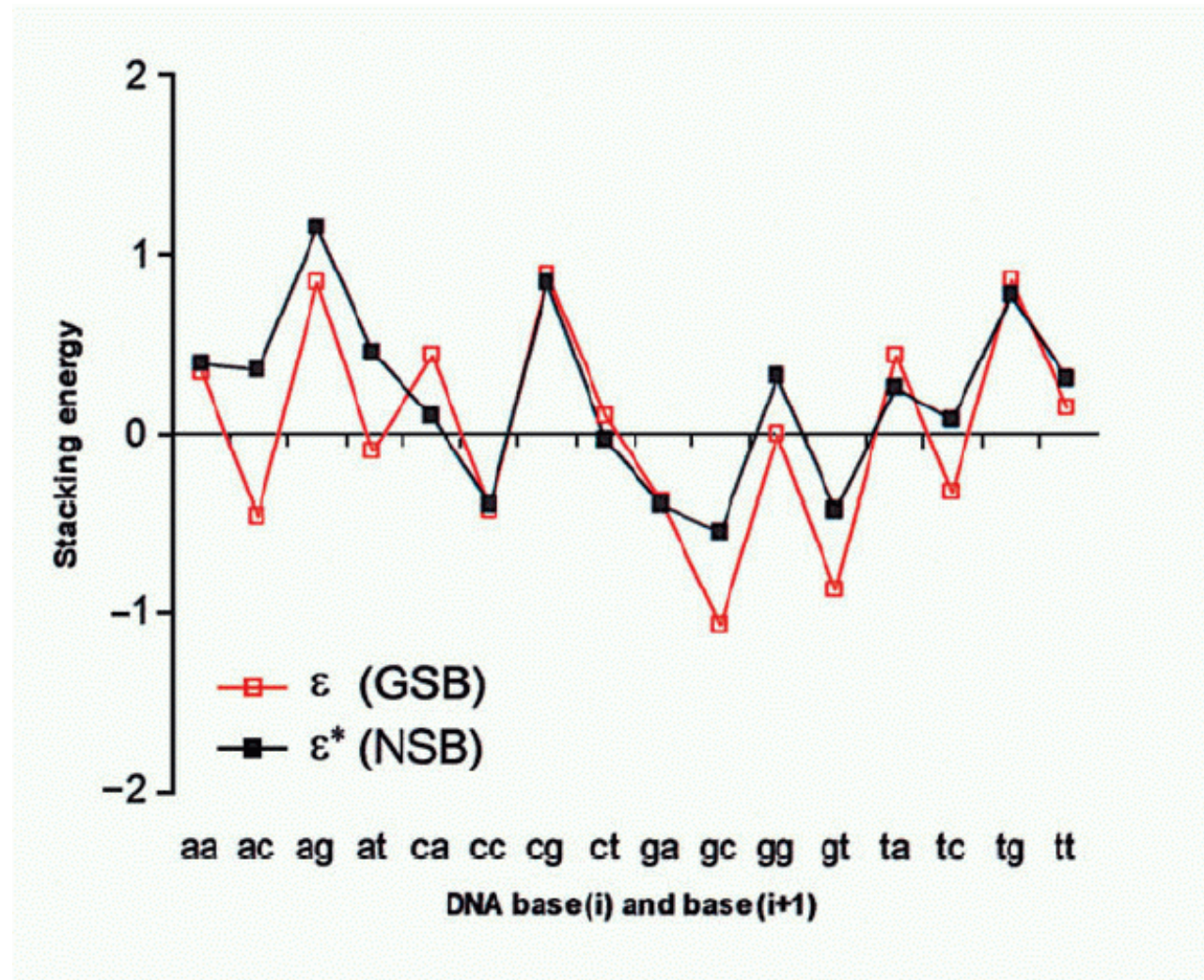
- The model is named Position Dependent Nearest Neighbour (PDNN)
- The hybridization is characterized by a energy, and the energy is approximated by pair-wised interaction between nearest neighbours.

PDNN Model

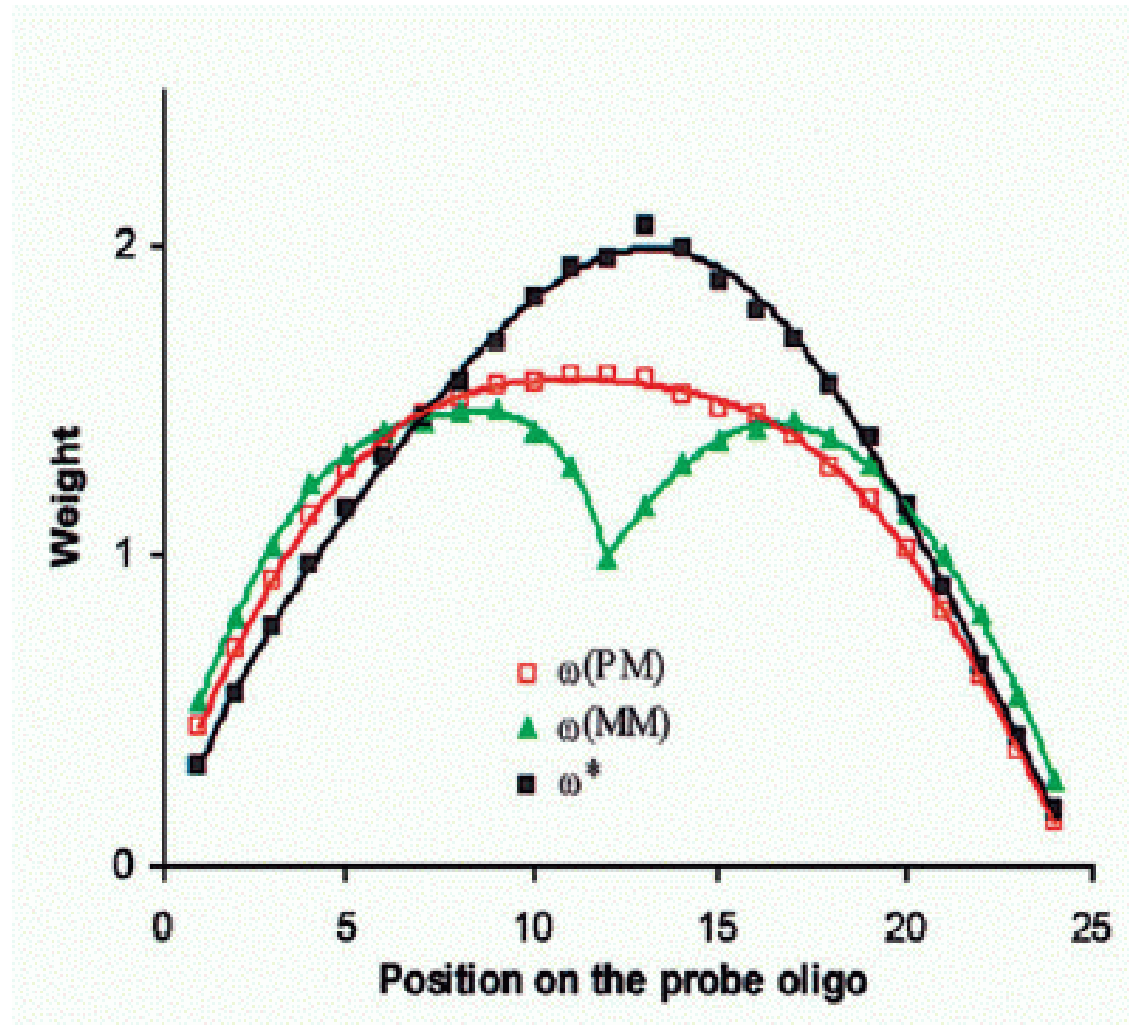
- Position Dependent Nearest Neighbor Model

$$\left\{ \begin{array}{l} E = \sum_{i=1}^{24} \omega_i \lambda(b_i, b_{i+1}) \\ I = \frac{N}{1 + \exp(E)} + b + \epsilon \end{array} \right.$$

Stack Energy



Positional Weights



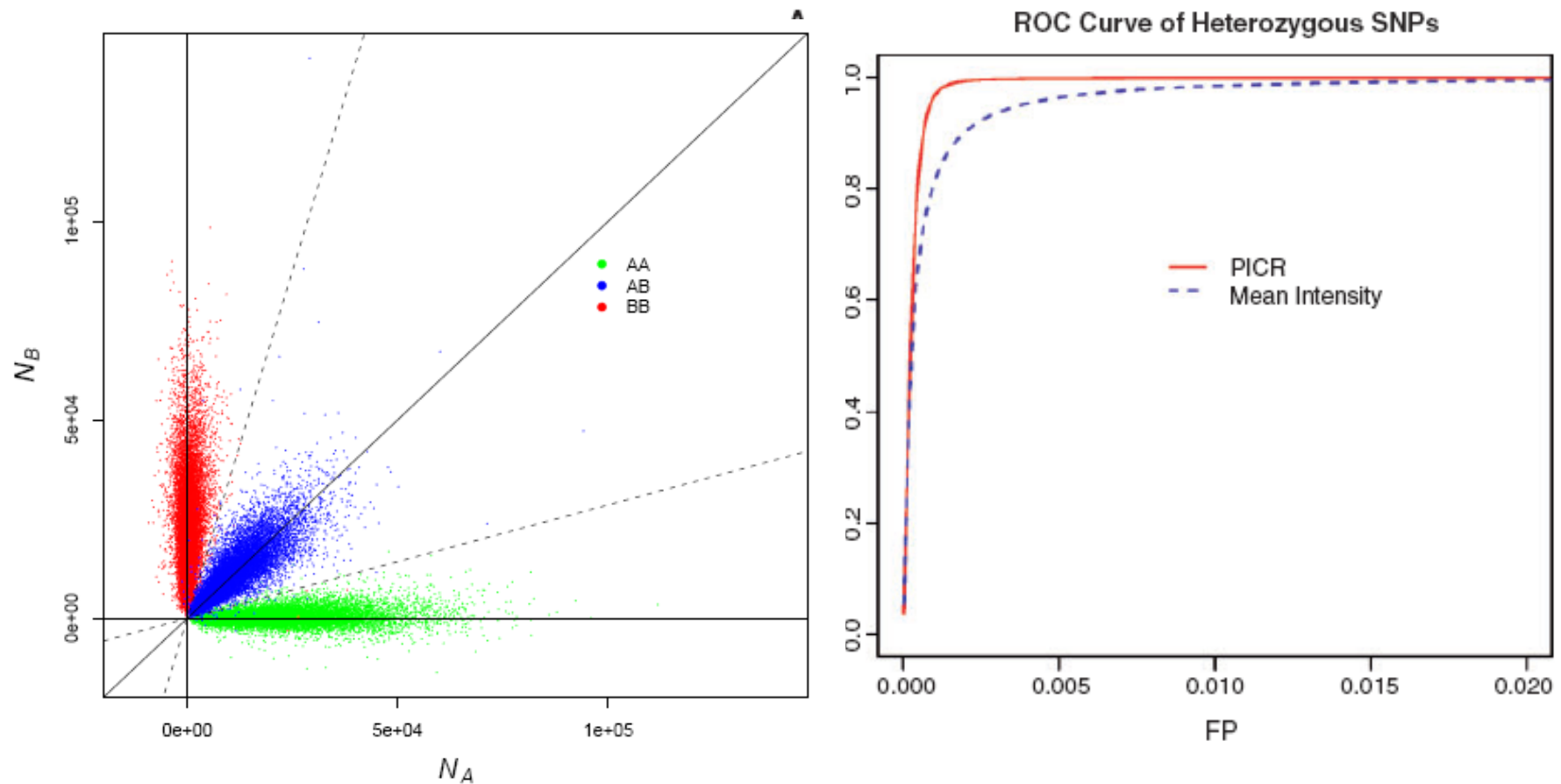
Extension : Generalized PDNN Model in SNP Array

$$\left\{ \begin{array}{l} E(S^P, S^T) = \sum_{i=1}^{24} \omega_i \lambda(S_i^P, S_{i+1}^P) \\ E_1(S^P, S^T) = \left\{ \sum_{l=1, l \neq 12+j, 13+j}^{24} \theta_l^j \lambda(S_l^P, S_{l+1}^P) \right\} \\ \quad + \kappa^j \delta(S_{12+j}^P S_{13+j}^P S_{14+j}^P, S_{12+j}^T S_{13+j}^T S_{14+j}^T) \\ E_2(S^P, S^T) = E_1^{S^P, S^T} + \xi^j (S_{12+j}^P S_{13+j}^P S_{14+j}^P, S_{12+j}^T S_{13+j}^T S_{14+j}^T) \end{array} \right.$$

Extension: Generalized PDNN Model in SNP Array

$$\left\{ \begin{array}{l} \phi(x) = \frac{1}{1 + \exp(x)} \\ I_{PA,ks} = N_A \phi(E(S^{PA,ks}, S^{TA})) + N_B \phi(E_1(S^{PA,ks}, S^{TB})) \\ \quad + b_{PA,ks} + \epsilon_{PA,ks} \\ I_{PB,ks} = N_A \phi(E_1(S^{PB,ks}, S^{TA})) + N_B \phi(E(S^{PB,ks}, S^{TB})) \\ \quad + b_{PB,ks} + \epsilon_{PB,ks} \\ I_{MA,ks} = N_A \phi(E_1(S^{MA,ks}, S^{TA})) + N_B \phi(E_{t_k}(S^{MA,ks}, S^{TB})) \\ \quad + b_{MA,ks} + \epsilon_{MA,ks} \\ I_{MB,ks} = N_A \phi(E_{t_k}(S^{MB,ks}, S^{TA})) + N_B \phi(E_1(S^{MB,ks}, S^{TB})) \\ \quad + b_{MB,ks} + \epsilon_{MB,ks} \end{array} \right.$$

More Accurate Genotyping



Wan L, Sun KL, Ding Q, Cui YH, Li M, Wen YL, Elston R, Qian MP and Fu WJ. Hybridization modeling of oligonucleotide SNP arrays for accurate DNA copy number estimation. Nucleic Acids Research, 37(17):e117.(2009)