



Design and exploitation of a versatile Arabidopsis whole-Genome Tiling Array

Tiling-array data: analysis and visualization

Caroline Bérard¹, Sandra Derozier², Sandrine Balzergue²,
Tristan Mary-Huard¹, Francois Roudier³, Stéphane Robin¹,
Alain Lecharny², Vincent Colot³, Michel Caboche²,
Sébastien Aubourg² and Marie-Laure Martin-Magniette^{1,2}

¹ UMR AgroParisTech/INRA MIA 518, Equipe statistique et génome, Paris.

² Unité de Recherche en Génomique Végétale (URGV), Evry.

³ Institut de Biologie de l'ENS, Equipe Epigénétique et Epigenomique Végétale.

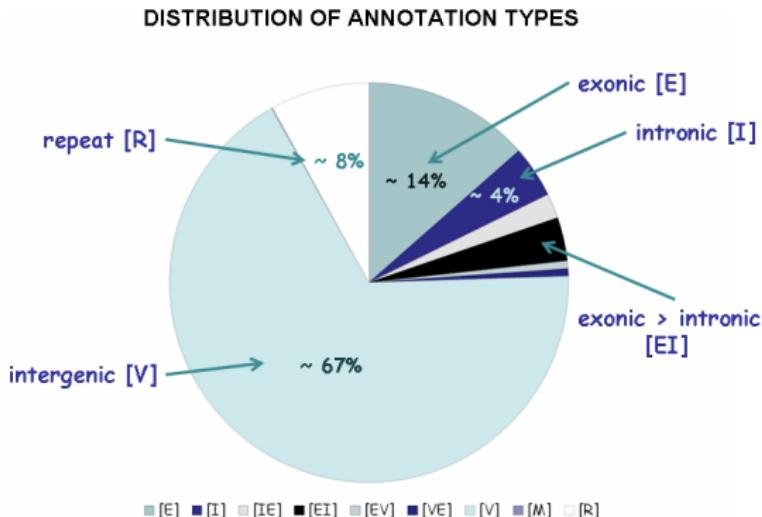
ANR/Genoplante *Tiling Array Genome (TAG) project*

Goal of the ANR Genoplante TAG Project

- Design of a *tiling-array* covering the *Arabidopsis thaliana* whole genome.
- Different types of application
 - ▶ **Transcriptome:** detection of transcripts, conditions of gene expression
 - ▶ **ChIP-chip:** study control mechanism of gene expression
(DNA methylation, histone modifications, transcription factor)
→ Development of adapted statistical methods
- Visualization of probe features and integration of the statistical results in the FLAGdb++ environment.

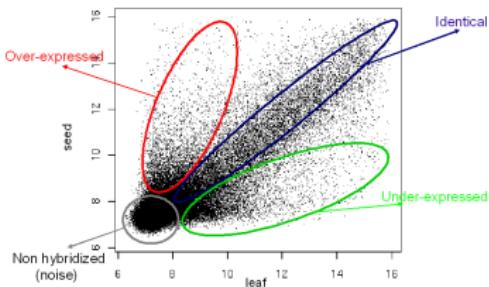
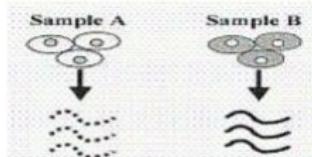
Tiling-array features

- 717.246 probes per strand, Resolution of 160 bp.
- NimbleGen probes
 - ▶ Variable length (between 50 and 75 nt)
 - ▶ Constant TM $\sim 76^\circ$



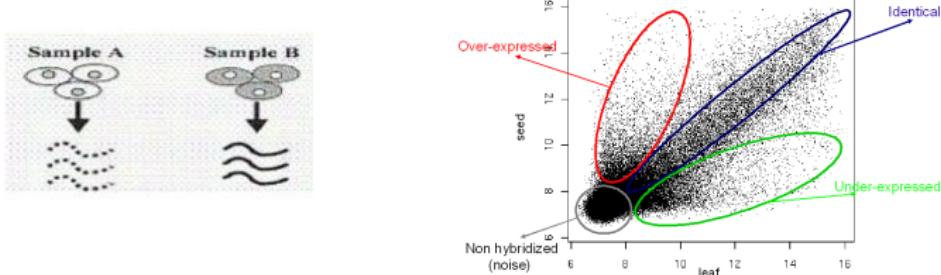
An unsupervised Classification Problem

- **Transcriptome:** noise, identical, over or under expressed (4 groups)

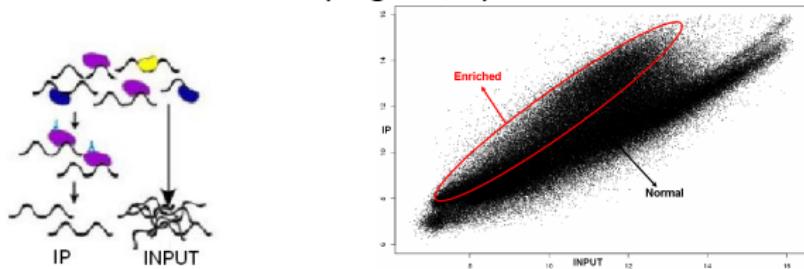


An unsupervised Classification Problem

- **Transcriptome:** noise, identical, over or under expressed (4 groups)

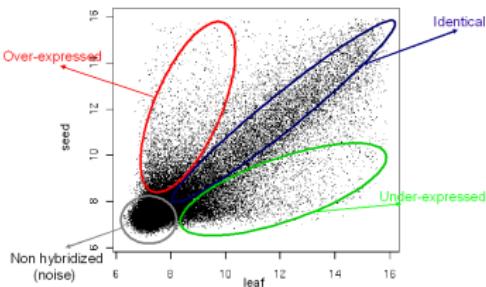
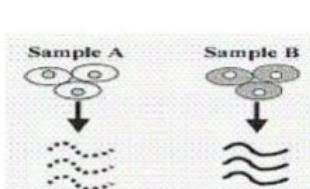


- **ChIP-chip:** enriched, normal (2 groups)

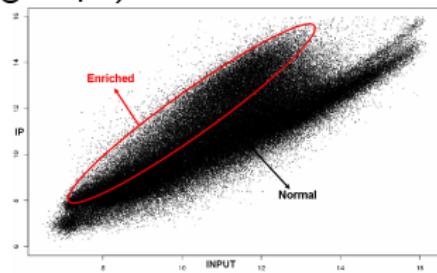
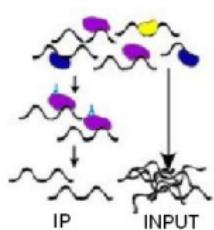


An unsupervised Classification Problem

- **Transcriptome:** noise, identical, over or under expressed (4 groups)



- **ChIP-chip:** enriched, normal (2 groups)



→ Find the status of the probe

Available informations

- Position of the probes along the genome $\rightsquigarrow t$
- Structural annotation $\rightsquigarrow C_t$

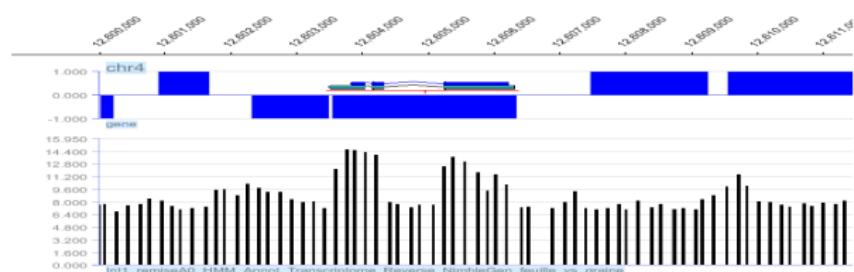


Available informations

- Position of the probes along the genome $\rightsquigarrow t$
- Structural annotation $\rightsquigarrow C_t$



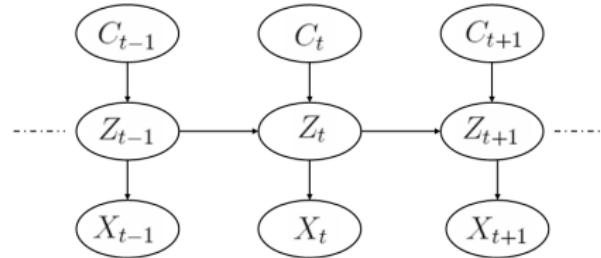
- Visualization of the signal intensity



→ Dependence between neighboring probes

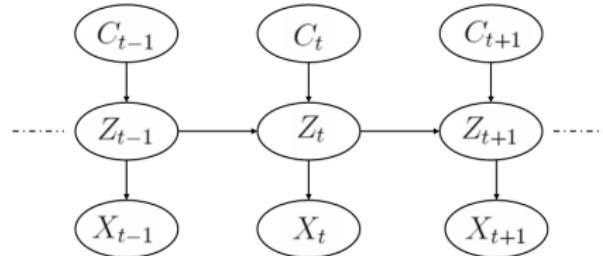
General Modeling

$X = \text{Data}$, $Z = \text{unknown status}$, $C = \text{annotation}$, $t = \text{position}$



General Modeling

$X = \text{Data}$, $Z = \text{unknown status}$, $C = \text{annotation}$, $t = \text{position}$



- **Transcriptome:** $X_t = (\text{Int}_1, \text{Int}_2)$
 - ▶ **Symmetrical:** Bidimensionnal Gaussian mixture
- **ChIP-chip:** $X_t = (IP, Input)$
 - ▶ **Non symmetrical:** Mixture of regressions

Transcriptome data: Model with HMM and Annotation

- Model

- ▶ C_t = annotation of the probe t (intron, exon, intergenic, ...)
- ▶ Z_t (status of the probe) \sim Markov chain
- ▶ $P(Z_t = l | Z_{t-1} = k, C_t = p) = \pi_{kl}^{C_t}$ → one transition matrix for each annotation category

Transcriptome data: Model with HMM and Annotation

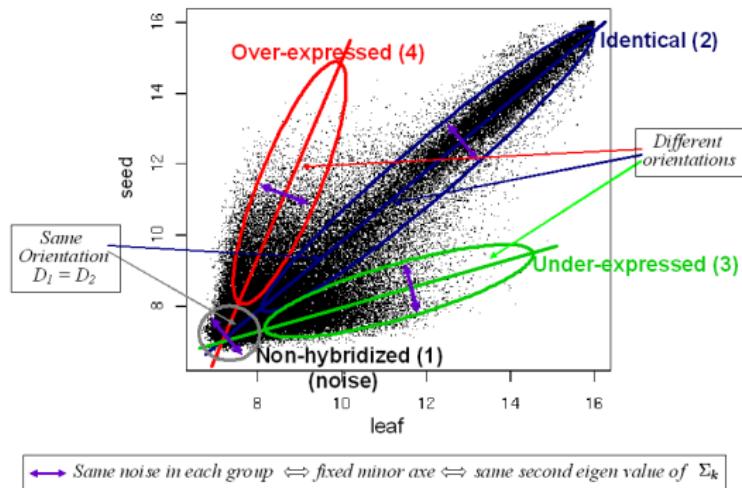
- Model

- ▶ C_t = annotation of the probe t (intron, exon, intergenic, ...)
- ▶ Z_t (status of the probe) \sim Markov chain
- ▶ $P(Z_t = l | Z_{t-1} = k, C_t = p) = \pi_{kl}^{C_t}$ → one transition matrix for each annotation category
- ▶ X : observed signal ($\text{Int}_1, \text{Int}_2$)
- ▶ $(X_t | Z_t = k) \sim \mathcal{N}(\mu_k, \Sigma_k)$
 - μ = mean vector \rightsquigarrow center of the ellipse
 - Σ = variance matrix \rightsquigarrow shape, orientation, volume of the ellipse
- ▶ **Bidimensional Gaussian Mixture:**

$$\begin{aligned} f(x, \psi) &= p_0 \phi(x | \mu_0, \Sigma_0) + p_1 \phi(x | \mu_1, \Sigma_1) \\ &\quad + p_+ \phi(x | \mu_+, \Sigma_+) + p_- \phi(x | \mu_-, \Sigma_-) \end{aligned}$$

- Assumptions on the ellipse form

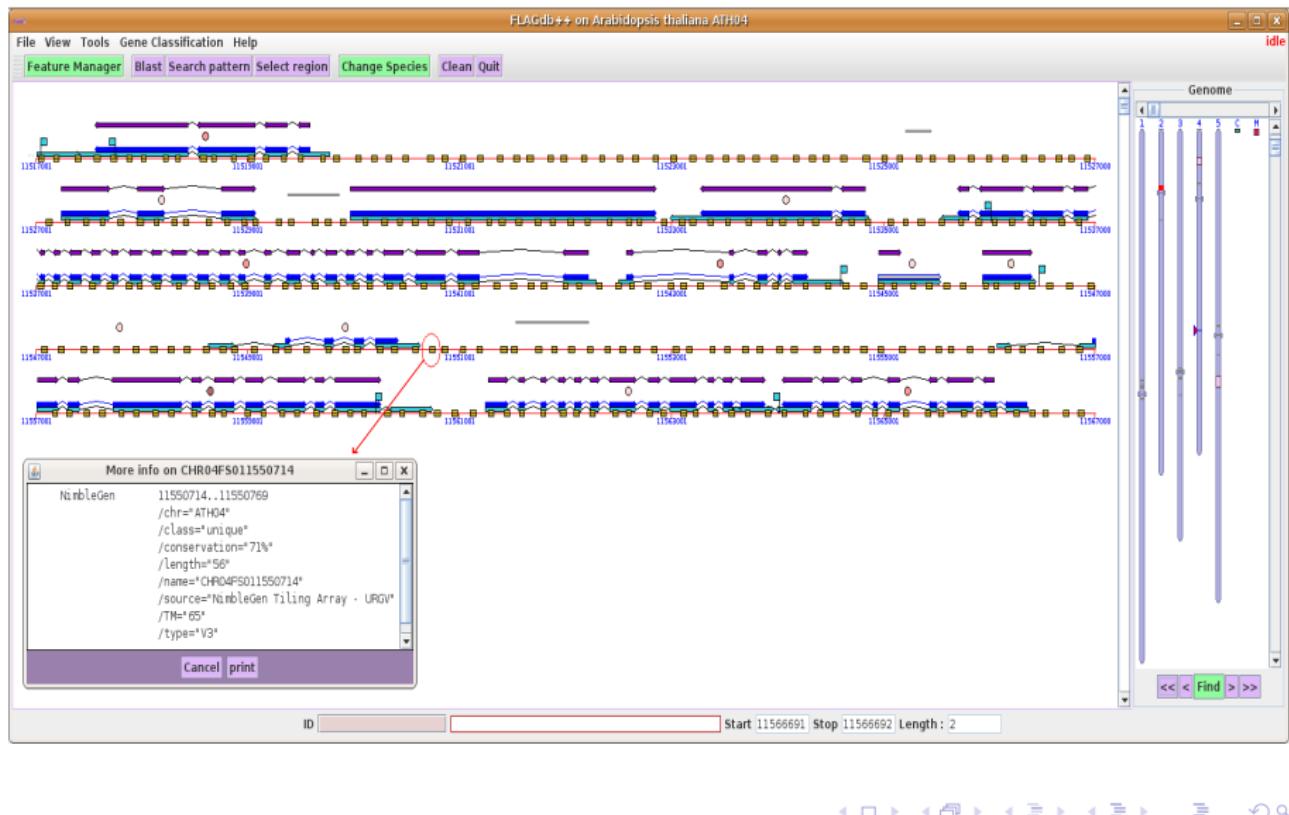
- Noise and **Identical** groups have same orientation
- Identical variance** in each group



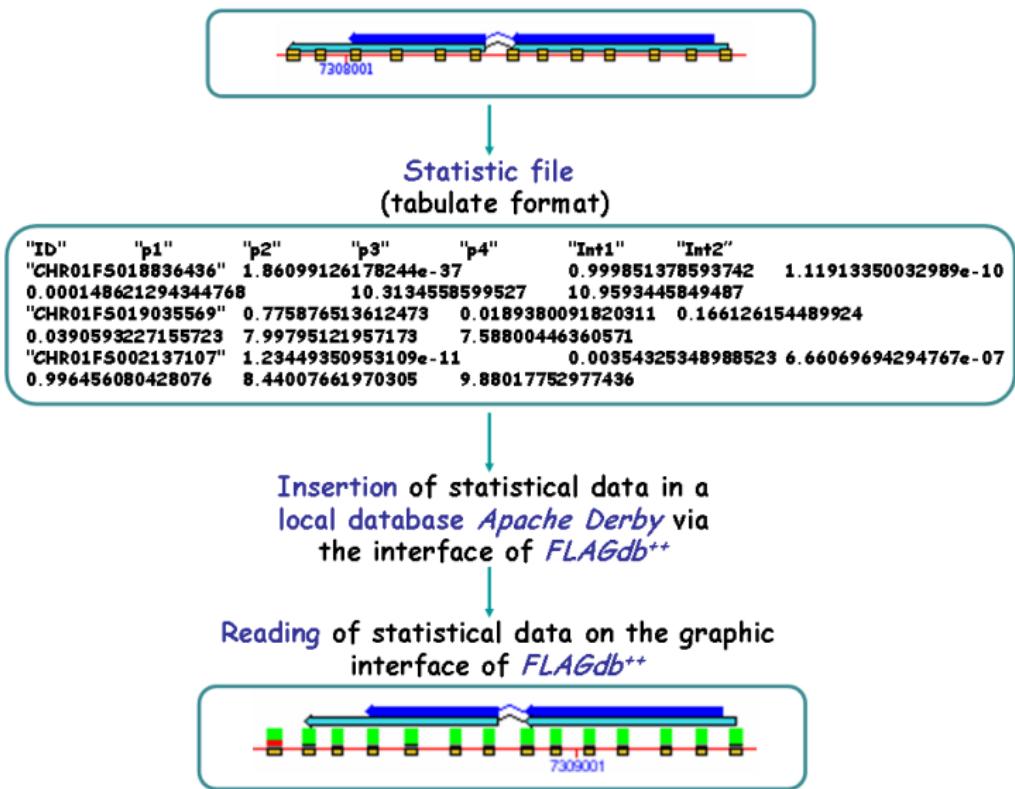
- Maximum likelihood estimation with constraints on the variance matrices using the EM algorithm

- Posterior probability: $\tau_t = \Pr\{Z_t = 1 | X\}$

View of NimbleGen probes in FLAGdb++

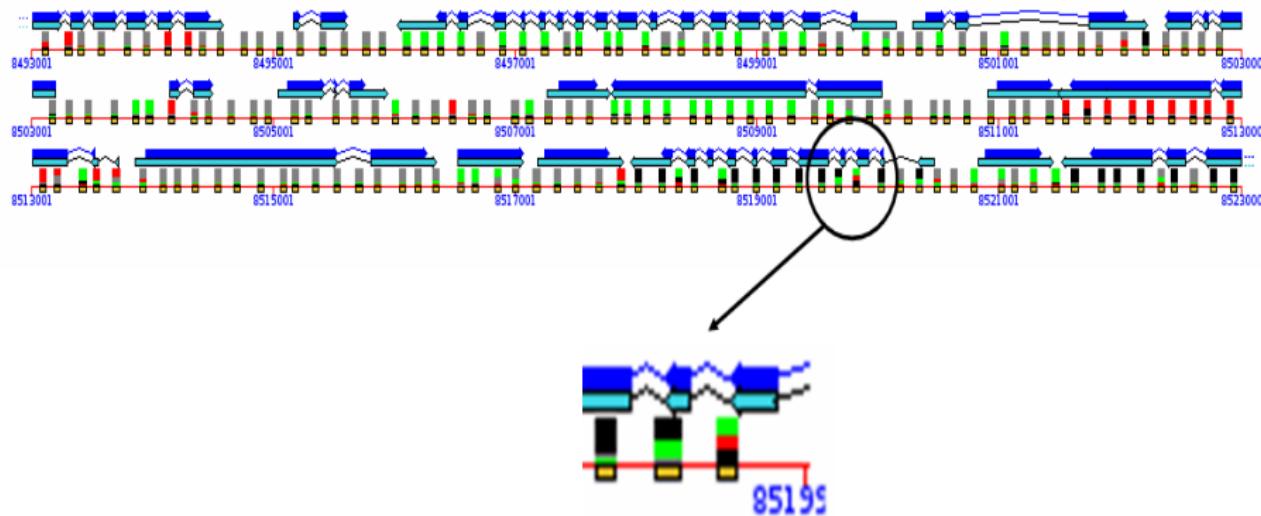


Integration of statistical results in FLAGdb++



Application on Transcriptome Data

Comparison of expression between seed 10 days after pollinisation and leaf of *Arabidopsis thaliana*



The colors of a probe depend on posterior probabilities.

Parameter estimation (in %)

Transition matrix of **intergenic** group:

	Noise	Ident.	Under-exp	Over-exp	Proportions
Noise	87	1	7	5	84
Ident.	95	3	1	1	1
Under-exp	77	1	19	3	9
Over-exp	75	2	5	18	6

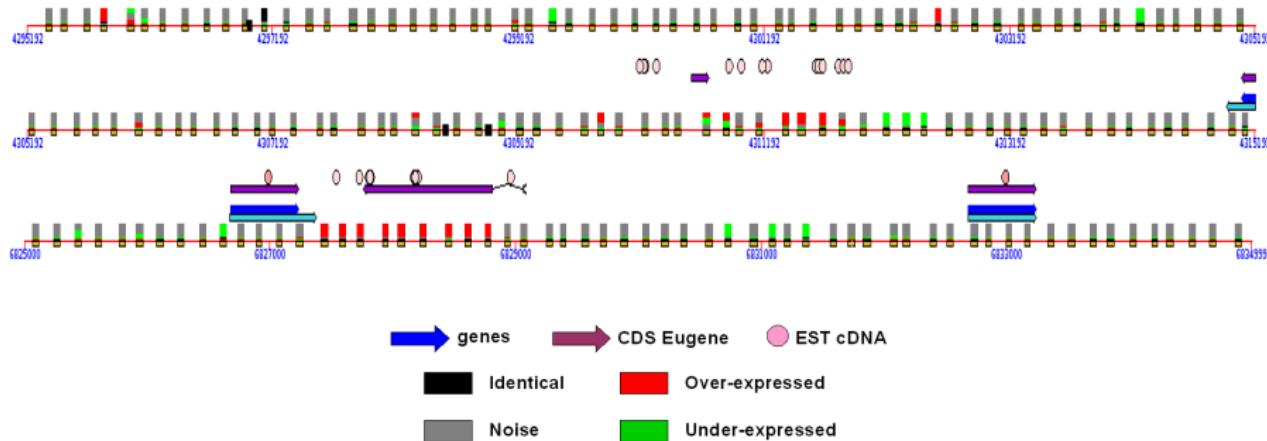
Transition matrix of **exonic** group:

	Noise	Ident.	Under-exp	Over-exp	Proportions
Noise	83	14	3	0	22
Ident.	2	90	6	2	41
Under-exp	7	5	87	1	23
Over-exp	8	6	1	85	14

Transition matrix of **intronic** group:

	Noise	Ident.	Under-exp	Over-exp	Proportions
Noise	87	2	8	3	60
Ident.	89	0	1	10	7
Under-exp	55	2	43	0	24
Over-exp	96	1	0	3	9

Genome exploration to improve structural annotation



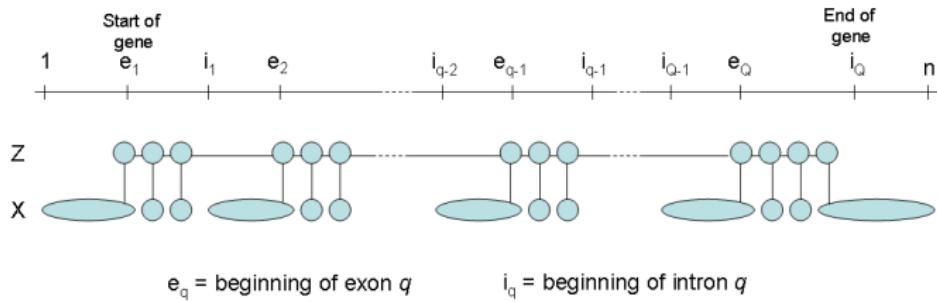
- Number of regions with expressed probes in intergenic

Nb probes	2	3	4	>5
Forward	798	167	70	88
Reverse	786	194	86	90

Probe classification → Gene classification

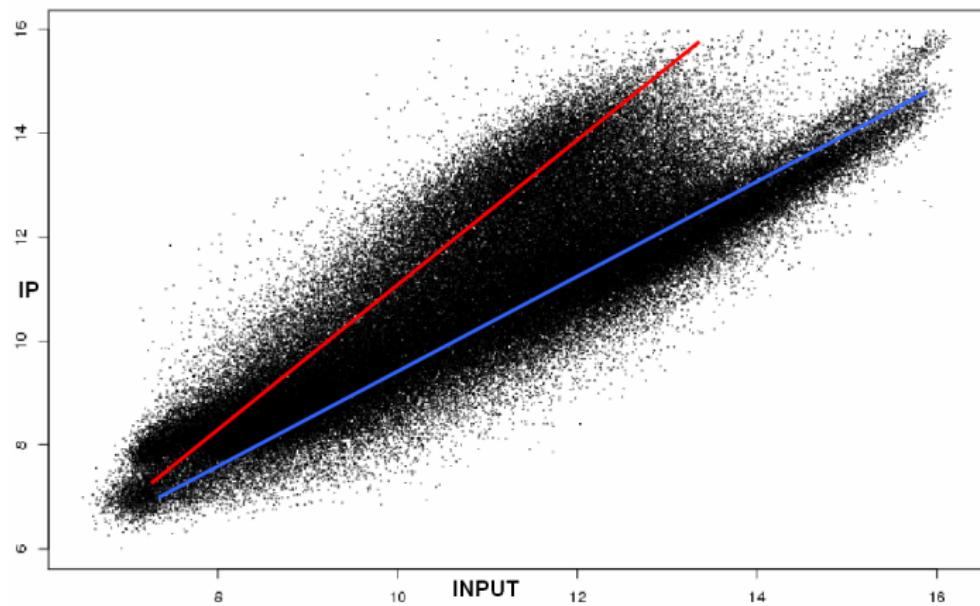
Gene = region covered by several probes \rightsquigarrow interesting biological entity

$$P(\text{for all exon probes } s \text{ for gene } g, Z_s = k | X)$$



\rightsquigarrow Calculations use the Forward step of the Forward/Backward algorithm.

ChIP-chip data



ChIP-chip data: Mixture of regressions

- Model

$$\Pr\{Z_t = 1\} = \pi, \quad \Pr\{Z_t = 0\} = 1 - \pi.$$

The relation between IP and Input depends on the status of the probe:

$$IP_t = \begin{cases} a_0 + b_0 Input_t + E_t & \text{if } Z_t = 0 \text{ (normal)} \\ a_1 + b_1 Input_t + E_i & \text{if } Z_t = 1 \text{ (enriched)} \end{cases} \quad V(IP_t) = \sigma^2$$

- Maximum likelihood estimation using the EM algorithm
- Posterior probability: $\tau_t = \Pr\{Z_t = 1|X\}$

ChIPmix: mixture model of regressions for two-color ChIP-chip analysis.

M-L. Martin-Magniette, T. Mary-Huard, C. Bérard and S. Robin. *Bioinformatics* (2008)

Classification: control of false positives

- ChIP-chip: 2 groups

Controlling false detections: We want to control the probability for the τ_t of a normal probe to fall above the classification threshold.
For a fixed risk α we calculate the threshold s such that

$$s : \quad \Pr\{\tau_t > s \mid t \text{ normal}, \log(\text{Input}) = X_t\} = \alpha$$

and if $\tau_t > s$ then the probe t is classified as 'enriched'.

The threshold s depends on both α and the log-Input X_t .
This control focuses on misclassifications in 'enriched' group.

Conclusions

- Versatile Arabidopsis whole-Genome Tiling Array.
- Protocols developed at URGV for Transcriptome, ChIP-chip and CGH.

~~> balzerg@evry.inra.fr

- General modeling of the hybridized signal.
 - ▶ Use the whole available information of the probes.
 - ▶ Comprises independent mixture (no spatial dependence) and homogeneous HMM (no annotation).
 - ▶ Can be also used with one-color *tiling arrays* to compare 2 samples.
 - ▶ Can be also used for ChIP-chip data to compare two IP samples.
 - ▶ Package R available upon request.

~~> caroline.berard@agroparistech.fr

- Visualization in FLAGdb++.

~~> <http://urgv.evry.inra.fr/FLAGdb>