

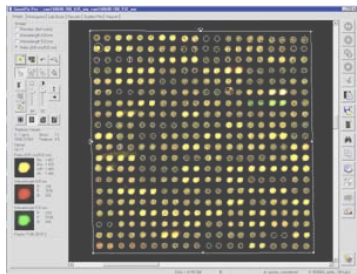
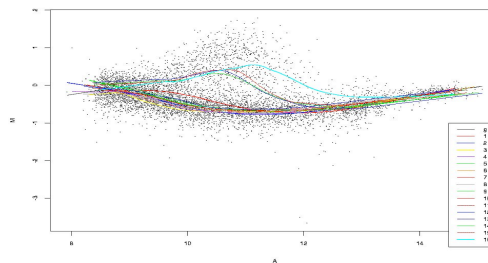
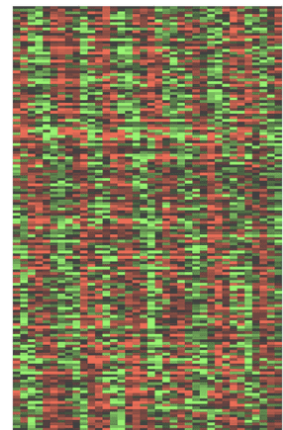
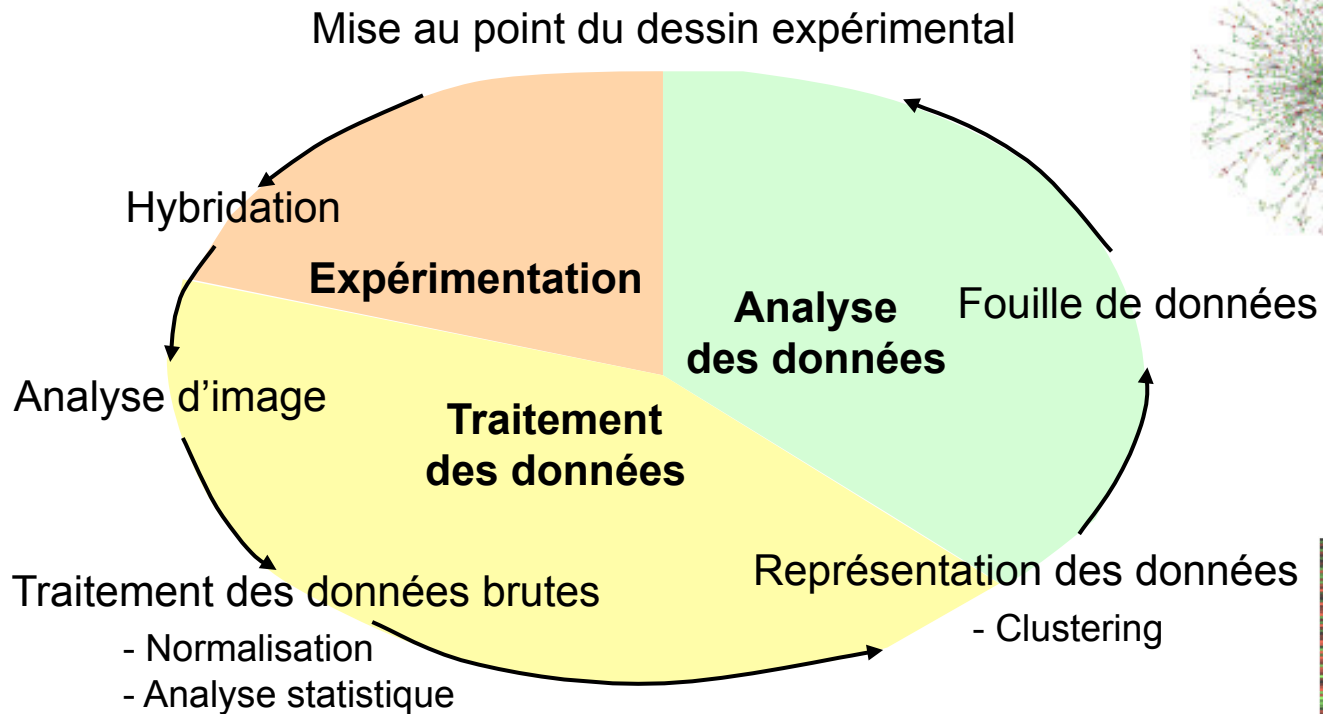
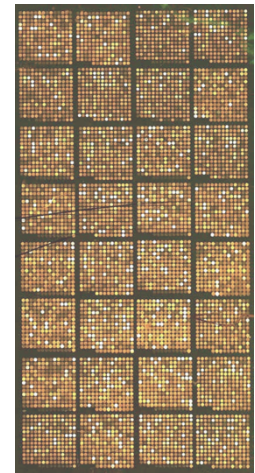
Master de génétique
UE génomique fonctionnelle
Université Denis Diderot – Janvier 2009

Normalisation

Stéphane Le Crom (lecrom@biologie.ens.fr)

Laboratoire de Génétique Moléculaire du Développement - INSERM U784
Plate-forme Transcriptome - IFR36
École normale supérieure

La bioinformatique dans une expérience biopuces

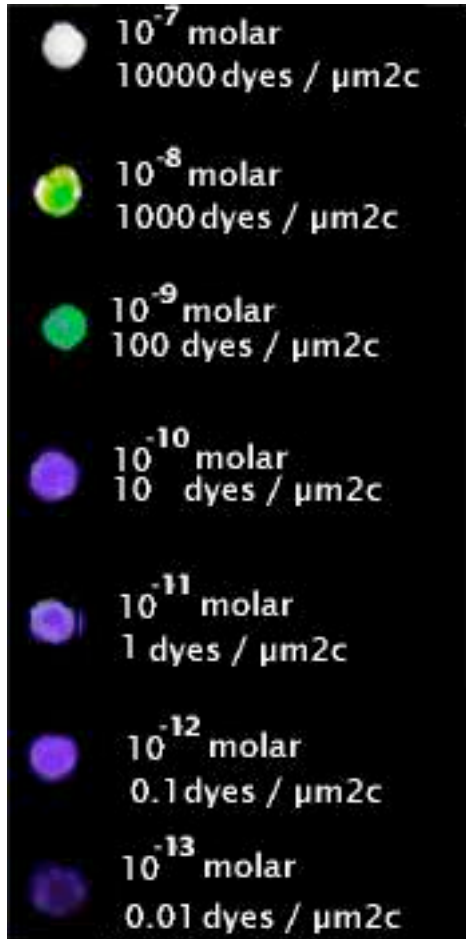


Analyse bioinformatique des puces à ADN

La segmentation du signal

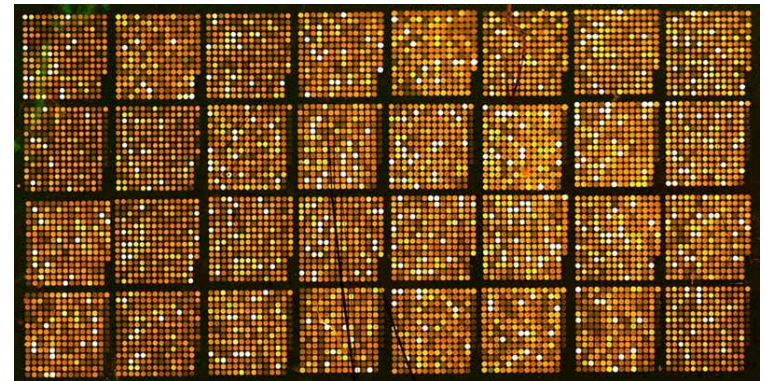
Les différents types d'images rencontrés

Echelle de couleur
=
Echelle Quantitative



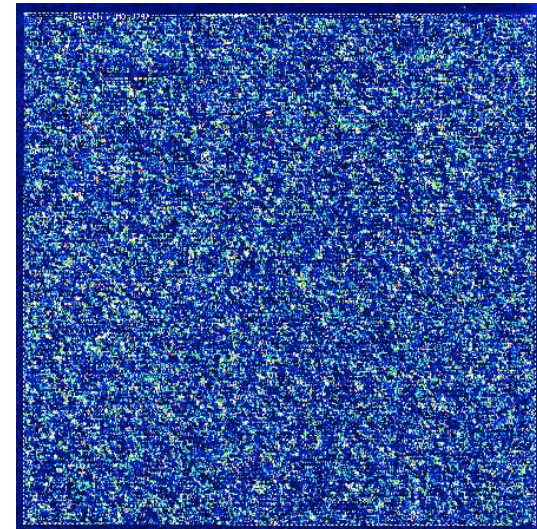
- **Puces à ADNc/oligonucleotides - ENS**

- 2 canaux
- Superposition (= Ratio)

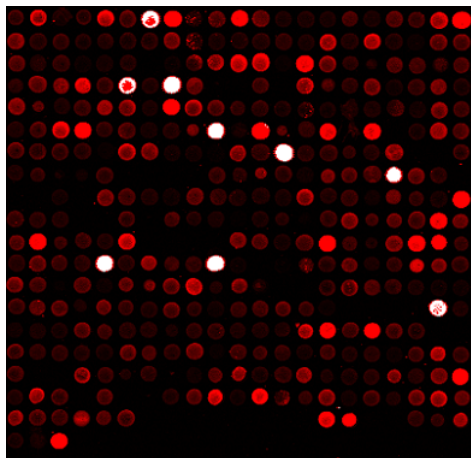


- **Puces à Oligonucleotides - Affymetrix**

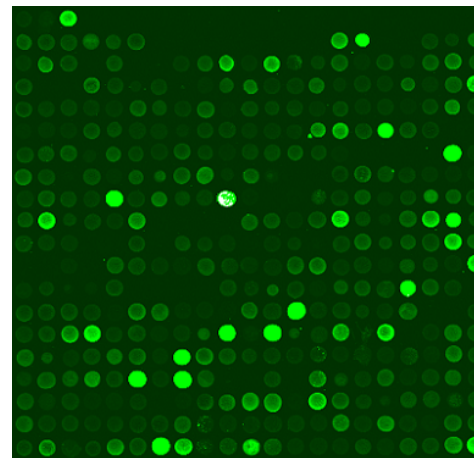
- 1 canal
- Intensité (= quantité d'ARN)



Obtention de l'image



Longueur d'onde
du Cy5



Longueur d'onde
du Cy3

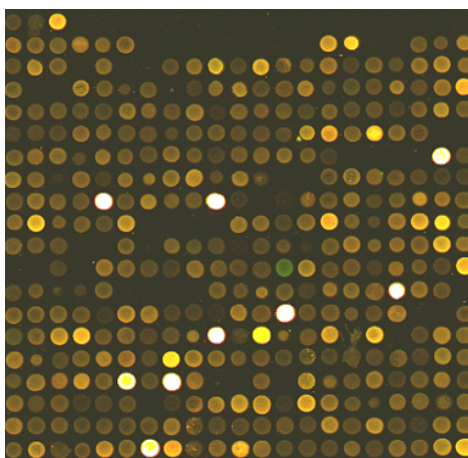
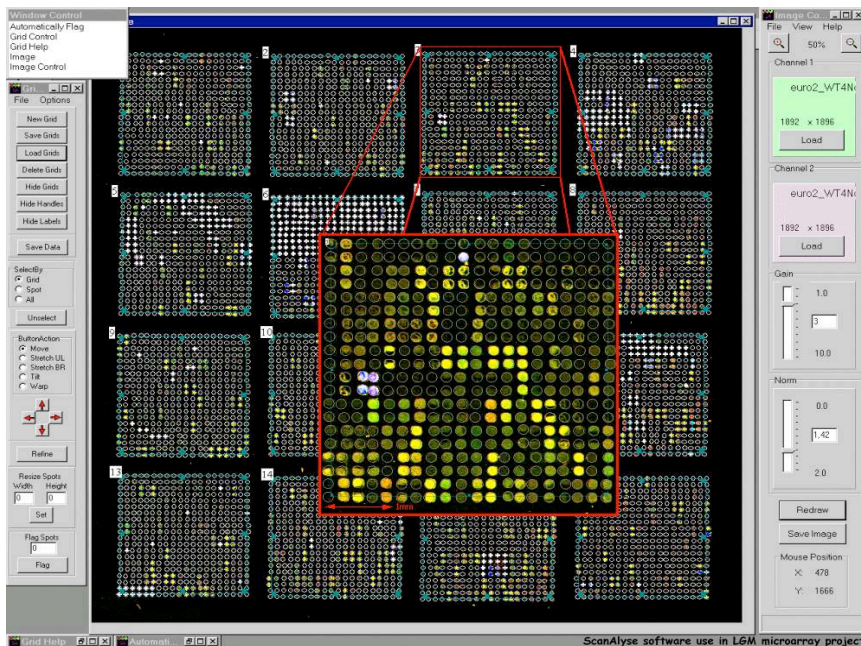


Image finale

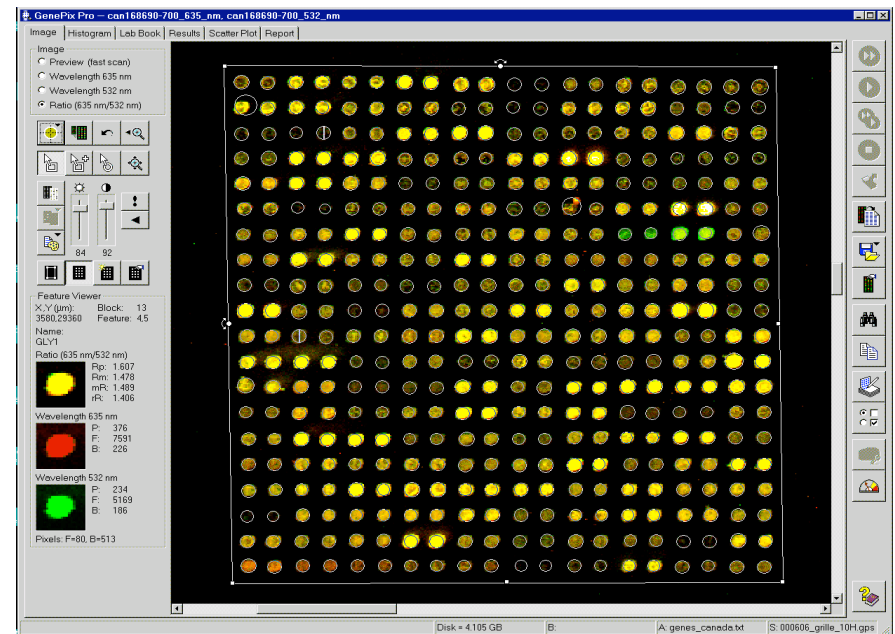
Principes généraux de l'analyse d'image

- Convertir l'image en valeurs numériques pour quantifier l'expression

Il existe des logiciels d'analyse d'image ...



ScanAlyze
(M. Eisen Stanford University)



Genepix Pro
(Axon software)

Les différentes étapes de l'analyse d'image

1 — Localisation des spots sur la lame

Pour chaque spot

2 — Délimitation des pixels correspondant à la zone d'hybridation

3 — Délimitation des pixels pour l'estimation du bruit de fond

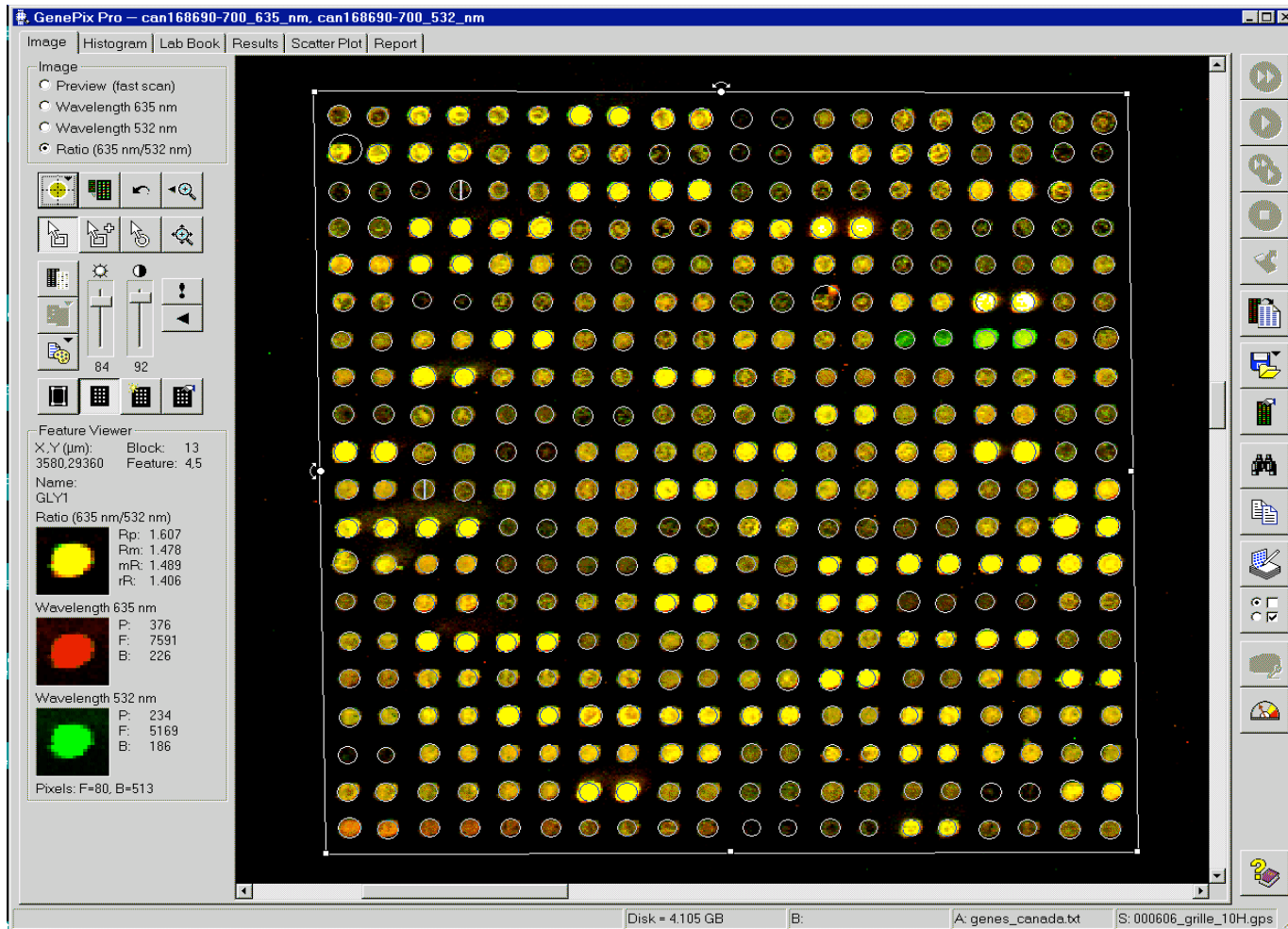
4 — Calcul de l'intensité globale de fluorescence

Sur l'ensemble de la lame

5 — Identification des spots déformés par des artéfacts

Localisation des spots sur la lame

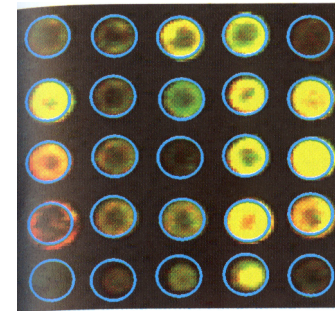
- Il est important d'assigner à chaque spot un identifiant correct !



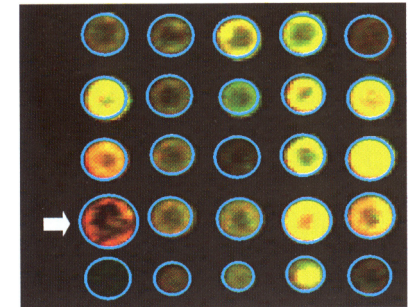
Délimitation des pixels de la zone d'hybridation

- **Différentes méthodes sont possibles (segmentation)**

- Cercles à diamètres fixes
- Cercles à diamètres variables
- Histogrammes
- Formes « adaptée »



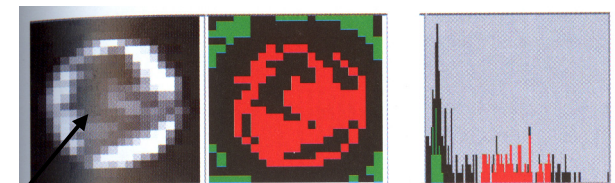
Cercles à diamètre fixe



Cercles à diamètre variable

- **Plusieurs logiciels sont disponibles**

- GenePix Pro
- ScanAlyze
- QuantArray
- ImaGene
- Dapple



Intensité variable

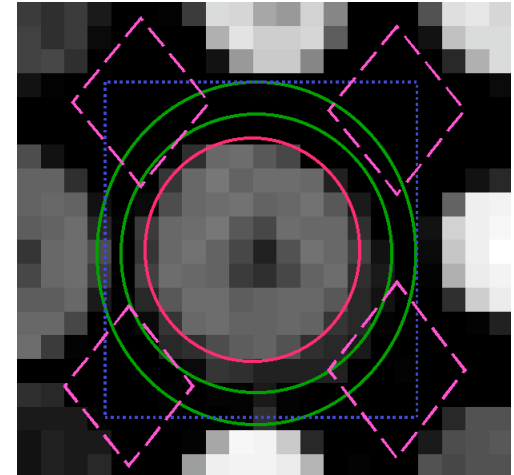
Méthode par histogramme

Délimitation des pixels du bruit de fond

Les signaux observés ont deux
composantes

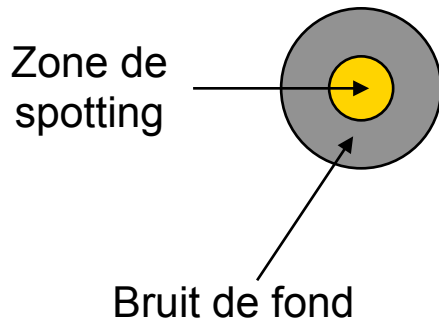
Fluorescence issue d'une
hybridation spécifique

Fluorescence non spécifique :
Bruit de fond

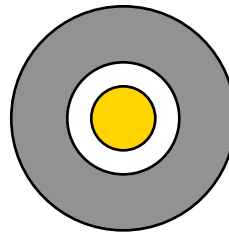


- **Analyse des pixels localisés à proximité de la région spottée**

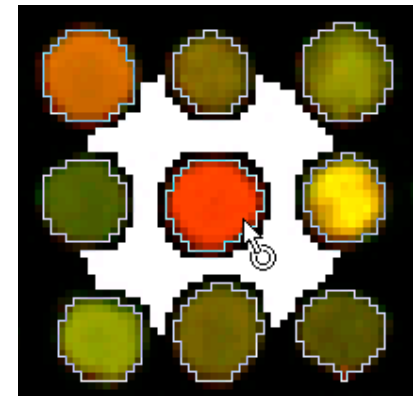
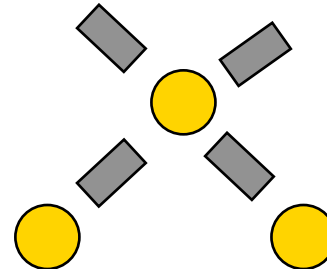
Méthode
« ScanAlyze »



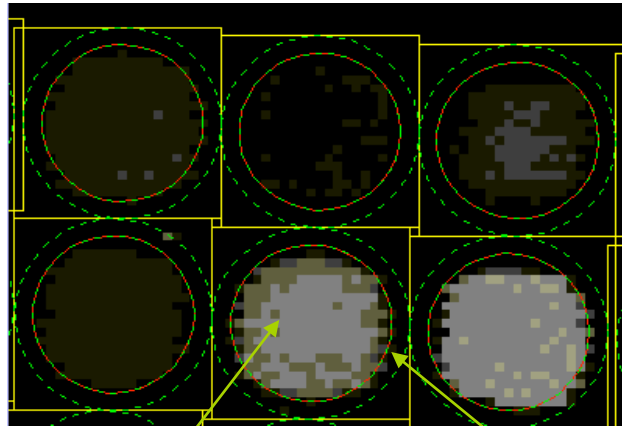
Méthode
« ImaGene »



Méthode
« GenePix »



Calcul de l'intensité globale de fluorescence



$$\text{Intensité net} = \text{Intensité brute} - \text{Bruit de fond}$$

Intensité à l'intérieur du spot

Mesure du bruit de fond

Intensité par pixel

Moyenne ou médiane ?



Intensité globale du spot

Critère qualité : taille du spot, déviation standard...

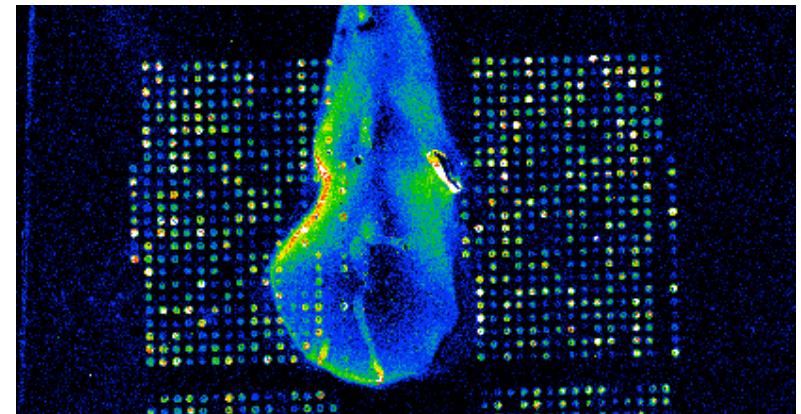
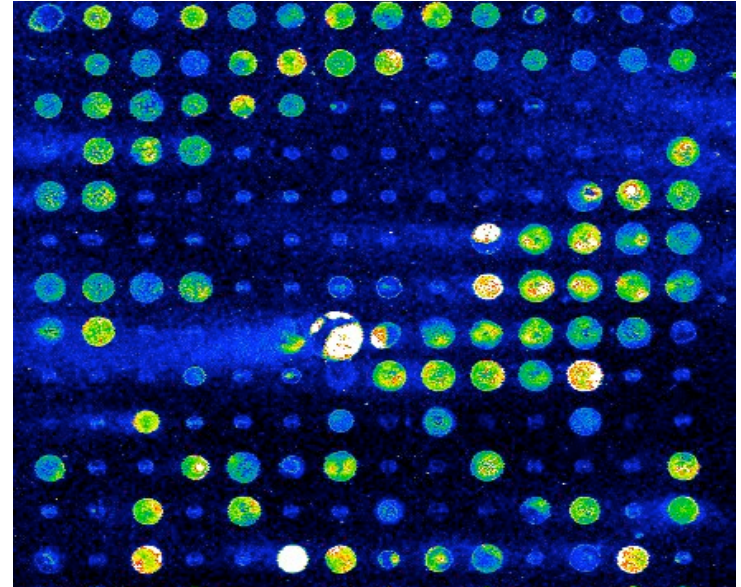
Annotation des spots non conformes

- **Exemples**

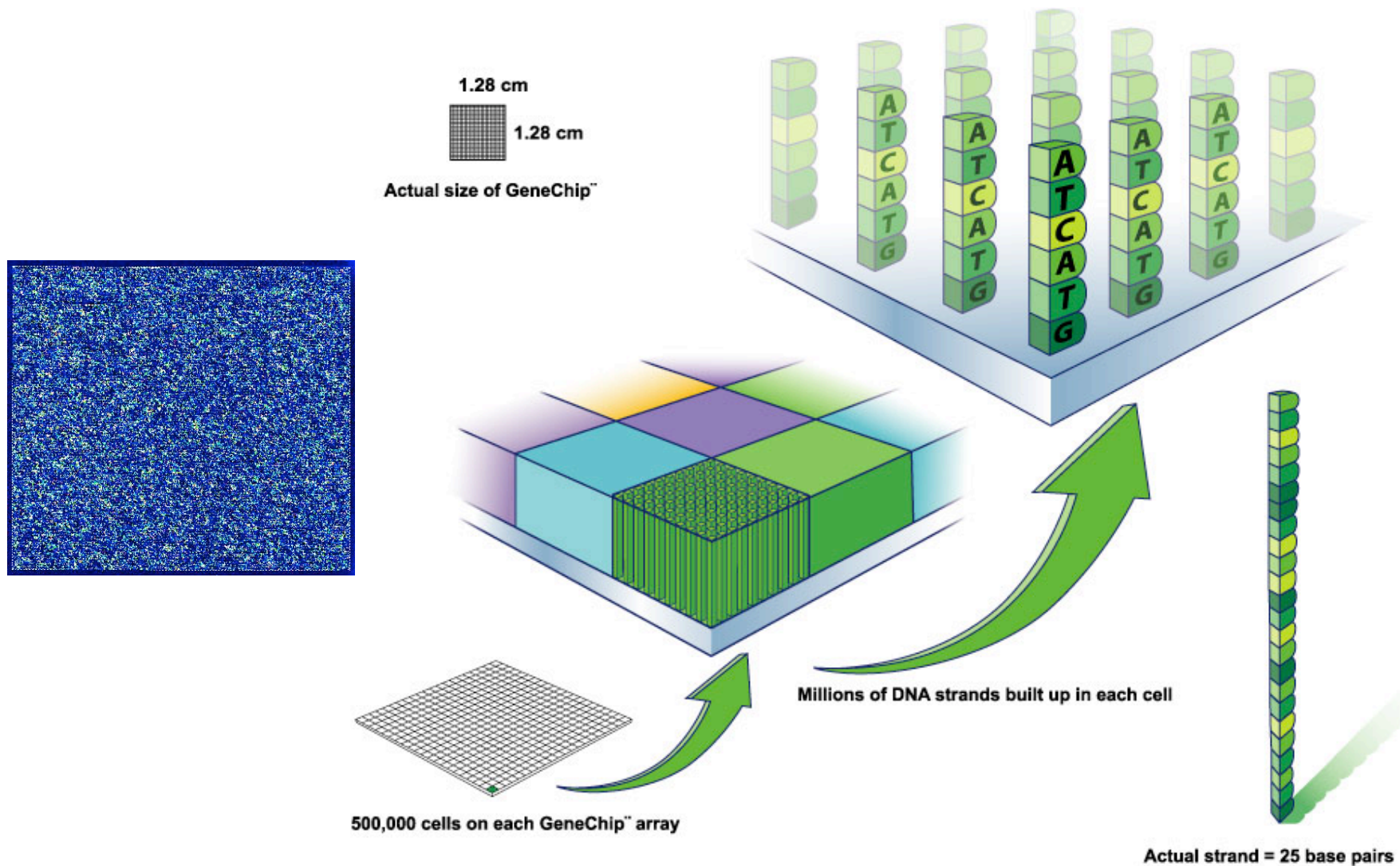
- Effet « comète »
- Tâches d'hybridation
- Problèmes de lavages

- **Solutions**

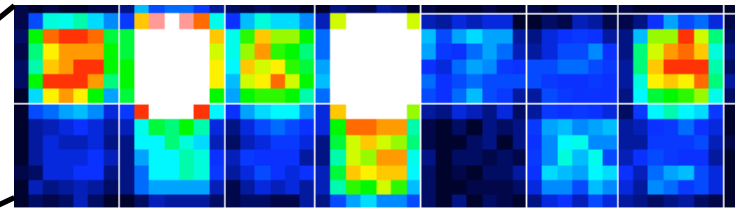
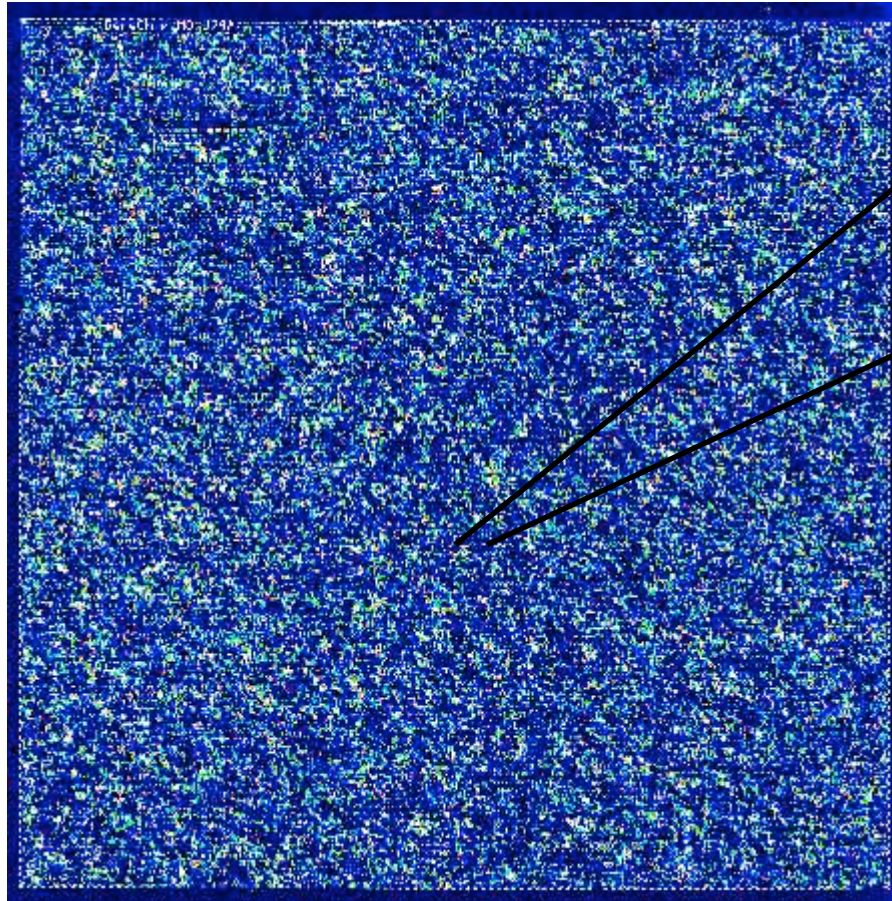
- Toujours regarder les contrôles qualité des lots de lames commandés
- Se référer aux aides en ligne
- Éliminer les spots artéfactuels à l'aide d'un filtre manuel ou automatique



Les puces Affymetrix (GeneChip)



Obtention de l'image (GeneChip Affymetrix)



PM

MM

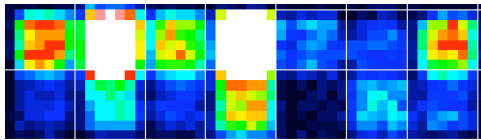
Des paires d'oligonucléotides sont créées pour chaque gènes :

- "perfect match" PM
- "mismatch" MM

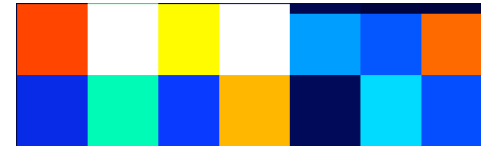
Sélection des spots corrects

- Estimation des valeurs moyennes

Image obtenue pour un gène

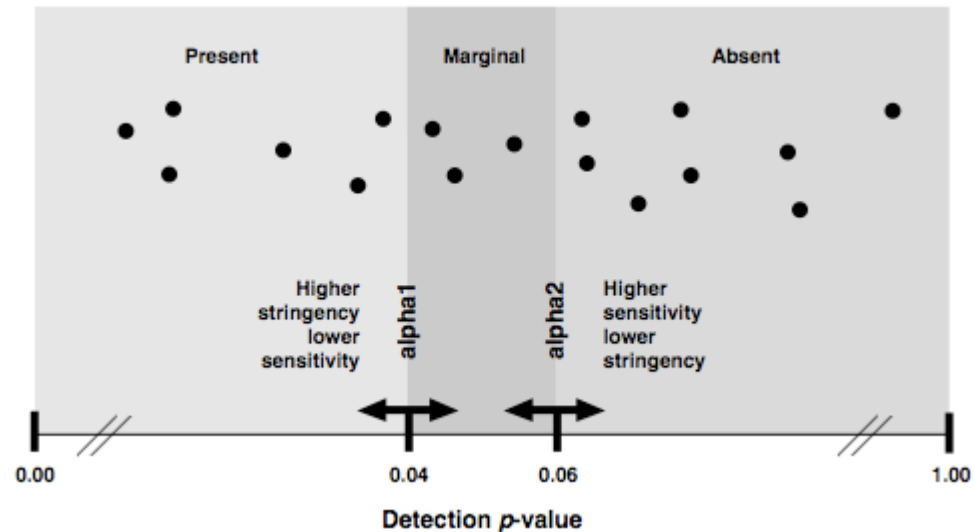
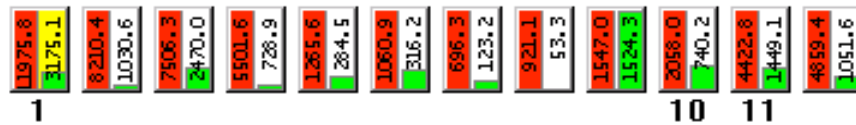


Valeurs moyennes calculées pour chaque cellule



96616_at (C3MG)
Intensity: 38 - 2174

■ Match
 ■ Mismatch
 ■ Image Mask
 ■ Probe Mask
 ■ Not in Average



- Calcul du nombre de paires positives et négatives (matrice de décision)
- Détermination du statut des gènes (présent, marginal ou absent)

Les données brutes en sortie d'analyse

- **Différents formats de fichiers**

- Le fichier d'expérience : .EXP
- Le fichier d'image équivalent aux fichiers TIFF de GenePix : .DAT
- Les données d'intensité pour chaque sonde sans traitement obtenues à partir du fichier .DAT : .CEL
- Les données traitées (par le programme d'Affymetrix, GCOS) et rassemblées par probeset : .CHP



Analyse bioinformatique des puces à ADN

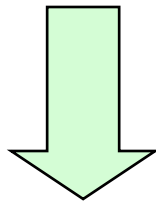
La normalisation d'une lame

Les sources de variations d'une puce à ADN

- Quantité d'ADN
- Efficacité des étapes expérimentales (extraction, RT, ...)
- Les fluorochromes (Cy3 brille plus que Cy5, accrochages ≠)

Erreur systématique

- Effets similaires sur plusieurs mesures
- Les corrections peuvent se mesurer à partir des données

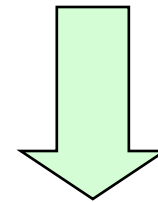


Calibration

- Rendement de la PCR
- Qualité des ADN déposés et du spotting
- Effets de cross-hybridation et hybridation non spécifique

Erreur stochastique

- Effets qui se produisent de façon trop aléatoire et qui du coup ne peuvent se mesurer comme du bruit

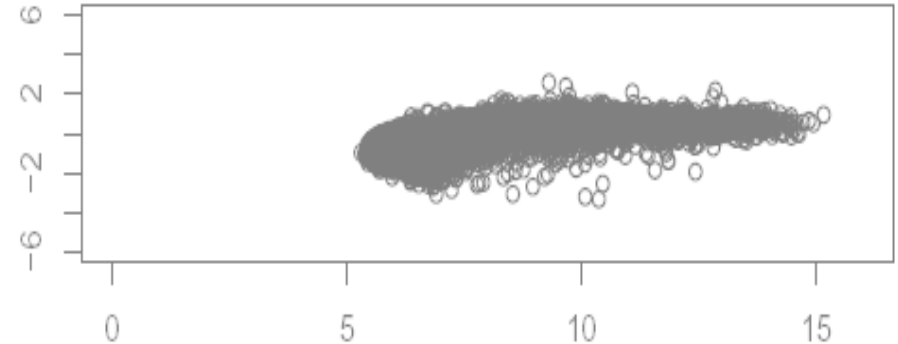


Modèle d'erreur

1^{ère} étape : le nettoyage des données

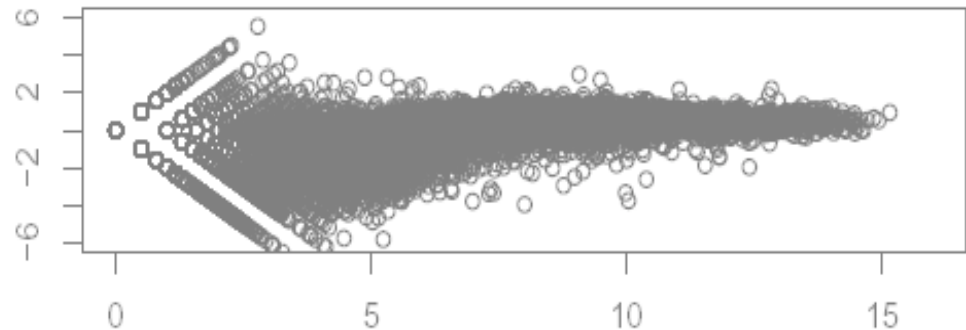
1/ Élimination des gènes dont l'annotation est différente de 0

- Éventuellement retourner voir l'image originale)



2/ Filtrage sur les intensités

- Saturation du scanner
- Écart avec le bruit de fond trop faible



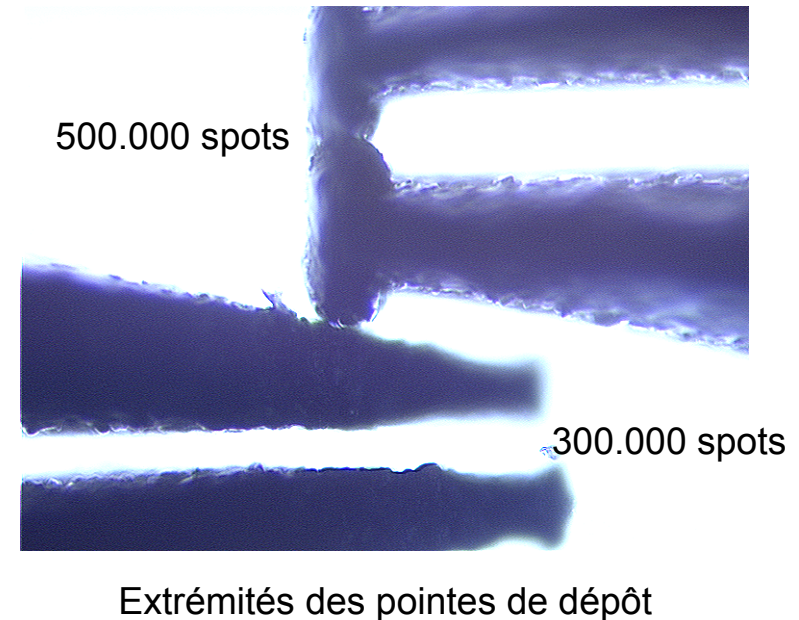
La normalisation

- **Pourquoi faire ?**

- Pour corriger les différences systématiques entre les mesures sur la même lame qui ne représentent pas de véritables variations biologiques.

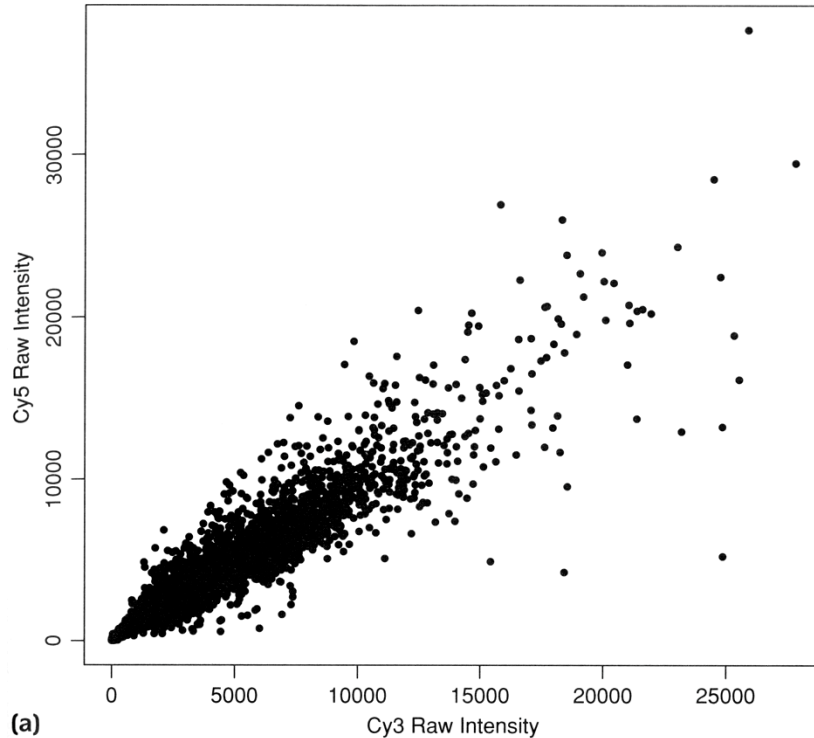
- **Pourquoi normaliser est nécessaire ?**

- En examinant les réplicats contre le même échantillons, où de vraies différences d'expression ne doivent pas apparaître.

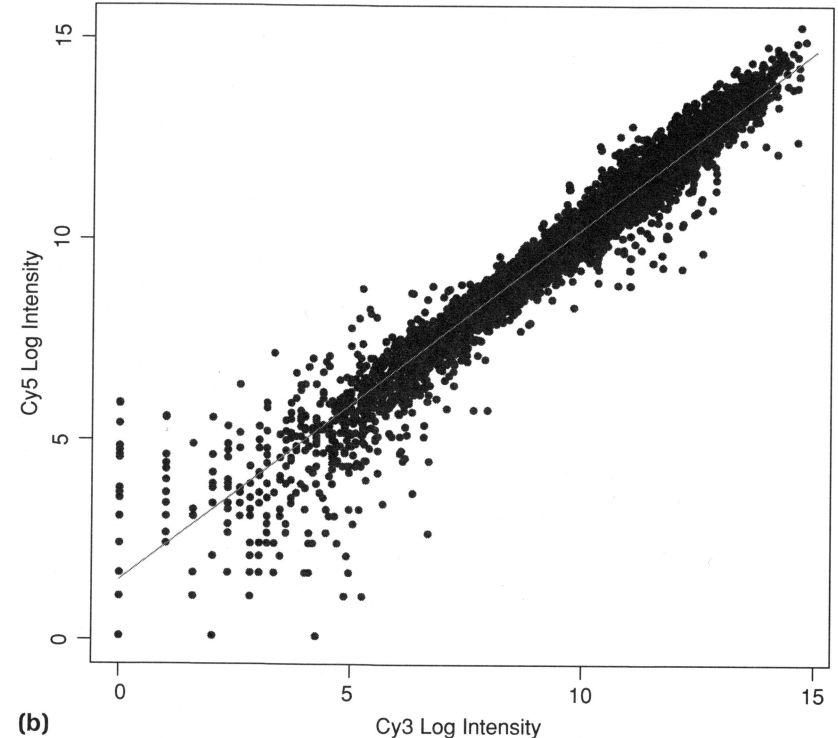


- La normalisation des résultats des puces à ADN permet la comparaison de plusieurs expériences (référence commune)
- La normalisation calibre les erreurs systématiques (et non stochastique)
- Il est indispensable d'effectuer deux transformations mathématiques sur les résultats bruts avant de normaliser les données

La transformation logarithmique



Les effectifs sont plus importants vers les faibles intensités



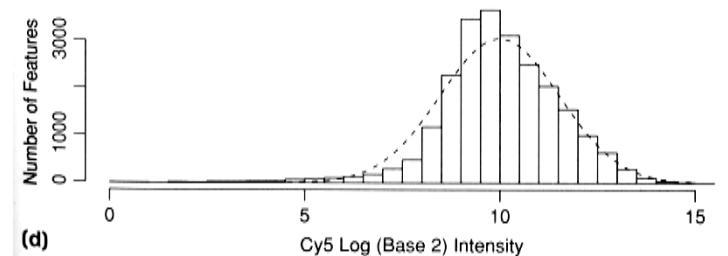
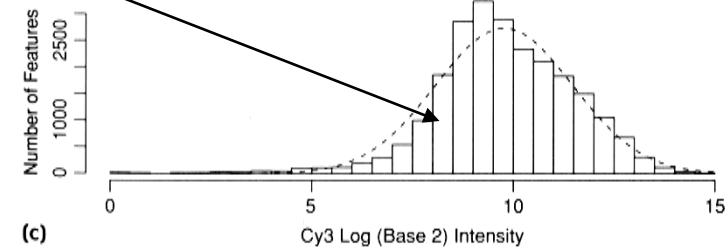
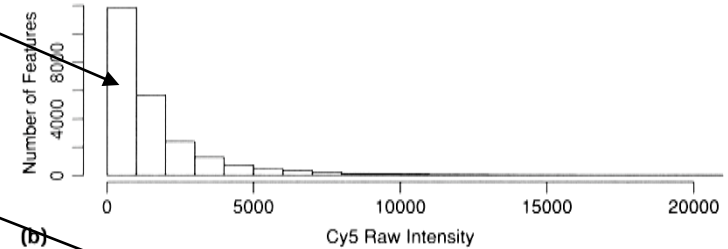
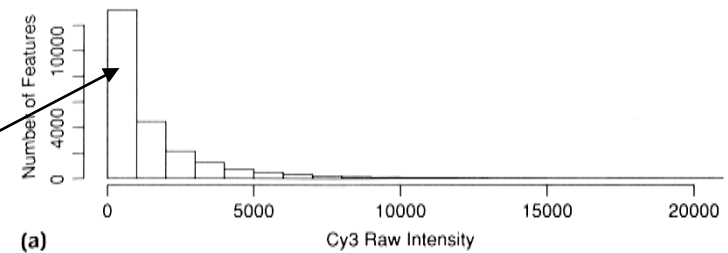
Les intensités sont distribuées de façon uniforme

Remarque : On choisira le logarithme en base 2 pour les analyses

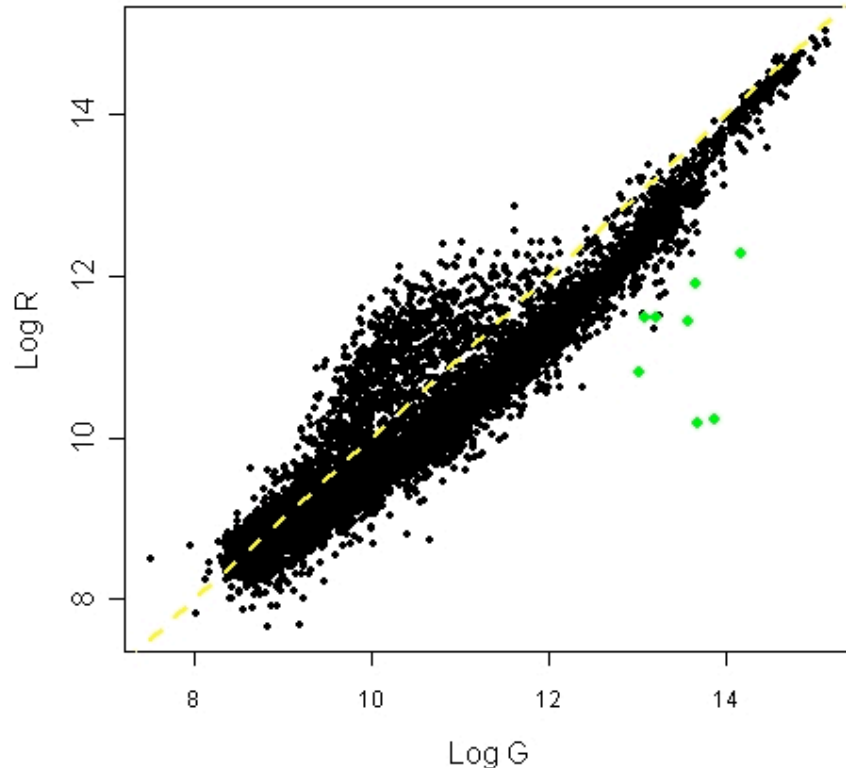
La transformation logarithmique

• Effet sur la distribution des intensités

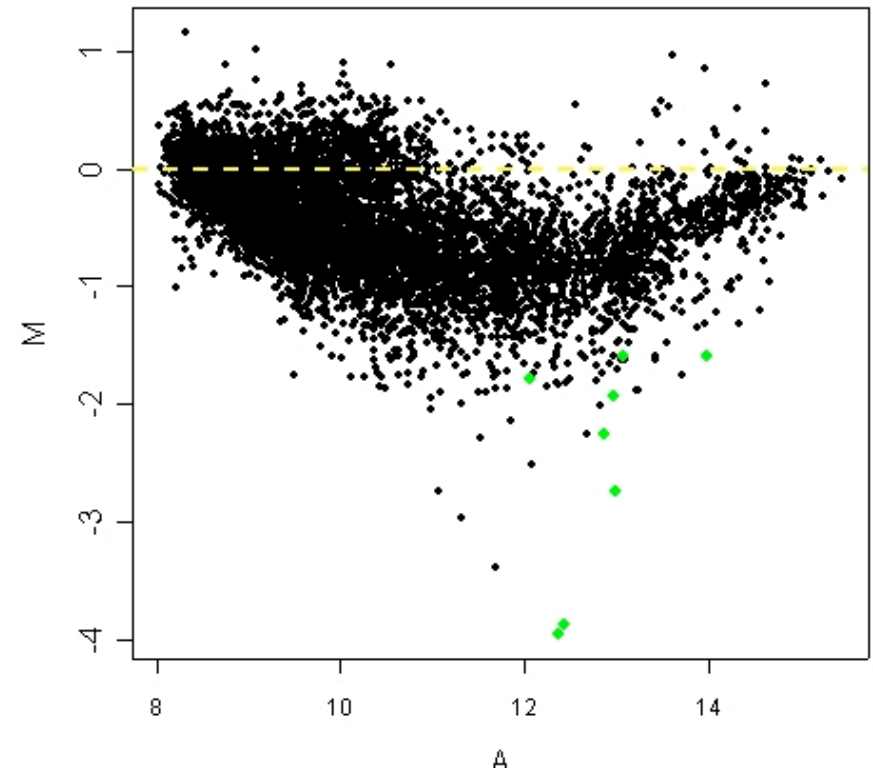
- La plupart des intensités mesurées sont faibles
- Distribution « en cloche »
- Recentrage de la distribution
- Rend symétrique les distribution
- Facilite l'utilisation des statistiques...



La rotation des graphiques (MA plot)



$\log_2 R$ vs $\log_2 G$



$M = \log_2 R/G$ vs $A = \log_2 \sqrt{RG}$

- Différences des intensités :

$$M = \log \text{ratio} = \log R/G = \log R - \log G$$

- Contre la moyenne des intensités :

$$A = \log \text{moyenne} = \log \sqrt{RG} = [\log R + \log G]/2$$

Quels spots utiliser pour normaliser ?

• Utilisation de contrôles positifs

- Prise en compte de spots contenant des gènes de ménages ou de l'ADN génomique
- Les contrôles positifs doivent être détectables, posséder une expression stable et tomber dans la gammes de détection du scanner

Avantage : peu de gènes sont nécessaires

Inconvénient : ces gènes subissent trop de fluctuations non contrôlées dans les systèmes biologiques

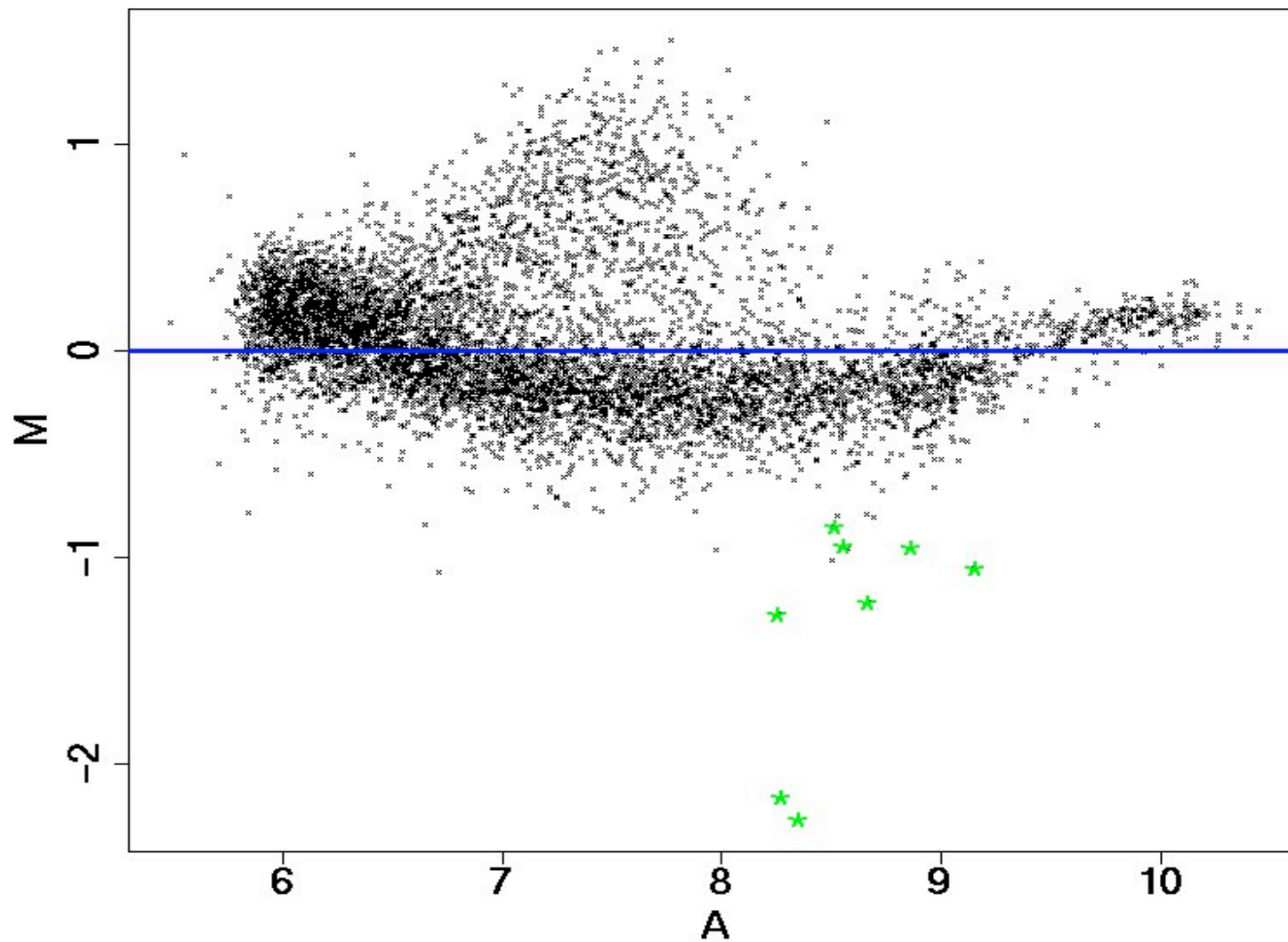
• Mesure de l'intensité globale

- Utilisation de l'intensité globale sur toute la membrane, mesurée pour tous les spots
- La mesure de l'intensité globale doit s'effectuer sur un nombre suffisant de spots et doit utiliser des valeurs homogènes

Avantage : mesure efficace sur un grand nombre de spots

Inconvénient : il est nécessaire que la majorité des gènes analysés n'ait pas une expression modifiée

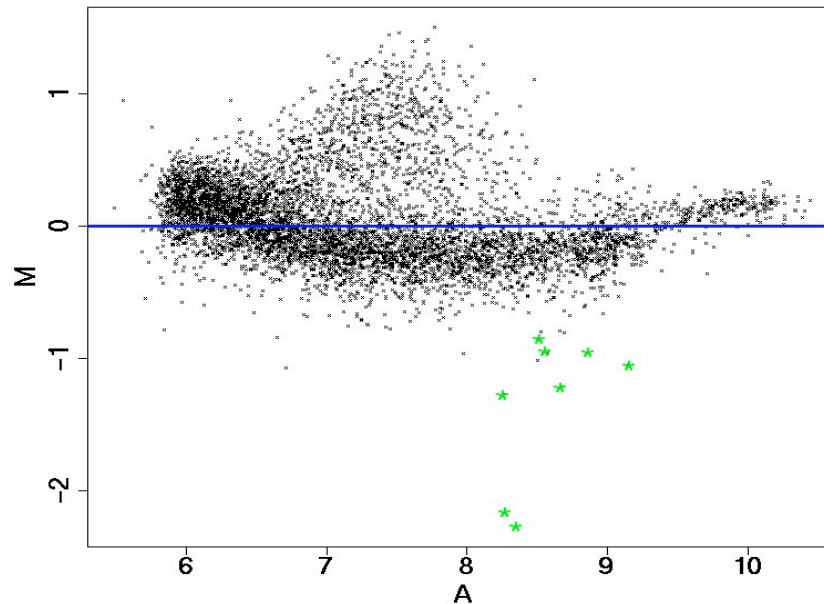
Les différentes méthodes de normalisation



Les différentes méthodes de normalisation

- Normalisation basée sur un ajustement global

- Médiane ou moyenne des log de ratios pour un gène particulier ou un ensemble de gènes (gènes de ménage)
- Normalisation en utilisant les intensités totales



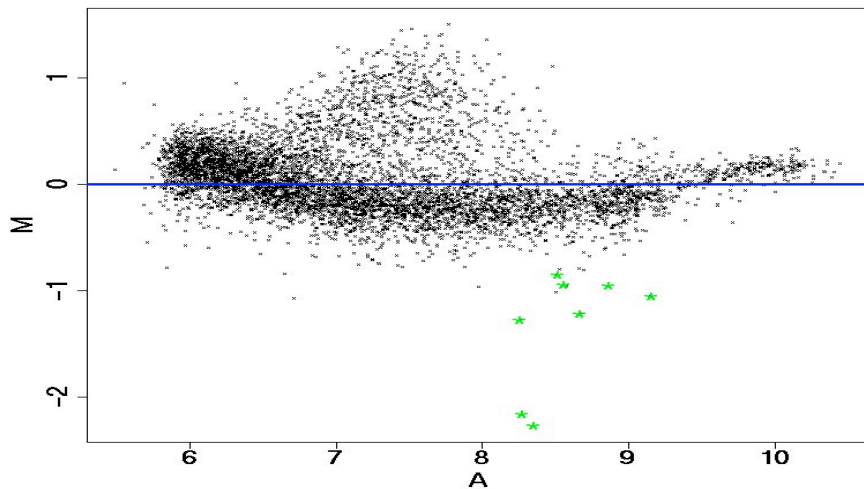
MA plot avec médianes ≈ 0

Les différentes méthodes de normalisation

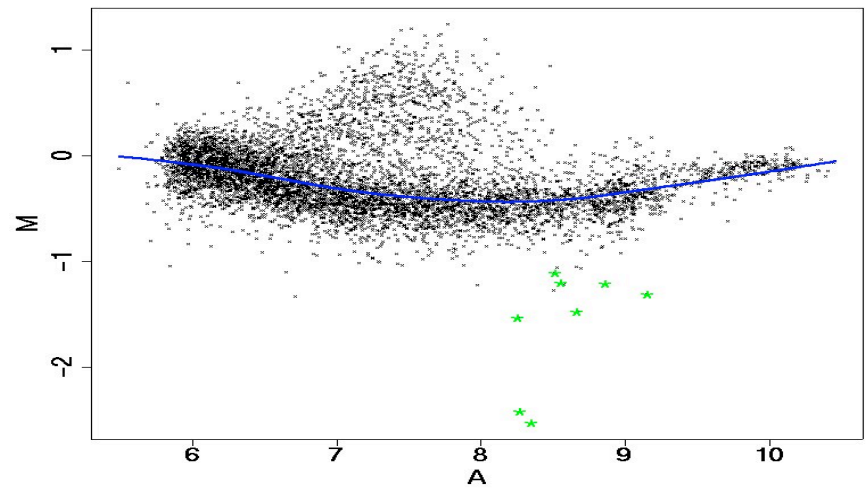
- Normalisation dépendante des intensités

- Utilisation d'une méthode de régression qui utilise la fonction lowess de Cleveland (1979) :

Lowess = LOcally WEighted Scatterplot SMOOTHing

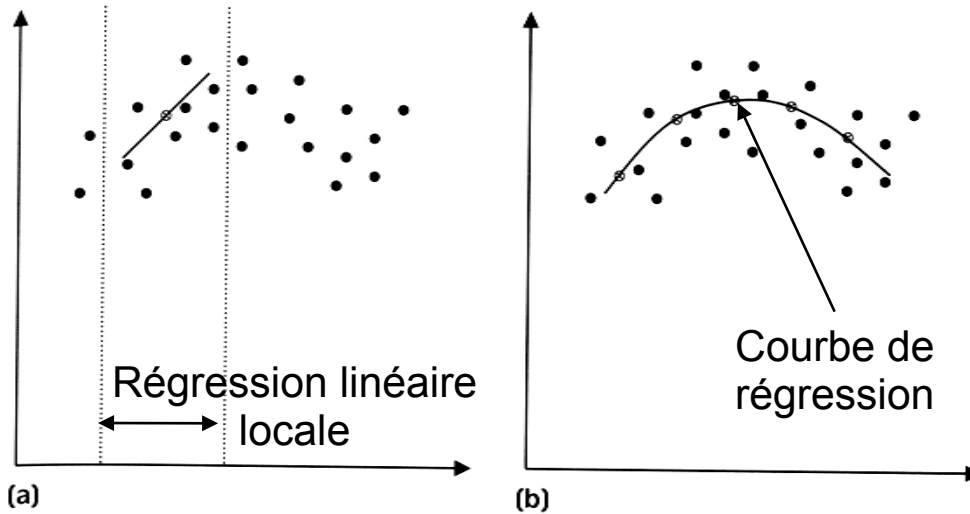


Normalisation sur les intensités globales



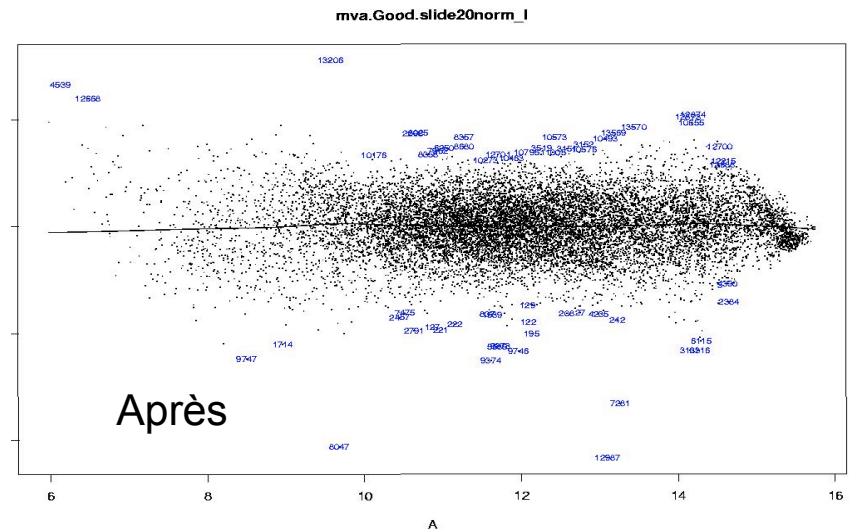
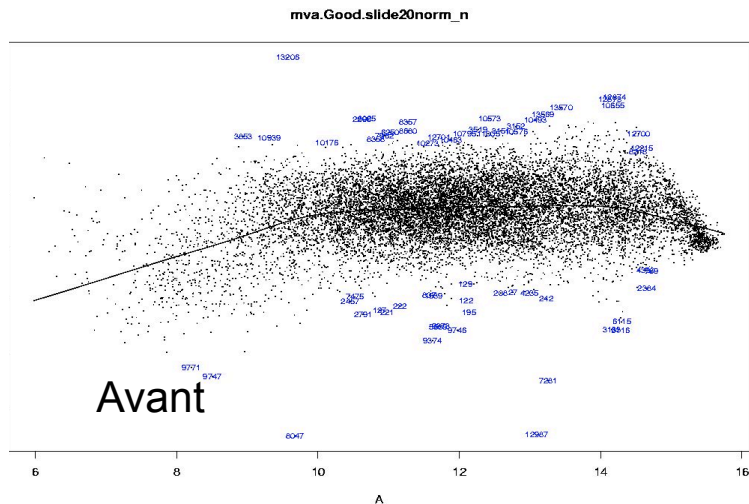
Normalisation dépendant des intensités

La normalisation par Lowess

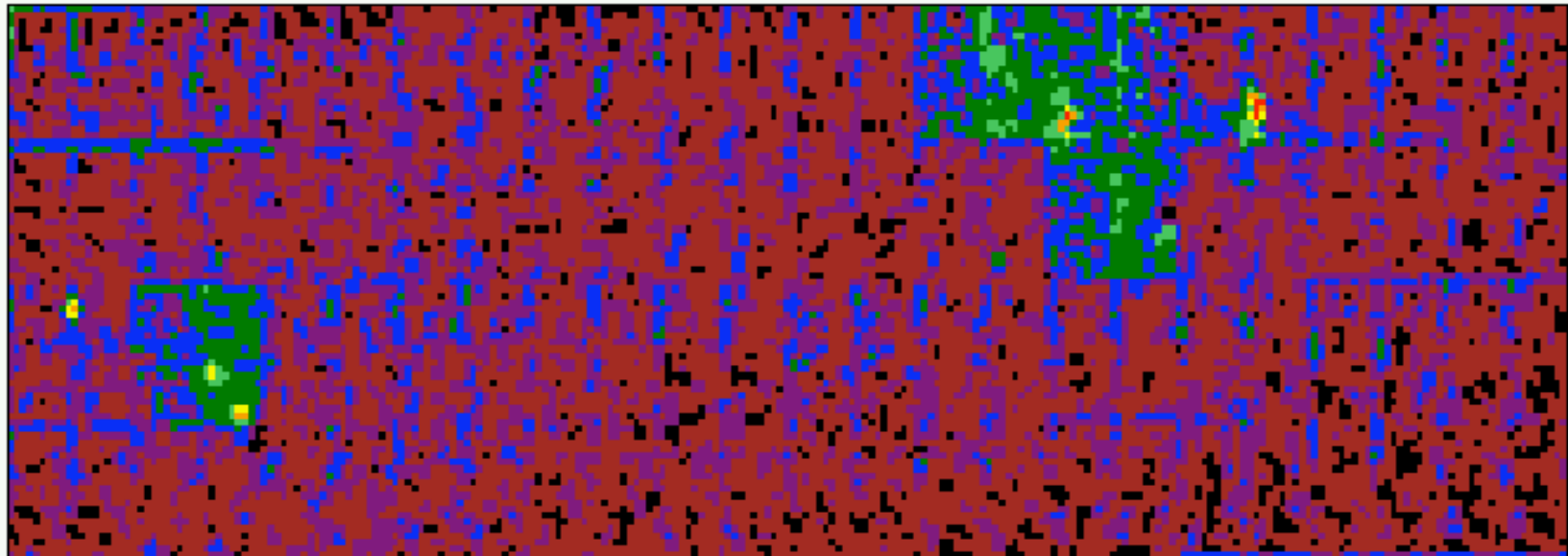


Les paramètres à prendre en compte :

- La taille des fenêtres
- Le chevauchement des fenêtres



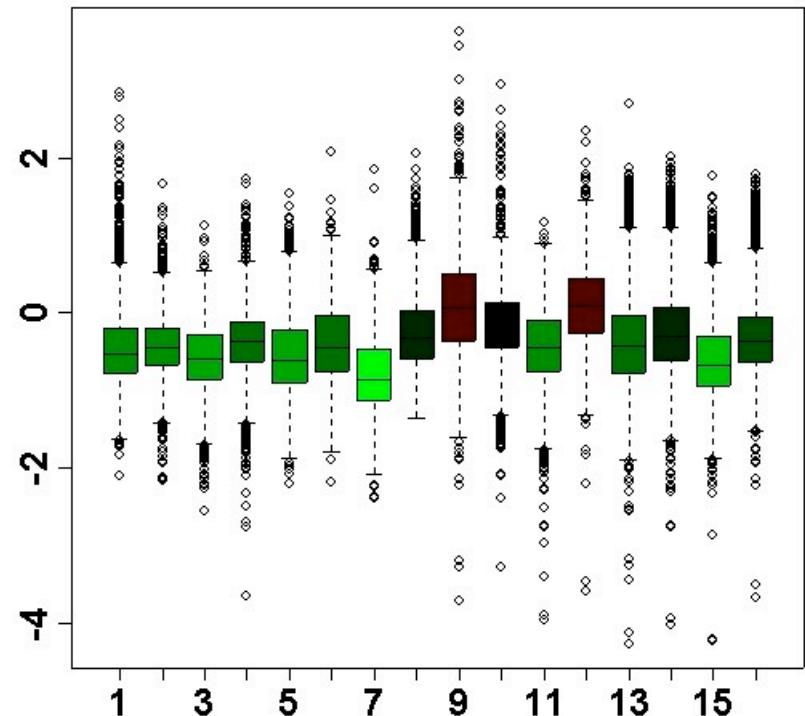
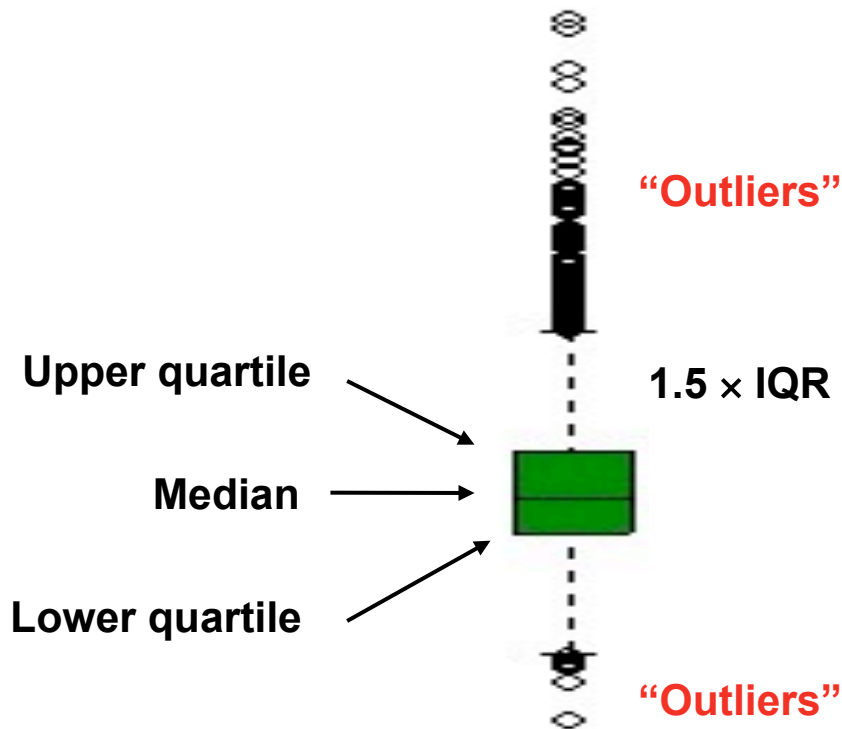
Comment s'affranchir des effets spatiaux ?



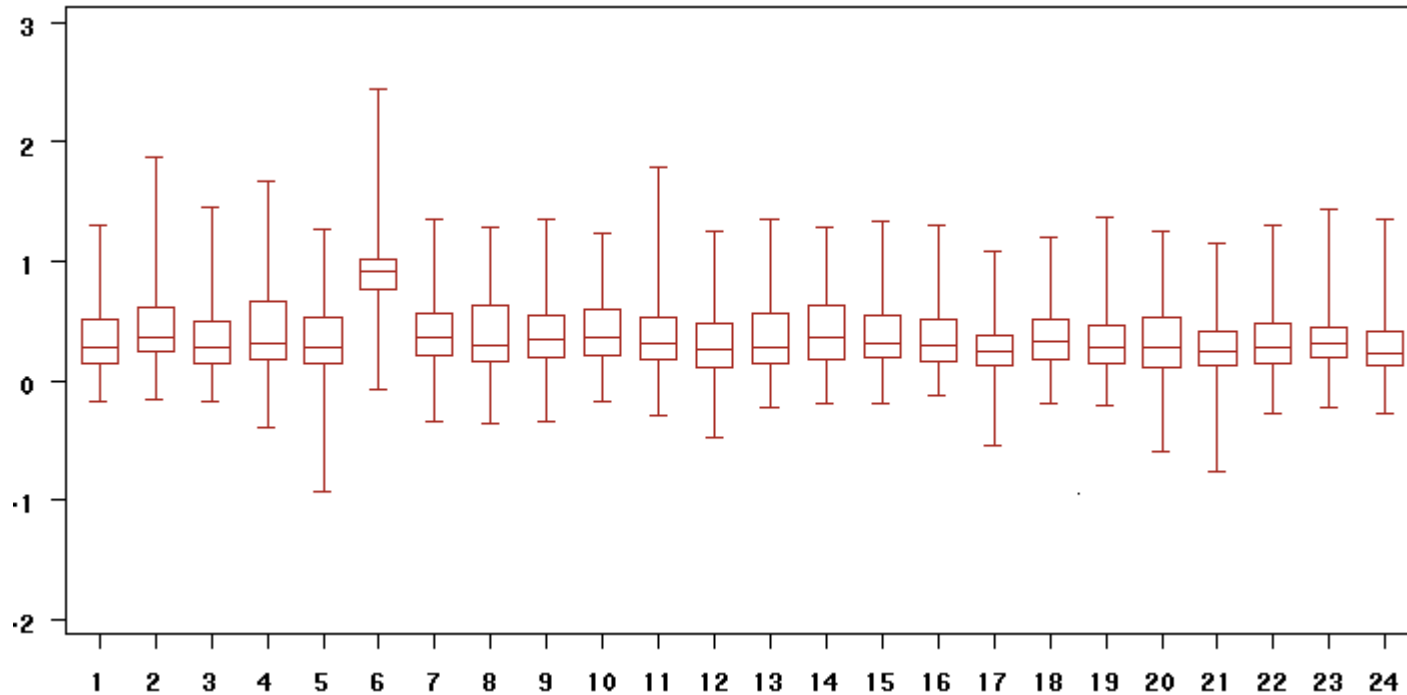
Effet de bloc ou de pointe de dépôt

Box plots

- Chaque distribution des ratios (M) d'un bloc est représentée par une boîte
- On peut visualiser directement la forme globale de la distribution (moyenne, écart-type) et comparer facilement et rapidement les $\log_2(\text{ratios})$



Visualisation de l'effet aiguille par box-plots



Les différentes méthodes de normalisation

• Normalisation en prenant en compte les pointes de spotting

• En plus des variations dépendantes de l'intensité, un biais spatial peut aussi être une source importante d'erreur systématique.

• Peu de méthodes de normalisation corrigent les effets spatiaux qui produisent des artéfacts d'hybridation comme les pointes de spotting ou les différentes plaques lors de la production des lames.

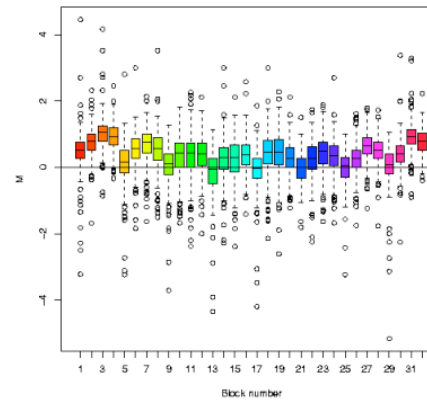
• Il est possible de corriger en même temps les biais dû à l'intensité et aux différentes pointes utilisées en effectuant une régression par Lowess sur les données à l'intérieur de chaque groupe de pointes soit :

$$\log_2 R/G \Rightarrow \log_2 R/G - c_i(A) = \log_2 R/(k_i(A)G)$$

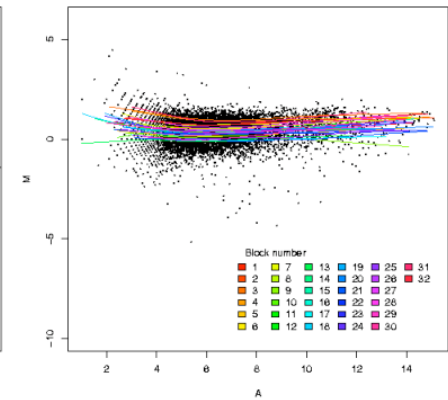
où $c_i(A)$ est le coefficient de régression Lowess sur le MA plot pour la grille i.

Avant normalisation

Boxplot of print-tip groups before normalization (excluding filtered spots)

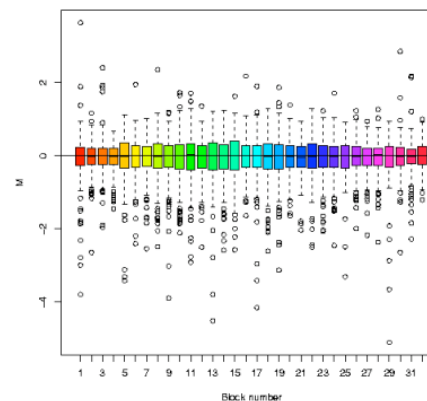


MA-plot before normalization (excluding filtered spots)

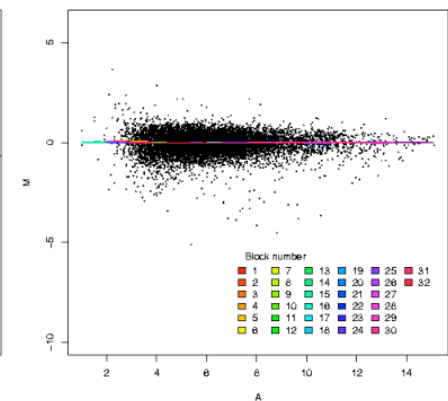


Après normalisation

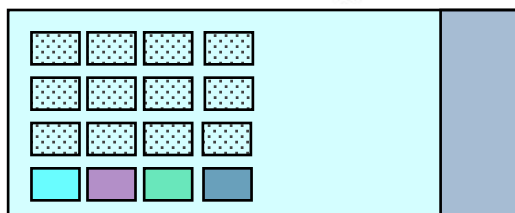
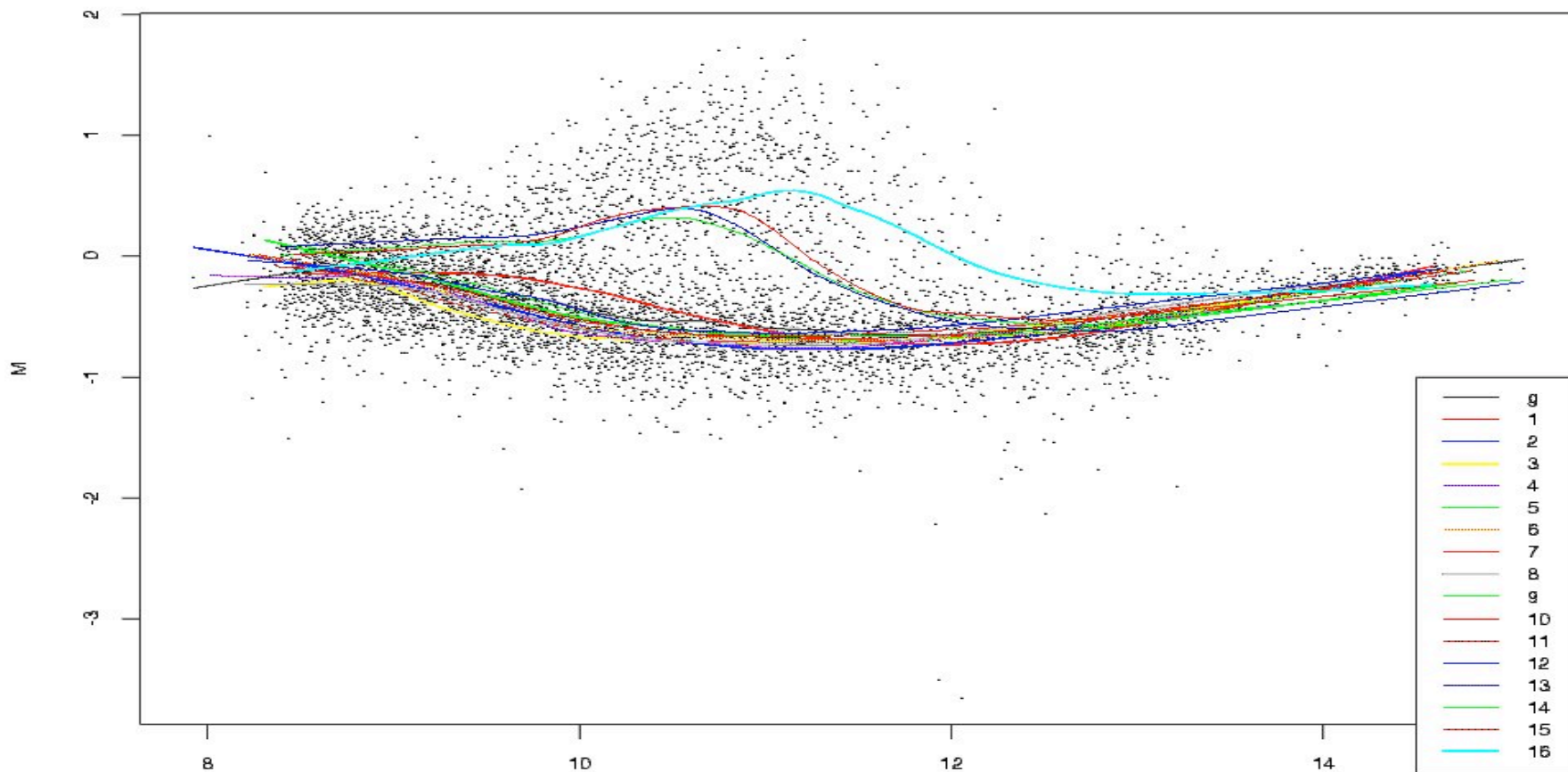
Boxplot of print-tip groups after normalization (excluding filtered spots)



MA-plot after normalization with lowess curves (excluding filtered spots)



Lowess par groupes de pointes



Lame de verre

Lames de verre avec ADNc

4x4 blocs = 16 groupes de pointes

Les différentes méthodes de normalisation

Avantages

Inconvénients

----- Régression linéaire -----

- Méthode très simple
- Alignement sur une horizontale

- Asymétrie entre Cy5 et Cy3
- Peu efficace sur les nuages déformés

----- Lowess -----

- Marche bien sur les nuages déformés
- Correction à la fois des intensités en Cy5 et Cy3

- Requiert l'utilisation de logiciels statistiques
- Nécessite de faire attention aux paramètres de la régression

----- Lowess par bloc -----

- Corrige les nuages déformés
- Prend en compte les effets locaux

- Risque de sur-corrrection du signal
- Il est nécessaire d'avoir assez de spots par bloc pour que la normalisation marche

La normalisation « deux couleurs »

- **Il y a quelques hypothèses à garder à l'esprit**
 - Utiliser la méthode Loess global implique qu'à l'échelle de l'abondance des ARNm :
 - seule une minorité des gènes est différentiellement exprimée
 - il y a un nombre égal de gènes différentiellement exprimés induits ou réprimés
 - Pour les méthodes spécifiques par bloc, il est nécessaire que les conditions précédentes soient respectées pour chaque bloc. D'un point de vue statistique, le nombre de spots concerné par la méthode ne doit pas être trop petit.
 - Utiliser un sous-ensemble de gènes spécifiques pour la normalisation (control, gènes de ménage) implique des hypothèses similaires.
- **Il y a des améliorations à apporter**
 - Utilisation d'une méthode de normalisation adaptée aux données utilisées.
 - Il est important de ne pas écraser les variations et de ne pas créer de faux positifs.

La normalisation « deux couleurs »

- **On recommande**

- D'effectuer la transformation en log₂ ratios (MA plot)
- D'utiliser la normalisation Loess globale pour corriger le biais de fluorochrome
- D'utiliser une normalisation par pointe (médiane) pour prendre en compte les biais spatiaux.
- De garder à l'esprit que la normalisation change les données brutes : il est donc nécessaire d'adapter la méthode de normalisation aux données.
- De garantir les mêmes conditions techniques pour toutes les lames qui seront utilisées dans vos expériences (même manipulateur, même scanner, même lot de lames ...)
- De ne pas hésiter à passer du temps sur les contrôles qualité (lames jaunes)

- **En principe**

- Les spots mauvais seront éliminés en faisant plusieurs réplicats

=> La chose la plus difficile c'est de corriger les biais techniques sans rien changer au signal étudié

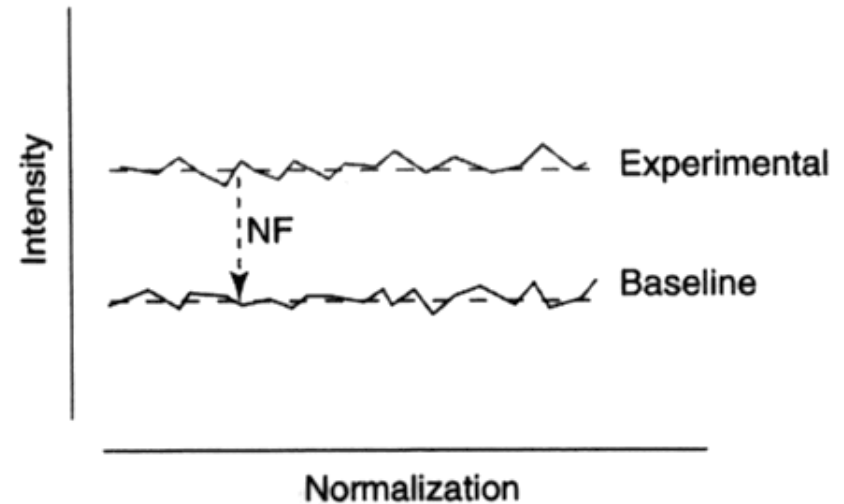
Analyse bioinformatique des puces à ADN

La normalisation entre lames

La normalisation avec les puces Affymetrix

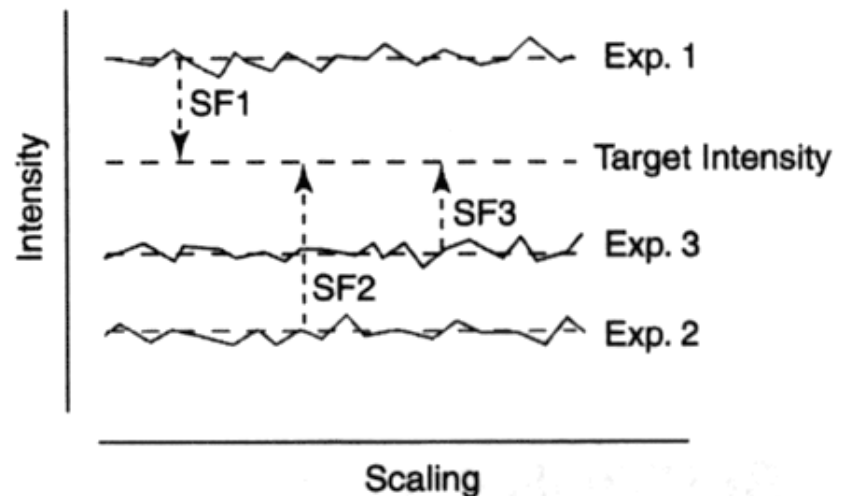
• Normalisation

- Comparaison entre deux expériences :
- Chaque condition est hybridée sur une seule membrane/lame.
- L'utilisation des intensités globales permet la normalisation de toutes les valeurs et donc la comparaison des deux expériences



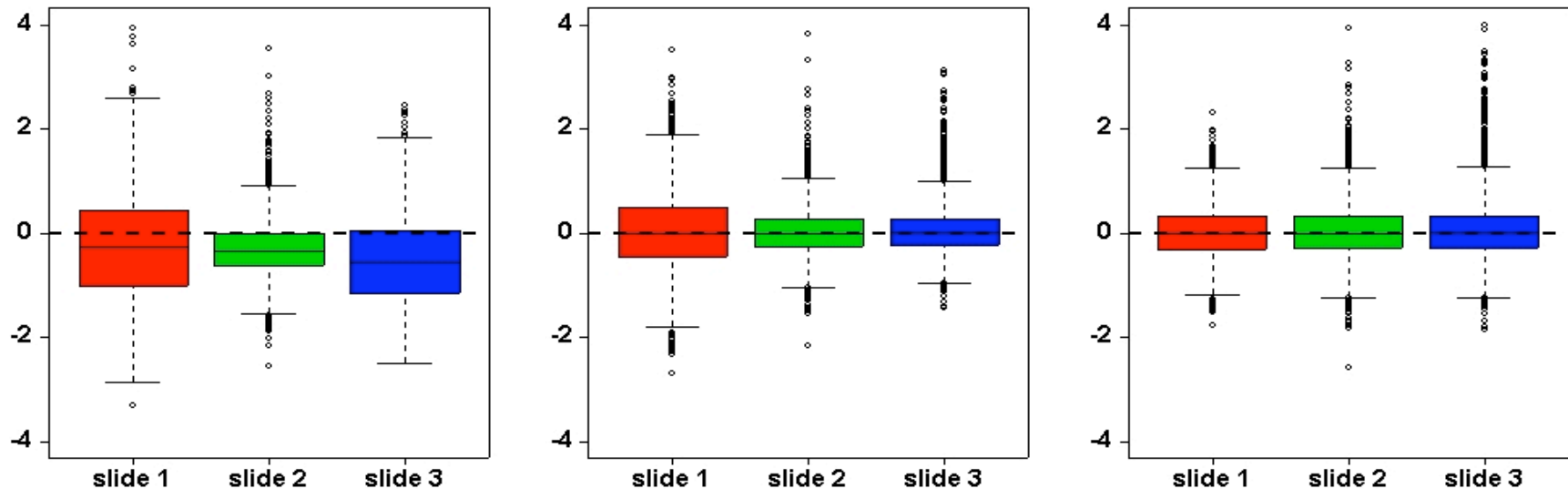
• Standardisation

- Comparaison entre plusieurs expériences :
- Toutes les expériences doivent être comparées à une condition contrôle.
- Il est nécessaire d'utiliser une intensité cible, fixée arbitrairement ou obtenue de façon absolue (gènes de ménage).



La normalisation entre lames

- **Hypothèse : les variations des distributions observées ne sont pas des changements biologiques réels**



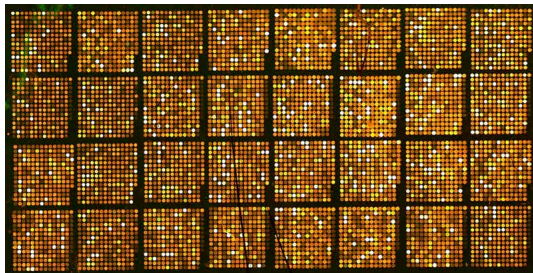
Box plot des distributions des $\log_2(\text{ratios})$ pour 3 hybridations identiques (réplicats) :

- Gauche : sans aucune normalisation
- Centre : après une normalisation Loess par bloc (centrage)
- Droite : après une normalisation entre lames (réduction)

Analyse bioinformatique des puces à ADN

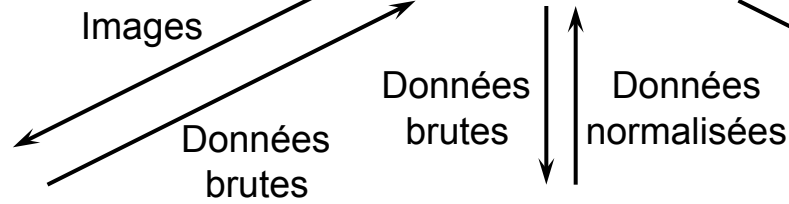
La gestion des données

La gestion du flux de données

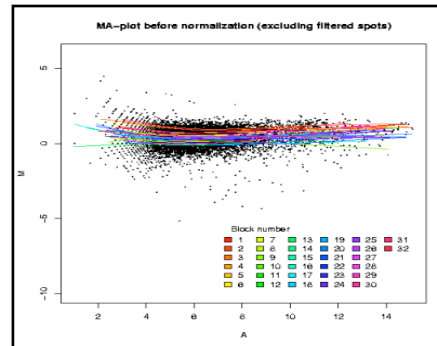
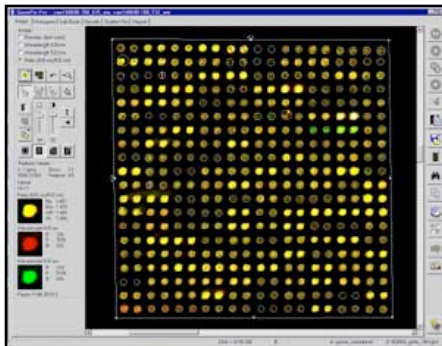


Images obtenues avec le scanner

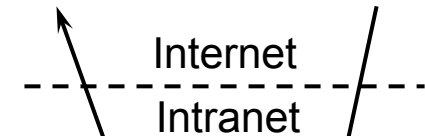
Serveur de Fichiers



Analyse d'images



Publication Web Base de données publiques



Données normalisées Données publiées

Interface Web



Les systèmes de stockage des données

- Il y a trois niveaux différents de gestion des données

- Les dépôts de données publics

Construits sur un schéma le plus flexible possible pour assurer le stockage de données hétérogènes comme les données provenant de différents organismes ou obtenues avec différents processus d'analyse

- Les bases de données institutionnelles

Construites afin d'aider un groupe d'utilisateurs sur une plate-forme technique dédiée ou pour répondre à un projet spécifique

- Les bases de données locales

Construites et installées pour un petit nombre d'utilisateurs et pour répondre à des questions très spécifiques et précises

NCBI - GEO

Gene Expression Omnibus

HOME SEARCH SITE MAP Handout NAR 2005 Paper NAR 2002 Paper FAQ MIAME Email GEO

NCBI - GEO Not logged in | Login

Gene Expression Omnibus: a gene expression/molecular abundance repository supporting MIAME compliant data submissions, and a curated, online resource for gene expression data browsing, query and retrieval.

Public data

GPL Platforms	2151
GSM Samples	76479
GSE Series	3361
Total	81991

Site contents

Documentation

- Overview | FAQ
- Web deposit guide
- Batch deposit guide
- Linking & citing
- Journal citations
- DataSet clusters
- GEO announce list
- Data disclaimer
- GEO staff

Query & Browse

- Repository browser
- Submitter contacts
- SAGEmap
- FTP site
- GEO Profiles
- GEO DataSets

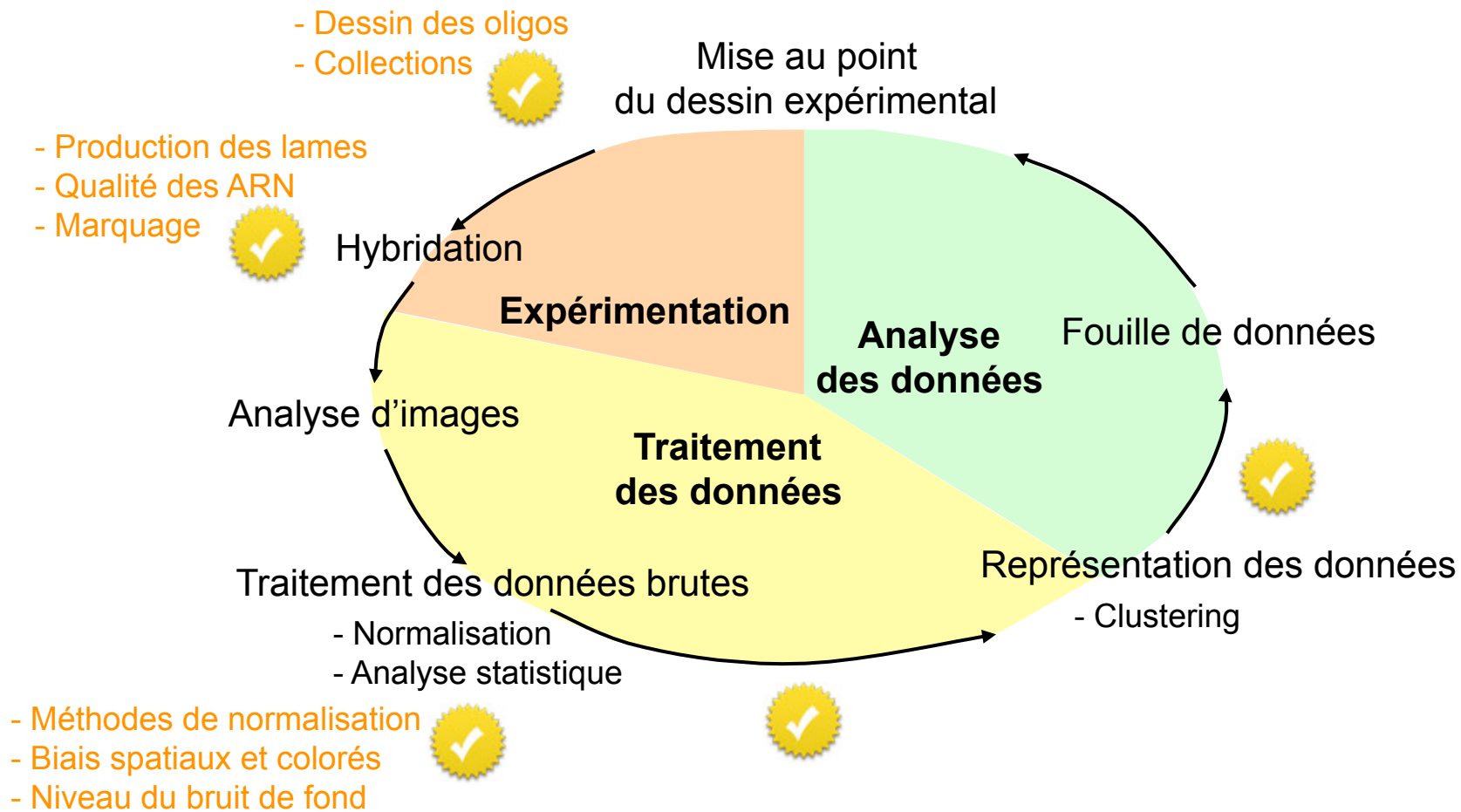
Deposit & Update

- Direct deposit
- Web deposit
- New account

Get GEO accession Scope: Self Format: HTML Amount: Quick go

Depositors only User: Password: LOGIN Recover a password

Être attentif à la qualité des données



Excel introduit des erreurs dans le nom des gènes

La conversion automatique

- La conversion automatique par défaut des dates introduit des erreurs. Par exemple, le gène suppresseur de tumeurs DEC1 est converti en « 1-DEC » (premier décembre).
- La conversion par défaut des nombres affecte les identifiants de clones de la forme nnnnnnnEnn, où n indique un chiffre. Par exemple, le clone RIKEN « 2310009E13 » est converti en nombre à virgule flottante « 2.31E+13 ». Une recherche a identifié plus de 2000 identifiants de ce type sur un total de 60770 clones RIKEN.
- Ces conversions sont irréversible, le nom de gène original ne peut plus être retrouvé.

The screenshot shows the NCBI LocusLink interface for the gene NEDD5. The search bar contains 'Hs NEDD5'. The gene name is 'NEDD5: neural precursor cell expressed, developmentally down-regulated 5'. In the 'Mouse Homology Maps' section, there is a table with columns for the comparison, distance, and a link. The link '2-Sep' is circled in red, indicating a mistaken identifier. Other links include 'Sept2' and 'AW208991'.

Comparison	Distance	Link	Species
NCBI vs. MGD	1 cM	2-Sep	Hs Mm
UCSC vs. MGD	1 cM	Sept2	Hs Mm
UCSC vs. Hudson et al.	1 1319.34 cR	AW208991	Hs Mm

Zeeberg BR et al. Mistaken identifiers: gene name errors can be introduced inadvertently when using Excel in bioinformatics. *BMC Bioinformatics*. 2004 5:80.

Quelques références

• Normalisation

- Quackenbush J. Microarray data normalization and transformation. *Nat Genet.* 2002 **32** Suppl:496-501.
- Yang YH, Dudoit S, Luu P, Lin DM, Peng V, Ngai J, Speed TP. Normalization for cDNA microarray data: a robust composite method addressing single and multiple slide systematic variation. *Nucleic Acids Res.* 2002 **30**(4):e15.
- Leung YF, Cavalieri D. Fundamentals of cDNA microarray data analysis. *Trends Genet.* 2003 **19**(11):649-59.

• Statistiques

- Saporta. Probabilités, analyse des données et statistiques. *Editions Technip*
- Daudin, Robin et Vuillet. Statistique inférentielle, idées, démarches, exemples. *Presses Universitaires de Rennes*
- Tassi. Méthodes statistiques. *Economica*

