# cDNA Microarray Analysis

## with BioConductor packages

Nolwenn Le Meur

Copyright 2006

# Data Analysis of Microarrays

Experimental Design

Image Analysis

Quality Assessment

Pre-processing

Background Correction

Normalization

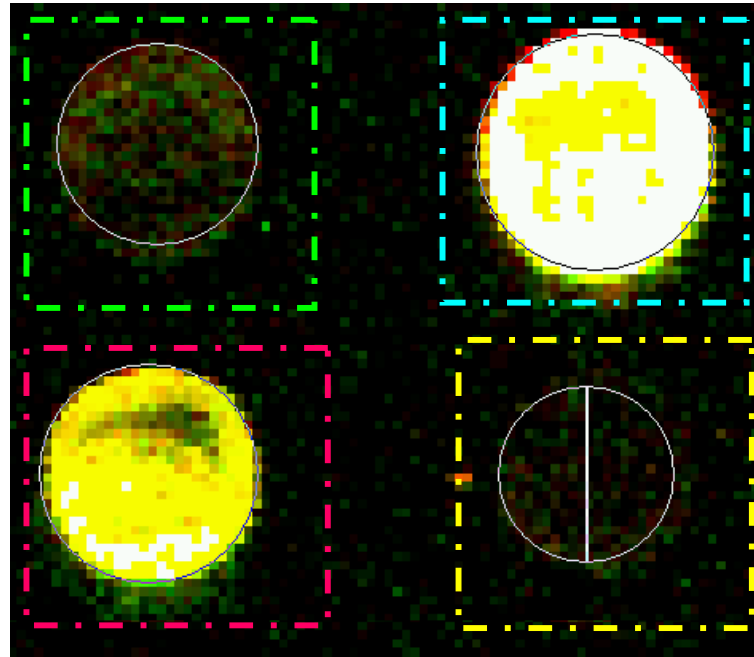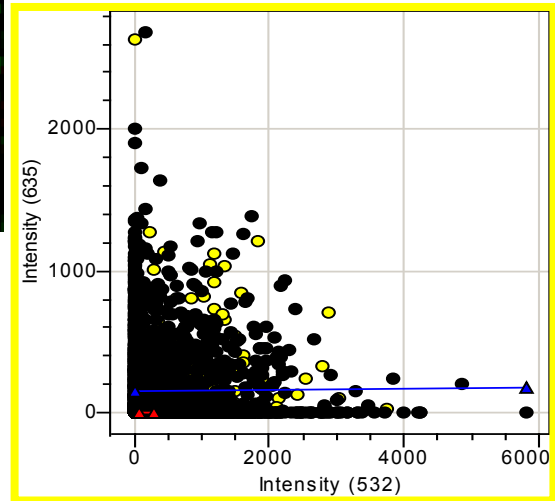Summarization

Analysis

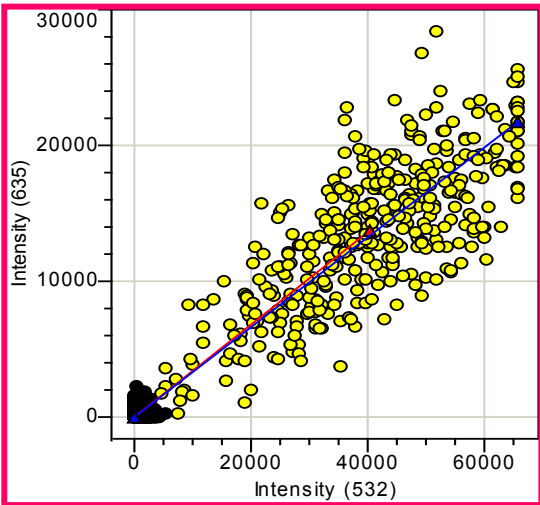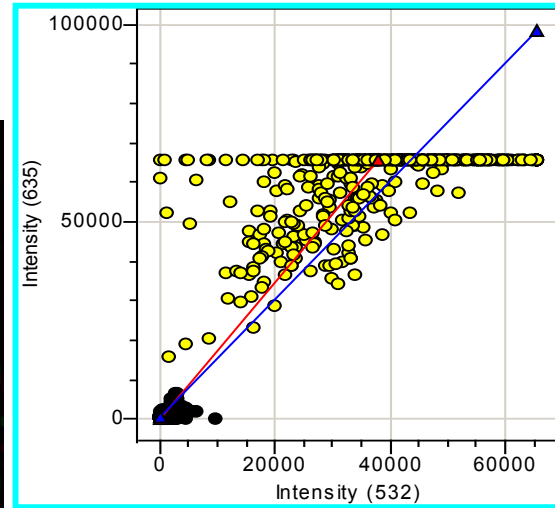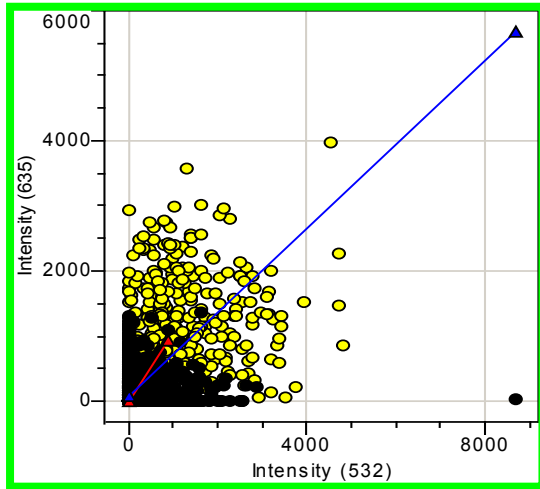Testing     Discovery     Prediction

# Outline

- **Data acquisition & Pre-processing (chap. 4)**
  - Image analysis
  - Quality assessment
  - Pre-processing
- **Differential expression (chap. 14, 15 & 23 )**

- **Lab : case studies (chap 4 & 23)**
  - marray & arrayQuality (Y.H Yang & A.C. Paquet)
  - limma (G.K Symth)

# Terminology

- **Target:** DNA hybridized to the array, mobile substrate.
- **Probe:** DNA spotted on the array (spot).
- **print-tip-group :** collection of spots printed using the same print-tip (or pin), aka. grid.

- **G, Gb:** Cy 3 signal and background intensities
- **R, Rb:** Cy5 signal and background intensities
- $M = \log_2(R) - \log_2(G)$
- $A = 1/2(\log_2(R) - \log_2(G))$

# Image Analysis

## 1. Location

## 2. Segmentation

## 3. Quantification
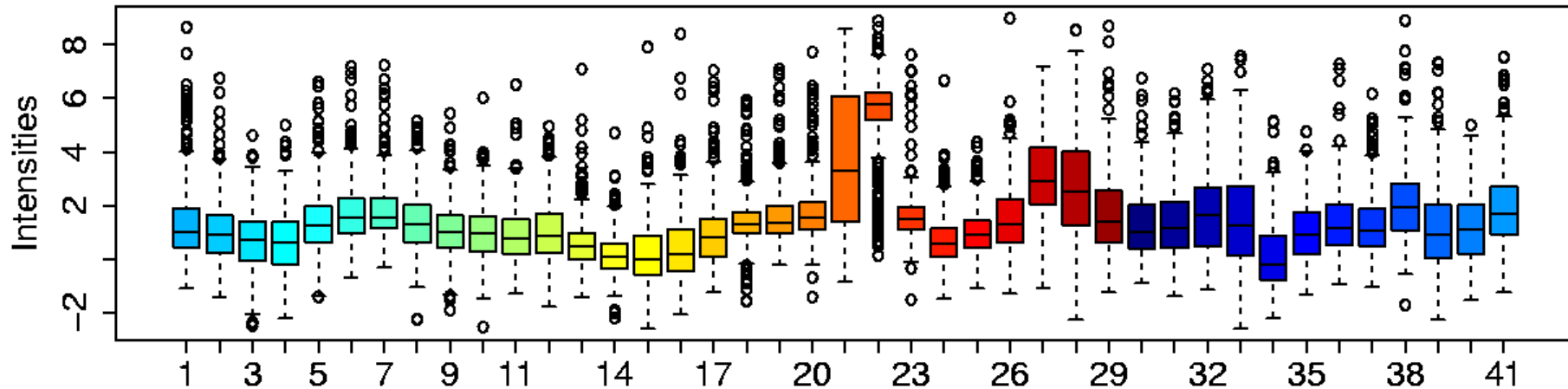
**Raw data**

# Quality Filtering



- ● Background
- ● Foreground

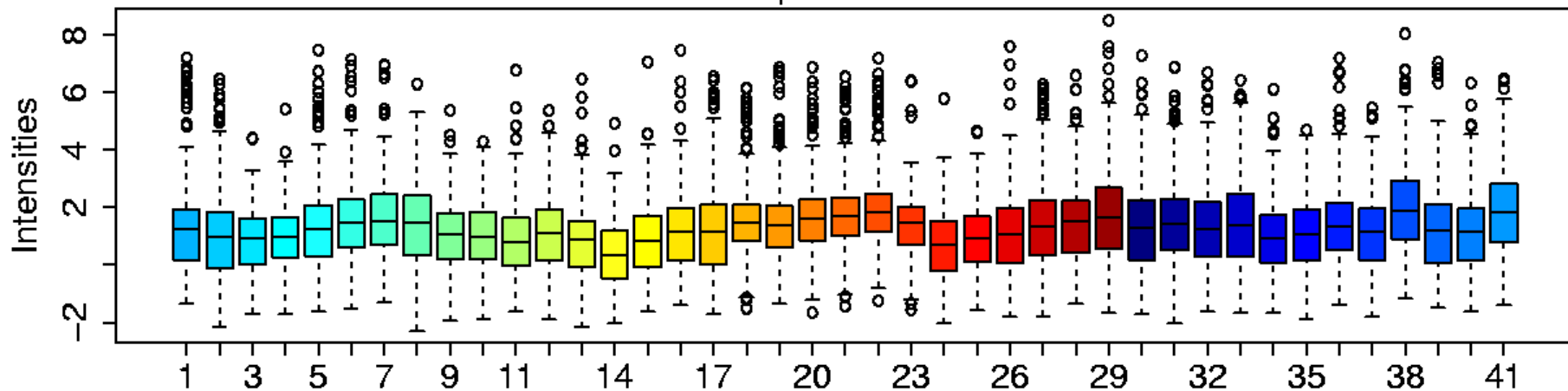# Quality Assessment

For each array:

- **Diagnostics plots** of spot statistics

  *e.g.* R and G log-intensities, M, A, spot area.

  - Boxplots;

  - 2D spatial images;

  - Scatter-plots, e.g. MA-plots;

  - Density plots.

- **Stratify** plots according to layout parameters, *e.g.* print-tip-group, plate.
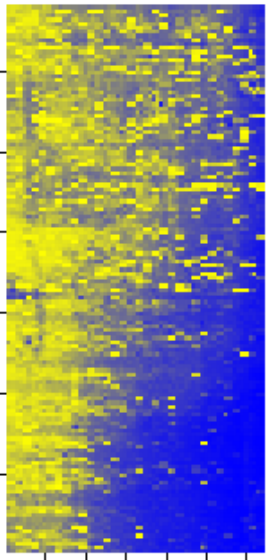
# PCR Plates - Boxplots

# Spatial Effects – Image Plots



R     Rb       R-Rb

**color scale by rank**

another array: print-tip

**color scale ~ log(G)**

**color scale ~ rank(G)**

max

min

# Spatial Effects



**1 pin** ➡ **1 block**

# Spotting Pin Quality Decline



SMP3 (0.25 ul uptake)  SMP3B (0.6 ul uptake)

after delivery of $5\times10^5$ spots

after delivery of $3\times10^5$ spots

# Print-tip Effects – ECDF plot

# Diagnostic plot with *arrayQuality*

# Data Exploration with *limma*



Example MA-Plot with Spot-Type Highlighting

(Limma user Guide)

# Quality Assessment: Summary

For each array:

- Diagnostics plots
- Stratify

BioC packages:

- *arrayQuality*
- *arrayMagic*
- *…*

# Outline

- Data acquisition & Pre-processing (chap. 4)
  - Image analysis
  - Quality assessment
  - **Pre-processing**
- Differential expression (chap. 14, 15 & 23 )

- Lab : case studies (chap 4 & 23)
  - marray & arrayQuality (Y.H Yang & A.C. Paquet)
  - limma (G.K Symth)

# Variance-Bias trade off

# Sources of Variation

- RNA extraction

- reverse transcription

- labeling efficiencies

- Scanner settings

**Systematic**

- similar effect on many measurements
- corrections can be estimated from data

**Calibration**

- PCR

- DNA concentration

- Printing or pin

- cross-hybridization

- …

**Stochastic**

- too random to be ex-plicitly accounted for
- "noise"

**Error Model**

# Background Correction

- none

- subtraction, movingmin

- *Minimun, edwards, normexp,…*

- More details … *limma*

  >?backgroundCorrect

# Background Correction



none                        substraction                        *normexp*

# Why Normalize?

# Normalization

Identify and remove the effects of systematic variation

- Normalization is closely related to quality assessment. In a ideal experiment, no normalization would be necessary, as the technical variations would have been avoided.

- Normalization is needed to ensure that differences in intensities are indeed due to differential expression, and not some printing, hybridization, or scanning artifact.

- Normalization is necessary before any analysis which involves within or between slide comparisons of intensities, e.g., clustering, testing.

# Data Transformation

**measured intensity  =  offset  +       gain   × true abundance**

$$Y_{ik} = B_{ik} + \alpha_{ik} S_k$$

- Intensity measurements adapt a distribution that is closer to the normal distribution
- Muliplicative noise becomes additive noise: variance more independent of intensity

Example: log transformation

# Normalization methods

- median
- loess
- 2D loess
- print-tip loess
- variance stabilisation

Two-channel

Separate-channel

Smyth, G. K., and Speed, T. P. (2003). In: *METHODS: Selecting Candidate Genes from DNA Array Screens: Application to Neuroscience*

BIOCONDUCTOR

FRED HUTCHINSON
CANCER RESEARCH CENTER

# Two channel normalization

- Location: centers log-ratios around zero using A and spatial dependent bias



Swirl 93 array: pre-normalization log-ratio M



Swirl 93 array: within-print-tip-group loess normalization log-ratio

# Two channel normalization

- **Location**: centers log-ratios around zero using A and spatial dependent bias

- **Scale**: adjust for different in scale between multiple arrays



median centered

Scaling

median centered & MAD scaled

# One channel normalization

- As technology improves the spot-to-spot varation is reduced

- Development of normalization techniques that work on the absolute intensities

Ex: quantile normalization (*limma*)

variance stabilization (*vsn*)

# Quantile Normalization

Before

After ₅



Bolstand *et al.*(2003)

# Variance Statibilizing Transformation



- Meaningful around 0
- Original intensities may be negatives

(Huber *et al.* 2004)

# Variance stabilization (*vsn*)



linear           log           arsinh

# Variance stabilization (*vsn*)

| log-ratio | $$\log \frac{x_i}{x_j}$$ |
|---|---|

| 'glog' (generalized log-ratio) | $$\log \frac{x_i + \sqrt{x_i^2 + c_i^2}}{x_j + \sqrt{x_j^2 + c_j^2}}$$ |
|---|---|

- interpretation as "fold change"

+ interpretation even in cases where genes are off in some conditions (negative values)

+ visualization

+ can use standard statistical methods (hypothesis testing, ANOVA, clustering, classification…) without the worries about low-level variability that are often warranted on the log-scale

# Preprocessing : Summary

For each array:

- Background correction or not
- Normalization: bias-variance trade-off
- Diagnostic plots

BioC pacakges:

- *marray*
- *limma*
- *…*

# Outline

- Data acquisition & Pre-processing (chap. 4)
  - Image analysis
  - Quality assessment
  - Pre-processing
- **Differential Expression (chap. 14, 15 & 23 )**

- Lab : case studies (chap 4 & 23)
  - marray & arrayQuality (Y.H Yang & A.C. Paquet)
  - limma (G.K Symth)

# Experimental Designs



**Reference design with dye swap**

**Loop**

**Avoid Confounding effect**

Yang, Y. H. et Speed, T. (2002). Design issues for cDNA microarray experiments. *Nat.Rev.Genet.*, **3**: 579-588.

# Experimental Designs

- Simple comparisons
- Technical replicates
- Dye swap
- Within array replicate spots
- Two groups
- Several groups
- Direct two color designs
- Factorial design
- Time Course
- …

**Case Studies Chap. 23**

# Differentially Express Genes

- Fold change

But no assessment of statistical significance

# Differentially Express Genes

Example: The two–sample t–statistic is used to test equality of the group means $\mu_1$, $\mu_2$.

The *p–value* $p_g$ is the probability under the null hypothesis (here: $\mu_1 = \mu_2$) that the test statistic is at least as extreme as the observed value $T_g$. Under the null hypothesis, $Pr(p_g < \alpha) = \alpha$.

# Differentially Express Genes

- Fold change
- Parametric test
  - standard $t$-test
  - Welch $t$-test
- Non parametric
  - Wilcoxon test
  - Mann-Whitney
- Permutation test

# Multiple testing

| Number of genes | Gene significance level | | | |
|---|---|---|---|---|
| | P-values < 0.01 | 0.05 | 0.1 | 0.15 |
| 10 | < 1 | < 1 | 1 | 1.5 |
| 20 | < 1 | 1 | 2 | 3 |
| 5000 | 50 | 250 | 500 | 750 |
| 10000 | 100 | 500 | 1000 | 1500 |

**Test of Thousands of hypotheses simultaneously!**

➢ **Increased chance of false positives**

Drăghici (Chapman &Hall 2003)

Individual p-values of 0.01 no longer correspond to significant findings.

**-> Adjust for multiple testing**

# Nonspecific filtering

- Remove genes :
  - Low intensities
  - Do not show sufficient variation across all samples
- Select genes :
  - Known to interact in a specific biological process, e.g. GO  (Chap 14.)

# Type of Error

|  | H$o$ is true | H$o$ is false |
|---|---|---|
| H$o$ not rejected | True negatives<br>1- $\alpha$ | False negatives<br>(Type II error)<br>$\beta$ |
| H$o$ rejected | False positives<br>(type I error)<br>$\alpha$ | True positives<br>(Power)<br>1-$\beta$ |

BIOCONDUCTOR

FRED HUTCHINSON CANCER RESEARCH CENTER

# Control of Error

- Type **II** error or Minimizing False negatives

    ->power of tests, sample size

- Type **I** error

    **->** Control false positive rate (FWER,FDR) *or p-value*

    – **F**amily **W**ise **E**rror **R**ate

       control probability of false positive on entire set of genes

    – **F**alse **D**iscovery **R**ate

       control false discovery rate on set of identified genes

# Control of Type Error I

| Control | Method | Pros/Cons |
|---------|--------|-----------|
| **FWER** | Bonferroni | Very conservative |
| | Šidák | Very conservative |
| | Holm | Assumption free, conservative |
| | Hochberg | Independent variables |
| | Modified Westfall & Young | Exploit *joint* distribution of test statistics, need replicates |
| **FDR** | Benjamini & Hochberg | Independent variables conservatives |
| | Benjamini & Yekutieli | |
| | Tusher | Sensitive to the number of replicates |

*Ge, Y & Dudoit, S. (2003) Technical report #633*

# FWER vs FDR

- **FWER** if high confidence in **all** selected genes is desired. Loss of power due to large number of tests: many differentially expressed genes may not appear significant.

- If a certain proportion of **false positives** is **tolerable**: Procedures based on **FDR** are more flexible; the researcher can decide how many genes to select, based on practical considerations

# Moderated t-statistics

- **t–test estimate** the variance of each gene individually.

    - > Ok if we have enough replicates,

    - but with few replicates (say 2–5 per group), these variance estimates are highly variable.

- **moderated t–statistic**, the estimated gene–specific variance $s_g^2$ is replaced by a weighted average of $s_g^2$ and $s_0^2$, which is a global variance estimator obtained from pooling all genes.

This gives an interpolation between the t–test and a fold–change criterion.

Examples: packages *limma, siggenes*

# *limma* moderated t-statistic

- complex experiments: linear models, contrasts

- empirical Bayes methods for differential expression: t-tests, F-tests, posterior odds

- inference methods for duplicate spots, technical replication

- control of FDR across genes and contrasts

# Differential Expr. : Summary

- Permutation tests
- Multiple testing
- Pre-filtering or subsetting
- Rank genes

BioC pacakges:

- *limma*
- *multtest*
- *…*

# BioC Task View: TwoChannel

## Subview of

- Microarray

## Packages in view

| Package | Maintainer | Title |
|---|---|---|
| arrayQuality | A. Paquet | Assessing array quality on spotted arrays |
| bridge | Raphael Gottardo | Bayesian Robust Inference for Differential Gene Expression |
| genArise | IFC Development Team | Microarray Analysis tool |
| GEOquery | Sean Davis | Get data from NCBI Gene Expression Omnibus (GEO) |
| limma | Gordon Smyth | Linear Models for Microarray Data |
| limmaGUI | Keith Satterley | GUI for limma package |
| maDB | Johannes Rainer | Microarray database and utility functions for microarray data analysis. |
| makePlatformDesign | Benilton Carvalho | Platform Design Package |
| marray | Yee Hwa (Jean) Yang | Exploratory analysis for two-color spotted microarray data |
| nnNorm | Tarca Laurentiu | Spatial and intensity based normalization of cDNA microarray data based on robust neural nets |
| nudge | N. Dean | Normal Uniform Differential Gene Expression detection |
| oligo | Benilton Carvalho | Oligonucleotide Arrays |
| OLIN | Matthias Futschik | Optimized local intensity-dependent normalisation of two-color microarrays |
| OLINgui | Matthias Futschik | Graphical user interface for OLIN |
| rama | Raphael Gottardo | Robust Analysis of MicroArrays |
| snapCGH | Mike Smith | Segmentation, normalisation and processing of aCGH data. |
| spotSegmentation | Chris Fraley | Microarray Spot Segmentation and Gridding for Blocks of Microarray Spots |
| vsn | Wolfgang Huber | Variance stabilization and calibration for microarray data |

BIOCONDUCTOR

FRED HUTCHINSON CANCER RESEARCH CENTER

# Outline

- Data acquisition & Pre-processing (chap. 4)
  - Image analysis
  - Quality assessment
  - **Pre-processing**
- Differential expression (chap. 14, 15 & 23 )

- **Lab : case studies (chap 4 & 23)**
  - marray & arrayQuality (Y.H Yang & A.C. Paquet)
  - limma (G.K Symth)

# Getting started

| | Preprocessing | | | |
|---|---|---|---|---|
| | *limma* package | | *marray* package | |
| Action | Function | Class - Object | marray | Class - Object |
| read target file | readTargets | *dataframe* | read.marrayInfo | marrayInfo |
| read image file | read.maimages | RGList | read.marrayRaw, read.GenePix, read.Spot, read.SMD, read.Agilent | marrayRaw |
| read gene list | readGAL | RGList$genes | read.Galfile | marrayInfo, marrayLayout |
| read spot type | readSpotTypes, controlStatus | RGList$genes$status | | |
| array layout | getLayout | RGList$printer | read.marrayLayout, Layout | marrayLayout |
| background correction | backgroundCorrect | | | |
| one array normalization | normalizeWithinArrays, MA.RG | MAList | maNormMain | marrayNorm |
| normalization between arrays | normalizeBetweenArrays | MAList | | |