# Analysing Illumina bead-based data using beadarray

Mark Dunning
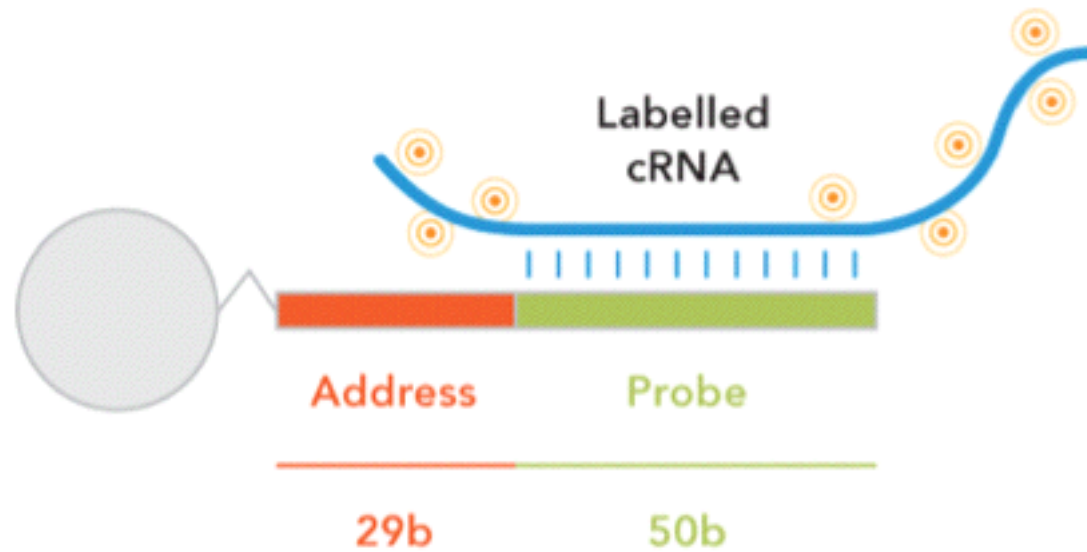
6th August 2007

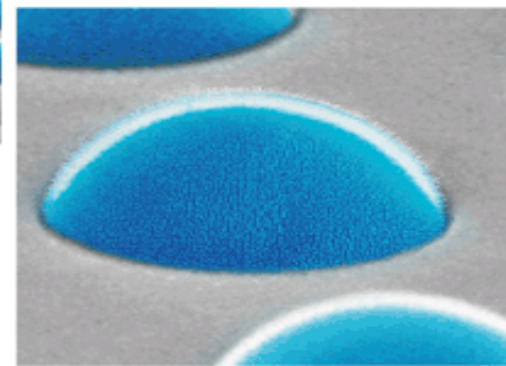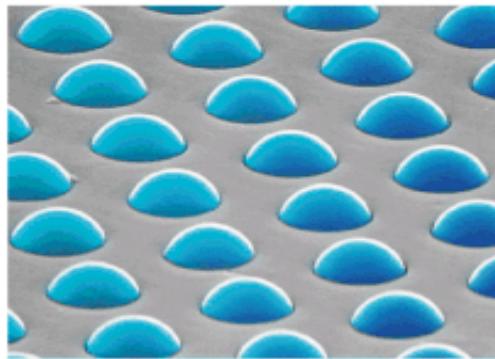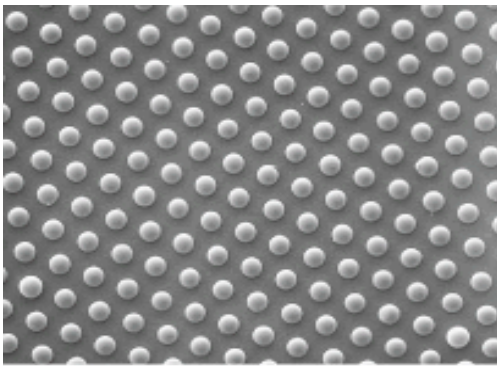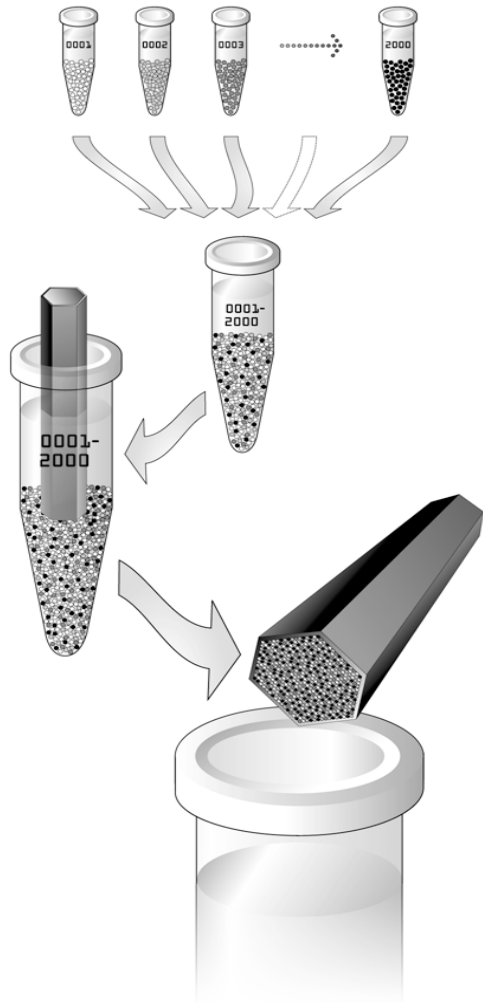UNIVERSITY OF CAMBRIDGE

CANCER RESEARCH UK

# The Bead



Each silica bead is 3 microns in diameter

700,000 copies of same probe sequence are covalently attached to each bead for hybridisation & decoding
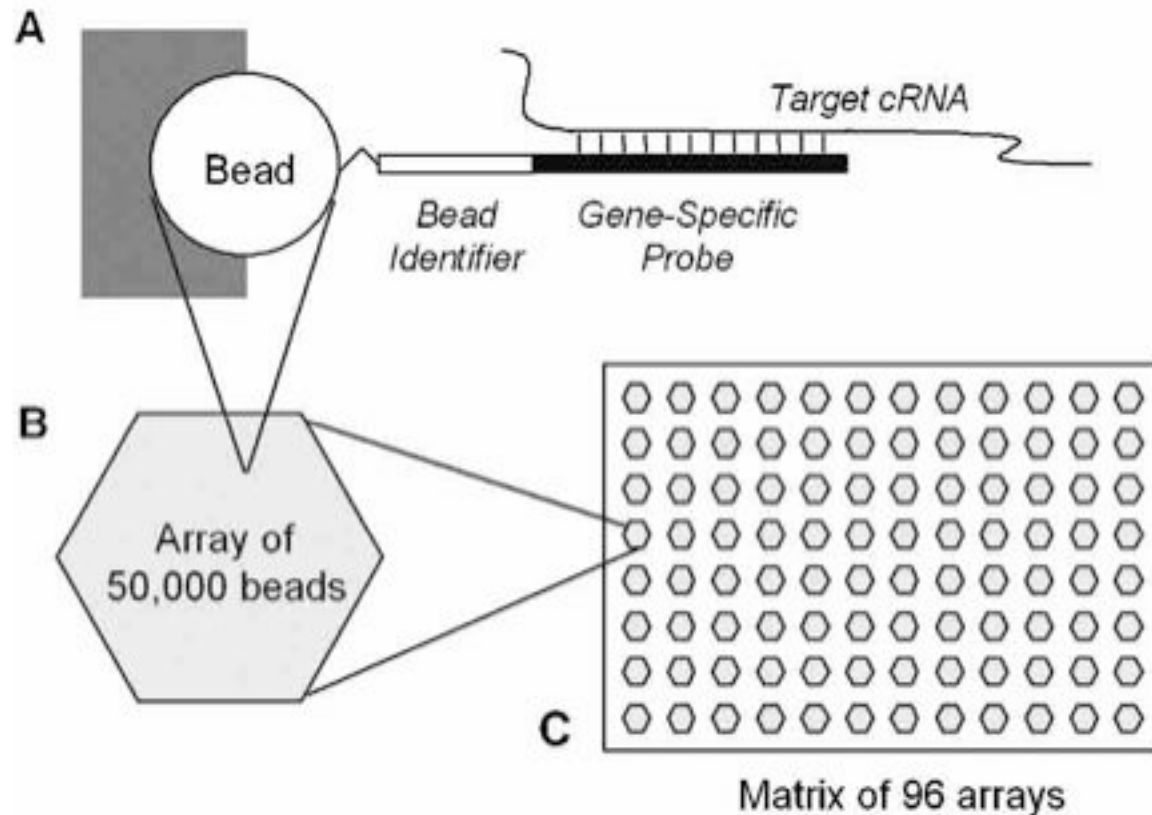
# Beads in Wells

# Bead Preparation and Array Production



- Bead pools produced containing 384 to 24,000 bead types

- Wells created in either fibre-optic bundle (hexagon) or chip (rectangle) & exposed to array

- Beads self-assemble into wells to form **randomly arranged** array of beads

- *Average of 30 beads of each type*

- Each array produced separately

# Combining Arrays - The SAM



~1500 bead types on array ~30 of each type
1 array = 1 sample or treatment
96 arrays processed in parallel - **High throughput**

# The SAM

# Combining Arrays - BeadChips



RefSeq BeadChip

8 arrays per chip 1 strip = 1 array

24,000 bead types from RefSeq database x 30 reps on each array

Whole Genome

6 arrays per chip: 2 strips = 1 array

48,000 bead types (24,000 RefSeq + 24,000 supplemental) on each array

# Raw Data

Illumina's scanning software (BeadScan) produces encrypted files (.idat, .locs etc) which are read by their proprietary analysis software (BeadStudio)

However, with modifications BeadScan you can also get more useful (readable) files for each array on a SAM or *strip* on a BeadChip

-Text file giving the identity and location of each individual bead -  with 50,000 rows for SAM ~ 1.1 million for BeadChip

-TIFF images (and not jpegs)

We refer to the TIFF and text files as the **bead level data** for an array

# Bead Level Text Files

Example of a bead level text file

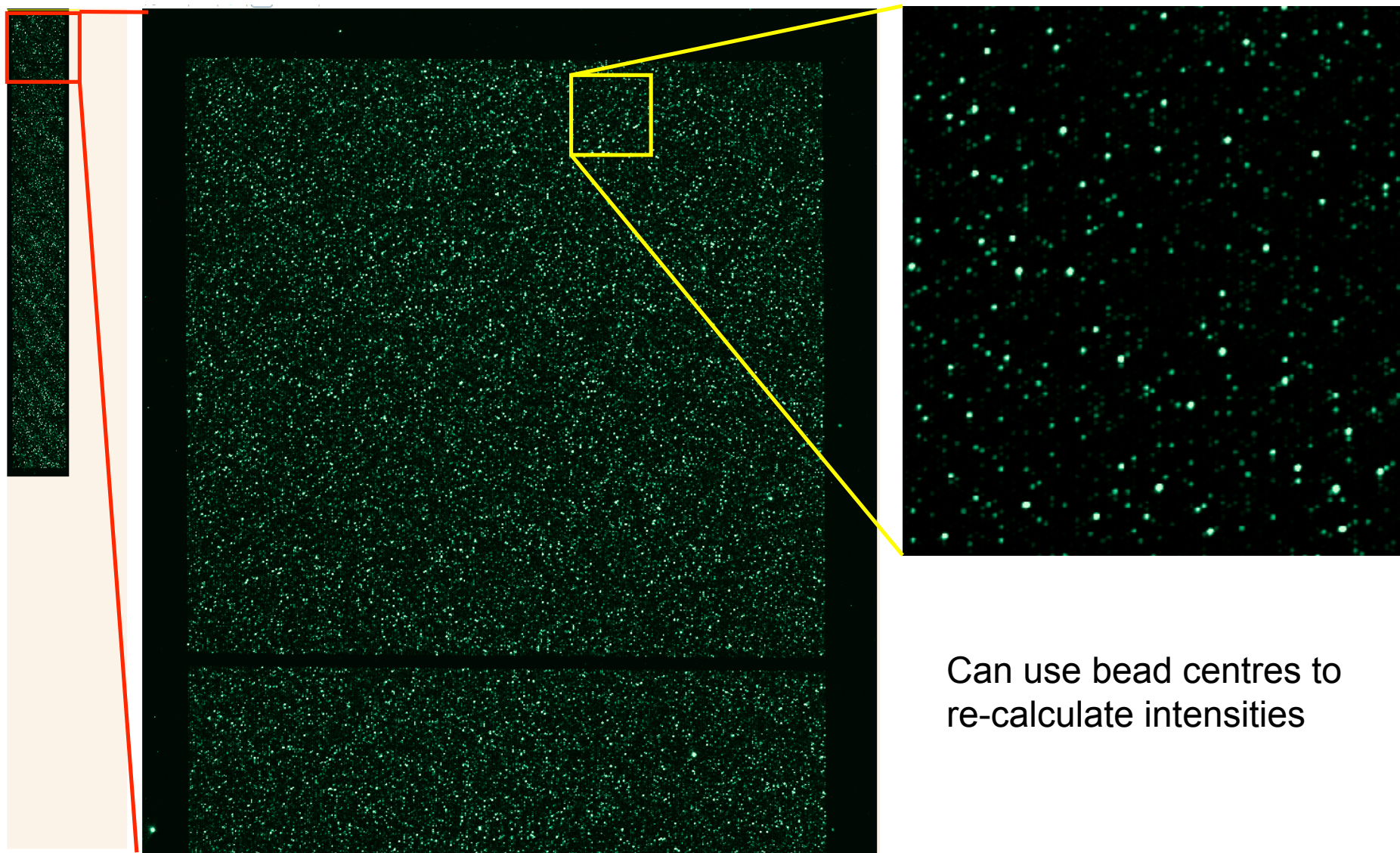| | A | B | C | D |
|---|---|---|---|---|
| | Code | Grn | GrnX | GrnY |
| 1 | Code | Grn | GrnX | GrnY |
| 2 | 2 | 1686 | 405.9445 | 994.7201 |
| 3 | 2 | 2148 | 1485.263 | 465.5954 |
| 4 | 2 | 2391 | 981.7433 | 710.9218 |
| 5 | 2 | 1961 | 414.4303 | 895.2175 |
| 6 | 2 | 2477 | 1026.212 | 942.4114 |
| 7 | 2 | 2659 | 720.4089 | 1370.215 |
| 8 | 2 | 1772 | 1139.226 | 816.4459 |
| 9 | 2 | 2737 | 1143.429 | 213.7267 |
| 10 | 2 | 2369 | 1110.516 | 203.423 |
| 11 | 2 | 2283 | 1483.378 | 548.7356 |
| 12 | 2 | 2371 | 895.504 | 976.541 |
| 13 | 2 | 2532 | 1667.515 | 864.9724 |
| 14 | 2 | 2558 | 1133.62 | 960.1776 |
| 15 | 2 | 1931 | 1127.286 | 1469.364 |
| 16 | 2 | 1760 | 279.3574 | 946.3187 |
| 17 | 2 | 2690 | 812.6176 | 803.8156 |
| 18 | 2 | 2583 | 1048.631 | 889.1783 |
| 19 | 2 | 2432 | 509.0219 | 1079.245 |
| 20 | 2 | 2538 | 929.3365 | 1226.301 |
| 21 | 2 | 2280 | 553.4136 | 885.7501 |
| 22 | 2 | 2077 | 714.496 | 250.4801 |
| 23 | 2 | 2551 | 536.4883 | 206.4698 |
| 24 | 2 | 1593 | 936.7546 | 543.4179 |
| 25 | 3 | 19868 | 1022.757 | 1404.977 |
| 26 | 3 | 20674 | 1398.915 | 971.864 |
| 27 | 3 | 21526 | 1333.79 | 1372.704 |

ProbeID

Corrected Intensity

Bead Centre

Information for **"all"** beads on an array (50,000 or 1 million rows)
Sometimes outliers or non-decoded beads are removed

# TIFF images



Can use bead centres to re-calculate intensities

# BeadStudio output



One set of observations (mean, se, detection etc) for each bead type. Local background correction was done and outliers removed before calculation of mean
All values are un-logged (1 - $2^{16}$)

# The 'beadarray' Library

Collection of BeadArray analysis functions written using **R**

Functions for reading SAM and BeadChip data in bead summary or bead level format

Options for image processing

Also quality control, diagnostic checks and normalisation

Compatible with Bioconductor & R packages (e.g. *limma*, *affy)*

beadarray has been part of the Bioconductor project since December 2005

Recently accepted for publication in Bioinformatics

# Why use beadarray?

Access to bead level data prior to processing by BeadStudio and re-visit image analysis

Quality control within arrays rather than just between arrays

Can be used to read expression / SNP / methylation and DASL data

Useful for those wishing to develop their own analysis methods (eg genotyping)

No need for Illumina (PC-based) analysis software

# Reading bead level data

Reading bead level data into beadarray is as easy as running the following

```
> BLData = readIllumina()
```

This reads all the bead level files that it finds in the **R** working directory and estimates the foreground and background intensities for each bead on each array using the images

Setting `useImages=FALSE` will take the corrected intensities from the text files

# Notes on `readIllumina`

Phenotypic information about the samples and metrics information provided by Illumina can also be read

`readIllumina` can read single or two-colour data from SAM or BeadChip experiments

Can take a lot of time and memory. Reading a BeadChip with image processing takes around 10 minutes and uses 2Gb RAM

Users can choose a smaller set of files to read, or choose not to repeat the image processing

# What is BLData?

BLData is a BeadLevelList object, which is an *environment* object

Information about BLData is organised into slots accessed by '@' - beadData, arrayInfo
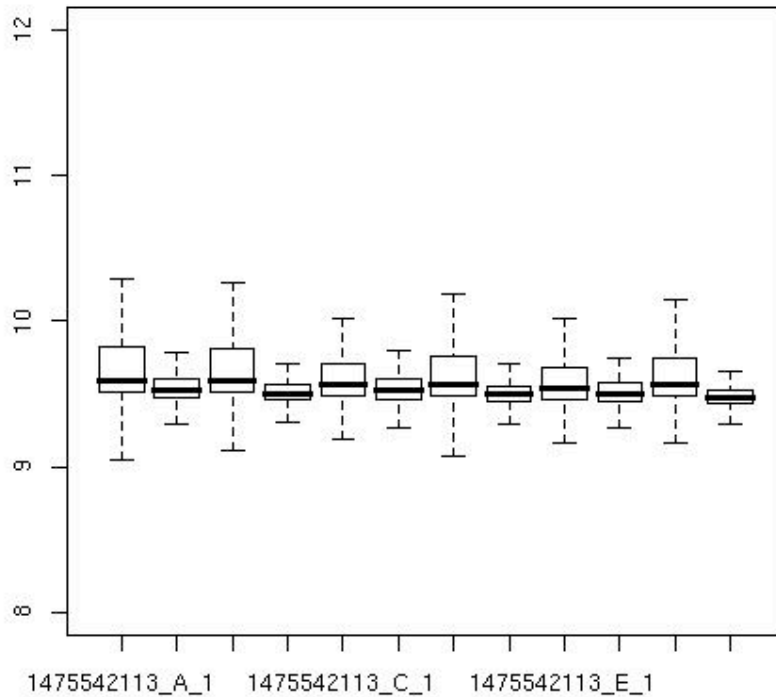
Arrays can be subset using '[['

The getArrayData function is also provided for convenience

See practical for examples

# Raw Foreground and Background

Foreground

Background



Raw foreground and background intensities from each *strip* on a BeadChip (BeadStudio merges the strips together)

The different strips can have different properties

# Compare with conventional arrays

# Background Correction

beadarray includes all the background correction methods available in limma

The default option is to simply subtract the background from the foreground for each bead



Not as many negative values as for conventional arrays (<0.01% of beads are negative with Illumina data compared with 20-30%)

# Spatial artefacts

Recall that spatial trends can be a cause for concern for microarray data

This should not be such a problem for BeadArrays due to the random positioning of beads and high number of replicates

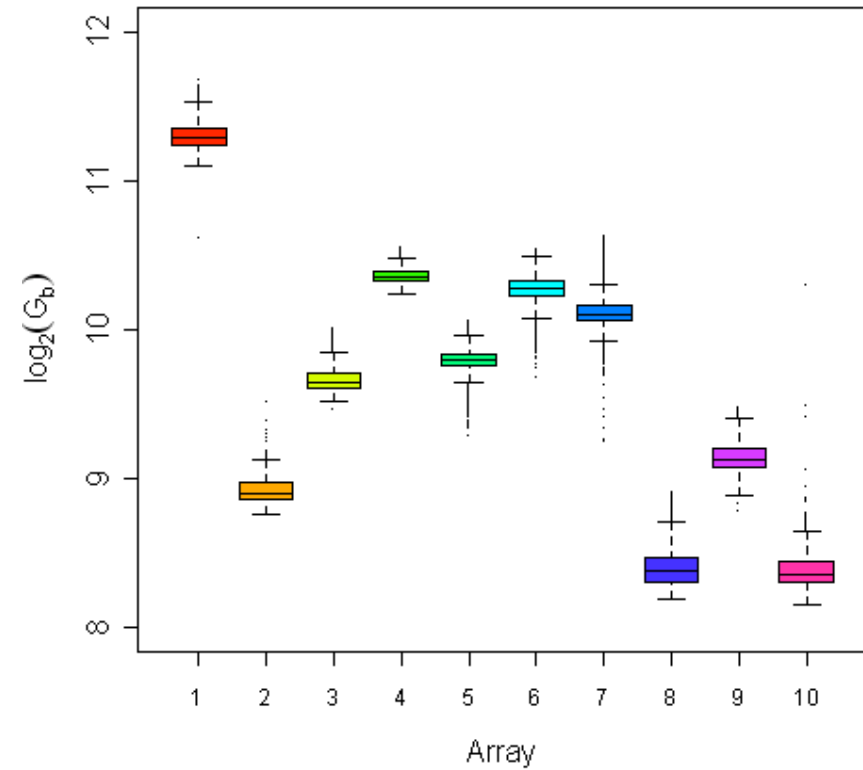beadarray includes functionality to check for serious spatial trends on arrays (as checking each array manually would be time-consuming)

The imageplot function can be used to investigate spatial trends (see practical)

# imageplot

```
>imageplot(BLData, what="G")
```

Useful things to plot include foreground, background, residuals

Spatial artefacts will often associate with outliers

This can give more detailed diagnostics for particular arrays

***All of which is not possible with summarised data***

Artefacts can be seen on original images with some effort

# Creating Bead Summary Data

We use the Illumina method to remove outliers using a 3 median absolute deviation cut-off from the median for each bead type

```
>BSData = createBeadSummaryData(BLData)
```

# Remarks

Can choose to summarise the data on the log2 scale

The resulting object BSData is an ExpressionSetIllumina object which extends an ExpressionSet. The expression matrix can be easily extracted for further analysis

If two-colour data is given, the two channels are summarised separately to give a SnpSetIllumina object

# Storing bead summary data

`BSData` is now an `ExpressionSetIllumina` object sharing many common properties with other Bioconductor objects

The expression values can be accessed using the `exprs` function

```
> E=exprs(BSData)
> dim(E)
[1] 47293    18
> E[1:3,1:3]
             AVG_Signal.IH.1 AVG_Signal.IC.1 AVG_Signal.IH.2
GI_10047089-S            87.8           131.8           231.9
GI_10047091-S           161.8           130.8           258.6
GI_10047093-S           481.2           401.4           499.4
```

For more details see the practical…

# readBeadSummaryData

```
>BSData = readBeadSummaryData(dataFile, qcInfo,
sampleSheet)
```
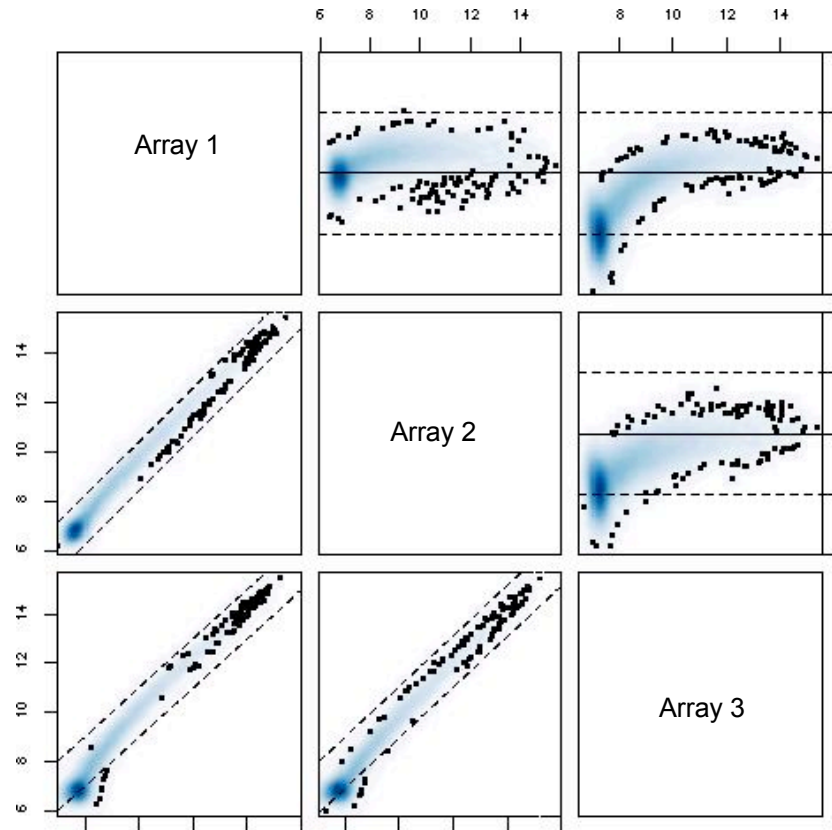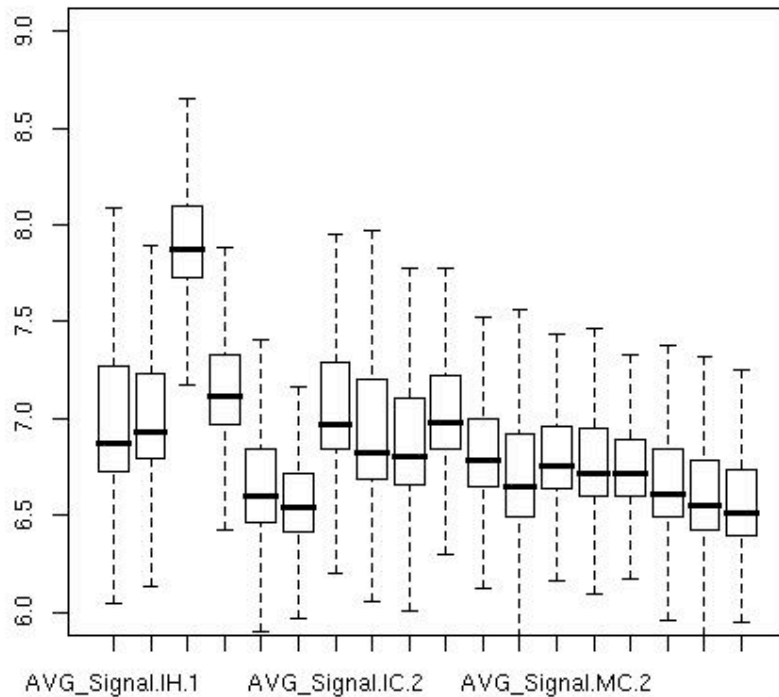
**Warning:** May need to change skip and sep parameters depending on version of BeadStudio

Eg

```
>BSData = readBeadSummaryData(dataFile, qcInfo,
sampleSheet, skip=7, sep=",")
```

Also, column headings in BeadStudio output sometimes change (BEAD_STDEV -> BEAD_STDERR)

# Quality Control



```
> boxplot(as.data.frame(log2(E)),outline=FALSE, ylim=c(6,9))


> plotMAXY(E, arrays=1:3)
```
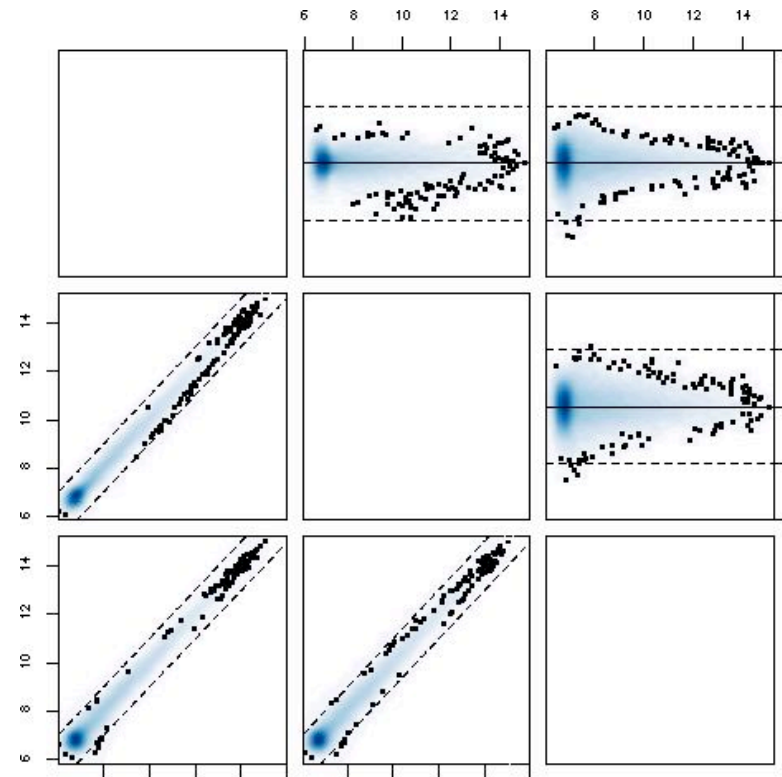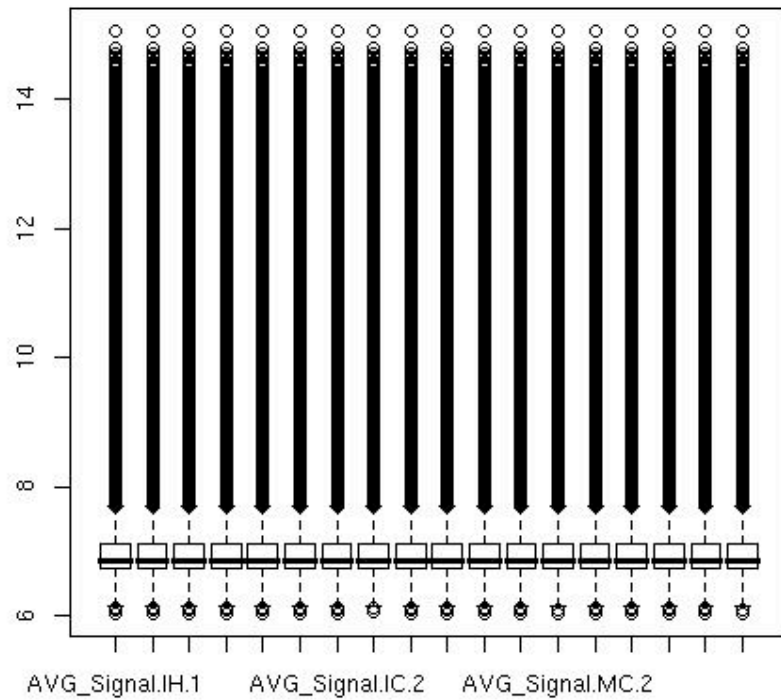
# Normalisation

Illumina data seems to be of good quality, however some trends can still be seen in the data (eg decrease in intensity across a chip)

Important not to remove any biological effects

Quantile normalisation seems to be effective

Or any other normalisation method from Bioconductor which can be used on an expression matrix. Many can be found in the affy package

```
>E = normaliseIllumina(BSData, method="quantile", transform="log2")
>boxplot(as.data.frame(log2(E)), outline=TRUE)
>plotMAXY(E, arrays=1:3)
```

And now over to the practical….

# Acknowledgements

Computational Biology Group (Cambridge)
**Matthew Ritchie**
Mike Smith
Julie Addison
Natalie Thorne
Andy Lynch
Nuno Barbosa-Morais
Simon Tavaré

Dermitzakis group (Sanger Institute - Cambridge)
Manolis Dermitzakis
Barbara Stranger
Matthew Forrest

Caldas Lab
Inma Spiteri
Anna Git

Illumina (San Diego)
Brenda Kahl
Semyon Kruglyak
Gary Nunn

UCSD(San Diego)
Roman Sasik

CANCER RESEARCH UK

MRC Medical Research Council