

Lecture 3: Microarray Technology

BIOINF3005/7160: Transcriptomics Applications

Dr Stephen Pederson

Bioinformatics Hub,
The University of Adelaide

March 23rd, 2020

Microarray Technology

Two Colour Microarrays

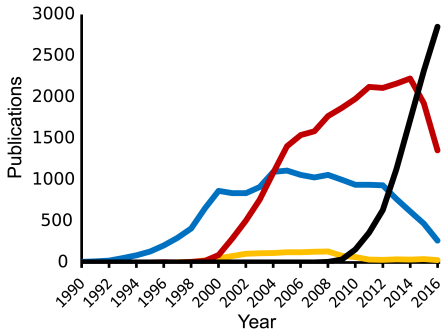
Single Channel Microarrays

Whole Transcript Arrays

Hypothesis Testing

Microarray Technology

Microarrays



EST (blue); SAGE / CAGE (yellow); Microarrays (red); RNA Seq (black)¹

¹Rohan Lowe et al. "Transcriptomics technologies". In: *PLOS Computational Biology* 13.5 (May 2017), pp. 1–23. DOI: 10.1371/journal.pcbi.1005457. URL: <https://doi.org/10.1371/journal.pcbi.1005457>.

Microarrays

- Microarrays effectively ushered in the modern era of transcriptomics
- Purely interested in *relative abundances*
- Could measure expression levels for 1000's of genes simultaneously, for *the first time*
- Were essentially glass slides with probes affixed to them

Microarrays

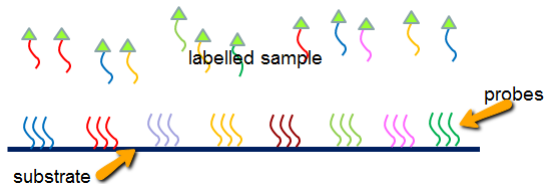
- Once again depends on reverse transcriptase for mRNA → cDNA
- **No reliance on Sanger Sequencing**
- Used probes (like a Northern blot) but the **cDNA is labelled and the probes are spatially fixed**
 - Probes must be designed beforehand
 - Probes are fixed to the array in *known locations*

Microarrays

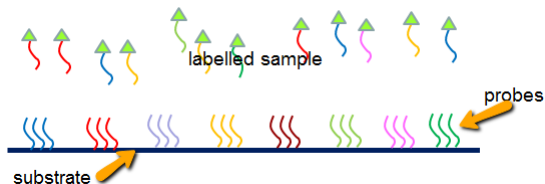
1. Fluorescent labelling during mRNA conversion to cDNA
2. Complimentary probes bind target sequences (hybridisation)
3. Fluorescence detection at each probe

Fluorescence Intensity \propto mRNA abundance

Microarrays



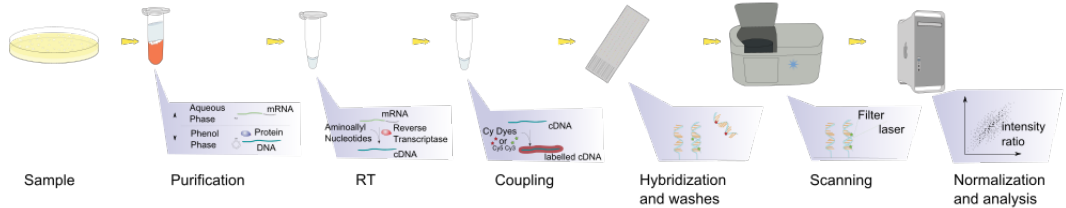
Microarrays



Highly abundant targets will yield more signal after hybridisation



Microarrays

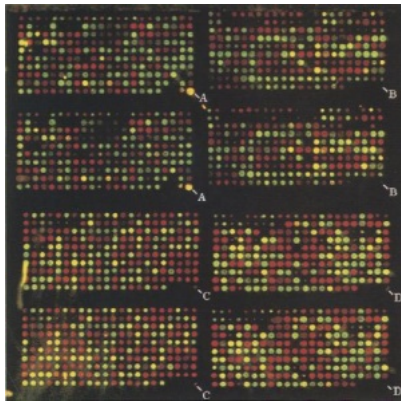


Two Colour Microarrays

Two Colour Microarrays

- Sometimes called “Low-Density Oligo Microarrays”
- Probes with known sequences are at known locations
 - Probes were 60-75mer complimentary cDNA
 - Originally printed in local facilities
- Samples are labelled with *either* Cy3 (Green @ 570nm) or Cy5 (Red @ 670nm)
- Two samples are hybridised to each array
 - Competitive hybridisation
 - Relative Red/Green intensities were of interest
 - Gave an estimate of logFC within each array

Two Colour Microarrays



A section of a two colour array²

²D Shalon, S J Smith, and P O Brown. "A DNA microarray system for analyzing complex DNA samples using two-color fluorescent probe hybridization." In: *Genome Research* 6.7 (1996), pp. 639–645. DOI: 10.1101/gr.6.7.639. eprint: <http://genome.cshlp.org/content/6/7/639.full.pdf+html>. URL: <http://genome.cshlp.org/content/6/7/639.abstract>.

Two Colour Microarrays

- Probes are “printed” to the array
 - Print tips can get clogged and be uneven
- Able to be customised for your own experiment
 - A mapping file for probe location to target sequence is required
- Both colours were scanned separately
 - One scan detects red only, the next detects green only
 - Each individual scan would have to be aligned spatially with the other

Two Colour Microarrays

- Spots were detected using astronomical software
 - Sizes were variable / irregular
- Detection of true signal above background (DABG)
 - Required “identified” (foreground) pixels and surrounding (background) pixels
 - Used surrounding pixels to estimate BG
 - Assumed BG was additive, e.g. $R = R_{bg} + R_{fg}$
- Dye bias was also noted \implies experiments often used dye swaps
 - A sample from “group 1” might be labelled with red on one array, then labelled with green on the next

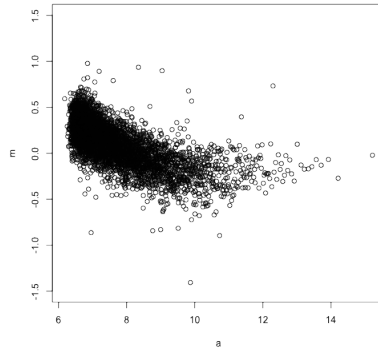
Two Colour Microarrays

- All intensities are transformed to the \log_2 scale
- Dye bias was checked using “MA Plots”
 - M was the *difference in intensity* across both channels
 - A was the *average intensity* across both channels

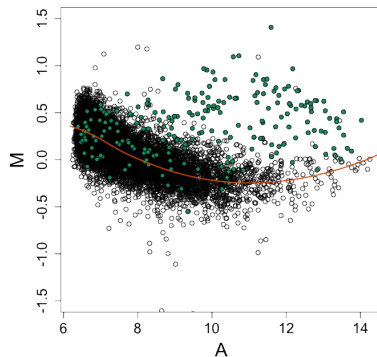
$$M = \log_2 R - \log_2 G$$

$$A = \frac{\log_2 R + \log_2 G}{2}$$

Two Colour Microarrays



Two Colour Microarrays

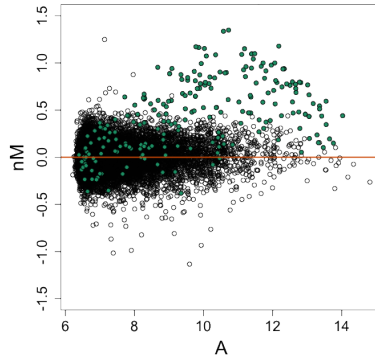


We can fit a **loess** curve through the data
(Here, spike-in controls are also highlighted)

Two Colour Microarrays

- loess: Locally estimated scatterplot smoothing
 - We use a sliding window and fit a polynomial line
 - Usually polynomial of order 1 (linear) or 2 (quadratic)
- Once we have the loess curve: we subtract it from the data
 - Explicitly assumes that the bulk of the difference is bias, i.e. *most genes are not differentially expressed*
 - No modification to the A values, or any R/G intensities

Two Colour Microarrays



No more dye bias ...

Two Colour Microarrays

- We use these normalised M values across arrays to estimate logFC
- Dye-swap complications \implies *Experimental Design*
- Robust suite of statistical tools developed from here
- The R package `limma` set the standard

Single Channel Microarrays

Single Channel Microarrays

- Affymetrix 3' Arrays became the dominant technology (until RNA seq)
- Probes target the 3' end of transcripts \implies intact transcripts
- Single channel (i.e. single colour) \implies one sample per array
- $\sim 1,000,000 \times 25$ -mer probes

Fluorescence Intensity \propto mRNA abundance

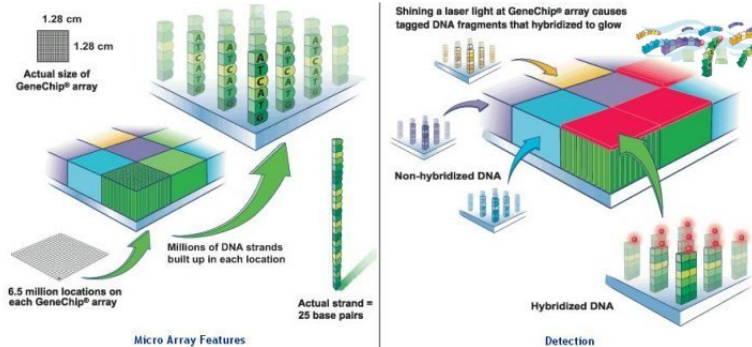
Single Channel Microarrays



Single Channel Microarrays

- Manufacture used photolithography
- Far greater density of probes than two-colour arrays
 - Shorter probes but far more of them
- Fixed array designs for each “model” and organism
- Probes designed based on known gene annotations at design-time
- Also need a mapping file from location to probe sequence

Single Channel Microarrays



3' Arrays

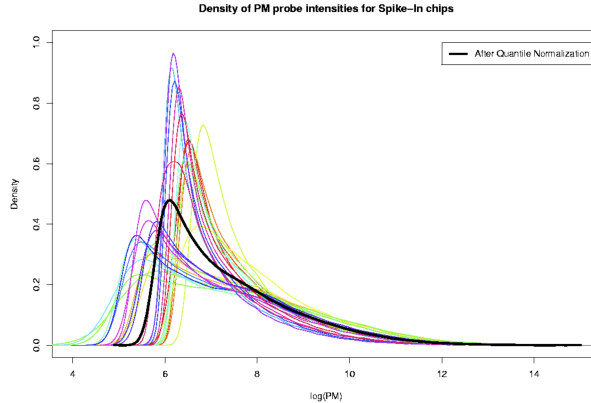
- Each 3' exon would be targeted by 11 unique probes
 - The set of **11 probes** would be collected together as a single **probeset**
- Alternate isoforms with different 3' exons could be detected easily as they would have distinct probesets
- Need a *Chip Description File* to map probes to array coordinates and probesets

3' Arrays

Key Technical Issues:

1. Differences between **arrays**
 - Hybridisation artefacts, cDNA/RNA concentration artefacts
2. Background Correction at the **probe** level
 - 25-mer probes \implies *non-specific binding*
 - Optical Background
3. Expression estimates at the **probeset** level
 - Some probes *unresponsive*, other probes *promiscuous*
 - Do you just *average them*?

Normalisation

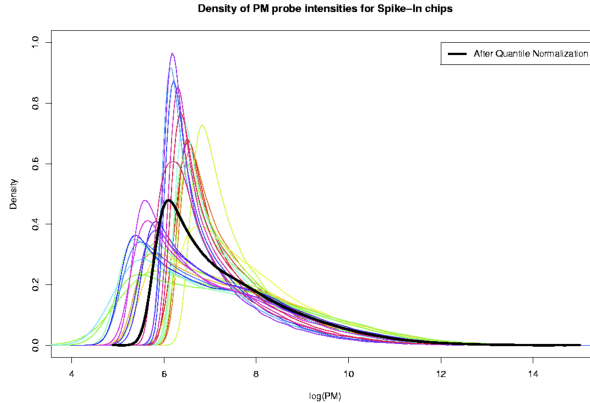


Quantile Normalisation

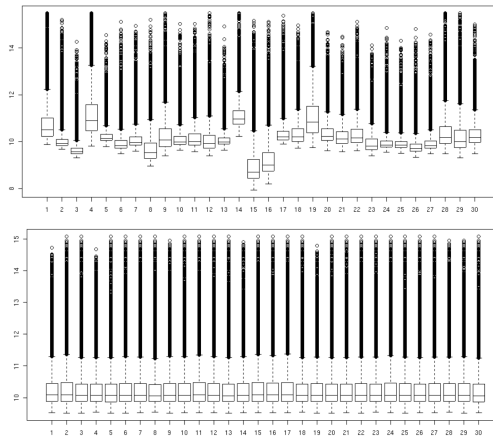
1. Find the probe with the lowest intensity on each array
 - This will be from different probesets and unrelated to each other
2. Find the average intensity across these probes
3. Assign this value to each probe
4. Repeat for the probes with the next lowest intensity until done
5. All arrays now have the same intensity distribution

Under this approach, **we are adjusting the raw intensities**

Quantile Normalisation



Quantile Normalisation



Background Correction

- Probes targeting 3' exons: Perfect Match (*PM*) probes
- Probes with middle base changes: MisMatch (*MM*) probes
- *MM* probes were expected to capture similar *NSB* behaviours to paired *PM* probe
 - Were often **brighter** than *PM* probes in pair
- Literally **half** of the array was *MM* probes

Background Correction

For a given PM/MM probe pair

$$PM = B + S$$

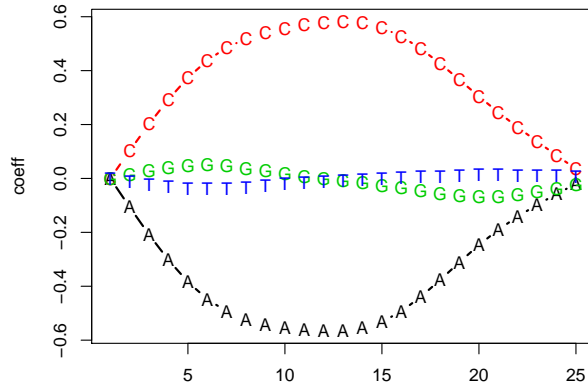
but $\dots MM \neq B$

- How do we estimate S ?
- $S \geq 0$

Background Correction

- Found $\hat{S} = E[S|PM]$ using a convolution of normal and exponential distributions (*RMA*)
- GC content and position in probe also impacted *NSB* \implies *GC-RMA*
- No need for the *MM* probe as a pair
 - *MM* probes still used in estimation of parameters

Background Correction



Probeset Summarisation

- Probes $j = 1, 2, \dots, 11$ need to be combined (summarised) within a **probeset**
 - This gives the **gene-level expression estimates** for **each array**
 - Poor performing probes were generally poor on all arrays
 - Promiscuous probes were general similar on all arrays
- Probe-level modelling gave μ_i for each array i
 - The model was fit robustly \implies outlier signal is down-weighted
 - Using $Y_{ij} = \log_2 \hat{S}_{ij}$:

$$Y_{ij} = \mu_i + \alpha_j + \varepsilon_{ij}$$

Now we have a single, gene-level estimate of expression for each array: $\hat{\mu}_i$

Analysis

- For each gene we take $\hat{\mu}_i$ and fit a linear model, conduct a t-test etc
- We will deal with the statistics very soon (FUN!)

Analysis

The basic process for single channel arrays:

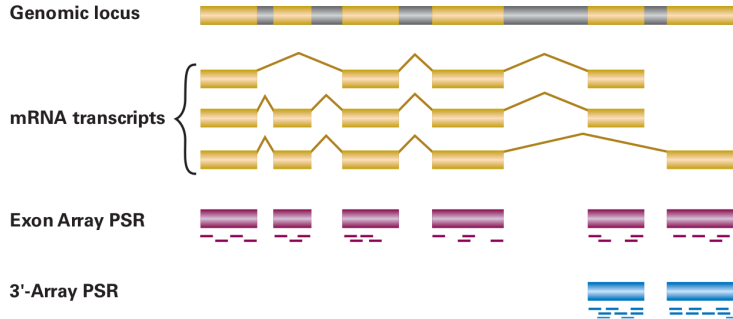
1. Normalise for technical differences
2. Find probe-level estimates of *true signal*
3. Obtain gene-level estimates of signal
4. Statistical Analysis across all genes

Whole Transcript Arrays

Whole Transcript Arrays

- The second generation of Affymetrix arrays were Gene/Exon Arrays
- Far greater density of probes (~ 5 -6 fold)
 - No *MM/PM* pairs
 - Antigenomic and MisMatch probe groups
- These target the whole transcript (WT), **NOT** just the 3' end
- How does RNA degradation impact this?
- How does alternate splicing impact this?

Whole Transcript Arrays



Whole Transcript Arrays

RNA Degradation

- 3' Arrays had 11 probes targeting the 3' end
- Easily comparable across genes to assess RNA quality
- Not the case for WT Arrays

Alternate Splicing

- Identifying the correct transcript remained largely unsolved
- Some exons may be missing
 - No true signal *implies* biases expression estimates down
 - Can appear as changes in expression, e.g. a short transcript in one condition will yield a lower expression estimate than a long transcript

Whole Transcript Arrays

- Before these technical problems were solved RNA Seq “exploded”
- How do we separate differential expression
 - i.e. changes in transcriptional activity and regulation
- from alternate isoform usage
 - e.g. changes in the dominant isoform, alternate promoter usage
- Many genes exist in *multiple isoforms in the same tissue*

These still remain (somewhat) unsolved in RNA Seq

Whole Transcript Arrays

- Exon Arrays disappeared very quickly
- Gene Arrays are still in active use (Cheap)
- Both are limited to genes/transcripts defined at time of array design
- Novel transcripts, retained introns etc **cannot** be detected

Hypothesis Testing

Hypothesis Testing

In biological research we often ask:

“Is something happening?” or “Is nothing happening?”

We might be comparing:

- Cell proliferation in response to antibiotics in media
- mRNA abundance in two related cell types
- Methylation levels across genomic regions
- Allele frequencies in two populations

Hypothesis Testing

How do we decide if our experimental results are “significant”?

- Is it normal variability?
- What would the data look like if our *experiment had no effect*?
- What would our data look like if there was *some kind of effect*?

Every experiment is considered as a random sample from all possible repeated experiments.

Sampling

Most experiments involve measuring something:

- Discrete values e.g. read counts, number of colonies
- Continuous values e.g. Ct values, fluorescence intensity

Every experiment is considered as a random sample from all possible repeated experiments.

Sampling

Many data collections can also be considered as experimental datasets

Example 1

In the 1000 Genomes Project a risk allele for T1D has a frequency of $\pi = 0.07$ in European Populations.

Does this mean, the allele occurs in exactly 7% of Europeans?

Sampling

Example 2

In our in vitro experiment, we found that 90% of HeLa cells were lysed by exposure to our drug.

- Does this mean that exactly 90% of HeLa cells will always be destroyed?
- What does this say about in vivo responses to the drug?

Population Parameters

- Experimentally-obtained values represent an **estimate** of the true effect
- More formally referred to as *population-level parameters*
- Every experiment is considered a *random sample of the complete population*
- Repeated experiments would give a **different** (*but similar*) estimate

Population Parameters

- Experimentally-obtained values represent an **estimate** of the true effect
- More formally referred to as *population-level parameters*
- Every experiment is considered a *random sample of the complete population*
- Repeated experiments would give a **different** (*but similar*) estimate

All population parameters are considered to be fixed values, e.g.

- Allele frequency (π) in a population
- The average difference in mRNA levels

The Null Hypothesis

All classical statistical testing involves:

1. a Null Hypothesis (H_0) and
2. an Alternative Hypothesis (H_A)

Why do we do this?

The Null Hypothesis

- We define H_0 so that we know what the data will look like if there is no effect
- The alternate (H_A) includes every other possibility besides H_0

An experimental hypothesis may be:

Example

$$H_0 : \mu = 0 \text{ Vs } H_A : \mu \neq 0$$

Where μ represents the *true average difference in a value* (e.g. mRNA expression levels)