# Microarray Informatics

## Donald Dunbar

MSc Seminar
3rd February 2010

---

## Aims

- To give a biologist's view of microarray experiments
- To explain some technologies involved
- To describe typical microarray experiments
- To show how to get the most from and experiment
- To show where the field is going

---

## Introduction

- Part 1
  - Microarrays in biological research
  - A typical microarray experiment
  - Experiment design, data pre-processing
- Part 2
  - Data analysis and mining
  - Microarray standards and resources
  - Recent advances

---

## Microarray Informatics

# Part 1

---

## Biological research

- Using a wide range of experimental and computational methods to answer biological questions
- Genetics, physiology, molecular biology…
- Biology and informatics → bioinformatics
- Genomic revolution
- What can we measure?

---

## The central dogma



promoter exon intron exon intron intron exon

30k
Gene: DNA

90k
Transcript: mRNA

100+k
Protein

**kinase, protease, structural receptor, ion channel…**

## Measuring RNA and proteins

- Proteins
  - Western blot
  - ELISA
  - Enzyme assay
- mRNA
  - Northern blot
  - RT-PCR

## Measuring RNA and proteins

- Protein levels/activities would be best
  - no real high throughput method
- mRNA levels will have to do
  - genome-wide physical microarrays
  - other 'array-like' technologies
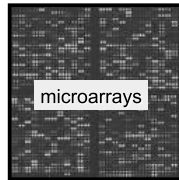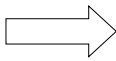  - sequencing (see later)

## Measuring transcripts

- Genome level sequencing
- New miniaturisation technologies
- Better bioinformatics

microarrays

## Microarrays: wish list

- Include all genes in the genome
- Include all splice variants
- Give reliable estimates of expression
- Easy to analyse
  - bioinformatics tools available
- Cost effective
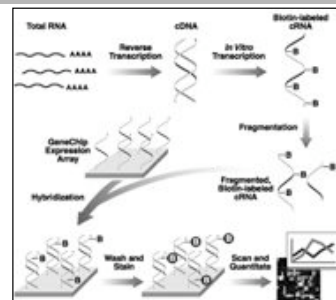
## Microarray technologies - 1

- Oligonucleotides - Affymetrix
- One chip all genes
- Chips for many species
- Several oligos per transcript
- Use of control, mismatch sequences
- One sample per chip
  - 'absolute quantification'
- Well established in research
- Expensive

AFFYMETRIX

## Microarray technologies - 1



AFFYMETRIX

## Microarray technologies - 2

- Illumina BeadChip
- Oligos on beads
- Hybridise in wells
- Compared to Affy
  - Higher throughput
  - Less RNA needed
  - Cheaper

## Problems with microarrays

- The gene might not be on the chip
- Can't differentiate splice variants
- The gene might be below detection limit
- Can't differentiate RNA synthesis and degradation
- Can't tell us about post translational events
- Bioinformatics can be difficult
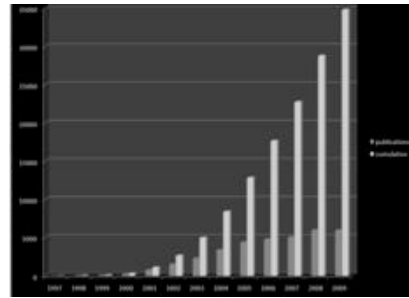- Relatively expensive

## History of Microarrays

- Developed in early 1990s after larger macro-arrays (100-1000 genes)
- Microarrays were spotted on glass slides
- Labs spotted their own (Southern, Brown)
- Then companies started (Affymetrix, Agilent)
- Some early papers:
  - *Nature* 1993 364(6437): 555-6 Multiplexed biochemical assays with biological chips. Fodor SP, et al
  - *Science* 1995 Oct 20;270(5235):467-70 Quantitative monitoring of gene expression patterns with a complementary DNA microarray. Schena M, et al

## Microarray publications

## Types of experiment

- Usually **control** v **test(s)**

| | |
|---|---|
| Placebo | Drug treatment     Drug 2… |
| Wild-type | Knockout |
| Healthy | Patient |
| Normal tissue | Cancerous tissue |
| Time = 0 | Time = 1     Time = 2… |

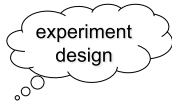## Types of experiment

- Usually **control v test(s)**
- But also **test v test(s)**
- Comparison:
  - placebo v drug treatment
  - drug 1 v drug 2
  - tissue 1 v tissue 2 v tissue 3 (pairwise)
  - time 0 v time 1, time 0 v time 2, time 0 v time 3
  - time 0 v time 1, time 1 v time 2, time 2 v time 3

## A typical experiment
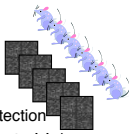
experiment design

## Experiment design: system

- What is your model?
  - animal, cell, tissue, drug, time…
- What comparison?
- What platform
  - microarray? oligo, cDNA?
- Record all information: see "standards"
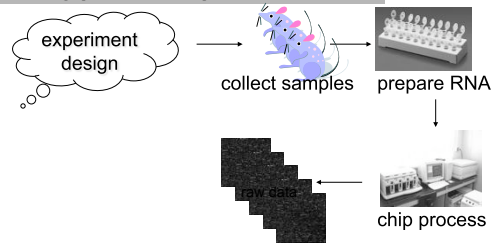
## Experiment design: replicates

- Microarrays are noisy: need extra confidence in the measurements
- We usually don't want to know about a specific individual
  - eg not an individual mouse, but the strain
  - although sometimes we do (eg people)
- Biological replicates needed
  - independent biological samples
  - number depends on variability and required detection
- Technical replicates (same sample, different chip) usually not needed

## A typical experiment

experiment design → collect samples → prepare RNA

raw data ← chip process

## Raw data

- Affymetrix GeneChip process generates:
  - DAT    image file
  - CEL    raw data file
  - CDF    chip definition file
- Processing then involves CEL and CDF

- Will use Bioconductor
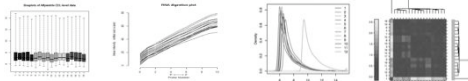
## Bioconductor (BioC)

- http://www.bioconductor.org/
- "Bioconductor is an open source software project for the analysis and comprehension of genomic data"
- Started 2001, developed by expert volunteers
- Built on statistical programming environment "R"
- Provides a wide range of powerful statistical and graphical tools

- Use BioC for most microarray processing and analysis
- Most platforms now have BioC packages
- Tutorial: manuals.bioinformatics.ucr.edu/home/R_BioCondManual

## Quality control (QC)

- Affymetrix gives data on QC
  - the microarray team will record these for you
  - scaling factor, % present, spiked probes, internal controls
- Bioconductor offers:
  - boxplots and histograms of raw and normalised data
  - RNA degradation plots
  - specialised quality control routines (eg arrayQualityMetrics)



February 3rd 2010          MSc Seminar: Donald Dunbar
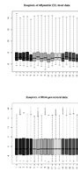
## Pre-processing: background

- Signal corresponds to expression…
  - plus a non-specific component (noise)
- Non specific binding of labelled target
- Need to exclude this background
- Several methods exist
  - eg Affy: PM-MM but many complications
  - eg RMA PM=B+S (don't use MM)

February 3rd 2010          MSc Seminar: Donald Dunbar
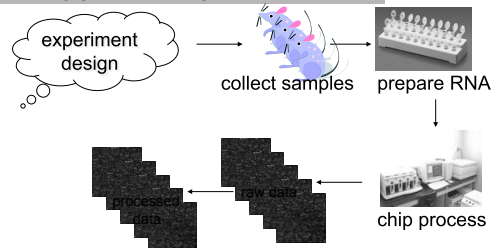
## Pre-processing: normalisation

- In addition to background corrections
  - chip, probe, spatial, intra and inter variation
  - need to remove to get at real expression differences
- Make use of statistics

combined with probeset summary:
get an expression value for the gene

- But seems to be non-linear dependency on intensity
  - additive and multiplicative errors
- Quantile normalisation often used
- Normalisation more complicated for 2-colour arrays
- Try to remove most noise at lab stage (ie control things well statistically)

*Now carry on with analysis!*

February 3rd 2010          MSc Seminar: Donald Dunbar

## A typical experiment



experiment design → collect samples → prepare RNA → chip process → raw data → processed data

February 3rd 2010          MSc Seminar: Donald Dunbar

## Part 1 Summary

- Microarrays in biological research
- Two types of microarray
- A typical microarray experiment
- Experiment design
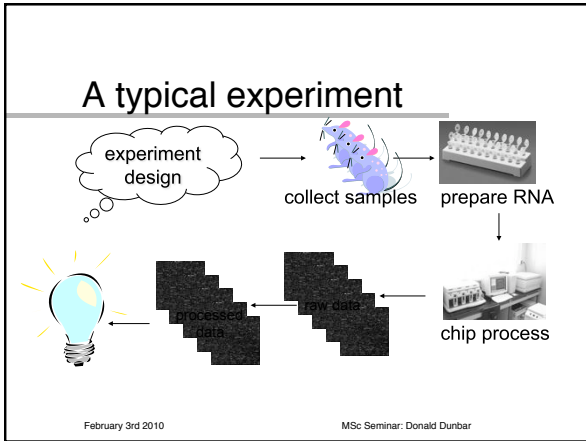- Data pre-processing

February 3rd 2010          MSc Seminar: Donald Dunbar

## Microarray Informatics

# Part 2

February 3rd 2010          MSc Seminar: Donald Dunbar

## A typical experiment

experiment design → collect samples → prepare RNA

processed data ← raw data ← chip process

## Data analysis

- Identifying differential expression
- Compare control and test(s)
  - t-test
  - ANOVA
  - SAM (FDR)
  - Limma
  - Rank Products
- Time series

control          treated

v

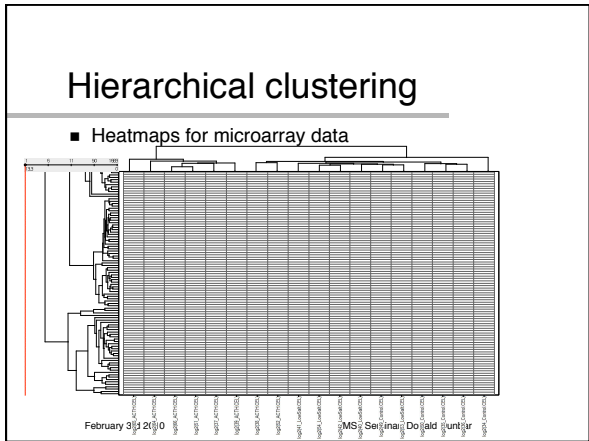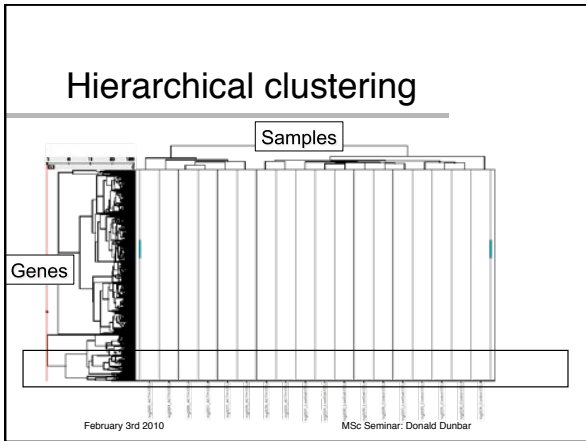0      1      2      3

v      v      v

## Multiple testing

- Problem:
  - statistical testing of 30,000 genes
  - at $\alpha = 0.05 \rightarrow 1500$ genes
- Need to correct this
  - Multiply p-value by number of observations
    - Bonferroni, too conservative
  - False discovery
    - defines a q value: expected false positive rate
    - Less conservative, but higher chance of type I error
    - Benjamini and Hochberg
- Then regard genes as differentially expressed
- Depends on follow-up procedure!

## Hierarchical clustering

- Look for structure within dataset
  - similarities between genes
- Compare gene expression profiles
  - Euclidian distance
  - Correlation
  - Cosine correlation
- Calculate with distance matrix
- Combine closest, recalculate, combine closest… (or split!)
- Draw dendrogram and heatmap

## Hierarchical clustering

Samples

Genes

## Hierarchical clustering

- Heatmaps for microarray data

## Hierarchical clustering

- Predicting association of known and novel genes
- Class discovery in samples: new subtypes
- Visualising structure in data (sample outliers)
- Classifying groups of genes
- Identifying trends and rhythms in gene expression
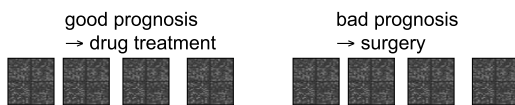- Caveat: you will always see clusters, even when they are not particularly meaningful (nb Ian Simpson)

## Sample classification

- Supervised or non-supervised
- Non-supervised
  - like hierarchical clustering of samples
- Supervised
  - have training (known) and test (unknown) datasets
  - use training sets to define robust classifier
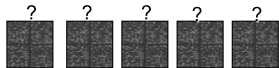  - apply to test set to classify new samples

## Sample classification

good prognosis
→ drug treatment

bad prognosis
→ surgery

Gene selection, training, cross validation →
classifier: gene x * 0.5 gene y * 0.25 gene z …

?   ?   ?   ?   ?

## Sample classification

good prognosis
→ drug treatment

bad prognosis
→ surgery

Apply classifier

## Sample classification

- Class prediction for new samples
  - cancer prognosis
  - pharmacogenomics (predict drug efficacy)
- Need to watch for overfitting
  - using too much of the data to classify
  - classifier loses specificity

## Annotation

- Big problem for microarrays
- Genome-wide chips need genome-wide annotation
- Good bioinformatics essential
  - use several resources (Affymetrix, Ensembl)
  - keep up to date (as annotation changes)
  - genes have many attributes
    - name, symbol, gene ontology, pathway…

## Data-mining

# Microarrays are a waste of time

# …unless you do something with the data

## Data-mining

- Once data are statistically analysed:
  - pull out genes of interest
  - pull out pathways of interest
  - mine data based on annotation
    - what are the expression patterns of these genes
    - what are the expression patterns in this pathway
  - mine genes based on expression pattern
    - what types of genes are up-regulated …
    - fold change, p-value, expression level, correlation
- Should be driven by the biological question

Gene set
Enriched regulatory sequences
Functional significance?

February 3rd 2010          MSc Seminar: Donald Dunbar



PubMatrix Results 13th September 2005

February 3rd 2010          MSc Seminar: Donald Dunbar



## Further data-mining

- Other tools available using
  - gene ontology (GO)
  - biological pathways (eg KEGG)
  - genomic localisation (Ensembl)
  - regulatory sequence data (Toucan, BioProspector)
  - literature (eg Pubmatrix, Ingenuity…)
- … to make sense of the data
- Links at: www.bioinf.mvm.ed.ac.uk/projects/analysis_tools.html

February 3rd 2010          MSc Seminar: Donald Dunbar

9

## Microarray Resources

- Microarray data repositories
  - Array express (EBI, UK)
  - Gene Expression Omnibus (NCBI, USA)
  - CIBEX (Japan)
- Annotation
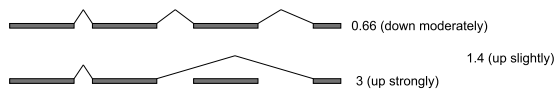  - NetAffx, Ensembl, TIGR, Stanford…

## Microarray Standards

- MIAME
  - Minimum annotation about a microarray experiment
  - Comprehensive description of experiment
  - Models experiments well, and allows replication
    - chips, samples, treatments, settings, comparisons
  - Required for most publications now
- MAGE-ML
  - Microarray gene expression markup language
  - Describes experiment (MIAME) and data
  - Tools available for processing

## Recent advances: Exon chips

- Affymetrix now have chips that allow us to measure expression of splice variants

0.66 (down moderately)

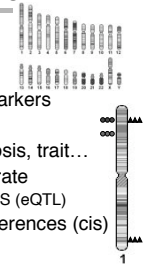1.4 (up slightly)

3 (up strongly)

New chips will give us much more information

## Recent advances: Genotyping chips

- All discussion on EXPRESSION chips
- Also can get chips looking at genotype
- Tell us the sequence for genome-wide markers
- Test 300,000 markers with one chip
- Look for association with disease, prognosis, trait…
- Combined with expression chips to generate
  - EXPRESSION QUANTITATIVE TRAIT LOCUS (eQTL)
  - Overlap of expression and genetic differences (cis)
  - Correlation at different locus (trans)

1

## Next Generation Sequencing

- Sequence rather than hybridisation
- Gene expression, genotyping, epigenetics
- New technologies: much cheaper than before
- Gene expression, genotyping, epigenetics
- Open ended (no previous knowledge required)
- Will take over in 2 years: the end of microarrays?

## Part 2 Summary

- Data analysis
- Data Mining
- Microarray Resources
- Microarray Standards
- Recent & future advances

## Seminar Summary

- Part 1
  - Microarrays in biological research
  - A typical microarray experiment
- Part 2
  - Data analysis and mining
  - Recent & future advances

## Contact



- Donald Dunbar
- QMRI Bioinformatics
- donald.dunbar@ed.ac.uk
- 0131 242 6700
- Room W3.01, QMRI, Little France
- www.bioinf.mvm.ed.ac.uk