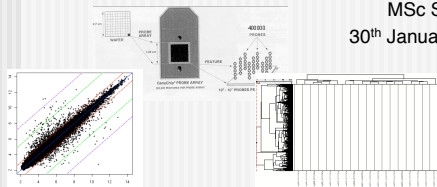


Microarray Informatics

Donald Dunbar

MSc Seminar
30th January 2008



Aims

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

January 30th 2008

MSc Seminar: Donald Dunbar

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

January 30th 2008

MSc Seminar: Donald Dunbar

Microarray Informatics

Part 1

January 30th 2008

MSc Seminar: Donald Dunbar

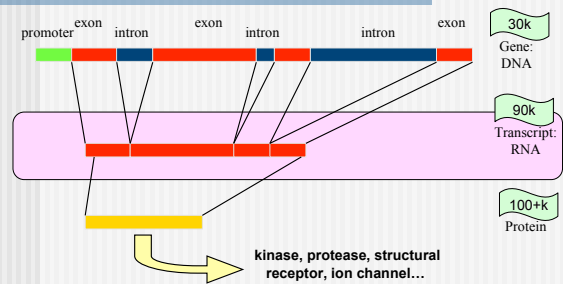
Biological research

- 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?

January 30th 2008

MSc Seminar: Donald Dunbar

The central dogma

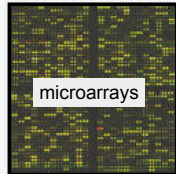


January 30th 2008

MSc Seminar: Donald Dunbar

Measuring transcripts

- Genome level sequencing
- New miniaturisation technologies
- Better bioinformatics

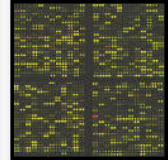


January 30th 2008

MSc Seminar: Donald Dunbar

Microarrays: wish list

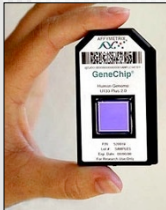
- Include all genes in the genome
- Include all splice variants
- Give reliable estimates of expression
- Easy to analyse
 - bioinformatics tools available
- Cost effective



January 30th 2008

MSc Seminar: Donald Dunbar

Microarray technologies - 1



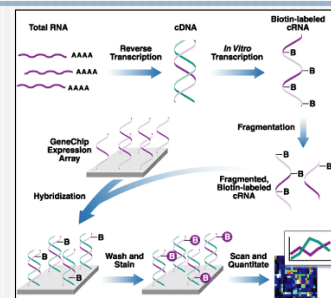
- Oligonucleotides - Affymetrix
- One chip all genes
- Chips for many species
- Several oligos per transcript
- Use of control, mismatch sequences
- One sample per chip
 - absolute quantification
- Well established in research
- Expensive



January 30th 2008

MSc Seminar: Donald Dunbar

Microarray technologies - 1



January 30th 2008

MSc Seminar: Donald Dunbar

Microarray technologies - 2

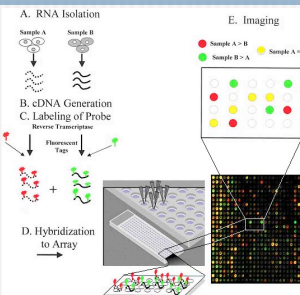
- cDNAs - Agilent
- One chip all genes
- Chips for many species
- One cDNA per transcript
 - quantification of ratios
- Established in research
- Expensive



January 30th 2008

MSc Seminar: Donald Dunbar

Microarray technologies - 2



January 30th 2008

MSc Seminar: Donald Dunbar

Problems with transcriptomics

- 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

January 30th 2008

MSc Seminar: Donald Dunbar

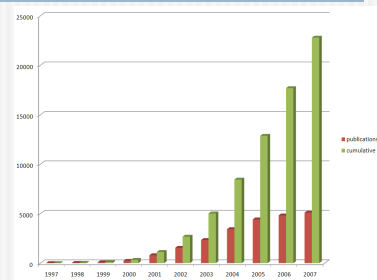
History of Microarrays

- 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:
 - *Int J Immunopathol Pharmacol.* 1990 19(4):905-914. Raloxifene covalently bonded to titanium implants by interfacing with (3-aminopropyl)-triethoxysilane affects osteoblast-like cell gene expression. Bambini et al
 - *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

January 30th 2008

MSc Seminar: Donald Dunbar

Microarray publications



January 30th 2008

MSc Seminar: Donald Dunbar

Types of experiment

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

Placebo

Wild-type

Healthy

Normal tissue

Time = 0

Drug treatment Drug 2...

Knockout

Patient

Cancerous tissue

Time = 1 Time = 2...

January 30th 2008

MSc Seminar: Donald Dunbar

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

January 30th 2008

MSc Seminar: Donald Dunbar

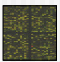

A typical experiment



January 30th 2008

MSc Seminar: Donald Dunbar

Experiment design: system

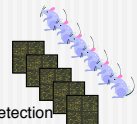
- What is your model?
 - animal, cell, tissue, drug, time...
- What comparison?  
- What platform
 - microarray? oligo, cDNA?
- Record all information: see "standards"

January 30th 2008

MSc Seminar: Donald Dunbar

Experiment design: replicates

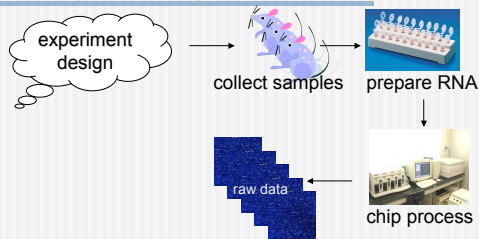
- Microarrays are noisy: need extra confidence in the measurements
- We usually don't want to know about a specific individual
 - eg not an individual mouse, but the strain
 - although sometimes we do (eg people)
- Biological replicates needed
 - independent biological samples
 - number depends on variability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



January 30th 2008

MSc Seminar: Donald Dunbar

A typical experiment



January 30th 2008

MSc Seminar: Donald Dunbar

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



January 30th 2008

MSc Seminar: Donald Dunbar

Bioconductor (BioC)



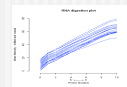
- <http://www.bioconductor.org/>
- "Bioconductor is an open source software project for the analysis and comprehension of genomic data"
- Started 2001, developed by expert volunteers
- Built on statistical programming environment "R"
- Provides a wide range of powerful statistical and graphical tools
- Use BioC for most microarray processing and analysis
- Make experiment design file and import data

January 30th 2008

MSc Seminar: Donald Dunbar

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 in simpleaffy)



January 30th 2008

MSc Seminar: Donald Dunbar

Pre-processing: background

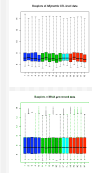
- Signal corresponds to expression...
 - plus a non-specific component (noise)
- Non specific binding of labelled target
- Need to exclude this background
- Several methods exist
 - eg Affy: PM-MM but many complications
 - eg RMA $PM=B+S$ (don't use MM)

January 30th 2008

MSc Seminar: Donald Dunbar

Pre-processing: normalisation

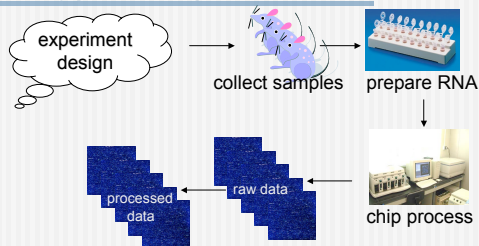
- In addition to background corrections
 - chip, probe, spatial, intra and inter array variation
 - need to remove to get at real expression differences
- Make use of statistics
 - combined with probe set summary: get an expression value for the gene
- But seems to be more dependency on intensity
 - additive and multiplicative errors
- Quantile normalisation often used
- Normalisation is complicated for 2-colour arrays
- Try to remove most noise at lab stage (ie control things well statistically)



January 30th 2008

MSc Seminar: Donald Dunbar

A typical experiment



January 30th 2008

MSc Seminar: Donald Dunbar

Part 1 Summary

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

January 30th 2008

MSc Seminar: Donald Dunbar

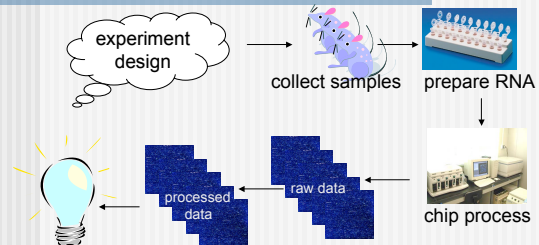
Microarray Informatics

Part 2

January 30th 2008

MSc Seminar: Donald Dunbar

A typical experiment

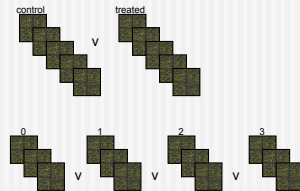


January 30th 2008

MSc Seminar: Donald Dunbar

Data analysis

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



January 30th 2008

MSc Seminar: Donald Dunbar

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

January 30th 2008

MSc Seminar: Donald Dunbar

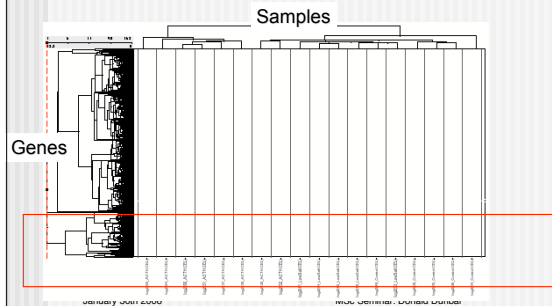
Hierarchical clustering

- 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

January 30th 2008

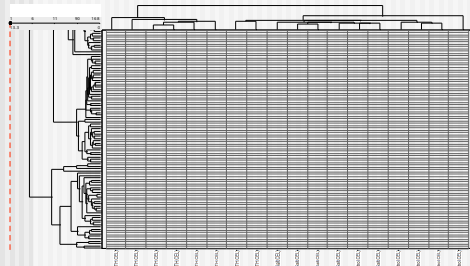
MSc Seminar: Donald Dunbar

Hierarchical clustering



Hierarchical clustering

- Heatmaps for microarray data



January 30th 2008

MSc Seminar: Donald Dunbar

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

January 30th 2008

MSc Seminar: Donald Dunbar

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

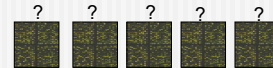
January 30th 2008

MSc Seminar: Donald Dunbar

Sample classification



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



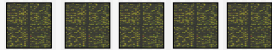
January 30th 2008

MSc Seminar: Donald Dunbar

Sample classification



Apply classifier



January 30th 2008

MSc Seminar: Donald Dunbar

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

January 30th 2008

MSc Seminar: Donald Dunbar

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...

January 30th 2008

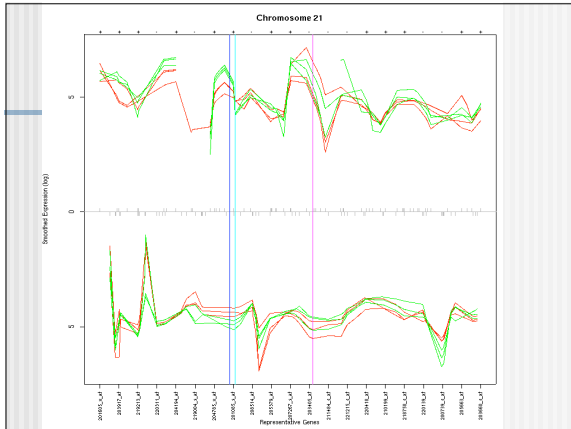
MSc Seminar: Donald Dunbar

Data-mining

Microarrays are a waste
of time
...unless you do
something with the data

January 30th 2008

MSc Seminar: Donald Dunbar



TOUCAN 2

[Launch Now](#) [Launch instructions](#) [Tutorial](#) [Manual](#)
[FAQ](#) [Version](#) [SOAP](#) [License](#)
[Terms](#) [Mailing](#) [Features](#) [Related](#)
[Service status](#) [Screenshots](#) [Acknowledgements](#) [FAQ](#)

Introduction

TOUCAN is a workflow for regulatory sequence analysis on mammalian genomes: comparative genomics, detection of significant transcription factor binding sites, and detection of cis-regulatory modules (combinations of binding sites) in sets of coexpressed/colocalized genes.

It is a platform independent, standalone Java application that is tightly linked with [Ensembl](#), and was built using the [BioJava](#) package. SOAP web services are used to remotely access [multiple algorithms](#) for comparative genomics, motif detection, and module detection.

Comments, suggestions, and bug reports can be sent to stijn.aertcaut@med.kuleuven.ac.be or toouan@latterer.cs.kuleuven.ac.be

Register

If you're using TOUCAN, please enter your email address:

Launch TOUCAN v. 2.2.5

To run TOUCAN you need to have **two things** installed:

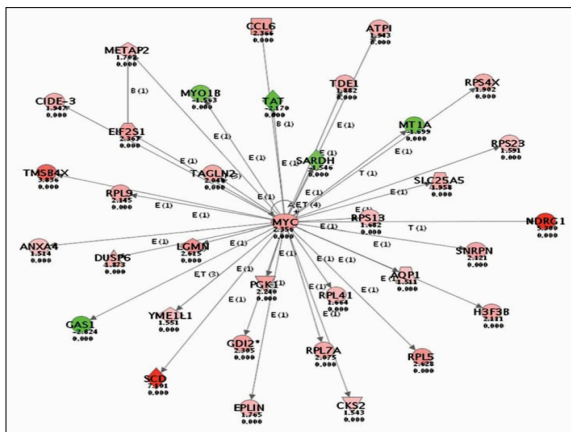
- Java 2 Platform, Standard Edition (J2SE), version 1.4.x, or 1.5.0
- Java Web Start which is shipped as part of J2SE. From J2SE 1.4.2 onwards it is installed together with the SDK/JRE (more info)

If you have fulfilled these requirements, then you can launch TOUCAN directly using this URL: <http://www.esat.kuleuven.ac.be/~aertca/stijn/software/toouan.jsp>

Alternatively you can type this command in a terminal window:
 javaws <http://www.esat.kuleuven.ac.be/~aertca/stijn/software/toouan.jsp>

We try to encourage to use Java Web Start because this way you will always have the latest version of the software. This is important because the properties change at least once a month to follow the newest Ensembl releases. If you are really unable to use Java Web Start, you can email us as we email, and we can give you the JAR file of TOUCAN.

January 30th 2008 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 or our text mining tool)
 - ... to make sense of the data

January 30th 2008 MSc Seminar: Donald Dunbar

Microarray Resources

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

January 30th 2008 MSc Seminar: Donald Dunbar

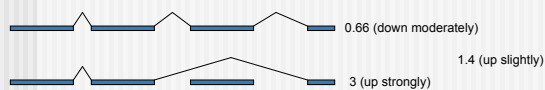
Microarray Standards

- 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

January 30th 2008 MSc Seminar: Donald Dunbar

Recent advances: Exon chips

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



New chips will give us much more information

January 30th 2008

MSc Seminar: Donald Dunbar

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 TRAITS LOCI (eQTL)
 - Overlap of expression and genetic differences (cis)
 - Correlation at different locus (trans)



January 30th 2008

MSc Seminar: Donald Dunbar

Part 2 Summary

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

January 30th 2008

MSc Seminar: Donald Dunbar

Seminar Summary

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

January 30th 2008

MSc Seminar: Donald Dunbar

Contact

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

January 30th 2008

MSc Seminar: Donald Dunbar