

Microarray Informatics

Donald Dunbar
MSc Seminar
4th February 2009

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

February 4th 2009 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

February 4th 2009 MSc Seminar: Donald Dunbar

Microarray Informatics

Part 1

February 4th 2009 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?

February 4th 2009 MSc Seminar: Donald Dunbar

The central dogma

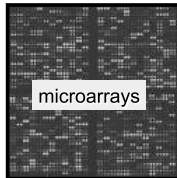
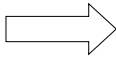
30k Gene: DNA
90k Transcript: RNA
100+k Protein

kinase, protease, structural receptor, ion channel...

February 4th 2009 MSc Seminar: Donald Dunbar

Measuring transcripts

- Genome level sequencing
- New miniaturisation technologies
- Better bioinformatics

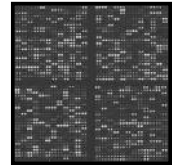


February 4th 2009

MSc Seminar: Donald Dunbar

Microarrays: wish list

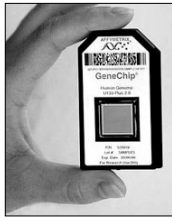
- Include all genes in the genome
- Include all splice variants
- Give reliable estimates of expression
- Easy to analyse
 - bioinformatics tools available
- Cost effective



February 4th 2009

MSc Seminar: Donald Dunbar

Microarray technologies - 1



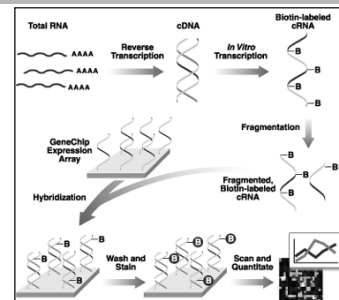
- Oligonucleotides - Affymetrix
- One chip all genes
- Chips for many species
- Several oligos per transcript
- Use of control, mismatch sequences
- One sample per chip
 - 'absolute quantification'
- Well established in research
- Expensive



February 4th 2009

MSc Seminar: Donald Dunbar

Microarray technologies - 1

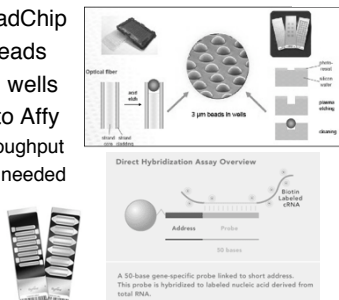


February 4th 2009

MSc Seminar: Donald Dunbar

Microarray technologies - 2

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



February 4th 2009

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

February 4th 2009

MSc Seminar: Donald Dunbar

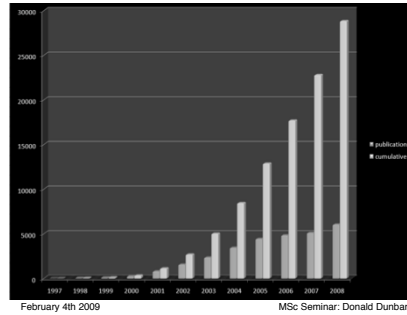
History of Microarrays

- 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

February 4th 2009

MSc Seminar: Donald Dunbar

Microarray publications

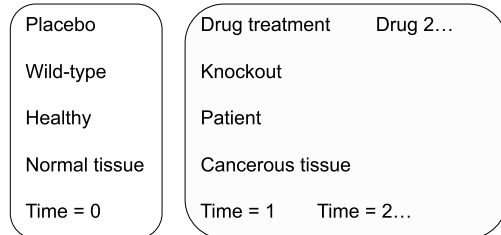


February 4th 2009

MSc Seminar: Donald Dunbar

Types of experiment

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



February 4th 2009

MSc Seminar: Donald Dunbar

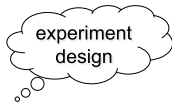
Types of experiment

- 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

February 4th 2009

MSc Seminar: Donald Dunbar




A typical experiment



February 4th 2009

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"

February 4th 2009

MSc Seminar: Donald Dunbar

Experiment design: replicates

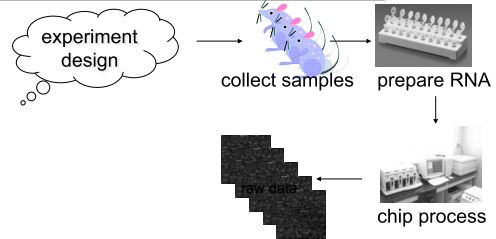
- Microarrays are noisy: need extra confidence in the measurements
- We usually don't want to know about a specific individual
 - eg not an individual mouse, but the strain
 - although sometimes we do (eg people)
- Biological replicates needed
 - independent biological samples
 - number depends on variability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



February 4th 2009

MSc Seminar: Donald Dunbar

A typical experiment



February 4th 2009

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



February 4th 2009

MSc Seminar: Donald Dunbar

Bioconductor (BioC)



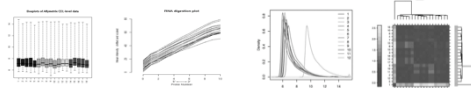
- <http://www.bioconductor.org/>
- "Bioconductor is an open source software project for the analysis and comprehension of genomic data"
- Started 2001, developed by expert volunteers
- Built on statistical programming environment "R"
- Provides a wide range of powerful statistical and graphical tools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

February 4th 2009

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 arrayQualityMetrics)



February 4th 2009

MSc Seminar: Donald Dunbar

Pre-processing: background

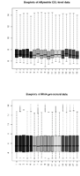
- Signal corresponds to expression...
 - plus a non-specific component (noise)
- Non specific binding of labelled target
- Need to exclude this background
- Several methods exist
 - eg Affy: PM-MM but many complications
 - eg RMA PM=B+S (don't use MM)

February 4th 2009

MSc Seminar: Donald Dunbar

Pre-processing: normalisation

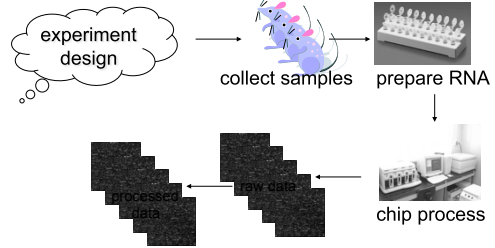
- In addition to background corrections
 - chip, probe, spatial, intra and inter-array variation
 - need to remove to get at real experimental differences
- Make use of statistics
 - combined with probe set summary: get an expression value for the gene
 - But seems to be no dependency on intensity
 - additive and multiplicative errors
- Quantile normalisation often used
- Normalisation is complicated for 2-colour arrays
- Try to reduce most noise at lab stage (ie control things well statistically)



February 4th 2009

MSc Seminar: Donald Dunbar

A typical experiment



February 4th 2009

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 4th 2009

MSc Seminar: Donald Dunbar

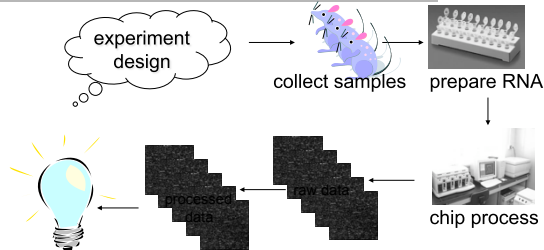
Microarray Informatics

Part 2

February 4th 2009

MSc Seminar: Donald Dunbar

A typical experiment

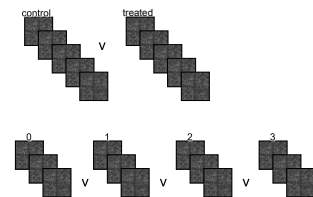


February 4th 2009

MSc Seminar: Donald Dunbar

Data analysis

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



February 4th 2009

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
- Depends on follow-up procedure!

February 4th 2009

MSc Seminar: Donald Dunbar

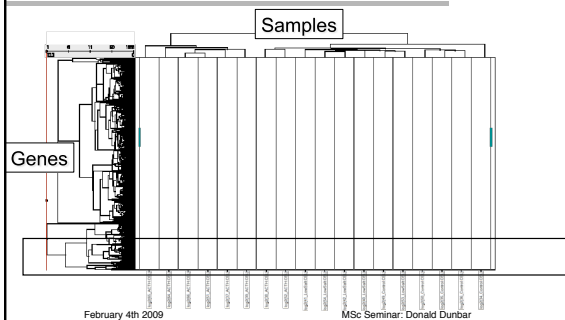
Hierarchical clustering

- 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

February 4th 2009

MSc Seminar: Donald Dunbar

Hierarchical clustering

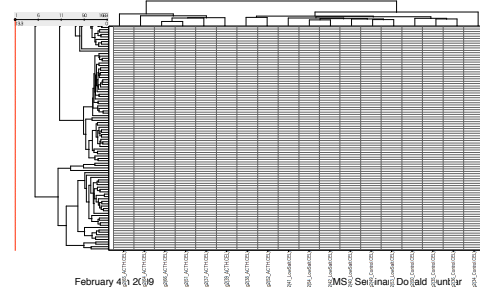


February 4th 2009

MSc Seminar: Donald Dunbar

Hierarchical clustering

- Heatmaps for microarray data



February 4th 2009

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

February 4th 2009

MSc Seminar: Donald Dunbar

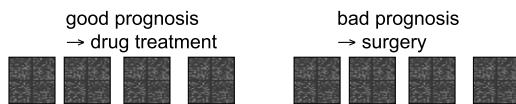
Sample classification

- 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

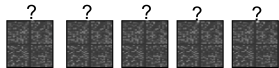
February 4th 2009

MSc Seminar: Donald Dunbar

Sample classification



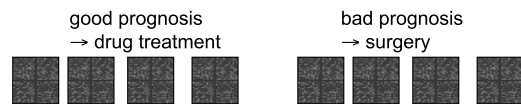
Gene selection, training, cross validation →
classifier: gene x * 0.5 gene y * 0.25 gene z ...



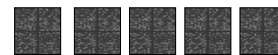
February 4th 2009

MSc Seminar: Donald Dunbar

Sample classification



Apply classifier



February 4th 2009

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
 - classifier loses specificity

February 4th 2009

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...

February 4th 2009

MSc Seminar: Donald Dunbar

Data-mining

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

February 4th 2009

MSc Seminar: Donald Dunbar

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

February 4th 2009

MSc Seminar: Donald Dunbar

Input a query or leave blank (that lists lots of data) for all data, then submit...

Annotation Queries
 Affymetrix ID
 Entrez Gene ID
 Gene Title
 Gene Symbol
 Gene Ontology Term
 Pathway
 Chromosome

Group
 Input or Output Any X

Comments
 Comments
 Submit Query or Reset

Expression maxima and minima
 BAT max/min < >
 BAT max - min
 WAT max/min
 WAT max - min
 Liver max/min
 Liver max - min

Correlation with circadian gene profiles
 Which gene? par1 Tissue: BAT Rank limit

Order
 Order output by: Gene symbol and Ascending

Submit Query or Reset

Home

February 4th 2009 MSc Seminar: Donald Dunbar

Entrez Gene ID	Gene title	Gene Symbol	Intestines										Group	Comments (links)																																																																																					
ADY ID			1A	2A	3A	4A	5A	6A	7A	1B	2B	3B	4B	5B	6B	7B																																																																																			
149851_at	period homolog 1 (Drosophila)	Per1 and PERCP	BAT193	331	393	1100	255	328	102	32	198	175	133	939	659	179	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100

February 4th 2009 MSc Seminar: Donald Dunbar

Detailed information: 1449851_at

Annotation Data for Per1
 Affymetrix ID: 1449851_at
 Entrez Gene ID: 13629
 Gene Title: period homolog 1 (Drosophila)
 Gene Symbol: Per1
 RefSeq ID: NM_0010002893

GO Biological Process information
 7677 / regulation of transcription, DNA-dependent / inferred from electronic annotation / 7622 / rhythmic behavior annotation / 7623 / circadian rhythm / inferred from electronic annotation

GO Molecular Function information
 4871 / signal transducer activity / inferred from electronic annotation

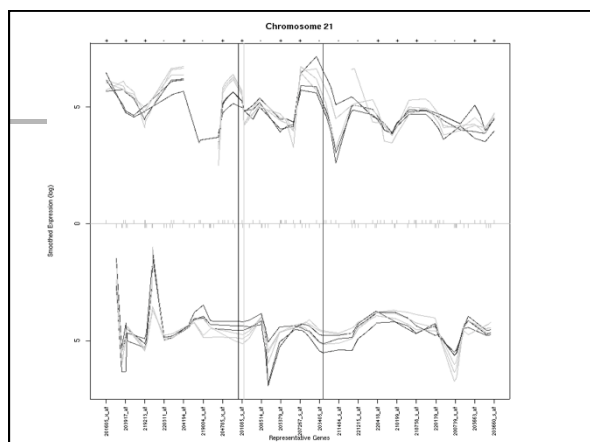
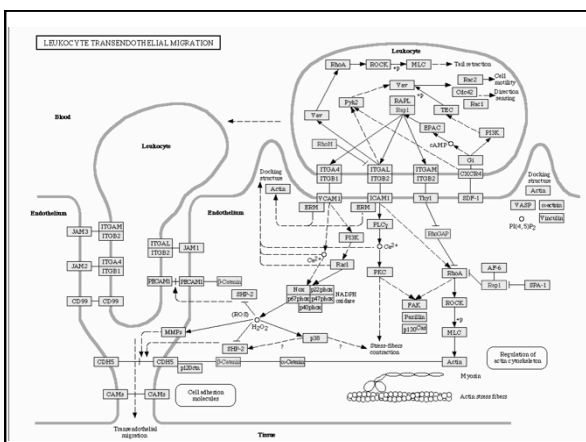
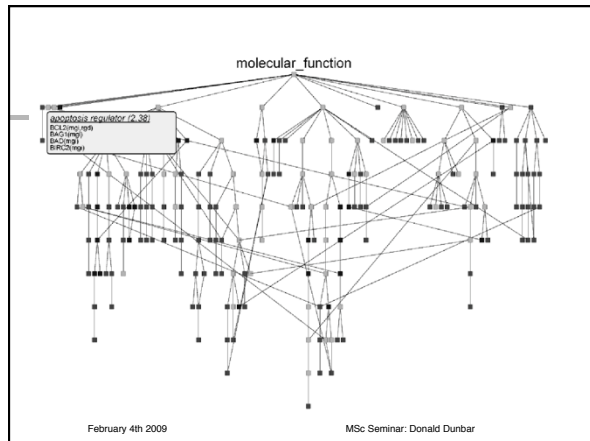
GO Cellular Component information
 5634 / nucleus / inferred from direct assay

Pathway
 Cytosol, Cytoplasm, Cytoplasm, Cytoplasm, Cytoplasm, Cytoplasm



Expression Data for Per1
 Intestines
 View these data as a [grid](#)

Entrez Gene ID	1A	2A	3A	4A	5A	6A	7A	1B	2B	3B	4B	5B	6B	7B
BAT193	331	393	1100	255	328	102	32	198	175	133	939	659	179	18
WAT192	796	963	3091	1127	974	412	612	593	1	1	1	1	1	1
Liver27	129	105	1241	651	444	348	32	79	253	109	170	220	22	22

February 4th 2009 MSc Seminar: Donald Dunbar



TOUCAN 2

LEUVEN  

Launch New: News, Version, Download, Service status
 Launch instructions: News, Version, Making, Screenshots
 Tutorials: Source, SOAP, Features, Acknowledgments
 Manual: License, Related, FAQ

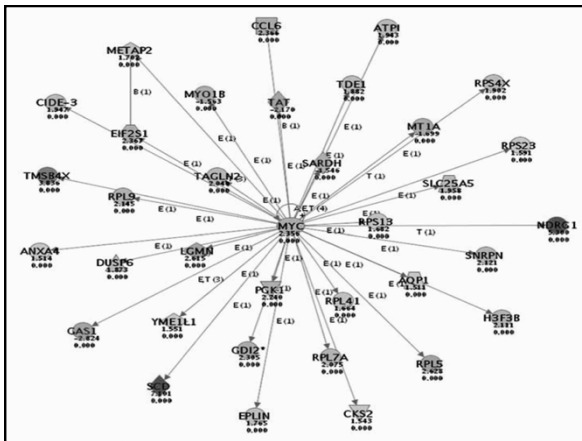
Gene set
Enriched regulatory sequences
Functional significance?

February 4th 2009 MSc Seminar: Donald Dunbar

PubMatrix Results 13th September 2005
 Genes (ACTH differentially expressed) Terms for DDUNBAR

PubMatrix	adrenal/hypertension/blood pressure	ACTH low salt/low sodium	hyperplasia/hypoplasia	steroid retention
(α-christinin and methylglucosylase domain 9 OR Adam1)	0	0	0	0
(α-christinin domain containing 2 OR ADAC2)	0	0	0	0
(arbitrating transcription factor 3 OR AAT3)	0	0	0	0
(ary1 Cys synthetase long-chain OR Acs4)	11	1	2	0
(axillin OR Axil)	22	14	118	5
(alcohol dehydrogenase 7 OR Adh7)	0	0	0	0
(aldehyde dehydrogenase 1 family, member 12 OR ALDH12)	0	0	0	0
(aldehyde dehydrogenase 18 family, member A1 OR ALDH18A1)	0	0	0	0
(aldol-keto reductase family 1, member C18 OR Aldk18)	0	0	0	0
(alkaline phosphatase 2, liver OR Alpl2)	52	53	145	8
(arginine vasopressin receptor 1 OR Avpr1)	0	120	124	7
(astrocytic permeability-increasing protein like 1 OR Bpl1)	0	0	0	0
(basic leucine zipper and W2 domain 1 OR Bvz1)	0	0	0	0
(BMP-binding endothelial regulator OR MGL1/2/4/8)	0	0	0	0
(branched-chain aminotransferase 1, cytosolic OR Bca1)	11	0	0	0
(CD5 antigen-like OR Cd5)	0	0	0	0
(CDC28 protein kinase regulatory subunit 2 OR Cks2)	0	0	0	0
(CCL4-related cell adhesion molecule 2 OR Ccr42)	11	0	0	0
(cell adhesion cycle 2 homolog OR Cdc4)	0	0	0	0
(cysteine, gamma 2 OR Cwrg2)	0	0	0	0
(chitinase 3 like 3 OR Ch3L3)	0	0	0	0
(cyclic-dependent kinase inhibitor 1C OR Cdk1c)	11	12	0	11
(cystatin B OR Cstb)	2	4	0	0
(cystathionase or oxidase OR Cth2)	128	141	108	14
(cystathionase receptor like factor 1 OR Cstf)	0	1	0	0
(DNA (cystine 5)-methyltransferase 3 like OR Dnm3)	0	0	0	0
(demon-regulated by Ccr4k1, 4 OR MGL1/2/4/8/9)	0	0	0	0
(dihydrobioin alpha OR Dha)	0	0	0	0
(enigma OR Emh)	11	5	11	11
(fatty acid desaturase 2 OR Fad2)	0	0	0	0
(fetal-related protein OR Frl)	11	0	0	0
(galactose-4-epimerase, alpha OR Gae)	53	141	108	14

February 4th 2009 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

February 4th 2009 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...

February 4th 2009 MSc Seminar: Donald Dunbar

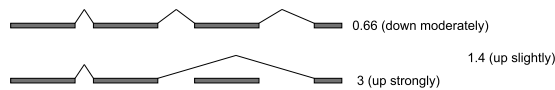
Microarray Standards

- 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

February 4th 2009 MSc Seminar: Donald Dunbar

Recent advances: Exon chips

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



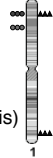
New chips will give us much more information

February 4th 2009

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



February 4th 2009

MSc Seminar: Donald Dunbar

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 5 years: the end of microarrays?

February 4th 2009

MSc Seminar: Donald Dunbar

Part 2 Summary

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

February 4th 2009

MSc Seminar: Donald Dunbar

Seminar Summary

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

February 4th 2009

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

February 4th 2009

MSc Seminar: Donald Dunbar

PhD opportunity

- Centre for Cardiovascular Science (Edinburgh)
- The Cellular and Molecular Basis of Cardiovascular Disease
- BHF funded PhDs
 - biologists (x4)
 - physical scientists (informatics, physics, maths....)
- Details on web:
 - <http://www.cvs.med.ed.ac.uk/Training/content.asp?SubCatID=44>

February 4th 2009

MSc Seminar: Donald Dunbar