Arrays & Expression

MMG 835, SPRING 2016 Eukaryotic Molecular Genetics

George I. Mias

Department of Biochemistry and Molecular Biology gmias@msu.edu

Arrays

- Spotted DNA arrays
 - Fluorescent (Cy3, Cy5)
- Long/short oligonucleotide arrays
- Glass slides
- Bead Arrays

affymetrix.com agilent.com illumina.com thermofisher.com







Arrays

- blocking compound 🌒 oxygen /afer adenine light source_ hydroxide thymine 25 nucleotides GeneChip® microarray LONG
- Affymetrix Genechip production
- 5-inch square quartz wafer
- 18-20 micron square windows
- nucleotide masks
- repeat for 25 nucleotides

The Science Creative Quarterly <u>http://www.scq.ubc.ca/spot-your-genes-an-overview-of-the-microarray/</u> (Jiang Long)

Array-based Comparative Genome Hybridization, CGH.

- genomic DNA, used as probe
- DNA copy-number variations
- Measuring relative copy number in the genome rather than in cellular mRNA.
- CGH resolution ~ 10-20 Mb

Array CGH



Reference and test DNA differentially Labelled with fluorescent tags (Cy5 and Cy3, respectively)





- COT-1 DNA.
- Samples are then hybridized to genomic arrays



- including BAC clones
- PCR fragments
- oligonucleotides







- CNVs in 39 unrelated healthy people and 16 individuals with known chromosomal abnormalities
- Control: pooled 55 normal individuals).
- Detected all the previously known abnormalities + 255 other CNVs (~12/person)
- 102 CNV occurred in more than one person
- 24 were found in >10%.



lafrate et al., Nature Genetics 36, 949 - 951 (2004)

relative loss case (6 gene signals) CGH

normal case (9 signals)



case with gain (12 signals)

Most common CNV at amylase gene cluster 1p13.3. Length of the polymorphic region varied from

~ 150 kb (a) to

lafrate et al., Nature Genetics 36, 949 - 951 (2004)

~ 425 kb (c).





Affymetrix

25-mer probes for both alleles

Location of the SNP locus varies from probe to probe.

Binding to both probes regardless of the allele.

More efficient when all 25 bases are complementary

Dimmer mismatched signal







Bead has 50-mer sequence attached complementary to sequence adjacent to SNP site.

BeadArray



Illumina

BeadArray



Illumina



Enzymatic single-base extension (T or G) complementary to the allele carried by the DNA (A or C).

Labeled/colored nucleotide signal.



illumina.com

SNP Arrays

DAY 3



Illumina

Resequencing Microarrays

Resequencing microarrays: probe provided for every possible single-base-pair mismatch

Sample ---TGCATGCTACCGTCACTGCGCGGTAT---



multiple probes (4 shown per position, can be more)

Resequencing Microarrays

Resequencing microarrays: probe provided for every possible single-base-pair mismatch Probe



multiple probes (4 shown per position, can be more)

Hybridization to Candidate Probes

Sample ---CCGGATGCTGCATGCTACCGTCACTTCGCGGTATATACCGAGCTGACAGTCAG---GGCCTACGACGTACGATGGCAGTGA TACGACGTACGATGGCAGTGACGCGC ACGTACGATGGCAGTGACGCGCCAT ACGATGGCAGTGACGCGCCATATAT TGGCAGTGACGCGCCATATATGGCT AGTGACGCGCCATATATGGCTCGACTGTC GCCATATATGGCTCGACTGTCAGTC

> In the absence of resequencing arrays, hybridization of the sample to candidate sequence probes can be used.

Hybridization to Candidate Probes



Mismatches from mutations in sample DNA lead to lower hybridization efficiency (compared to complete complementary sequence matching).

Hybridization to Candidate Probes



- Advantage: fewer probes may be enough to detect variation.
- If mismatch inferred need small-scale sequencing to identify the variant nucleotide.

Resequencing Microarrays

The effect of SNP hybridization as a function of its position in a probe.



- More central positions result in greatest decrease
- SNPs positioned at the end of probes much less likely to result in a significant decrease in hybridization.

Mutation Detection



- Drug resistant Saccharomyces cerevisiae (budding yeast). Candidate SNPs are identified by a positive log likelihood value.
- Small number of candidate SNPs detected throughout the genome

Feuk, Nature Reviews Genetics 7, 85-97 (2006). Gresham, D. et al., Science 311, 1932–1936 (2006).

Mutation Detection



Feuk, Nature Reviews Genetics 7, 85-97 (2006). Gresham, D. et al., Science 311, 1932–1936 (2006).

- Transcriptomics Collective of all expressed transcripts in a given cell .
- \cdot Gene "expressed" when transcribed to RNA
- Level of RNA produced from a gene is controlled by:

Transcription

Stability/Degradation

- Dynamic (Vs. static genome)
- Tissue specific

Gene expression may also refer to levels of protein instead

Gene Expression Time Period Key Facts Gene Model Phenotype

 1900s
 Cells contain genes. (location, composition and structure unknown)
 →
 →

 1950-1960s
 Gene located on specific regions of chromosomes. (composed of DNA transcribed into RNA)
 →
 →



1970 2000s
 Genes clustered together in generic regions separated by large intergenic regions of "junk DNA". (mostly transcribed into protein-coding RNAs which can be spliced into multiple isoforms)





Spliced RNA isoforms

2001present Genes have distal 5' TSSs and regulatory regions and share their genomic sequences with long and short noncoding RNAs encoded on both strands.

Interleaved RNAs



plus other more subtle phenotypes due to mutations or interactions with interleaved ncRNAs

Gingeras, Genome Res. 17(6):682-90 (2007)

Gene Expression Transcriptome has mainly

- mRNA (messanger RNA)
- rRNA (ribosomal)
- tRNA (transfer)

But also

- ncRNA (non-coding)
 - lincRNA (long intergenic ncRNA);
 - snRNA (small nuclear)
 - snoRNA (small nucleoral)
 - miRNA (micro)
 - siRNA (small interfering)
 - · piRNAs (PIWI-interacting) RNAs





PASRs, promoter-associated short RNAs

tiRNAs, transcription initiation RNAs

spliRNAs, splice site RNAs

Morris & Mattick, Nature Reviews Genetics 15, 423–437 (2014)

ncRNA*	No. of known transcripts	Transcript lengths (nucleotides; nt)	Functions		
Precursors to short RNAs					
miRNA	1,756	>1,000	Precursors to short (21–23 nt) regulatory RNAs		
snoRNA	1,521	>100	Precursors to short (60–300 nt) RNAs that help to chemically modify other RNAs		
snRNA	1,944	1,000	Precursors to short (150 nt) RNAs that assist in RNA splicing		
piRNA	89	Unknown	Precursors to short (25–33 nt) RNAs that repress retrotransposition of repeat		
tRNA	497	>100	Precursors to short (73–93 nt) transfer RNAs		

Kowalczyk & Higgs, Nature 482, 310–311 (2012)

	ncRNA*	No. of known transcripts	Transcript lengths (nucleotides; nt)	Functions
	Long ncRNAs			
	Antisense ncRNA	5,446	100->1,000	Mostly unknown, but some are involved in gene regulation through RNA interference
	Enhancer ncRNA (eRNA)	>2,000	>1,000	Unknown
	Enhancer ncRNA (meRNA)	Not fully documented	As variable as the length of mRNAs	Unknown, but they resemble alternative gene transcripts
	Intergenic ncRNA	6,742	10 ² —10 ⁵	Mostly unknown, but some are involved in gene regulation
	Pseudogene ncRNA	680	10 ² —10 ⁴	Mostly unknown, but some are involved in regulation of miRNA
	3' UTR ncRNA	12	>100	Unknown
Κ	owalczyk & Higg	s,Nature 482, 3 ⁻	10–311 (2012)	



NHGRI, https://www.genome.gov/Images/EdKit/bio2j_large.gif

% Tissue-Total Number Both % Tissue-Number Alternative events detected isoforms tissueregulated regulated transcript events $(\times 10^3)$ $(\times 10^{3})$ detected regulated (observed) (estimated) 37 35 Skipped exon 10,436 6,822 65 72 Retained intron 1 1 167 96 57 71 Alternative 5' splice 15 15 2,168 1,386 64 72 site (A5SS) Alternative 3' splice 17 16 64 74 4,181 2,655 site (A3SS) Mutually exclusive 4 167 95 57 4 66 exon (MXE) Alternative first 14 13 10,281 5,311 52 63 exon (AFE) Alternative last 9 8 5,246 47 2,491 52 exon (ALE) 7 Tandem 3' UTRs 7 3,801 80 5,136 74 Total 105 100 37,782 22,657 68 60 Constitutive exon or region Body read Junction read pA Polyadenylation site Alternative exon or extension Inclusive/extended isoform Exclusive isoform Both isoforms

Pervasive tissuespecific regulation of alternative mRNA isoforms.

Wang et al., Nature 456, 470-476 (27 November 2008)

Gene Expression Methods

- Northern Blots
 - Measure RNA levels by hybridization of a labeled probe to total RNA after gel electrophoresis and transfer to a membrane.
- Reporter Genes
 - Use of activity of a foreign protein to measure the amount of transcription from a promoter.
- Quantitative real-time RT-PCR.
 - Uses amplification by sequence-specific primers
- Microarrays
 - high throughput
- RNA-Sequencing

1-Color Expression Arrays

mRNA extracted

reverse transcribe to cDNA (complementary DNA),

Scan

Relative Expression

The Science Creative Quarterly, <u>scq.ubc.ca</u> (Jiang Long)

2-Color Expression Arrays

mRNA extracted

reverse transcribe to cDNA (complementary DNA),

Fluorescent labels (Cy3, Cy5)

Reference Relative ratios (control for amounts & conditions)

The Science Creative Quarterly, <u>scq.ubc.ca</u> (Jiang Long)

Expression BeadChips

Up to 30 beads per probe average

Expression BeadChips

Expression BeadChips

Probes	Description	Human HT-12 v4.0*	Mouse WG-6 v2.0	Mouse Ref- 8 v2.0	Rat Ref-12	Human WG DASL HT*
		12-sample	6-sample	8-sample	12-sample	12-sample
RefSeq Con	itent					
NM	Coding transcript, well-established annotation	28,688	26,766	24,854	6,277	27,253
XM	Coding transcript, provisional annotation	11,121	6,856	796	15,983	426
NR	Non-coding transcript, well-established annotation	1,752	56	47	1	1,580
XR	Non-coding transcript, provisional annotation	2,209			12	26
Source	RefSeq source release	Human RefSeq Rel 38	Mouse Re	RefSeq I 22	Rat RefSeq Rel 16	Human RefSeq Rel 38
Supplemen	tary Content					
UniGene	Experimentally confirmed mRNA sequences that align to EST clusters	3,461			250	
RIKEN FANTOM2	Exemplar protein-coding sequences from the RIKEN FANTOM2 database		5,659			
RefSeq Release 5	Transcripts with NM and XM annotation in RefSeq Release 5 (Build 33.1)		3,573			
MEEBO	Probes to transcripts that do not align with 100% accuracy to RefSeq, but are confirmed as valid mRNA mapping to clusters in Expressed Sequence Tag databases ⁶		2,371			
Total		47,231	45,281	25,697	22,523	29,285

*> 99.98% of the bead types are present on any HumanHT-12 array

Box plots

Box plots

median confidence interval

mean confidence interval

Wolfram Research

MA Plots

MA Plots

Define:

 $M = \log_2(\mathbf{R}/\mathbf{G})$

 $A = (1/2) \times \log_2(\mathbf{R} \times \mathbf{G}).$

Can fit with a lowess curve, which is used $_{0}$ to normalize the gene-expression $^{\geq 0}$ measurements.

Lowess, locally weighted scatterplot smoothing

5 -5 10 15 5 $\mathbf{0}$ A

Volcano Plots

Statistical significance Vs. *fold changes* for all genes on the array.

Plotted as

 $-\log_{10}(p-values)$ Vs.

log₂ (fold change).

Volcano Plots

Upper corners represent genes that show both statistical significance and large fold changes.

p-value Histograms

p-values for a chosen test of differential expression for each gene.

Can be used for falsediscovery rates estimations.

Heatmaps & Dendrograms

- Each cell represents a value (scale important)
- Often generated following hierarchical clustering analyses.
- Dendrograms shown to the sides indicating *clustering* results.
- Clustering can be both vertical and horizontal.

Clustering techniques compare genes pairwise to compute *similarity*.

Euclidean distance (compare to usual vectors and using Pythagorean theorem) Expression level in sample 2

Expression level in sample 1

Pearson correlation coefficient as a distance

Calculated from the distances of each point from the linear-regression line (known as residuals). Expression level of gene 2

Expression level of gene 1

Mutual information

Model-free measurement of the degree of information content in one gene known from another gene.

Highest when genes are randomly distributed separately, but show a non-random joint distribution.

 \sim Expression level of gene

Expression level of gene 1

Manhattan distance (city-block distance, L1 norm)	$d_{fg} = \sum_{c} \left e_{fc} - e_{gc} \right $
Euclidean distance (L2 norm)	$d_{fg} = \sqrt{\sum_{c} (e_{fc} - e_{gc})^2}$
Mahalanobis distance	$d_{fg} = (e_f - e_g)' \Sigma^{-1} (e_f - e_g)$, where Σ is the (full or within-cluster) covariance matrix of the data
Pearson correlation (centered correlation)	$d_{fg} = 1 - r_{fg}$, with $r_{fg} = \frac{\sum_{c} (e_{fc} - \bar{e}_{f})(e_{gc} - \bar{e}_{g})}{\sqrt{\sum_{c} (e_{fc} - \bar{e}_{f})^{2} \sum_{c} (e_{gc} - \bar{e}_{g})^{2}}}$
Uncentered correlation (angular separation, cosine angle)	$d_{fg} = 1 - r_{fg}$, with $r_{fg} = \frac{\sum_{c} e_{fc} e_{gc}}{\sqrt{\sum_{c} e_{fc}^2 \sum_{c} e_{gc}^2}}$
Spellman rank correlation	As Pearson correlation, but replace e_{gc} with the rank of e_{gc} within the expression values of gene g across all conditions $c = 1C$
Absolute or squared correlation	$d_{fg} = 1 - r_{fg} \text{ or } d_{fg} = 1 - r_{fg}^{2}$

 $d_{fg'}$ distance between expression patterns for genes f and g. $e_{gc'}$, expression level of gene g under condition c.

Hierarchical Clustering

Sorting so similar genes appear near each other.

The length of the branch is inversely proportional to the degree of similarity.

Shades of red indicate increased relative expression; shades of green indicate decreased relative expression.

Self Organizing Maps (SOMs)

Self-organizing maps find variable-sized clusters of genes that are similar to each other, given the input number of clusters to find.

Expression level in biological sample 1

Relevance Networks

Find and display pairs of genes with strong positive and negative correlations, then construct networks from these gene pairs.

Strength of correlation is proportional to the thickness of the edges, and red indicates a negative correlation.

Principal-components analysis

Shows clustering or scatter of genes (or samples) when viewed along two or three principal components

Principal component can be thought of as a 'metabiological sample', which combines all the biological samples so as to capture the most variation in gene expression. Principal component 2

Principal component 1

The nearest-neighbour supervised method

- First construct hypothetical genes that best fit the desired patterns (e.g.a gene with high expression in disease 1 and low expression in disease 2, or vice versa).
- 2. Find individual genes that are most similar to the hypothetical genes.

Support Vector Machines

Instead of restricting to individual genes, find the line (or plane) that best separates groups of biological samples.

40 genes measured under two different conditions.

4 clusters of different sizes, shapes and numbers of genes.. Euclidean distance, which corresponds to the straight-line distance between points in this graph, was used for clustering. Right: the standard red-green representation of the data and corresponding cluster identities

Hierarchical clustering finds an entire hierarchy of clusters. The tree was cut at the level indicated to yield four clusters.

k-means (with k = 4) partitions the space into four subspaces, depending on which of the four cluster centroids (stars) is closest.

SOM finds clusters, which are organized into a grid structure (in this case a 2 × 2 grid).

Patterns in Data

- Foreskin fibroblasts grown in culture and were deprived of serum for 48 hr.
- Serum was added back and samples taken at time 0, 15 min, 30 min, 1 hr, 2 hr, 3 hr, 4 hr, 8 hr, 12 hr, 16 hr, 20 hr, 24 hr - final datapoint was from unsynchronized sample.
- cDNA microarray with 8,600 distinct human genes.
- Measurements are relative to time 0.
- Genes were selected if their expression level deviated from time 0 by at least a factor of 3.0 in at least 2 time points. the corresponding region of the dendrogram.

Eisen et al. PNAS 95:14863-14868 (1998) Iyer et al., Science 283(5398), pp. 83-87 (1999)

Clusters with genes in: (A) cholesterol biosynthesis (B) the cell cycle

- (C) the immediate—early response
- (D) signaling andangiogenesis,
- (E) wound healing and tissue remodeling

Patterns in Data

Permutations:

random 1, within rowsrandom 2, within columnsrandom 3, both directions

Eisen et al. PNAS 95:14863-14868 (1998)