



MAPS: a microarray project system for gene expression experiment information and data validation

Pierre R. Bushel¹, Hisham Hamadeh¹, Lee Bennett¹, Stella Sieber², Karla Martin¹, Emile F. Nuwaysir^{1,3}, Kate Johnson^{1,3}, Kelli Reynolds¹, Richard S. Paules² and Cynthia A. Afshari^{1,*}

¹Laboratory of Molecular Carcinogenesis and ²Laboratory of Environmental Carcinogenesis/Mutagenesis, National Institute of Environmental Health Sciences, PO Box 12233, Research Triangle Park, NC 27709, USA

Received on November 29, 2000; revised and accepted on January 21, 2001

ABSTRACT

Summary: MAPS is a MicroArray Project System for management and interpretation of microarray gene expression experiment information and data. Microarray project information is organized to track experiments and results that are: (1) validated by performing analysis on stored replicate gene expression data; and (2) queried according to the biological classifications of genes deposited on microarray chips.

Availability: MAPS is accessible at <http://www.dir.niehs.nih.gov/microarray/software/maps/>

Contact: bushel@niehs.nih.gov

INTRODUCTION

DNA microarray technology has risen to the forefront of gene expression analysis as a tool by which researchers can detect genome-wide differential expression of thousands of genes. Interest in identifying transcription signature patterns has led to an explosion in the use of microarray technology to simultaneously analyze changes in gene expression profiles (Hughes *et al.*, 2000).

Statistical methods and information systems have been developed to facilitate the analysis and management of complex and large-scale microarray gene expression results. For example, methods have been developed where computed confidence intervals are used to determine significantly changed genes from a ratio distribution of gene expression data (Chen *et al.*, 1997). Differentially expressed genes can be validated by comparing replicate measurements of microarrays and performing subsequent biological assays to confirm the biological significance of

altered gene expression (Wittes and Friedman, 1999). In addition, the microarray database ArrayDB was designed to store, analyze and associate gene expression data with information from remote biological resources (Ermolaeva *et al.*, 1998).

MicroArray Project System (MAPS) is a relational database management system with a web interface to cDNA microarray project information, validate replicate gene expression experimental results, and query gene expression data based on gene classifications of interest. In cDNA microarray analysis, gene expression ratios are computed from pixel intensity values acquired from 16 bit gray scale images. Differentially expressed genes are identified based on confidence intervals determined from a distribution of the ratios of intensity values, with significantly changed genes being those that fall outside of the confidence limits (Chen *et al.*, 1997). For a given confidence level the probability of chance occurrences in replicate experiments is determined using a binomial distribution model (Casella and Berger, 1990). Loading microarray information and data into MAPS allows a user to query for gene expression results which meet defined search parameters, validation criteria, and gene classifications of interest which are defined in the database by grouping genes into functional categories. MAPS, in contrast to ArrayDB, provides the advantage of associating experimental and biological microarray information with statistically validated gene expression data.

IMPLEMENTATION AND WEB INTERFACE

MAPS is implemented with a Microsoft Access database jet engine (Microsoft). Over 40 tables are defined in the database to manage microarray project information, detail experimental design, track clones, sample preparation,

*To whom correspondence should be addressed.

³Present address: NimbleGen Systems, LLC, University Research Park, 505 South Rosa Road, Madison, WI 53719, USA.

Please enter your query

(a) Conf Level/Std Dev Used: Experiment: (Select multiple) ExpID: 60 - Rat 99-04406 - - D-Mannitol ExpID: 61 - Rat 99-04408 - - D-Mannitol ExpID: 62 - Rat 99-04410 - - D-Mannitol ExpID: 63 - Rat 99-04412 - - D-Mannitol Gene ID(s): (Enter a return after each CloneID i.e. One per line) AA963928 AA997797 AA957359 AA998164 Regulatory Flags: (Select multiple) All p53 Validation Ratio: Hits in Replicate Experiments Validate using the above ratio: Sort by: GeneID Mean Ratio Value # of Hits Then: GeneID Mean Ratio Value # of Hits Then: GeneID Mean Ratio Value # of Hits Display Ratio Outlier Statistics (Warning: may require a long display time): Indicate CVs Greater Than: Indicate log₂ Intensity Products Less Than: Scale of Individual Ratio Values: Linear Log₂

(b) Validated Outliers @ 95.00 Confidence Level Binomial Probability of chance occurrences at Ratio outlier frequency queried:
$$P(k \text{ out of } n) = \frac{n!}{k!(n-k)!} (p^k)(q^{n-k})$$
 n = the number of expts k = the number of times an outlier occurred p = the probability of random occurrence q = the probability of an outlier not occurring P: exactly 3 out of 4 0.00047500

Experiment: Rat 99-04406 Treatment: D-Mannitol @ 5000 mg/Kg/day X 14 Day(s) ExpID: 60 Chips Omitted: P30S88 P31S21
 Experiment: Rat 99-04408 Treatment: D-Mannitol @ 5000 mg/Kg/day X 14 Day(s) ExpID: 61 Chips Omitted: P30S90 P31S22
 Experiment: Rat 99-04410 Treatment: D-Mannitol @ 5000 mg/Kg/day X 14 Day(s) ExpID: 62 Chips Omitted: P30S92 P31S23
 Experiment: Rat 99-04412 Treatment: D-Mannitol @ 5000 mg/Kg/day X 14 Day(s) ExpID: 63 Chips Omitted: P31S29 P31S30

Found 4

GenBank Acc	Gene ID	Replicate Exp Hit Ratio	Flag	PRC	Description	log ₂ Cal. Ratio	Mean(log ₂ Cal. Ratio)	Standard Deviation	Standard Error	CV
AA963928	AA963928	3 out of 4	N/A	9G6	"Rattus norvegicus syndecan mRNA, complete cds"	0.89...174.4...605.10/299.70 1.83...15.9...557.20/115.80 0.83...18.3...851.30/378.80	1.18	0.560	0.323	0.473
AA997797	AA997797	4 out of 4	N/A	11G12	"Rattus norvegicus brain digoxin carrier protein mRNA, complete cds"	1.08...21.9...3022.30/1294.40 1.05...24.3...7042.00/3130.90 0.84...19.5...1379.50/569.00 0.82...26.3...13564.90/6108.50	0.95	0.137	0.068	0.145
AA957359	AA957359	3 out of 4	N/A	9F11	"Rattus norvegicus p55CDC mRNA, complete cds"	-1.18...23.8...2720.30/5563.10 -1.29...24.5...3247.50/7289.00 -1.06...22.5...2023.00/3076.90	-1.18	0.114	0.066	0.097
AA998164	AA998164	3 out of 4	N/A	12B11	Rat mRNA for cyclin B	-2.64...16.6...135.80/754.60 -3.06...17.2...140.00/1093.90 -2.18...14.9...97.90/525.30	-2.63	0.437	0.253	0.166

Select a link to see current UniGene build information for a clone

(c) 1 records satisfy the query in UniGene for terms aa997797 and for the organism Rattus norvegicus

UniGene	Description	Symbol
Rn.5641	Rattus norvegicus brain digoxin carrier protein mRNA, complete cds	

Fig. 1. MAPS differentially expressed gene validation. Screenshots from the MAPS gene expression validation pages illustrating the query form for detecting valid differentially expressed genes across multiple experiments (a), results of genes meeting specified search criteria with probability of chance occurrence (b) and display of the hyperlink to the cDNA clone current UniGene cluster (c).

labeling and hybridization, and survey the quality control and assurance of processed cDNA chips. Connection to the database is accomplished by open database connectivity, a ColdFusion application server (Allaire), and an IIS web server (Microsoft). Structured query language, proprietary ColdFusion tagging, and JavaScript are integrated in platform-independent web browser forms to provide users with seamless access to the database. Queried information and data are displayed in hypertext markup language tables and formatted to allow users to drill-down to detailed record information as well as hyperlink to interface with remote biological resources (Figure 1).

In the coming age of high-throughput genomics analysis using microarray technology scientists will rely heavily on robust database systems and statistical procedures to assist in the organization and interpretation of results. MAPS is a platform for scientists to not only manage information and analyze microarray gene expression data but also proceed towards discovering heuristic knowledge from biological information and data.

ACKNOWLEDGEMENTS

The authors would like to acknowledge our collaborators Raymond Stoll, Kerry Blanchard, and Supriya Jayadev for the study of rat liver gene expression, and thank Yidong Chen for ArraySuite data analysis software, Joseph Haseman for statistical consulting, Richard Lowry for JavaScript binomial probability calculations and Alex Merrick and Leping Li for review of the manuscript.

REFERENCES

Casella,G. and Berger,R.L. (1990) *Statistical Inference*. Duxbury Press, Belmont, CA.
 Chen,Y., Dougherty,E.R. and Bittner,M.L. (1997) Ratio-based decisions and the quantitative analysis of cDNA microarray images. *J. Biomedical Optics*, **2**, 364–374.
 Ermolaeva,O. *et al.* (1998) Data management and analysis for gene expression arrays. *Nature Genet.*, **20**, 19–23.
 Hughes,T.R. *et al.* (2000) Functional discovery via a compendium of expression profiles. *Cell*, **102**, 109–126.
 Wittes,J. and Friedman,H.P. (1999) Searching for evidence of altered gene expression: a comment on statistical analysis of microarray data. *J. Natl Cancer Inst.*, **91**, 400–401.