

ArrayXPath II: mapping and visualizing microarray gene-expression data with biomedical ontologies and integrated biological pathway resources using Scalable Vector Graphics

Hee-Joon Chung¹, Chan Hee Park¹, Mi Ryung Han¹, Seokho Lee¹, Jung Hun Ohn¹, Jihoon Kim¹, Jihun Kim¹ and Ju Han Kim^{1,2,*}

¹Seoul National University Biomedical Informatics (SNUBI) and ²Human Genome Research Institute, Seoul National University College of Medicine, Seoul 110-799, Korea

Received February 13, 2005; Revised and Accepted March 31, 2005

ABSTRACT

Summary: ArrayXPath (<http://www.snubi.org/software/ArrayXPath/>) is a web-based service for mapping and visualizing microarray gene-expression data with integrated biological pathway resources using Scalable Vector Graphics (SVG). Deciphering the crosstalk among pathways and integrating biomedical ontologies and knowledge bases may help biological interpretation of microarray data. ArrayXPath is empowered by integrating gene-pathway, disease-pathway, drug-pathway and pathway-pathway correlations with integrated Gene Ontology, Medical Subject Headings and OMIM Morbid Map-based annotations. We applied Fisher's exact test and relative risk to evaluate the statistical significance of the correlations. ArrayXPath produces Javascript-enabled SVGs for web-enabled interactive visualization of gene-expression profiles integrated with gene-pathway-disease interactions enriched by biomedical ontologies.

INTRODUCTION

Cluster analysis is one of the most powerful methods for the exploratory analysis of gene-expression data. Gene-expression clusters based on similarity measures between expression profiles have positional associations along the chromosomes (1,2), exhibit common *cis*-regulatory elements in their upstream regions (3) and are coordinated by shared sets of regulators (4). Gene-expression clusters can be assigned to the well-known functional categories of the MIPS classification

(5), the Gene Ontology (GO) terms (6) or pathway resources (7) using annotations from public databases (3,8,9).

ArrayXPath (7) is a web-based application that (i) receives a clustered gene-expression profile of any microarray platform in a tab-delimited text format; (ii) automatically resolves the microarray probe identifiers (i.e. GenBank accession number, UniGene ID, LocusLink ID, official gene symbol, SwissProt ID or TrEMBL ID); (iii) searches major public pathway resources (i.e. GenMAPP, KEGG, BioCarta and PharmGKB Pathways); (iv) maps the different identifier sets between microarray probes and pathway nodes; (v) tests the statistical significance of the associations between gene-expression clusters and pathways (hence providing an automated annotation of clusters with the ranked pathways); (vi) visualizes expression levels onto pathways and (vii) allows web-based user navigation through multiple clusters and pathways enriched with animation features, using Javascript-enabled Scalable Vector Graphics (SVG).

Although biological pathways can provide key information about the organization of biological systems, relatively small number (i.e. ~3000) of genes compared with the estimated number (i.e. >30 000) of genes for our species, as reported by our previous work (7), do appear in major pathway resources, resulting in low coverage for genome-wide expression data analysis. Although GO-based annotations give lesser information than pathway-based ones, increasingly more gene products are being annotated by GO terms, resulting in much higher coverage. As of February 2005, we found that 13 949 LocusLink IDs have at least one GO annotations. Therefore, integrating lesser-knowledge-higher-coverage GO-based annotations can complement more-knowledge-lower-coverage pathway-based annotations for microarray data analysis.

Gene-pathway correlation alone may not be sufficient for deciphering the genomic secret of normal and pathological

*To whom correspondence should be addressed. Tel: +82 2 740 8320; Fax: +82 2 742 5947; Email: juhan@snu.ac.kr

physiology. Integrating not only biological (i.e. GO) but also clinical ontologies like the disease nomenclature system supported by Medical Subject Headings (MeSH) can provide further information for genotype-to-phenotype associations. OMIM Morbid Map provides valuable gene-disease correlations. Integrating drug-pathway correlations from PharmGKB Pathways (10) can also improve ArrayXPath.

We found that pathways had significant and informative crosstalk. Many genes appear in multiple pathways. Systematic analysis and interactive visualization of the complex crosstalk structures among pathways, pathway nodes (i.e. gene products) and gene-expression clusters may help understanding gene-pathway correlations. Figure 1a shows the concept diagram of new ArrayXPath that integrates the quinta-partite graph structure of cluster, gene, disease, pathway and GO-term associations from multiple resources.

Here we present an improved version of ArrayXPath. In addition to the functionalities described above, this is a software that (i) tests the statistical significance of the associations between gene-expression clusters and GO-based annotations to complement pathway-based microarray data analysis; (ii) allows users to search disease-related pathways; (iii) visualizes the global crosstalk of biological pathways by measuring and mapping the similarity distances superimposed by the local crosstalk of the subset of pathways matched to input gene-expression clusters and (iv) visualizes the detailed local crosstalk through gene-cluster, gene-pathway, gene-disease and gene-GO associations using interactive multi-partite graph representations in SVG. OMIM knowledge base and drug-pathway correlations from PharmGKB Pathways are also tightly integrated to ArrayXPath.

INPUT AND OUTPUT

Input

Input to ArrayXPath is a common tab-delimited text file for a clustered gene expression profile: <Probe ID>-<Cluster ID>-<Expression level at condition>]. The first column must contain either GenBank accession number, UniGene ID, LocusLink ID, SwissProt ID, TrEMBL ID or an official gene symbol. The second column contains the cluster ID. The third to *i*th columns are optional and contain expression levels. ArrayXPath does not perform cluster analysis *per se*. The input format is designed primarily for a partitioned clustering algorithm (i.e. *K*-means or Self-Organizing Maps) but a clustering result from a hierarchical algorithm (i.e. dendrogram) may be applied by carefully choosing a threshold. One can search disease-related pathways and their correlations by entering a disease name.

Output

ArrayXPath produces a list of the best-matching pathways and GO terms for each cluster with statistical significance scores of non-random association. Relevant pathways are listed in ascending order of *P*-values (and multiple-comparison corrected *Q*-values) (11). ArrayXPath provides a summary statistic for the overall mapping between input clusters and all pathways and GO terms matched.

If one chooses a pathway among the list, ArrayXPath outputs a Javascript-enabled SVG file, color-coded both by expression level and by cluster membership at each pathway-node level. If one chooses a cluster, ArrayXPath outputs cluster-pathway-disease diagram with significantly associated GO terms and OMIM information (Figure 1c). The cluster-centric view visualizes the related genes, pathways and diseases by measuring the shared membership of gene products. The whole quinta-partite associations (Figure 1a) can be interactively navigated by choosing one of the cluster, pathway or disease nodes from the graph in SVG.

Each node in the pathway graph and the correlation multi-partite graph is enriched with a hyperlink to an automated summary page for the corresponding gene product(s) provided by our integrated database: GRIP (Genome Research Informatics Pipeline, <http://grip.snubi.org/>) (7).

METHODS

Pathway integration and resolving diverse identifiers

ArrayXPath searches publicly available major pathway resources including KEGG, GenMAPP, BioCarta and PharmGKB Pathways. We have created a repository of meta-information by parsing SBML files for KEGG (<http://www.systems-biology.org/001/001.html>) and HTML files for GenMAPP (http://www.genmapp.org/HTML_MAPPs/Human/MAPPIndex_Hs_Contributed.htm) and BioCarta (<http://www.biocarta.com/genes/allPathways.asp>), and by manually encoding PharmGKB pathways (<http://www.pharmgkb.org/search/pathway/pathway.jsp>). A variety of gene-product identifiers including GenBank accession number, UniGene ID, LocusLink ID, EC number, official gene symbol, SwissProt ID and TrEMBL ID are inconsistently used for the pathway nodes as well as microarray probes, resulting in enormous ambiguity in integrating data from different resources. By integrating major databases including GenBank, UniGene, LocusLink, Homologene, SwissProt, Ensemble, UCSC Golden Path and NetAffyx (<http://www.affymetrix.com/analysis/index.affx>), ArrayXPath automatically matches the probe identifiers of microarray data to the identifiers of pathway nodes. When a pathway node is a composite type, i.e. consists of more than one element, ArrayXPath separately matches and visualizes each probe identifier to the corresponding individual element of the composite object.

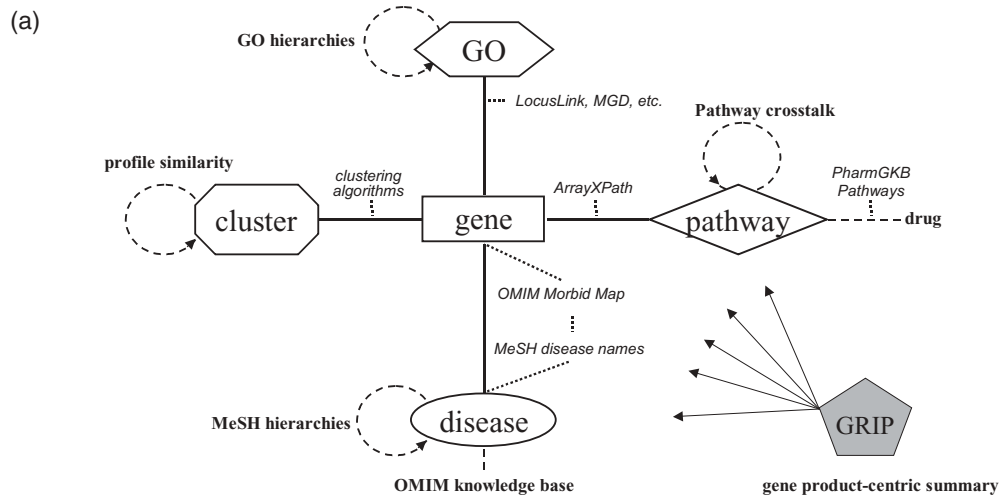
Table 1 shows the distribution of the pathway nodes identified from KEGG, GenMAPP, BioCarta and PharmGKB Pathways for *Homo sapiens*, *Mus musculus* and *Rattus norvegicus*. We found 1942 redundant nodes representing genes and proteins for the 45 GenMAPP pathways. Among the 1454 non-redundant elements, ArrayXPath successfully assigned 1391 gene products (95.7%) to either official gene symbols ($n = 1329$; 91.4%), LocusLink ID ($n = 39$; 2.7%) or SwissProt ID ($n = 23$; 1.6%). It is trivial to resolve LocusLink ID once we obtain official gene symbol. The number 39 for the LocusLink-ID column in Table 1 means that we could resolve LocusLink IDs for the 1368 (= 1329 + 39) elements. Similarly, there were far more than 23 elements having SwissProt IDs. The number 23 for SwissProt-ID column means that we could resolve SwissProt IDs for the 23 elements for which we failed to resolve official gene symbol or

LocusLink ID. Only 63 (4.3%) elements in GenMAPP remain unresolved because of intractable ambiguity. KEGG has 256 non-composite (i.e. simple) and 121 composite elements (i.e. enzymes), containing 256 and 505 gene products, respectively. Among the 256 simple-type elements, 21 appear as members of composite type elements. Overall, KEGG has 740 unique elements and ArrayXPath successfully assigned all of them (100%) either to official gene symbol ($n = 720$, 97.3%) or LocusLink ($n = 20$, 2.7%). PharmGKB Pathways added 133 official gene symbols and 1 LocusLink ID.

Overall, ArrayXPath identified 3151 gene products for the four major pathways in our species. We created a pre-computed association table of these elements to all resolvable IDs and to official gene symbols for reliable mapping of incoming microarray-probe identifiers.

Search pathways by disease name in MeSH (PathMeSH)

ArrayXPath allows one to search disease-related pathways. The OMIM Morbid Map (<http://www.ncbi.nlm.nih.gov/Omim/getmorbid.cgi>) contains official gene symbol, alias gene symbol and cytogenetic location of disease-related genes with OMIM ID and the associated disease name. The C category among the 15 branches of MeSH contains disease names and their entry terms with hierarchical structure. We extracted the gene-related and disease-related information from OMIM and MeSH. We mapped the disease names by using exact keyword match method provided by MeSH. We mapped the disease-related genes onto the integrated pathway resources by using our integrated database, GRIP, as described above. Among the 3259 official gene symbols resolvable from



(b) Screenshot of the ArrayXPath web interface showing search results for 'Breast Neoplasms'.

Search keyword is Breast Neoplasms that is disease name.
 PathMeSH find 32 pathways and search 18 genes.
 This view is pathway-based gene information. Another view is [gene-based pathway information](#).

Gene list
 [AR] [ATM] [BRCA1] [BRCA2] [CDH1] [CHEK2] [ESR1] [FBP1CC1] [TP53] [TSG101] [BACH1] [BRIP1] [CDS1] [PHB] [PPM1D] [RAD54L] [SLC22A1L] [XRCC3]

Select ordering Pvalue Relative Risk

Pathway	Gene Symbol	Drug	OMIM	Cytoband	P-value	Relative Risk
BioCarta//Hs_Role of BRCA1, BRCA2 and ATR in Cancer Susceptibility	ATM		607585	11q22.3	0.000001	55.0183150
	BRCA1		113705	17q21		
	BRCA2		600185	13q12.3		
	CHEK2		604373	22q12.1		
	TP53		191170	17p13.1		
BioCarta//Hs_ATM Signaling Pathway	ATM		607585	11q22.3	0.000030	45.2030075
	BRCA1		113705	17q21		
	CHEK2		604373	22q12.1		
	TP53		191170	17p13.1		
BioCarta//Hs_Cell Cycle: G2/M Checkpoint	ATM		607585	11q22.3	0.000080	35.7261505
	BRCA1		113705	17q21		
	CHEK2		604373	22q12.1		
	TP53		191170	17p13.1		
BioCarta//Hs_Regulation of cell cycle progression by Plk3	ATM		607585	11q22.3	0.000089	75.4250000
	CHEK2		604373	22q12.1		
	TP53		191170	17p13.1		
BioCarta//Hs_PB Tumor Suppressor/Checkpoint Signaling in response to					0.0022818	29.9515385

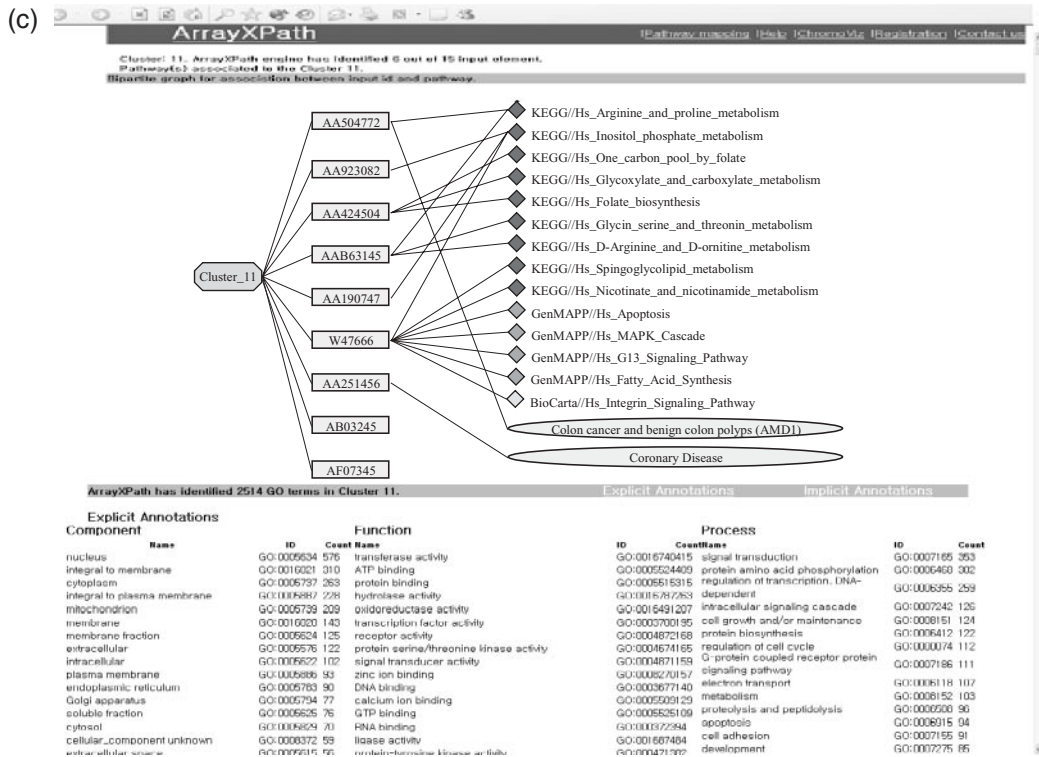


Figure 1. ArrayXPath functions. (a) Gene-cluster, gene-pathway, gene-disease and gene-GO associations (solid lines) are the building blocks of the quinta-partite graph representation used by ArrayXPath integration. Dotted lines explain how the associations are determined. GO and MeSH have their own hierarchical organizations, clusters can be organized by profile similarity measures, and pathways by crosstalks (broken circular arrows). (b) PathMeSH returns a list of disease-related pathways with statistical significance scores by integrating pathway resources, MeSH disease names and OMIM Morbid Map. (c) When one chooses a cluster, ArrayXPath outputs the cluster-centric view of the associations of related genes, pathways and diseases through the shared membership of gene products. The whole quinta-partite associations can be interactively navigated by choosing cluster, pathway or disease node from the graph in SVG. ArrayXPath also provides GO-based annotation and OMIM information to complement pathway-based analysis of gene expression clusters.

Table 1. Distribution of pathway-node identifiers among the major pathway resources

	Pathway	Gene/protein		ID resolution				Metabolite		Embedded pathway	Free text description				
		Simple	Complex	Redundant	Total	OGS	LL	SP	UR						
<i>H.sapiens</i>															
KEGG	70	(256) ^a (505) ^a 740	(121)	(469) (637) (1106)	740	720	20	0	0	1896	(2624)	0	(0)	121	(275)
GenMAPP	45	1454		(1942)	1391	1329	39	23	63	83	(97)	4	(4)	130	(372)
BioCarta	346	1584		(8976)	1584	1580	4	0	0	0	(0)	50	(141)	18	(53)
PharmGKB	9	134		(189)	134	133	1	0	0	11	(25)	1	(1)	23	(26)
Overall	470	3088		(12 900)	3025	2938	64	23	63	1990	(2746)	55	(146)	55	(146)
<i>M.musculus</i>															
BioCarta	277	1260		(7646)	1260	1224	36	0	0	1	(1)	37	(113)	75	(311)
<i>R.norvegicus</i>															
KEGG	63	(160) (277) 451	(72)	(527) (753) (1280)	451	282	169	0	0	1774	(2435)	0	(0)	167	(72)

Numbers in parentheses are redundant counts. KEGG has 121 composite elements containing 505 identifiable gene products. OGS, official gene symbol; LL, LocusLink; SP, SwissProt; and UR, unresolved.

^aThere were 21 elements redundant in the simple (256) and composite (505) elements so that 740 unique elements were found in the 70 KEGG human pathways.

the Morbid Map, we found that 2395 genes had disease names that could be mapped to MeSH disease names through headings or entry terms. We successfully mapped 1928 genes onto both MeSH disease names and pathway nodes (i.e. gene products). It means that about 64% (i.e. 1928/3025) (Table 1) of the non-redundant nodes in all pathway resources of our

species have at least one link to human pathophysiology through standard disease name(s) in MeSH.

If the input disease name is matched to the corresponding MeSH heading or entry term, ArrayXPath outputs the list of the pathways containing the disease-related gene product (Figure 1b). ArrayXPath determines the statistical significance

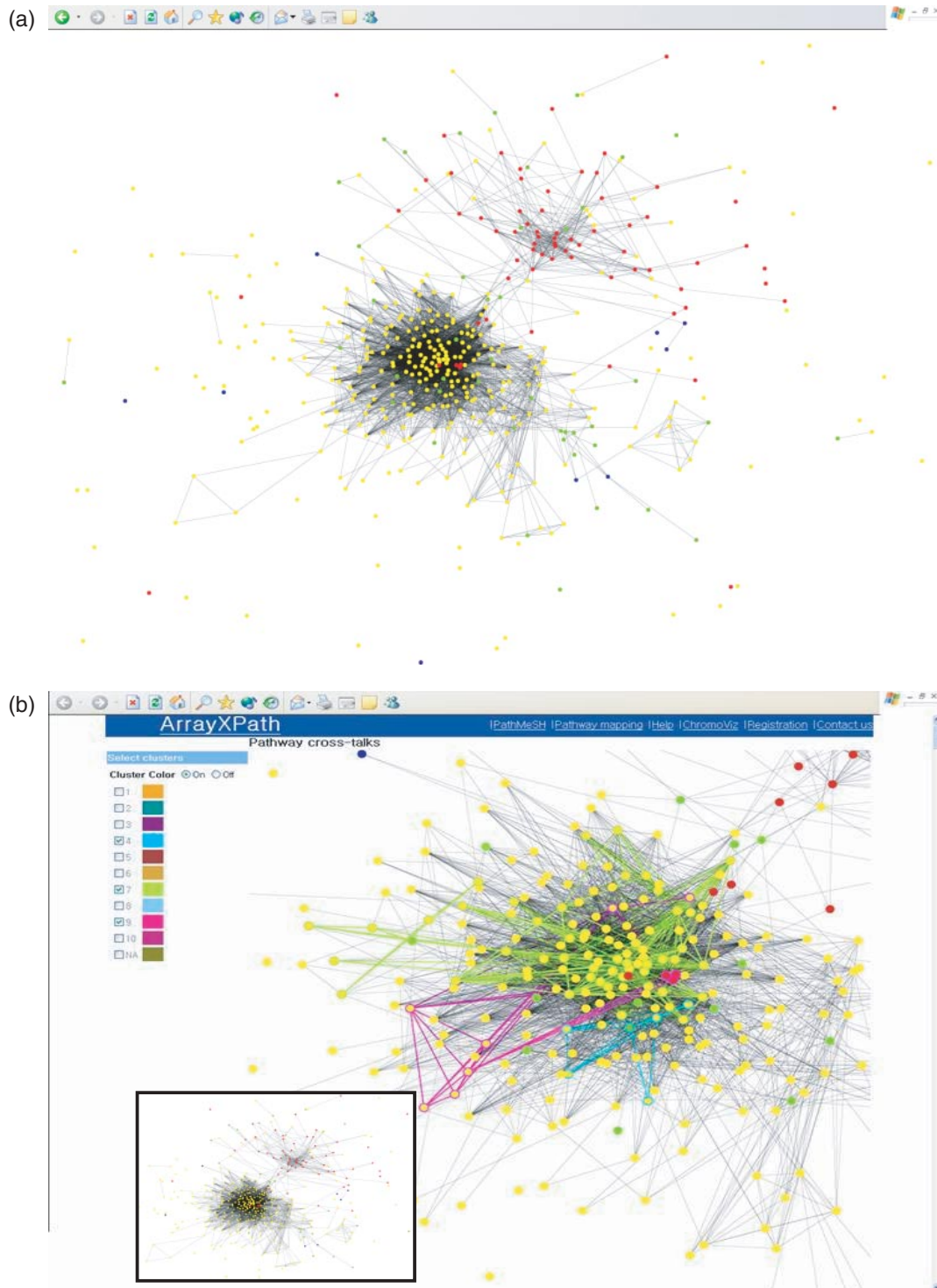


Figure 2. Pathway crosstalk. (a) Calculating pairwise similarity matrix between each pair of pathways and applying multi-dimensional scaling method created the global crosstalk graph of major biological pathways. Yellow nodes represent BioCarta, green nodes GenMAPP, red nodes KEGG and blue nodes PharmGKB Pathways (see Methods). (b) ArrayXPath interactively visualizes the local crosstalk (i.e. sky-blue, light-green and purple lines) of the pathways associated with the selected clusters (i.e. clusters 4, 7 and 9 respectively), superimposed on the global pathway crosstalk graph.

of the association between a pathway and a disease name in terms of the non-random proportion of matched entities. ArrayXPath applies Fisher's exact test by constructing a 2×2 contingency table containing the two pathway memberships (within and without the pathway) as column variables and

the disease memberships (within and without the disease) as row variables. We used Fisher's exact test because a large sample approximation is inappropriate in the pathway case (a 2×2 table often contains a cell with expected values < 5).

Visualization of the correlational structure of biological pathways

Visualizing the crosstalk among pathways may reveal important biological understanding (12). We created a pairwise similarity matrix of pathway distances by calculating the ratio of the number of genes in the intersection to that in the union of each pair of pathways. Multi-dimensional scaling of the similarity matrix and drawing the edges of the pathway pairs above certain similarity thresholds create a global crosstalk graph among all biological pathways (Figure 2a). It is demonstrated that metabolic pathways of KEGG (i.e. red nodes) and signal transduction pathways of BioCarta (i.e. yellow nodes) tend to form separate clusters.

ArrayXPath interactively visualizes the local crosstalk of the pathways associated with the user-selected clusters, superimposed on the global pathway crosstalk graph (Figure 2b). Figure 2b shows an interesting example where the pathways mapped onto the selected three gene-expression clusters form distinct clusters.

ArrayXPath interactively visualizes the detailed local crosstalk. The shared membership of gene products in gene-expression clusters, pathways and disease names can be captured by multi-partite graph representations in SVG (Figure 1a and c). By selecting a node from one view, one can interactively navigate the different views of the whole associations.

Mapping GO-based annotations

Among the 33 108 LocusLink IDs (as of February 2005), we identified 3025 gene products (i.e. ~9%) in major pathway resources (Table 1). Integrating GO-based annotations covering 13 949 LocusLink IDs (as of 2005 February) can help filling the gaps for pathway-based microarray data analysis. ArrayXPath provides both implicit and explicit GO annotations (Figure 1c) (13). While explicit annotation provides the GO terms directly mapped onto the members of gene-expression clusters, implicit annotation considers all ancestor terms from the GO hierarchical tree structure, providing general understanding. We applied hypergeometric distribution to evaluate the statistical significance of the associations.

DISCUSSION

ArrayXPath is a web-based service for mapping and visualizing microarray gene expression clusters with biomedical ontologies and major biological pathway resources using SVG. It permits one to input a clustered gene expression data in a tab-delimited text format via an Internet connection.

We found that integrating biomedical ontologies including GO-based annotations, disease names supported by MeSH, and genotype-to-phenotype information from OMIM Morbid Map greatly improve the capability of ArrayXPath to interpret gene-expression profiles. Integrated analysis and interactive visualization of the global and local crosstalks among pathways can facilitate system-level understanding of microarray gene-expression data. Although we evaluated the statistical significance of each association in the present study, combined analysis of the quinta-partite relationship may improve our

inference. Future studies are required to develop a computational method to reconstruct the whole correlational structure and extract more biology from gene-expression microarray data. Standard web-based integration of a wide range of bioinformatics modules and heterogeneous genomic data will obviously help advance biological science.

ACKNOWLEDGEMENTS

The authors thank all SNUBI members and external testers for their efforts to input and evaluate the pathway diagrams. This study was supported by a grant from Korea Health 21 R&D Project, Ministry of Health and Welfare, Republic of Korea (0412-MI01-0416-0002). Funding to pay the Open Access publication charges for this article was provided by the above mentioned grant.

Conflict of interest statement. None declared.

REFERENCES

- Roy,P.J., Stuart,J.M., Lurd,J. and Kim,S.K. (2002) Chromosomal clustering of muscle-expressed genes in *Caenorhabditis elegans*. *Nature*, **418**, 975–979.
- Lercher,M.J., Urrutia,A. and Hurst,L.D. (2002) Clustering of housekeeping genes provides a unified model of gene order in the human genome. *Nature Genet.*, **31**, 180–183.
- Tavazoie,S., Hughes,J.D., Campbell,M.J., Cho,R.J. and Church,G.M. (1999) Systematic determination of genetic network architecture. *Nature Genet.*, **22**, 281–285.
- Segal,E., Shapira,M., Regev,A., Pe'er,D., Botstein,D., Koller,D. and Friedman,N. (2003) Module networks: identifying regulatory modules and their condition-specific regulators from gene expression data. *Nature Genet.*, **34**, 166–176.
- Mewes,H.W., Frishman,D., Guldener,U., Mannhaupt,G., Mayer,K., Mokrejs,M., Morgenstern,B., Munsterkotter,M., Rudd,S. and Weil,B. (1999) MIPS: a database for genomes and protein sequences. *Nucleic Acids Res.*, **30**, 31–34.
- Ashburner,M., Ball,C.A., Blake,J.A., Botstein,D., Butler,H., Cherry,J.M., Davis,A.P., Dolinski,K., Dwight,S.S., Eppig,J.T. *et al.* (2000) Gene Ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nature Genet.*, **25**, 25–29.
- Chung,H.J., Kim,M., Park,C.H., Kim,J. and Kim,J.H. (2004) ArrayXPath: mapping and visualizing microarray gene expression data with integrated biological pathway resources using Scalable Vector Graphics. *Nucleic Acids Res.*, **32**, W464–W464.
- Dennis,G., Jr, Sherman,B.T., Hosack,D.A., Yang,J., Gao,W., Lane,H.C. and Lempicki,R.A. (2003) DAVID: Database for Annotation, Visualization, and Integrated Discovery. *Genome Biol.*, **4**, R60.
- Zeeberg,B.R., Feng,W., Wang,G., Wang,M.D., Fojo,A.T., Sunshine,M., Narasimhan,S., Kane,D.W., Reinhold,W.C., Lababidi,S. *et al.* (2003) GoMiner: a resource for biological interpretation of genomic and proteomic data. *Genome Biol.*, **4**, R28.
- Hewett,M., Oliver,D.E., Rubin,D.L., Easton,K.L., Stuart,J.M., Altman,R.B. and Klein,T.E. (2002) PharmGKB: the Pharmacogenetics Knowledge Base. *Nucleic Acids Res.*, **32**, 163–165.
- Storey,J.D. and Tibshirani,R. (2003) Statistical significance for genomewide studies. *Proc. Natl Acad. Sci. USA*, **100**, 9440–9445.
- Zheng,C.J., Zhou,H., Xie,B., Han,L.Y., Yap,C.W. and Chen,Y.Z. (2004) TRMP: a database of therapeutically relevant multiple pathway. *Bioinformatics*, **20**, 2236–2241.
- Robinson,P.N., Wollstein,A., Böhme,U. and Beattie,B. (2004) Ontologizing gene-expression microarray data: characterizing clusters with Gene Ontology. *Bioinformatics*, **20**, 979–981.