

Can we interpret biological data by computational means?

Ju Han Kim, M.D., Ph.D.
juhan@snu.ac.kr

Integrative biochip informatics

- Exploration: mapping and clustering
- Pattern recognition: finding gene – drug assoc.
- Prognostic subgroup prediction
- Interpretation
 - ✓ GRIP: automatic probe annotation
 - ✓ Xperanto: databasing expression
 - ✓ BioCANDI: database-integrated analysis tool
 - ✓ ChromoViz: visualizing regional regulation
 - ✓ GOTree and GOChase: GO-based annotation
 - ✓ ArrayXPath: pathway-based interpretation
 - ✓ PathMeSH: disease - gene - pathway
- BioEMR: towards real-time clinical trials

Integrative biochip informatics

- Exploration: mapping and clustering
- Pattern recognition: finding gene – drug assoc.
- Prognostic subgroup prediction
- Interpretation
 - ✓ GRIP: automatic probe annotation
 - ✓ Xperanto: databasing expression
 - ✓ BioCANDI: database-integrated analysis tool
 - ✓ ChromoViz: visualizing regional regulation
 - ✓ GOTree and GOChase: GO-based annotation
 - ✓ ArrayXPath: pathway-based interpretation
 - ✓ PathMeSH: disease - gene - pathway
- BioEMR: towards real-time clinical trials

Multi-dimensional Scaling

Table 1 Flying Mileages Between 10 American Cities

	Atlanta	Chicago	Denver	Houston	Los Angeles	Manhattan	New York	San Francisco	Seattle	Washington, DC
0	587	1212	701	1936	604	748	2139	2132	547	Atlanta
587	0	928	540	1181	1133	1489	1937	1937	547	Chicago
1212	928	0	879	831	1726	1631	949	1021	1494	Denver
701	540	879	0	1374	968	1420	1645	1891	1250	Houston
1936	1181	1133	1374	0	2350	2451	1489	1937	1043	Los Angeles
604	1188	1733	968	2359	0	1092	2584	2734	923	Miami
748	713	1631	1420	2451	1092	0	2571	2408	205	New York
2139	1937	949	1869	2451	2591	2351	0	678	2445	San Fr.
2132	1937	1021	1391	2591	2734	2469	678	0	1785	Seattle
547	547	1494	1250	2734	2469	2442	2329	0	1494	Wash.

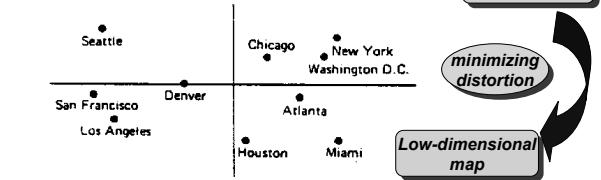
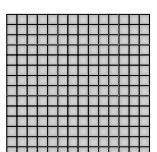


Figure 1 CMDS of flying mileages between 10 American cities
Young, 1985. Encyclopedia of Statistical Sciences

Subjective similarity between colors



2-d map

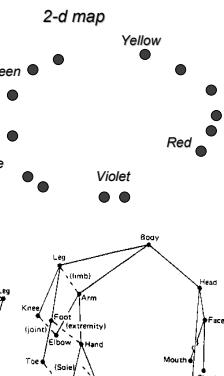
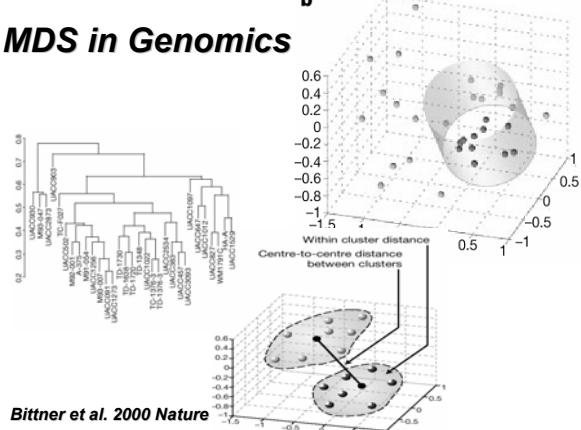


Figure 2 (a) RMDS of children's similarity judgments about 15 body parts; (b) RMDS of adults' similarity judgments about 15 body parts.

MDS in Genomics



Bittner et al. 2000 Nature

Temporal patterns revealed by MDS

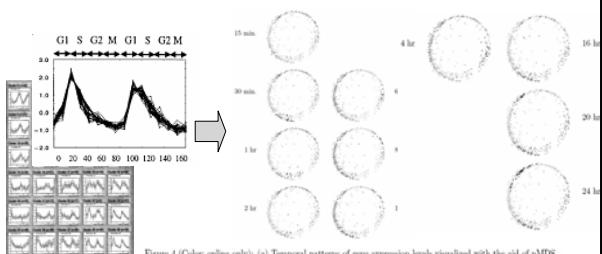
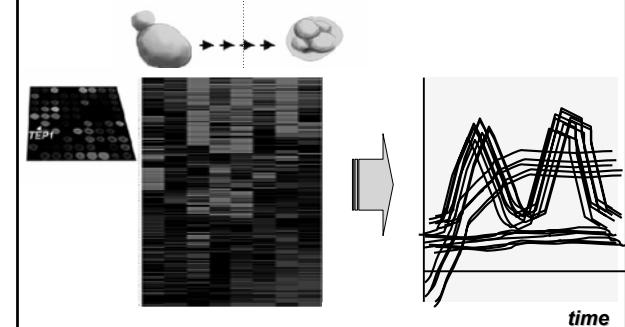


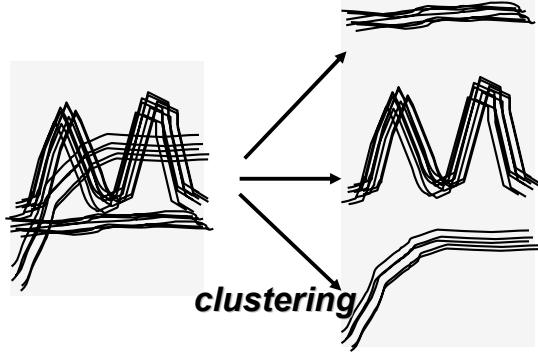
Figure 4 (Color online only): (a) Temporal patterns of gene expression levels visualized with the aid of tMDS. Colors indicate relative intensity of experimental values normalized so that $\sum_i t_{ij} = 0$ and $\sum_i t_{ij}^2 = 1$ where t_{ij} is experimental variable of j th genes at time t . (red > 1.6, yellow > 1.2, green > 0.8, pale blue > 0.4, gray < 0.4). Time sequences are the same as explained at Fig. 1. (b) Gene expression data as a function of the angle measured from the vertical axis in (a). The horizontal axis corresponds to t . The color convention is the same as in (a). The figures are arranged in two columns, but this is solely for the layout purpose; there is no distinction between two columns.

Taguchi et al., Bioinformatics 2004

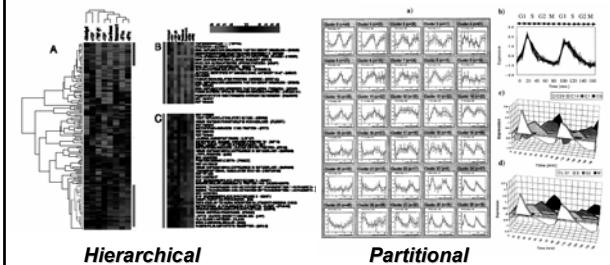
Biochip informatics: clustering



Biochip informatics: clustering

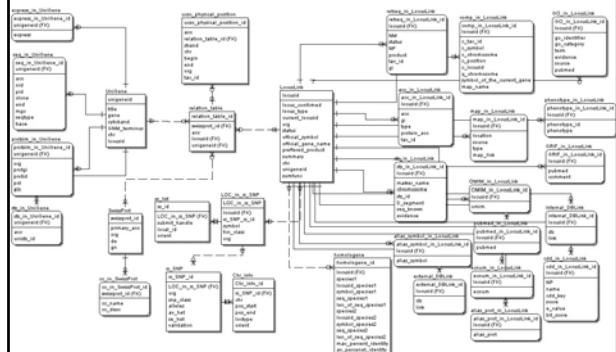


Hierarchical & Partitional Clustering



GRIP
Genome Research Informatics Pipeline
Integrating biological databases

Part of GRIP integration



GRIP: gene / protein information

The GRIP (Gene Research Information) system provides detailed information about genes and proteins. It includes a search interface, a genome browser, and various analysis tools. A screenshot shows a search result for a gene, displaying its physical location, sequence, and associated publications.

GRIP: SNP information

The GRIP SNP information interface displays details for a specific SNP (rs214906). It shows the SNP ID, chromosome, position, alleles, and other relevant data. A screenshot also shows a search result for the SNP, listing its characteristics and associations.

Xperanto

Expressionist's Esperanto in XML

Databasing expression with standards

nature

19 April 2001 Volume 410 Issue no 6831

Free and public expression

After a slow start, progress towards developing public repositories for gene expression data is poised to accelerate. For the many biologists working with DNA microarrays, that should be welcome news.

With a single format for gene expression data, databases should be able to 'talk' to one another and exchange data. The existence of a standard language should also spur development of software tools to query the databases, and to manage and display gene expression data.

nature

26 September 2002 Volume 419 Issue no 6905

Microarray standards at last

Not a moment too soon, the microarray community has issued guidelines that will make their data much more useful and accessible. *Nature* and the *Nature* research journals will respond accordingly.

You read a paper with a fascinating conclusion about the expression of several genes. You decide to use some of the same experiments on your system of choice. But when you wade through hundreds of pages of supplementary information, it's clear that the data won't be compatible. The authors that you want to compare with have used different array designs.

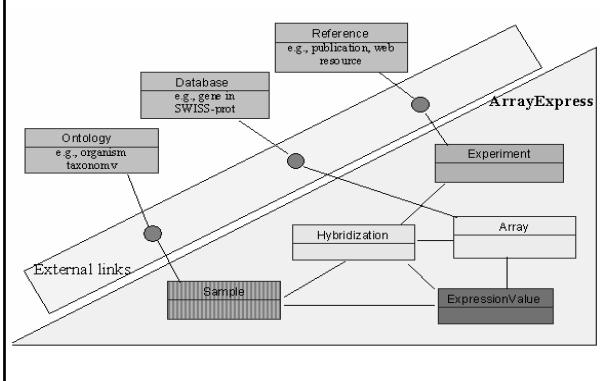
Welcome to the exciting but frustrating world of DNA microarray research. Microarrays are plastic or glass chips spotted with tiny amounts of thousands of probes, used to query the activity levels of that many genes in any tissue or organism at one time. Variables in every step of the experiment often make cross-paper comparison virtually impossible. Microarray papers also pose a considerable strain on the refereeing process: the vast amounts of data mean that critical review is a monumental task.

Now, however, the field is finally getting some standards. Not enough, though. The new guidelines are not given enough teeth, forcing editors to review them to think that they must transcribe the primary data set. In other cases, the primary data provided in proprietary software and so are impossible to comment on. Many journals allowed authors to put the huge data files on their own websites for the review process, until it became clear that unscrupulous authors compensated the anonymity of referees by tracking who had visited the website.

For authors, the proposal provides a checklist of variables that should be included in every microarray publication, at http://www.mged.org/Workgroups/MIAAME/miaame_checklist.html. This checklist, with all variables completed, would be supplied as supplementary material to manuscripts. The MGED proposal suggests that journals require submission of microarray data to one of two databases emerging as the main public repositories: GEO (www.ncbi.nlm.nih.gov/geo/) or ArrayExpress (www.ebi.ac.uk/arrayexpress).

Harrid editors can rejoice that, at last, the community is taming the unruly beast that is microarray information. Therefore, all submissions to *Nature* and the *Nature* family of journals received on or after 1 December containing new microarray experiments must include the mailing of five compact discs to the editor. These discs should contain the raw data in a standard, open-access MAGE standard. The information must be supplied in a format that could be read by widely available software packages. Data integral to the paper's conclusions should be submitted to the ArrayExpress or GEO databases, with accession numbers where available, supplied at or before acceptance for publication.

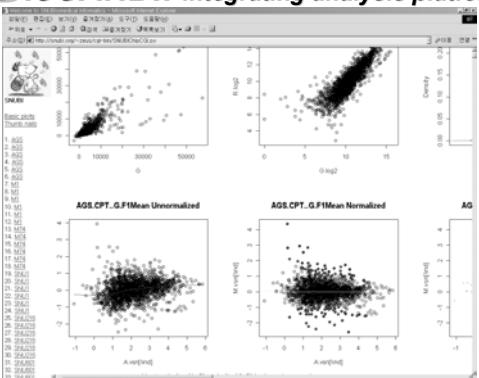
How much data should authors provide to the community?



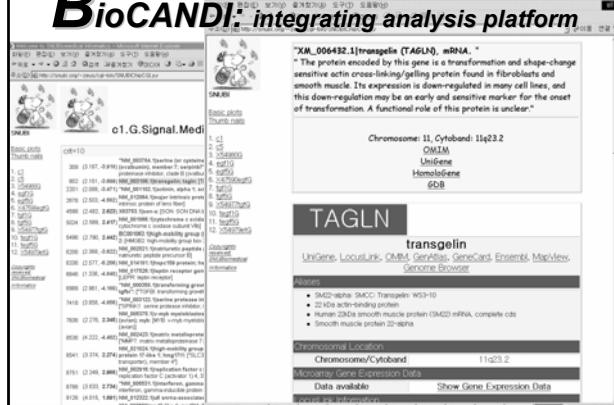
BioCANDI
BioChip Analysis N Data Integration

Database-backed microarray analysis

BioCANDI: integrating analysis platform



 BioCANDI: integrating analysis platform
bioRxiv preprint doi: <https://doi.org/10.1101/2023.09.21.553711>; this version posted September 21, 2023. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).



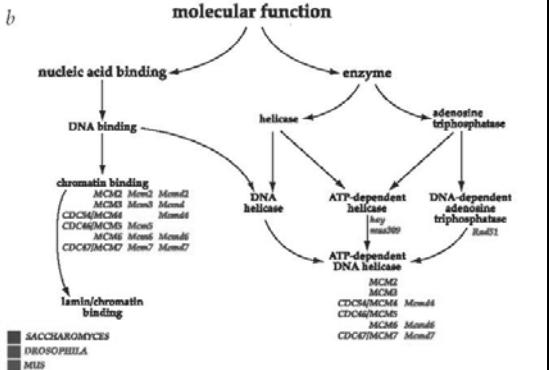
BioCANDI: integrating analysis platform



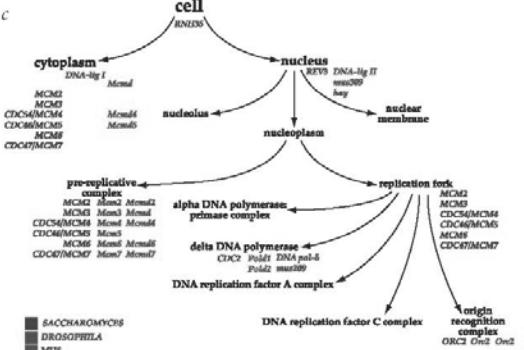
GOChase and GOTree

Interpreting gene expression clusters

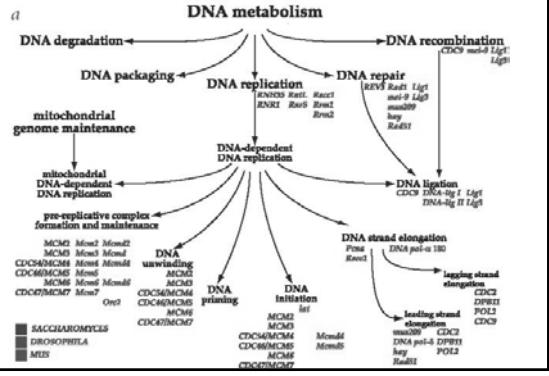
Gene Ontology



Gene Ontology



Gene Ontology



Upstream(Motif / Regulon) search

Structure & Function: upstream search

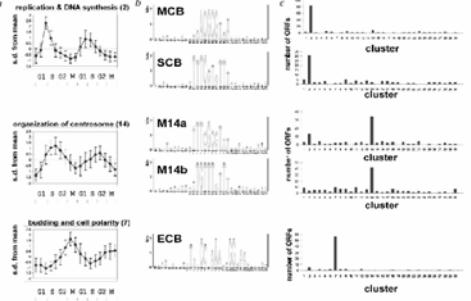


Fig. 5 Top: Prominent clusters, their motifs and overall distribution in all clusters. **a**, Mean temporal representation of a cluster, measured according to the biological functions. The letters represent motif frequency. Each letter is the mean of the motif frequencies of all clusters containing that motif. Error bars indicate the standard deviation of each cluster about the mean of the particular time point. **b**, Sequence logos representation of the motif(s) distributed within the cluster. The height of each letter is proportional to its frequency. The letters are sorted with the most frequent one on top. The overall height of the stack signifies information content of the sequence motif.

GOA-based interpretation of clusters

Table I: Comparison of GOTM with related software^a

	FatiGO	Onto-Express	GOSurfer	GoMiner	GOTM
Interface/OS Input Identifier	Web Unigene ID, Gene symbol, Swiss-Prot ID, Ensembl ID, GenBank ID	Web GenBank ID, Affymetrix probe set ID, LocusID, Unigene ID	Windows Affymetrix probe set ID, LocusID, Unigene ID	Windows/Mac HUGO gene names	Web LocusID, Gene symbol, Affymetrix probe set ID, UnigeneID, Swiss-Prot ID, Ensembl ID
Multi-level analysis Visualization of classification Statistical Analysis	No Bar chart, Table, Fixed tree	Yes Bar chart, Table	Yes Fixed tree	Yes Expandable tree, DAG	Yes Expandable tree, Bar chart, Fixed tree
Correction for multiple tests Visualization of enriched GO categories	Yes Bar chart, Table	Yes Bar chart, Table	No Highlight in the full GOTree	No Highlight in the full DAG	No Sub-tree and DAG of enriched GO categories, Highlight in the full GOTree and bar chart
Availability	Free	Partially free	Free	Free	Free

Gene Ontology based interpretations



Every error that has a beginning has an end.

GOChase: correcting errors from gene ontology-based annotations for gene products

Yu Kang Park, Chan Hee Park, Jin Han Kim

The Gene Ontology (GO) is a controlled biological vocabulary that provides three structured networks of terms to describe biological processes, cellular components, and molecular functions. Many databases of gene products are annotated using the GO vocabularies. We found that some GO-updating operations are not easily traceable by the current biological databases and GO browsers. Consequently, numerous annotation errors arise and are propagated throughout biological databases and GO-based bioinformatic analyses. GOChaser is a set of web-based utilities to detect and correct the errors in GO-based annotations.

A availability: <http://www.ncbi.nlm.nih.gov/Software/GCG/Chloro.html>

GOChase:

Park et al., Bioinformatics 2004

The Gene Ontology (GO) is a rapidly growing hierarchy of controlled vocabularies for describing biological processes, cellular components, and molecular functions.

GOTree and GOChase

ArrayXPath and ChromoViz

Interpreting gene expression clusters

ChromoViz

Pathway map with expression prof.

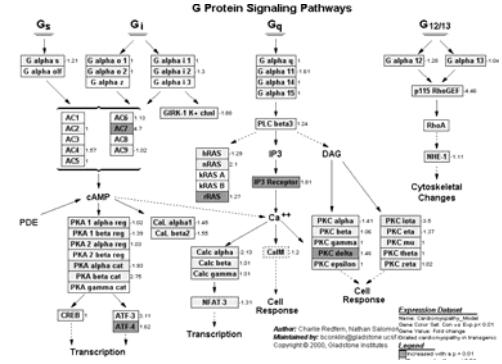
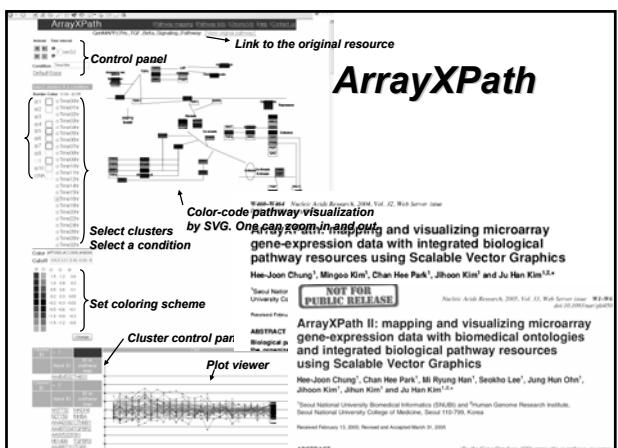


Table 1. Distribution of pathway-node identifiers among the major pathway resources

Pathway	Gene / Protein			ID resolution				Metaboplite	Embedded pathway	Free text description
	simple	compl	redundant	Total	OGS	LL	SP			



File Edit View Go Bookmarks Tools Help
Welcome to ArrayXPath
ArrayXPath | Pathway mapping | Help | ChromoViz | Registration | Contact us

E-Mail: Password: Log in

ArrayXPath has searched

Source	# pathways	# nodes (non-redundant)	# nodes (redundant)
GemMAPP	45	1391	1875
PharmGKB	9	154	180
KEGG	70	740	1788
BioCarta	346	1594	8969

ArrayXPath has identified 490 (27.62%) out of 1774 input elements in the following pathways

Source	# pathways	# identified nodes	# total nodes
GemMAPP	41	204	747
PharmGKB	5	14	110
KEGG	58	124	1474
BioCarta	288	290	4843

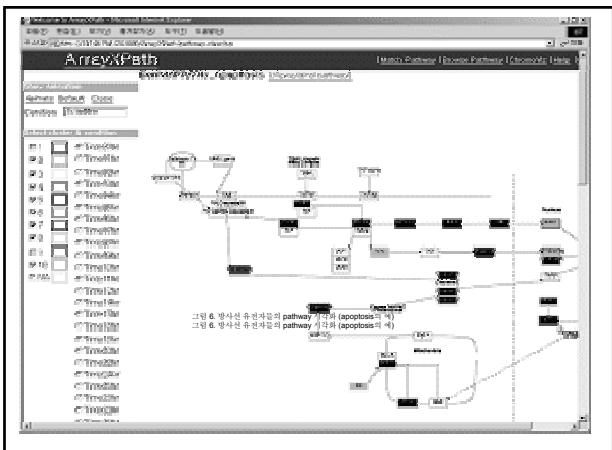
Set P-value cut off 0.05 Change

Cluster 1 has 29 gene products [Show Bipartite Graph]
ArrayXPath has identified 8 out of 29 input elements in 7 out of 45 GemMAPP pathways.

Clusters

GenMAPP 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 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

Find Next Find Previous Highlight Match Case



Current Molecular Medicine 2005, 5, 11-21

11

Pathway and Ontology Analysis: Emerging Approaches Connecting Transcriptome Data and Clinical Endpoints

L. Yue and W.C. Reisdorf

Bioinformatics, GlaxoSmithKline, USA

Abstract: The increasing use of gene expression profiling offers great promise in clinical research into disease biology and its treatment. Along with the ability to measure changing expression levels in thousands of genes at once, comes the challenge of analyzing and interpreting the vast sets of data generated. Analysis tools are evolving rapidly to meet such challenges. The next step is to interpret observed changes in terms of the biological properties or relationships underlying them. One powerful approach is to make associations between the genes that are under investigation and well-known biochemical or signaling pathways, and further to assess the significance of such associations.



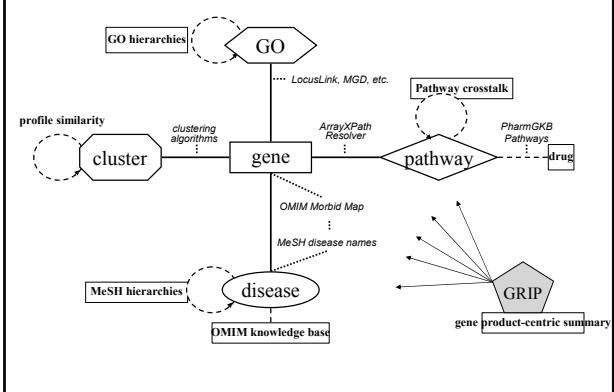
Table 1. Tools used to analyze pathway and biological relationship data.

Tool	Pathway/Ontology	Statistical Methods	Visualization	Refs	URL
ArrayXPath	Pathway	F, M	Y	[38]	http://www.snuib.org/software/ArrayXPath/
Pathway Miner	Pathway	F	Y	[39]	http://www.biorag.org/pathway.html
Knowledge Editor	Both	None	Y	[53]	http://scope.gsc.riken.jp/
EASE	Pathway	O	N	[47]	http://www.DAVID.ncbi.nlm.nih.gov
GeneMerge	Both	H	N	[52]	http://www.oeb.harvard.edu/freelab/publications/GeneMerge.html
MAPPFinder	Both	O	Y		
DAVID	Both	None	Y		
GFINDer	Both	H, F, C	N		
OntoTools	Both	H, F, C, O	N		
GOSurfer/ClipInfo	Ontology	C, M	Y		
GOFish	Ontology	None	Y		
GOGet/GoView	Ontology	None	Y		
GOTree Machine	Ontology	H	Y		
EntrezGO	Ontology	F, M	N		
FuncAssociate	Ontology	F, M	N		

One complication when analyzing such large data sets is that one is effectively doing thousands of comparisons of genes and pathways, and one should also include a correction factor for multiple hypothesis testing. Some of the tools listed in the table incorporate such a correction, but a general discussion of the methods and their relevant merits is beyond the scope of this review (and the expertise of the authors).

Of the tools primarily focused on pathways, ArrayXPath [38] appears to be the most comprehensive. The tool accepts a range of gene identifiers (Swissprot, GenBank, UniGene, LocusLink, etc.) and includes KEGG, BioCarta and GenMAPP as its data source. ArrayXPath uses Fisher's exact test along with the False Discovery Rate method for multiple hypothesis testing correction and provides a very clear graphical display of results. Pathway Miner [39] can map the genes from up to four different experiments onto KEGG, BioCarta or GenMAPP pathways, and it can

Where to go from here? Understanding gene expression clusters



Welcome to ArrayXPath - Microsoft Internet Explorer

Search keyword is Breast Neoplasms that is disease name.
PathMeSH find 32 pathways and search 13 genes.
This view is pathway-based gene information. Another view is gene-based pathway information.

Gene list
 [AM] [ATM] [BRCA1] [BRCA2] [CDH1] [CMM12] [ESR1] [ERBBCC1] [TP53] [TP53BP1] [TP53IP1] [TP53IPI] [TP53M] [TP53N] [TP53NL] [TP53NL1] [TP53NL2]

Select ordering P-value Relative Risk

Pathway	Gene Symbol	Drug	CMM	Cytoband	P-value	Relative Risk
BioCarta/Hs_Role of BRCA1, BRCA2 and ATM in Cancer Susceptibility	ATM	60705	11q22.3	17q12	0.000000	55.039190
	BRCA1	61335	17q11			
	BRCA2	60205	13q12.3			
	CHEK2	60473	22q12.1			
	TP53	181170	17q12.1			
BioCarta/Hs_ATM Signaling Pathway	ATM	60705	11q22.3	17q11	0.000000	45.320000
	BRCA1	61335	17q11			
	CHEK2	60473	22q12.1			
	TP53	181170	17q12.1			
BioCarta/Hs_Cell Cycle G2M Checkpoint	ATM	60705	11q22.3	17q11	0.000000	35.720000
	BRCA1	61335	17q11			
	CHEK2	60473	22q12.1			
	TP53	181170	17q12.1			
BioCarta/Hs_Regulation of cell cycle progression by Pho	ATM	60705	11q22.3	17q11	0.000000	75.420000
	CHEK2	60473	22q12.1			
	TP53	181170	17q12.1			
BioCarta/Hs_Tumor Suppressor/Checkpoint Signaling in response to	TP53	60705	11q22.3	17q11	0.000000	25.961500

ArrayXPath

Cluster: H1. ArrayXPath engine has identified 6 out of 15 input element.

Biocarta graph for association between input list and pathway.

Cluster: H1
 AA504772
 AA923082
 AA424504
 AA631425
 AA190747
 AA747666
 AA251456
 AB07245
 AF07345
 Coronary Disease
 Colon cancer and benign colon polyps (AMDD1)
 BioCarta/Hs_Integrin_Signaling_Pathway
 GenMAP/Hs_Apoptosis
 GenMAP/Hs_MAPK_Cascade
 GenMAP/Hs_G1_3_Signaling_Pathway
 GenMAP/Hs_Fatty_Acid_Synthesis
 GenMAP/Hs_Nicotinate_and_nicotinamide_metabolism

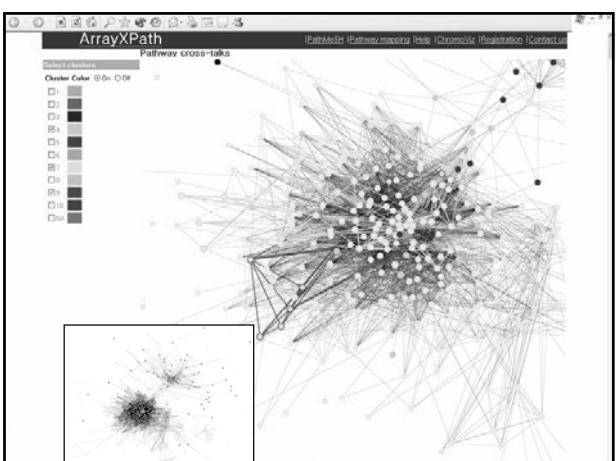
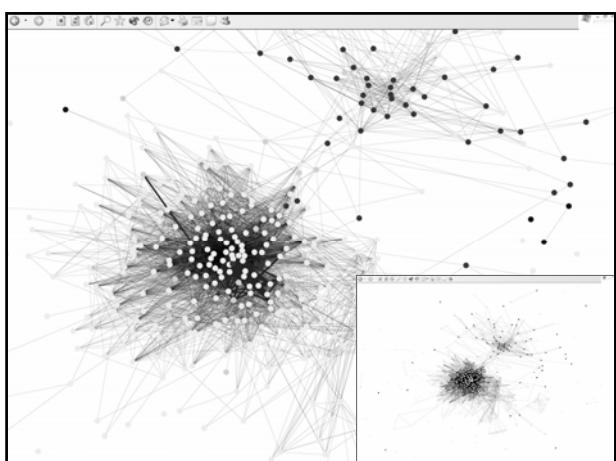
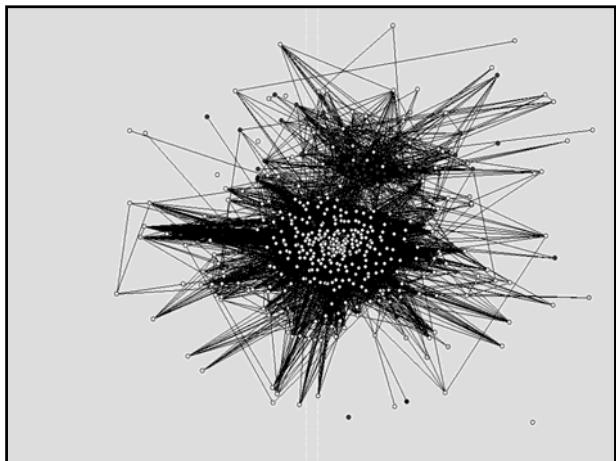
ArrayXPath has identified 2614 GO terms in Cluster: H1.

Explicit Annotations

Component	Function	Process
nucleus	G0009546 nucleic acid metabolism	G00095768 nucleic acid biosynthesis
Integral to membrane	G0006032 310 ATP binding	G0006460 protein amino acid phosphorylation
cytoskeleton	G0006032 263 protein binding	G0006951 505 regulation of transcription, DNA-templated
integral to plasma membrane	G0006951 263 protein binding	G0006959 259
mitochondrion	G0005079 239 oxidoreductase activity	G0007042 126
membrane	G0005079 239 oxidoreductase activity	G0007042 126
membrane fraction	G0005079 125 receptor activity	G000670188 protein biosynthesis
cellular	G0005079 125 receptor activity	G00068412 120
intracellular	G0005079 102 signal transducer activity	G000687119 G-protein coupled receptor protein
plasma membrane	G0005079 90 pho ion binding	G000697157 electron transport
extracellular matrix	G0005079 71 calcium ion binding	G000699129 metabolism
Goldapexaurus	G0005079 70 phosphorus binding	G000699136 metabolism
cellular fraction	G0005079 70 phosphorus binding	G000699136 metabolism
cytosol	G0005079 70 RNA binding	G000699154 synthesis
cellular component unknown	G0005079 69 kinase activity	G000699154 synthesis

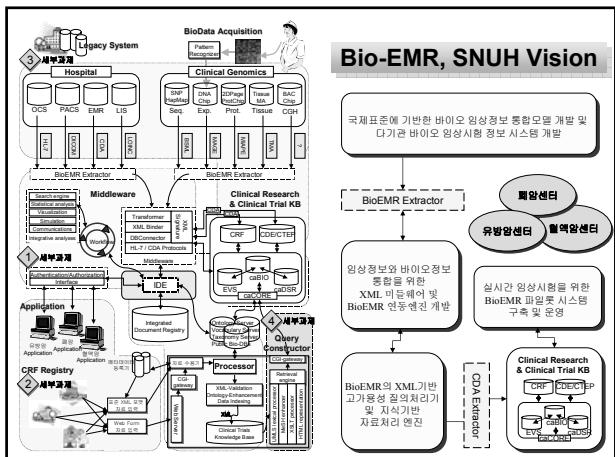
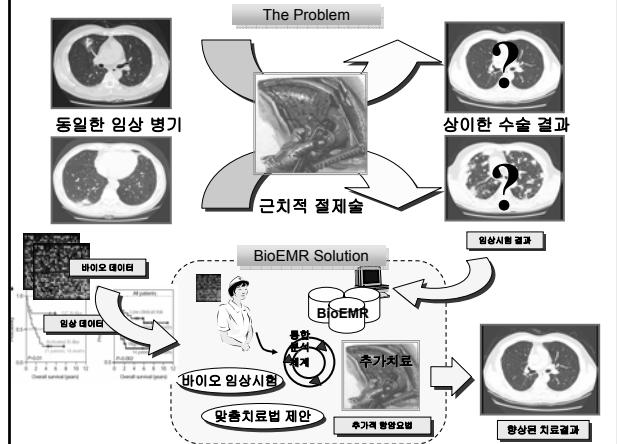
Logical Annotations

Implicit Annotations



Integrative biochip informatics

- Exploration: mapping and clustering
- Pattern recognition: finding gene – drug assoc.
- Prognostic subgroup prediction
- Interpretation
 - ✓ GRIP: automatic probe annotation
 - ✓ Xperanto: databasing expression
 - ✓ BioCANDI: database-integrated analysis tool
 - ✓ ChromoViz: visualizing regional regulation
 - ✓ GOTree and GOChase: GO-based annotation
 - ✓ ArrayXPath: pathway-based interpretation
 - ✓ PathMeSH: disease - gene - pathway
- BioEMR: towards real-time clinical trials



Then, isn't it just a tool?
Life

