

[S-3]

Development of Toxicogenomics Knowledge Base by Integrating High-throughput Genomic Data

Ju Han Kim, M.D., Ph.D., M.S.

Seoul National University Biomedical Informatics (SNUBI)

Bioinformatics is a rapidly emerging field of biomedical research. A flood of large-scale genomic and postgenomic data means that many of the challenges in biomedical research are now challenges in computational sciences. Postgenome informatics, powered by high throughput technologies and genomic-scale databases, is likely to transform our biomedical understanding forever much the same way that biochemistry did a generation ago. In this talk, I will describe how these technologies will impact toxicology research, introducing recent advances in databasing gene expression profiles with the emphasis on the necessity of tight integration of private and public databases and intelligent analysis toolkits. I will introduce some of our research efforts for toxicogenomics knowledge base. Xperanto (Expressionist's Esperanto in XML) integrates major data models for DNA microarray, tissue microarray and array CGH data with extended clinical and histo-pathological information models and supports analysis tools in an effort to establish a comprehensive knowledge base for toxicogenomics research. Each step will be given with real examples from ongoing research activities in the context of clinical relevance.

Development of Toxicogenomics Knowledgebase by Integrating High-throughput Genomic Technologies

Pharmaceutical side effects and Drug-drug interactions

- One in five new drugs have some serious side effects (by Lasser KE, *et al.* JAMA 2002 287:2215)
- Need ways to predict potential interactions and side effects
 - impact of environmental elements that have the potential for biological damage
 - the relationship of the activities of toxins to genetic makeup
- Previous tools for measuring toxic events
 - in vivo animal models such as rat, mouse, and dog
 - in vitro assays
 - Ames testing
 - micronucleus assays
 - unscheduled DNA synthesis measurement

New way to measure toxic events through gene expression profiles

- Human Genome Project + Microarray technology
 - interrogate the expression of tens of thousands of genes simultaneously in response to a toxic agent
 - investigate effects of human genetic variation on drug toxicity and efficacy

3

cDNA microarray schema

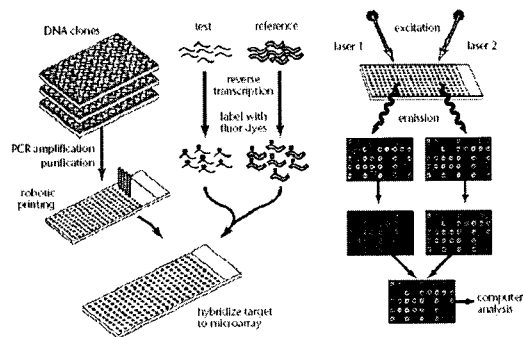


Figure from *Nature Genetics* (1999)
21:11

4

Gene expression profiling of chemotherapeutic drugs

Cluster analysis of temporal changes in gene expression of MCF-7 cells treated with doxorubicin

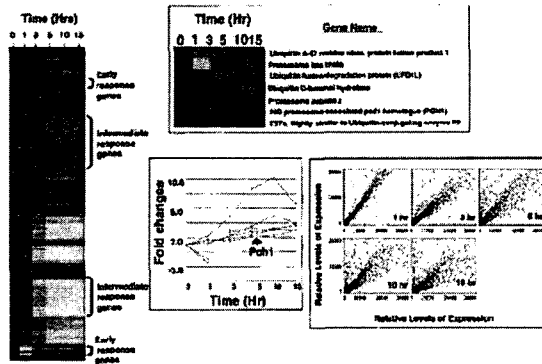


Figure from *Pharmaceutical Research* (2002) 19:1776

5

DNA arrays for toxicogenomics

DNA arrays allow quantitative measurements of gene expression to be made for tens of thousands of genes in parallel

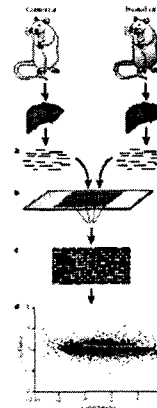


Figure from *Nature Reviews* (2001) 1:84

6

Gene expression profiling of toxic compounds

A compendium showing gene-expression change induced by 48 different hepatotoxic compounds in the rats

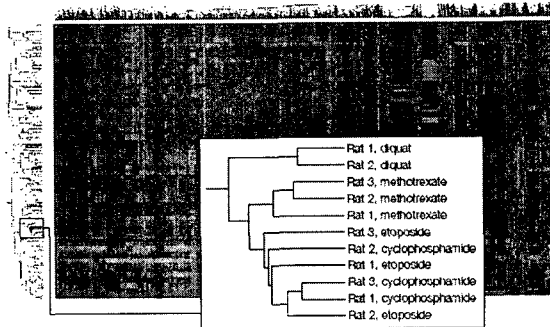
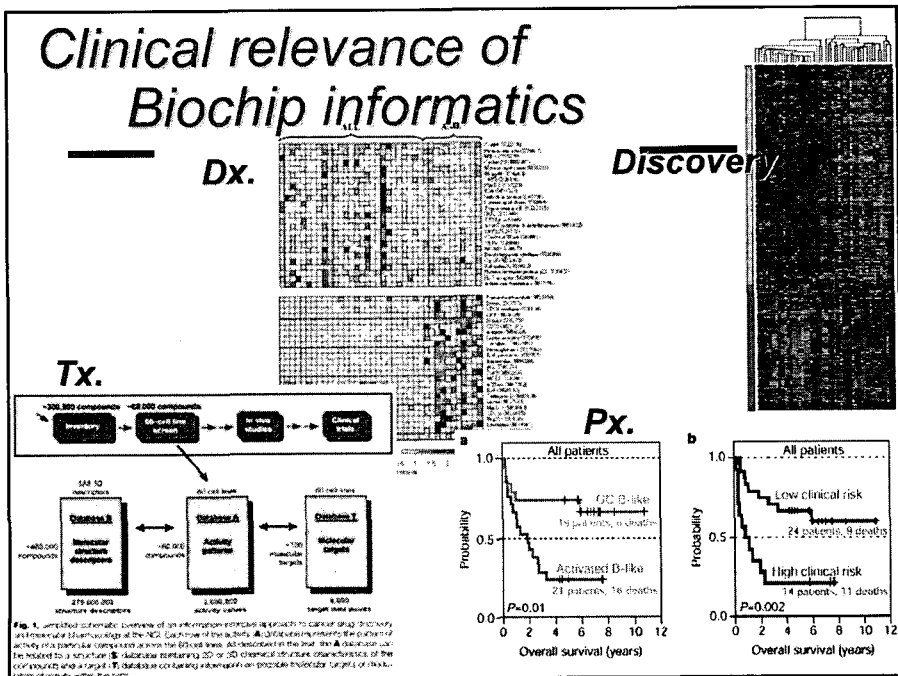
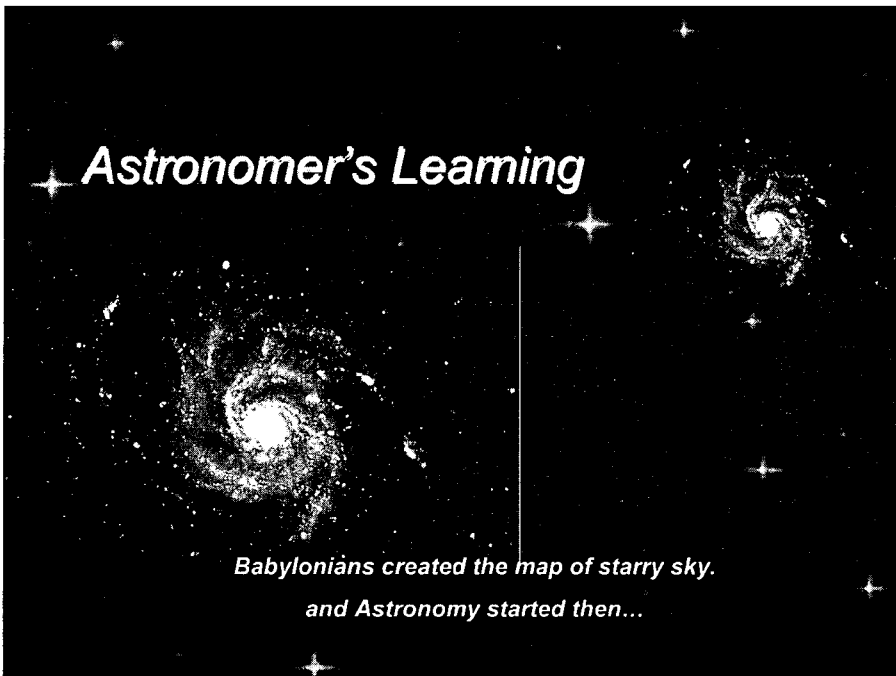
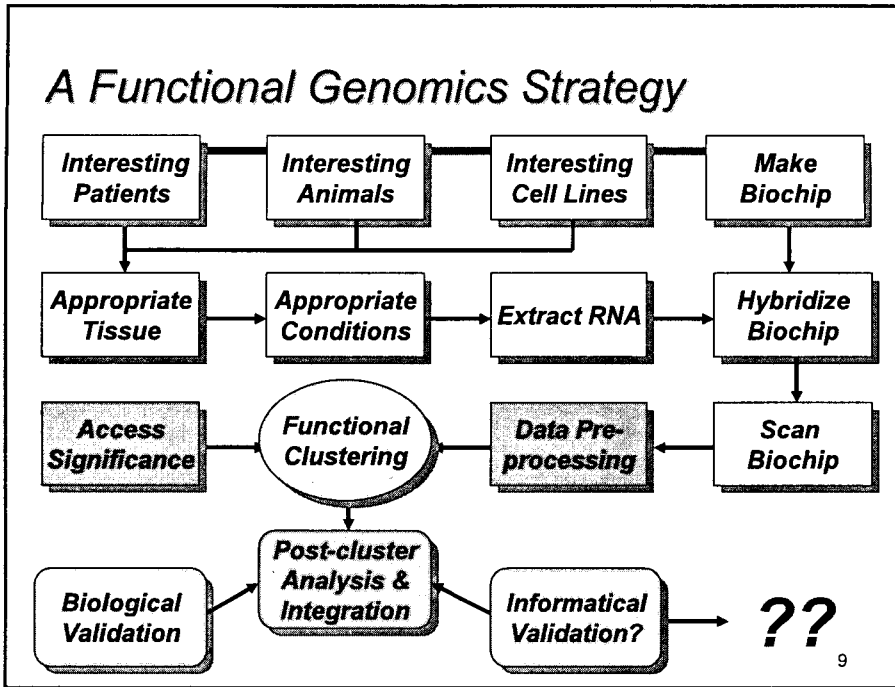


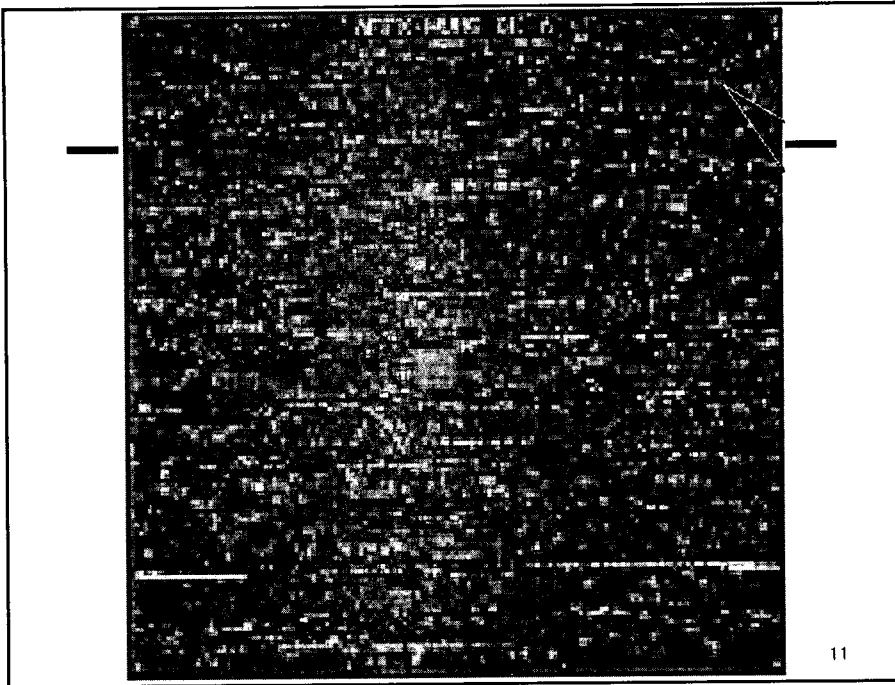
Figure from *Nature Reviews* (2001)
1:84

7

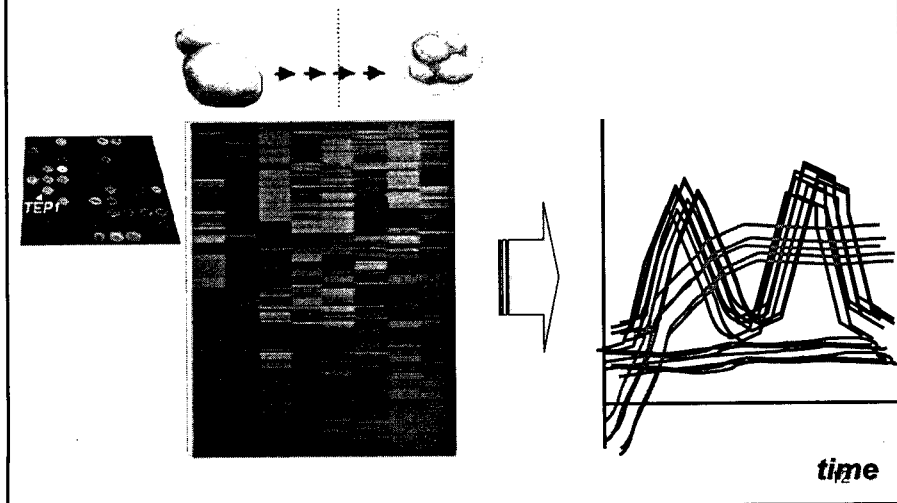
Clinical relevance of Biochip informatics



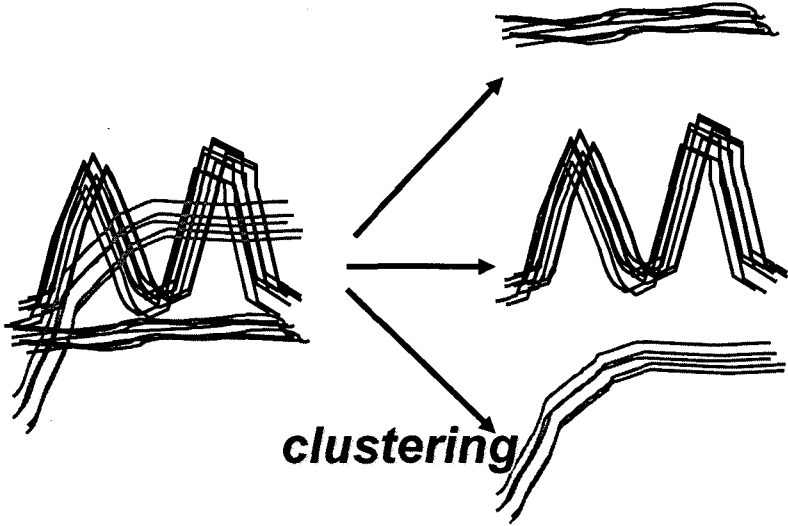




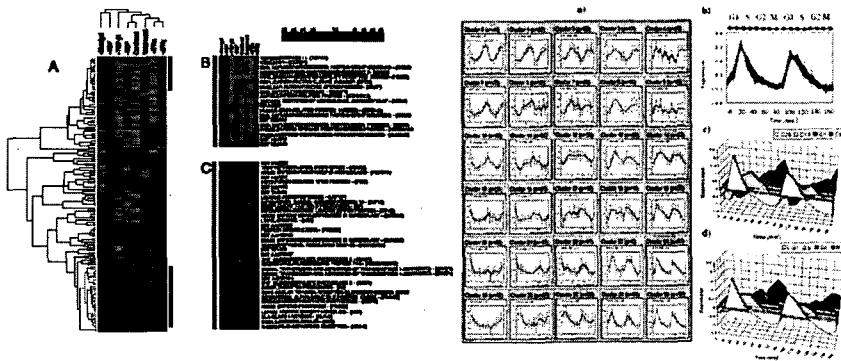
Biochip informatics: clustering



Biochip informatics: clustering

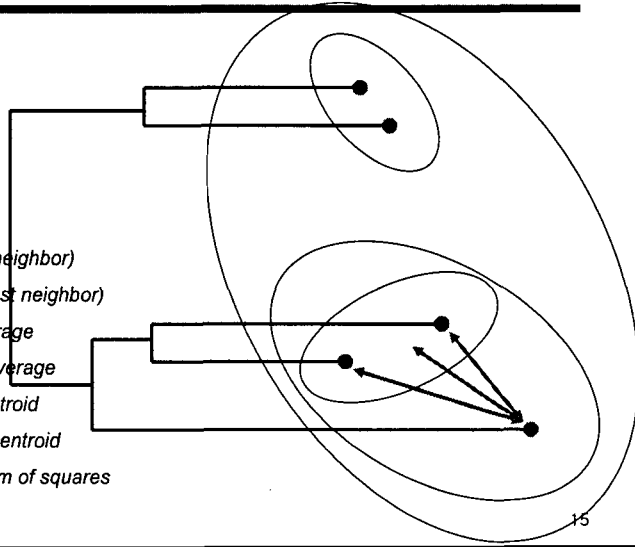


~~Hierarchical & Partitional Clustering~~



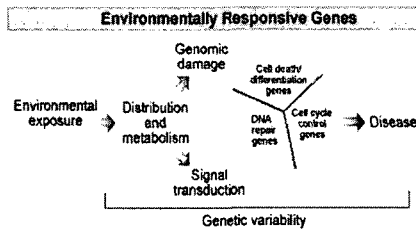
Hierarchical clustering in Genomics

- single-linkage (nearest neighbor)
- complete-linkage (farthest neighbor)
- weighed pair-group average
- unweighed pair-group average
- weighted pair-group centroid
- unweighted pair-group centroid
- Ward's method: min. sum of squares



Toxicogenomics

- the application of genomic tools to the study of biological responses to toxic substances



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.

17

nature

26 September 2002 Volume 419 Issue no 6906

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, you find that crucial details needed for replication are missing.

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 an monumental task.

Yet referees sometimes feel they are not given enough details, leading cautious reviewers to think that they must reanalyse the primary data set. In other cases, the primary data provided are 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 compromised 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.nsgf.org/Workgroups/MIAME/miame_checklist.html. This checklist, with all variables completed, would be supplied as supplementary information at the time of submission. The MGED group suggests that journals require submission of microarray data to either of two databases emerging as the main public repositories: GEO (www.ncbi.nlm.nih.gov/geo/) or ArrayExpress (www.ebi.ac.uk/arrayexpress/).

Harried 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 disks to the editor. These disks should include necessary information compliant with the MIAME 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?

Microarray standards are needed to facilitate moving data around

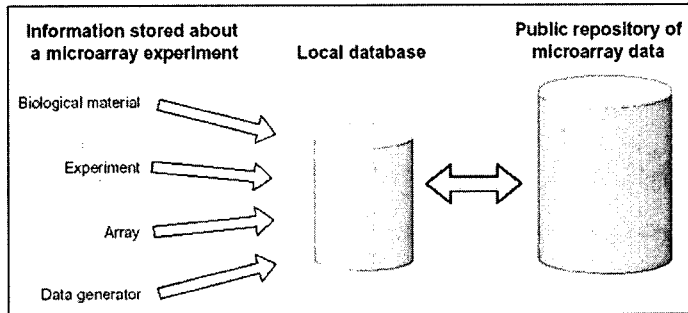


Figure from *Nature Genetics* (2002) 32: 460

19

MGED standards

The screenshot shows the 'MGED Home' website. At the top, it reads 'Microarray Gene Expression Data Society - MGED Society'. Below this, there is a paragraph describing the society's mission: 'The Microarray Gene Expression Data (MGED) Society is an international organization of biologists, computer scientists, and data analysts that aims to facilitate the sharing of microarray data generated by functional genomics and proteomics experiments.' It also mentions the current focus on establishing standards for microarray data annotation and packaging. A 'Latest News' section is visible at the bottom, with entries for 'MGED 6 Entry Registration Details' dated 21/06/2003 and 'MIAME-TOX Document' dated 14/02/2003.

<http://www.mged.org>

20

Microarray Standards

- **MIAME**

- *Minimum Information About a Microarray Experiment*
- *Experimental Design, Array Design, Hybridization, Samples, Measurements and Normalization*

- **MAGE-ML**

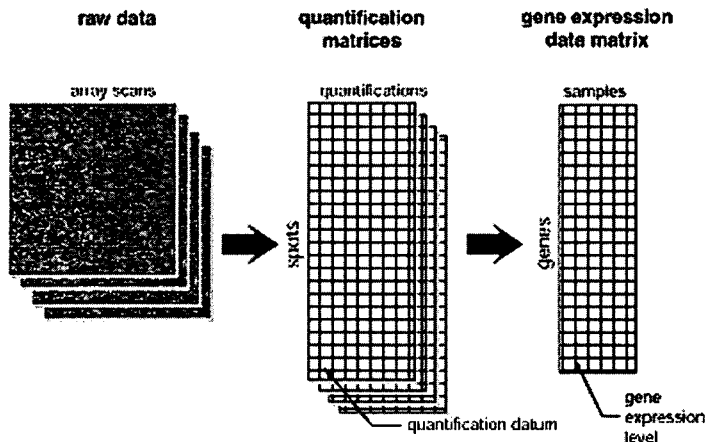
- *XML Implementation of the MIAME Standard*
- *De Facto Widespread Industry Support*

- **MAGE-OM**

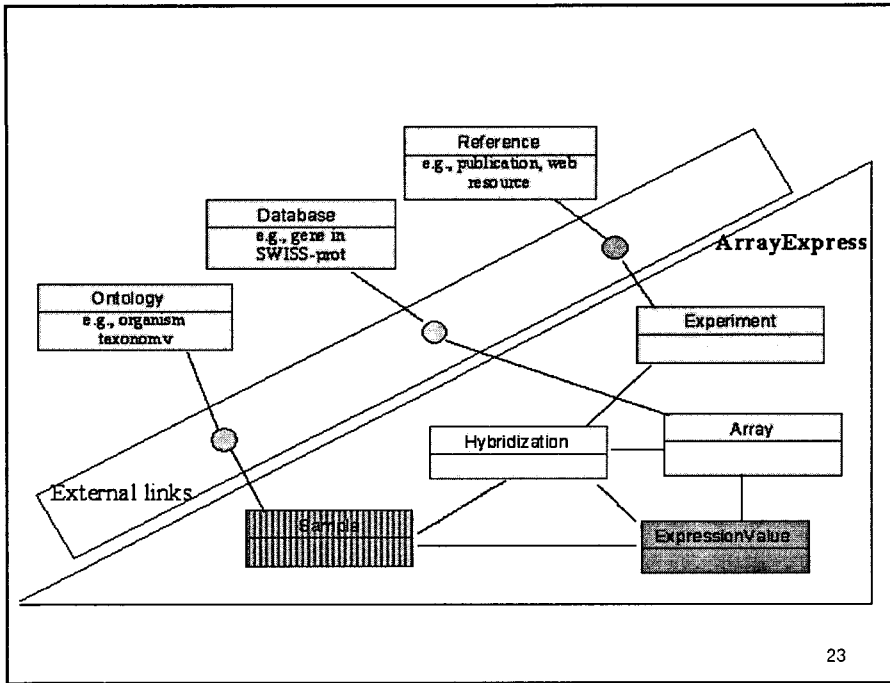
- *Object Model as a framework for developing MAGE*
- *OMG specifications are developed in UML*

21

Three levels of microarray gene expression data processing



22



Goals of MGED Efforts

- Elucidate information needed for experiments
 - MIAME
- Provide a means to share this information
 - MAGE
- Provide a common language for experiments
 - MGED Ontology
- Provide standard operating procedures for analysis
 - Data transformation

Standards will aid data integration in toxicogenomics efforts

HESI, NIEHS,
FDA, EBI,
academics
- new MIAME-
Tox

Minimum Information About a Microarray Experiment – MIAME for Toxicogenomics (MIAME/Tox)

DRAFT – Based on MIAME 1.1 (February 12, 2003)

Background: MIAME and MIAME/Tox

The MIAME/Tox document is based on the MIAME 1.1 document¹ produced by the MIAME (Minimum Information About a Microarray Experiment) Society². The goal of MIAME (minimum information about microarray experiments)³ is to outline the minimum information required to interpret unambiguously and potentially reproduce and verify an array-based gene expression microarray experiment. Although details for particular experiments may be different, MIAME aims to define the core that is common to most experiments. MIAME is not a formal specification, but a set of guidelines. More information about the MIAME rationale can be found in Minimum information about a microarray experiment (MIAME)—toward standards for microarray data, A. Brazma, et al., *Nature Genetics*, vol. 29, (September 2001), pp 335-337⁴. Although MIAME concentrates on information content and should not be confused with a data format, it also tries to provide a conceptual structure for microarray experiment descriptions. Similarly, MIAME/Tox seeks to provide such a conceptual structure in the context of toxicogenomics.

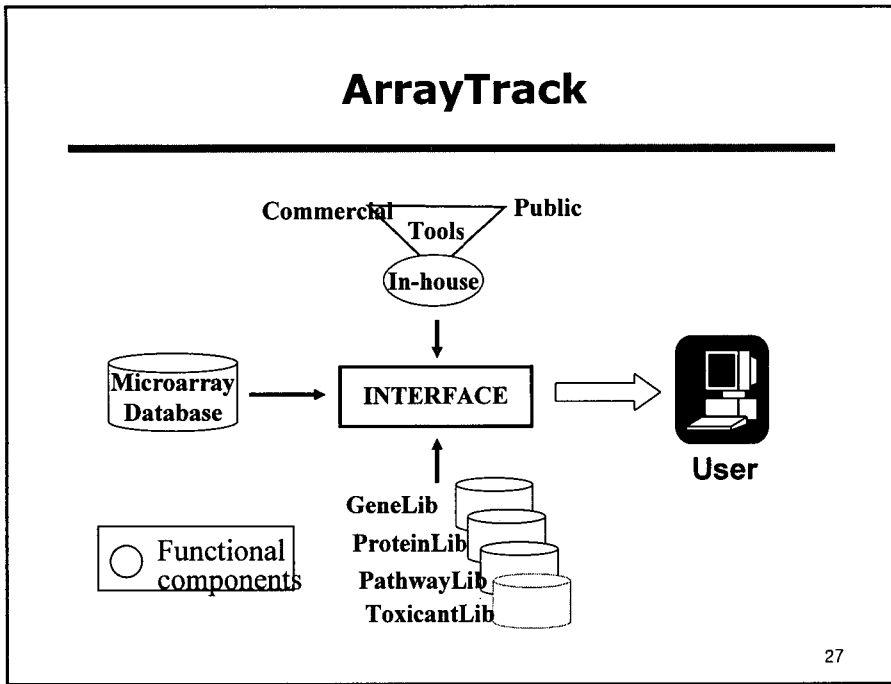
In addition to MIAME, a standard microarray data model and exchange format (MIAME/Tox) which is able to capture information specified by MIAME, has been suggested by EBI (the European Bioinformatics Institute) and recently became an Advisory Specification of the OBI (Ontology for Biological Investigation) (see <http://www.ncbi.nlm.nih.gov/obo/owl/MIAME/Tox.owl>). Many organizations (including EBI, Genentech, and others), have contributed ideas to MACE. MACE collectively refers to the MACE-XML (extensible markup language) format and MACE-ML (markup language) derived from the model. MACE-XML is able to capture information specified by MIAME and will be the standard microarray data model, while MACE-ML is the standard exchange format.

25

Minimum information to be recorded about toxicogenomics experiments

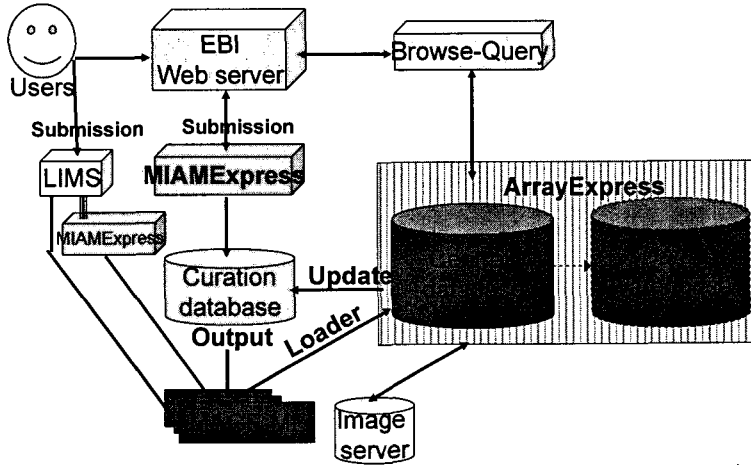
- Experimental design parameters, animal husbandry information or cell line and culture information, exposure parameters, dosing regimen, dose groups, and in-life observations.
- Microarray data, specifying the number and details of replicate array bioassays associated with particular samples, and including PCR transcript analysis if available.
- Numerical biological endpoint data, including necropsy weights or cell counts and doubling times, clinical chemistry and enzyme assays, hematology, urinalysis, other.
- Textual endpoint information such as gross observations, pathology and microscopy findings.

26



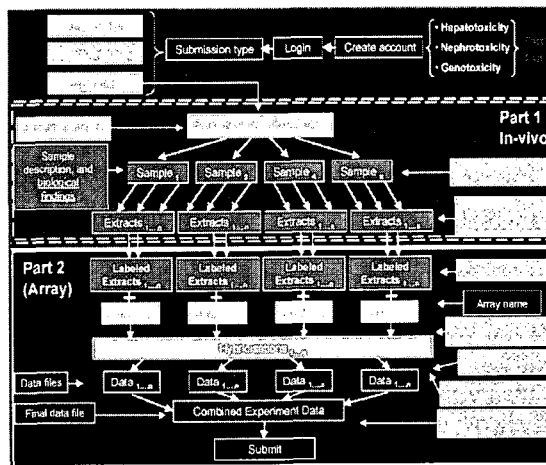
- ## dbZach
-
- Platform independent
 - Minimum Information About Microarray Experiment (MIAME) Supportive
 - Database (Oracle 9i)
 - Data Modeling
 - Modular Subsystem Architecture
 - Data Management
 - Tools
 - Graphical
 - User-friendly
 - Java 2
 - Data Mining
 - Data Visualization
- 28

ArrayExpress – data flow



29

Toxicogenomics specific MIAMExpress (ILSI- HESI) - in development



30

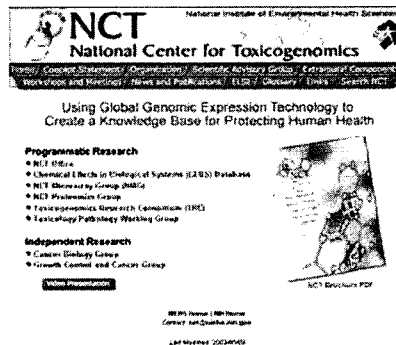
toxicogenomics project by ILSI

- International Life Sciences Institute (ILSI/HESI) toxicogenomics database
- cross-platform gene expression data on the effects of various toxic compounds

31

Government efforts: National Center for Toxicogenomics

- Goal
 - support research in the field of toxicogenomics
 - compile, analyze, and publish the resulting data
 - identify genes that are regulated in response to toxicans
 - develop a "Toxchip" to monitor changes



National Institute of Environmental Health Sciences
NCT
National Center for Toxicogenomics

History: Genetic Characterization, Organizational, Scientific Advisory Council, Administrative Components, Workshops and Meetings, News and Publications, ELSI, Outreach, Links, Search, etc.

Using Global Genomic Expression Technology to Create a Knowledge Base for Protecting Human Health

Programmatic Research

- NIEHS
- Chemical & Toxicologic Systems (CTS) Database
- NIEHS Research Group (NIEHS)
- NIEHS Publications Group
- Toxicogenomics Research (Genotoxicity) (TRG)
- Toxicology Pathology Working Group

Independent Research

- Cancer Biology Center
- Growth Control and Cancer Group

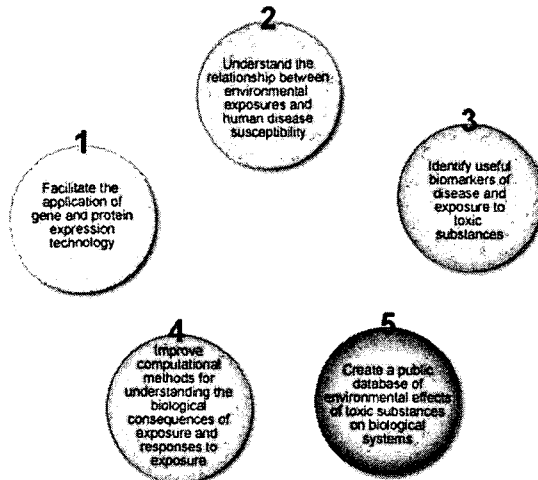
Other Resources

NCT Structure PDF

NIEHS Home | NIEHS Home
Contact: 404-242-3400, 404-242-3400
Last Modified: 2003/06/03

32

Goals of National Center for Toxicogenomics



33

Toxicogenomics Infrastructure

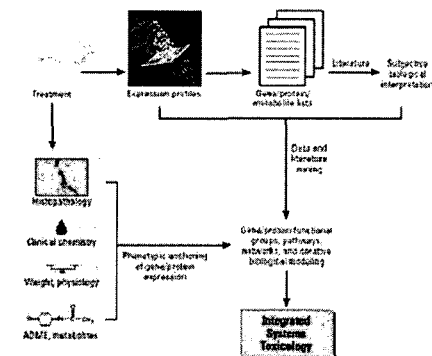
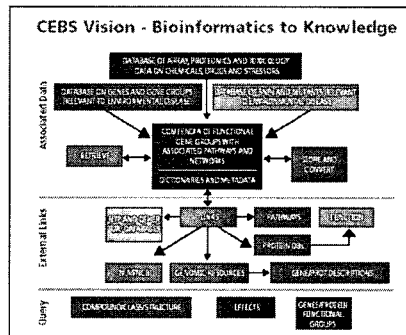


Figure 2. Integration of microarray expression profiles with literature mining, phenotypic anchoring, and iterative biological modeling for systems toxicology. ADME refers to absorption, distribution, metabolism, and excretion.

Figure from *EHP Toxicogenomics* (2003) 111:15

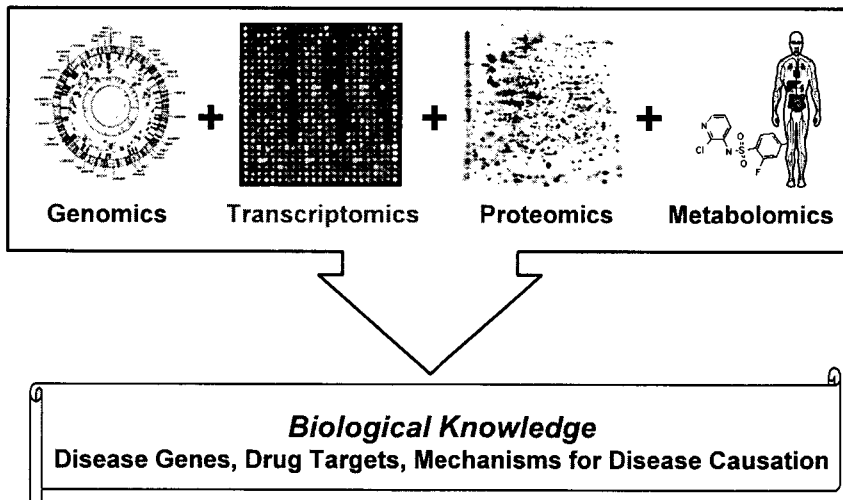
34

Chemical Effects in Biological Systems (CEBS)

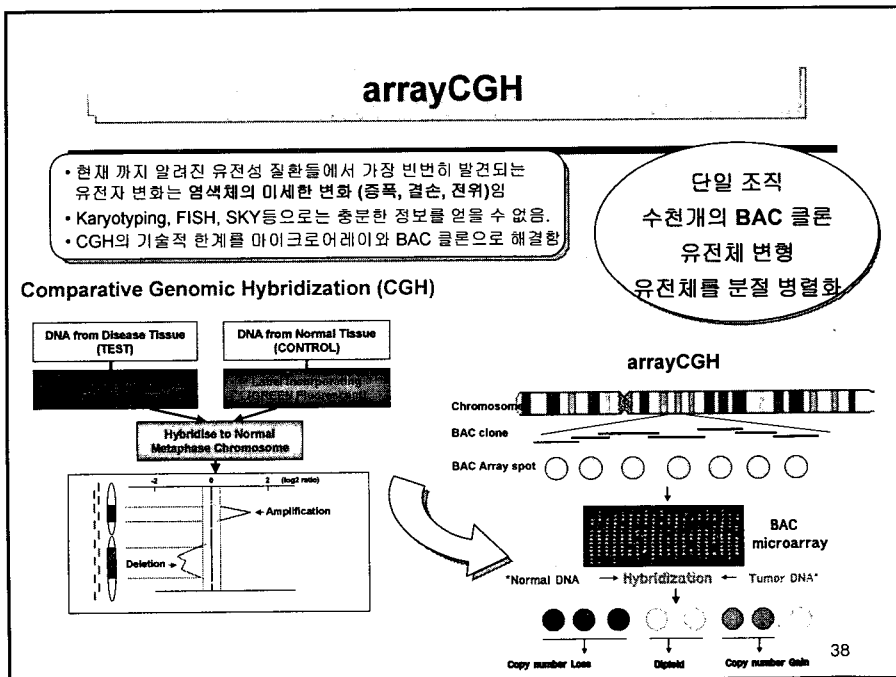
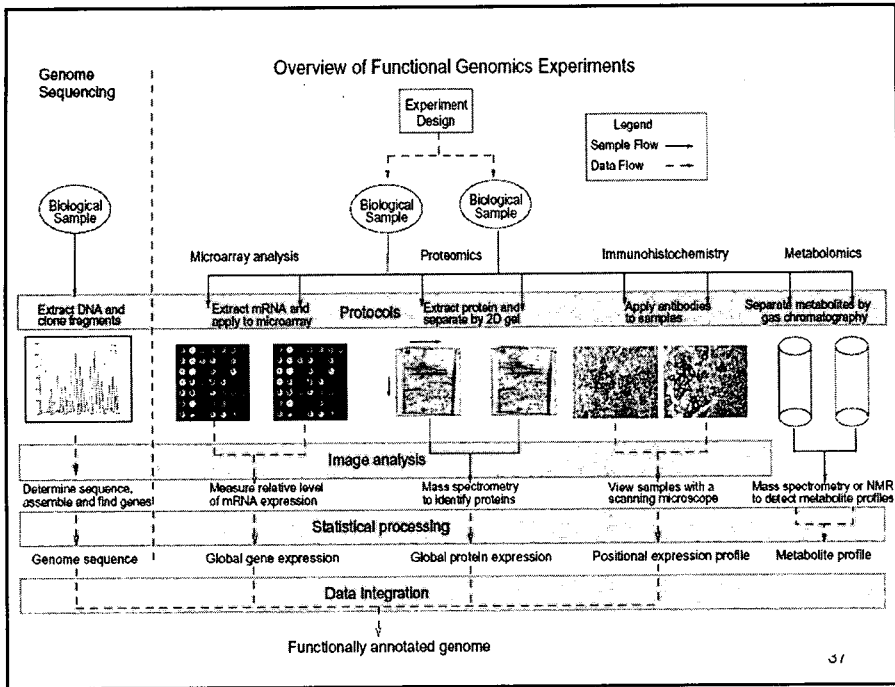


35

Proteomics in the "-omics" world.....



36



Tissue microarray

a

조직 표본에서 원기둥 모양의 Core 획득 (직경 0.6mm)

b

Recipient Block에 Core 이식 (수백회 반복)

c

검적성 필름으로 조직 마이크로어레이 제작 > 1,000 장 (두께 0.1mm)

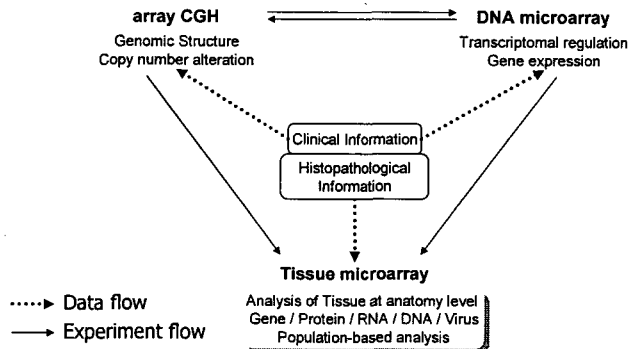
Target	Probe	방법
단백질	항체	면역조직화학법
DNA	DNA, fluorescent labeled	FISH
RNA	RNA, biotin labeled	in situ hybridization
Virus	primers	in situ PCR
DNA	primers	methylation specific in situ PCR

복수 조직
조직정보 포함
단일 검출 (다양한)
병리조직을 병렬화

Tissue microarray

4.0 mm 24 specimens	2.0 mm 60 specimens	1.0 mm 140 specimens	0.6 mm 240 specimens
<i>SuperBioChips</i>	<i>SuperBioChips</i>	<i>SuperBioChips</i>	<i>SuperBioChips</i>

Toxicogenomics research using high throughput technologies



41

Xperanto: *Expressionist's Esperanto in XML*

- MGED MIAME, MAGE-OM Standards
- Linux, MySQL, XML : open source
- Oligo and cDNA microarray
- Histopathology and Laboratory Data Model
- Batch upload and error check
- Diagnostic plots
- Visualization : genome browser, pathway browser)
- Automated gene / protein annotation
- Multi-user and group Environment

Park et al. Genomics and Informatics 2005

Xperanto: Expressionist's Esperanto in XML
 Now 2 Experiments

cell_type_comparison_design, MAGIC Oligo - Human 10K[M], 03-05-16

Now 2 Experiments

E-MANP-1 compound_treatment_design compound APDES-41MBL-1 03-03-07

HELa cells grown to subconfluent density and Calcein cells grown to high density were either treated with an iron source (control) or an iron chelator (desferal). Cells were harvested and RNA was prepared by RNeasy (Qiagen) and the RNA was subjected to a subsequent clean up using the RNeasy (Qiagen) clean up procedure. [ass](#)

Name	Labeled Extracts	Date	Files	Lab
035Mm_hybridization	S: heli2304mfraklabel S: F030205mfraklabel	03-00-00	[Icons]	X (MIDHE)
035Mm_hybridization	S: cec02295mfraklabel S: cec02295mfraklabel	03-00-00	[Icons]	X (MIDHE)
ass1		03-00-00	[Icons]	X (MIDHE)

Logout [Logout] [Logout] [Logout]

tst4 RNA_stability_design cell type MAGIC Oligo - Human 10K[M] 03-05-12

Name	Labeled Extracts	Date	Files	Lab
tst4-1	S: cec02235mfraklabel	03-00-00	[Icons]	X (MIDHE)
47	S: heli2304mfraklabel	01-02-13	[Icons]	X (MIDHE)
70	S: heli2304mfraklabel	03-05-16	[Icons]	X (MIDHE)
44	S: heli2304mfraklabel	03-05-19	[Icons]	X (MIDHE)
45-1	S: cec02235mfraklabel	03-05-14	[Icons]	X (MIDHE)
45-2	S: heli2304mfraklabel	03-00-00	[Icons]	X (MIDHE)
42-1	S: heli2304mfraklabel	03-00-00	[Icons]	X (MIDHE)
40-4	S: cec02235mfraklabel	03-00-00	[Icons]	X (MIDHE)
40-5	S: cec02235mfraklabel	03-00-00	[Icons]	X (MIDHE)
40-6	S: cec02235mfraklabel	03-00-00	[Icons]	X (MIDHE)

Logout [Logout] [Logout] [Logout]

Park et al. Genomics and Informatics 2005

Core module

Data management: log-in and main page

Microarray Data Input - Microsoft Internet Explorer

http://xperanto.studip.org/

experiment has 18 Experiments [add new exp.]

E-MANP-1 compound_treatment_design compound APDES-41MBL-1 03-03-07

MAGIC Oligo - Human 10K[M] MIDHE

Register: login
 Available Experiments: 58
 Available Experiments Results: 14754
 Total Quantification Files: 54
 Total Image Files: 9

Logout [Logout] [Logout] [Logout]

tst4 RNA_stability_design cell type MAGIC Oligo - Human 10K[M] 03-05-12

MAGIC Oligo - Human 10K[M] MIDHE

Register: login
 Available Experiments: 58
 Available Experiments Results: 15216
 Total Quantification Files: 16
 Total Image Files: 15

Logout [Logout] [Logout] [Logout]

Winstan MAGIC Oligo - Human 10K[M] 03-07-25

MAGIC Oligo - Human 10K[M] MIDHE

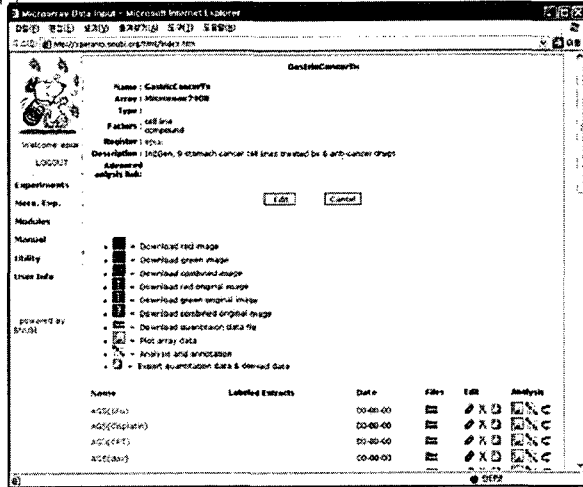
Register: login
 Available Experiments: 78
 Available Experiments Results: 27228
 Total Quantification Files: 20
 Total Image Files: 84

Logout [Logout] [Logout] [Logout]

Logout [Logout] [Logout] [Logout]

Core module

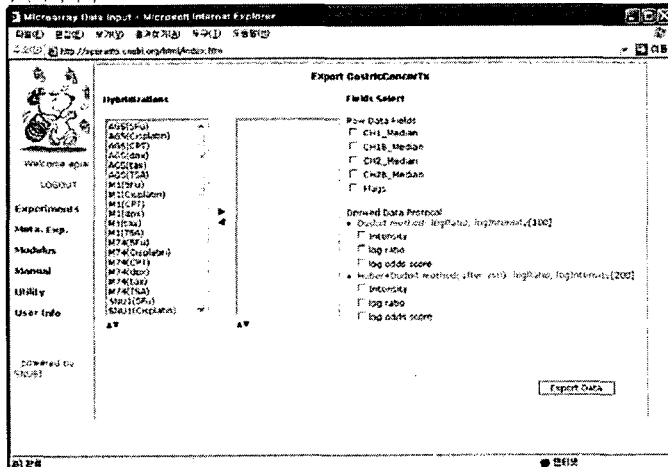
Data management: experiment



45

Core module

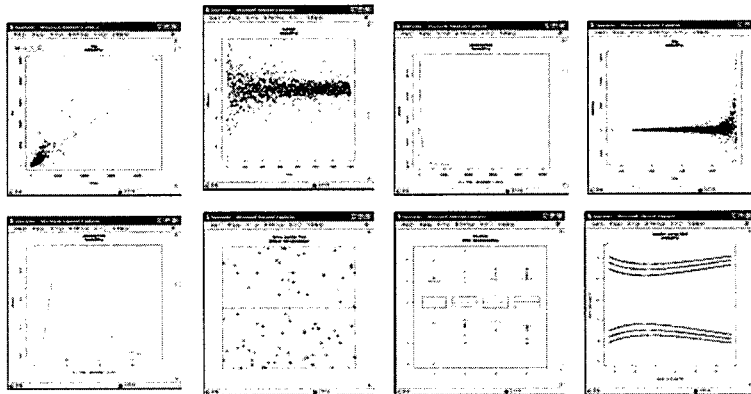
Data management: data export



46

Core module

Data Analysis: plots in normalization process



47

Web interface Submit Data

The screenshot shows a web interface for submitting data. It features a navigation bar at the top with various icons and a search bar. The main content area is divided into several sections:

- Equipment and Data:** A table with columns for Equipment, Quantity, and Date. The table contains several rows of data.
- Specimen and Data:** A table with columns for Specimen, Quantity, and Date. The table contains several rows of data.
- Clinical Information:** A section with various input fields and checkboxes for clinical data.
- Organ Specific Information:** A section with various input fields and checkboxes for organ-specific data.
- Additional Clinical Information:** A section with various input fields and checkboxes for additional clinical data.

The interface is designed for data entry and includes various validation and submission options.

48

CPCTR tissue array 1-2_1-1

Layout
quest Name: CPCTR tissue array 1-2_1-1 Date: num. of rows: 5, num. of cols: 5
 Case Diameter: 0.5 mm
 Block construction protocol: cpctr_block
 Array construction protocol: cpctr_arraying_protocol
 Experiment: CPCTR TMA 1-2 block with 12 substrates and 250 genes. This block contains the first substrate: prostate cancer tissue microarray with outcomes data created by Drs. Andre Sefla and Moun-Datta for the
 RibonArray: CPCTR prostate cancer, prostatic adenocarcinoma, adenocarcinoma of prostate, carcinoma of prostate, tma, tissue micro array, tissue micro-array, tissue microarray, tissue repository, tma, microarray cancer metastasis, marker, cancer test, validation, cpctr, cooperative prostate cancer tissue resource
 Reporter: Internal link: <http://www.prostatewiki.org>
 Person: Description: CPCTR prostate cancer, prostatic adenocarcinoma, adenocarcinoma of prostate, carcinoma of prostate, tma, tissue micro array, tissue micro-array, tissue microarray, tissue repository, tma, microarray cancer metastasis, marker, cancer test, validation, cpctr, cooperative prostate cancer tissue resource

Export (CSV) [XLSX]				
	add	Modify		
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

Specimen and Donor block information

Specimen ID: 3154147577 Diagnostic Description: Location: N32
Repository Institution Name: PC16 168
Tissue fixation Fixative Type: Protocol
Donor block ID: DR3 168
Additional URL

Clinical Information

Demography
 Patient ID: 3154149577 Age: 74 Sex: Male Race: Caucasian Responsible Physician:
 Description: Year_of_Birth: 1931
Diagnosis
 Diagnosis: Prostate Cancer Date: 1991-00-00 Responsible Physician:
 Description:
Operation
 Surgery: Prostatectomy Date: 1992-00-00 Responsible Physician:
 Description:
Molecular Analysis
 Type: Protocol: IRT Responsible Physician:
 Description:

Additional Clinical Information

Category	Attribute	Value	Description

Xperanto: Expressionist's Esperanto in XML

AGS.CPT..G.FMean Unnormalized

AGS.CPT..G.FMean Normalized

AG

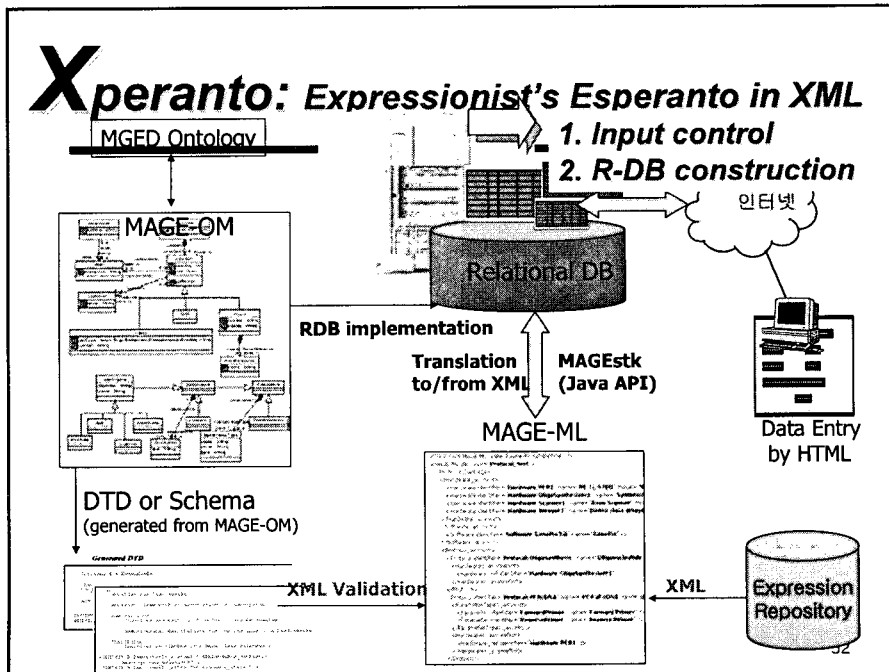
50

Xperanto: Expressionist's Esperanto in XML

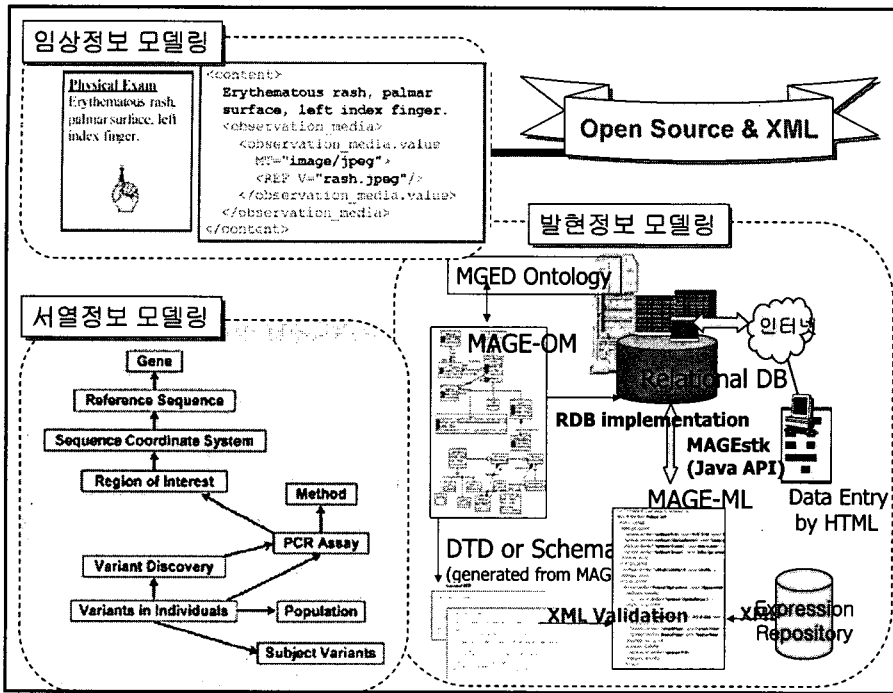
Gene Name	Accession	Expression Value
PCDH19	PM20510	1.0
PCDH19	PM20511	1.0
PCDH19	PM20512	1.0
PCDH19	PM20513	1.0
PCDH19	PM20514	1.0
PCDH19	PM20515	1.0
PCDH19	PM20516	1.0
PCDH19	PM20517	1.0
PCDH19	PM20518	1.0
PCDH19	PM20519	1.0
PCDH19	PM20520	1.0
PCDH19	PM20521	1.0
PCDH19	PM20522	1.0
PCDH19	PM20523	1.0
PCDH19	PM20524	1.0
PCDH19	PM20525	1.0
PCDH19	PM20526	1.0
PCDH19	PM20527	1.0
PCDH19	PM20528	1.0
PCDH19	PM20529	1.0
PCDH19	PM20530	1.0
PCDH19	PM20531	1.0
PCDH19	PM20532	1.0
PCDH19	PM20533	1.0
PCDH19	PM20534	1.0
PCDH19	PM20535	1.0
PCDH19	PM20536	1.0
PCDH19	PM20537	1.0
PCDH19	PM20538	1.0
PCDH19	PM20539	1.0
PCDH19	PM20540	1.0
PCDH19	PM20541	1.0
PCDH19	PM20542	1.0
PCDH19	PM20543	1.0
PCDH19	PM20544	1.0
PCDH19	PM20545	1.0
PCDH19	PM20546	1.0
PCDH19	PM20547	1.0
PCDH19	PM20548	1.0
PCDH19	PM20549	1.0
PCDH19	PM20550	1.0
PCDH19	PM20551	1.0
PCDH19	PM20552	1.0
PCDH19	PM20553	1.0
PCDH19	PM20554	1.0
PCDH19	PM20555	1.0
PCDH19	PM20556	1.0
PCDH19	PM20557	1.0
PCDH19	PM20558	1.0
PCDH19	PM20559	1.0
PCDH19	PM20560	1.0
PCDH19	PM20561	1.0
PCDH19	PM20562	1.0
PCDH19	PM20563	1.0
PCDH19	PM20564	1.0
PCDH19	PM20565	1.0
PCDH19	PM20566	1.0
PCDH19	PM20567	1.0
PCDH19	PM20568	1.0
PCDH19	PM20569	1.0
PCDH19	PM20570	1.0
PCDH19	PM20571	1.0
PCDH19	PM20572	1.0
PCDH19	PM20573	1.0
PCDH19	PM20574	1.0
PCDH19	PM20575	1.0
PCDH19	PM20576	1.0
PCDH19	PM20577	1.0
PCDH19	PM20578	1.0
PCDH19	PM20579	1.0
PCDH19	PM20580	1.0
PCDH19	PM20581	1.0
PCDH19	PM20582	1.0
PCDH19	PM20583	1.0
PCDH19	PM20584	1.0
PCDH19	PM20585	1.0
PCDH19	PM20586	1.0
PCDH19	PM20587	1.0
PCDH19	PM20588	1.0
PCDH19	PM20589	1.0
PCDH19	PM20590	1.0
PCDH19	PM20591	1.0
PCDH19	PM20592	1.0
PCDH19	PM20593	1.0
PCDH19	PM20594	1.0
PCDH19	PM20595	1.0
PCDH19	PM20596	1.0
PCDH19	PM20597	1.0
PCDH19	PM20598	1.0
PCDH19	PM20599	1.0
PCDH19	PM20600	1.0

51

Xperanto: Expressionist's Esperanto in XML



52



Knowledge base for Drugs

Genes - Drugs - Phenotypes

Phenotype Data Submission

Phenotype Data submission

Phenotype ID	Name	Description	Submit
1	Submit
2	Submit

Browse by Gene, Drug & Phenotype

KPRN knowledgebase currently has

- Genes: 17,037 HGNC Symbols
- Drugs: 4,653 Drug Names
- Diseases: 4,053 Disease Name

