

# **Comparative evaluation of microarray-based gene expression databases**

**Hong-Hai Do,  
Toralf Kirsten,  
Erhard Rahm**

**University of Leipzig, Germany  
*www.izbi.de, dbs.uni-leipzig.de***

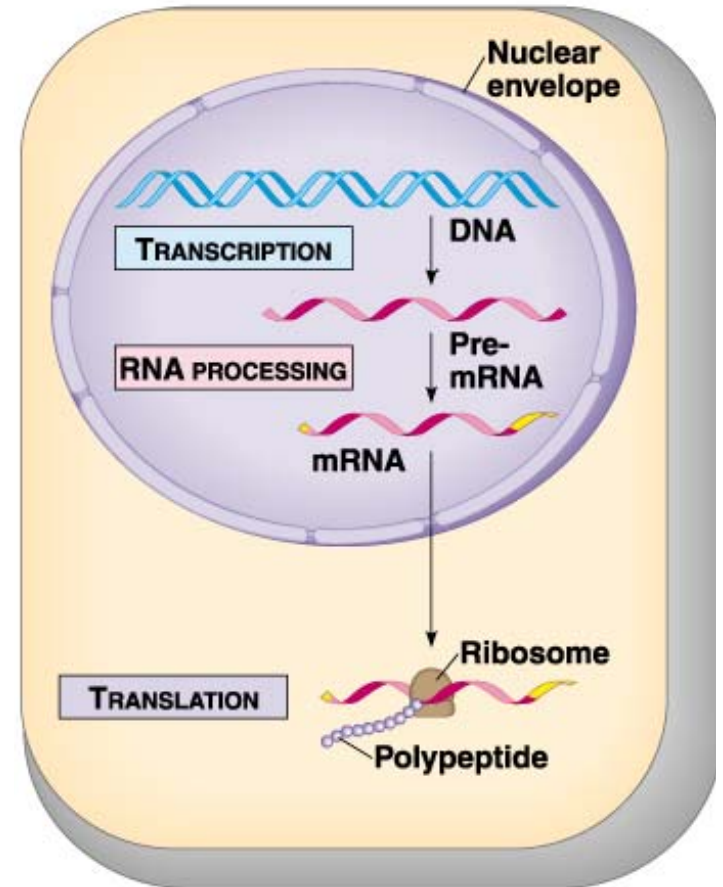
# Content

---

- Gene expression analysis
- Microarray database requirements
- Evaluation of 8 database solutions
- The GeWare project in Leipzig
- Conclusions

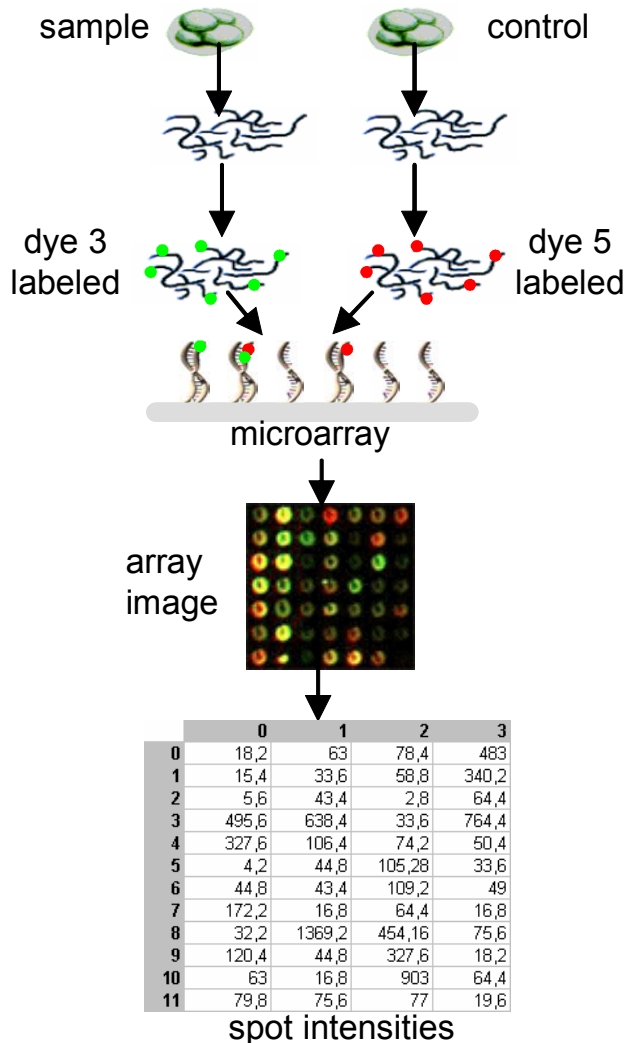
# Gene Expression Analysis

- Goal: Characterization of functions of genes and their mutual influence in the regulatory network
- Measuring mRNA amount in cells under different conditions
- Microarrays
  - Measuring expression of thousands of genes simultaneously
  - Large amounts of data with every experiment



# Microarray Experiment

## cDNA Arrays



## Oligonucleotide Arrays

(1) Cell selection

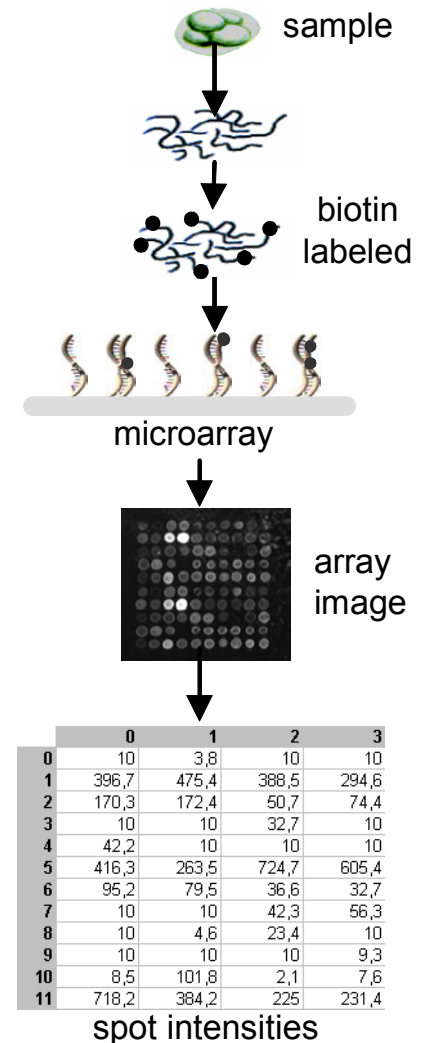
(2) RNA/DNA preparation

(3) Hybridization

(4) Array scan

(5) Image analysis

(6) Expression analysis



# Local Situation

---

- In Leipzig: ~15 different user groups:
  - Comparative primate genomics (human vs. chimpanzee)
  - Change detection in signal transduction in thyroid pathologies
  - Gene expression profiling of brain tumors
  - ...
- Affymetrix oligonucleotide microarrays
  - About 300-500 experiment series / year (trend ↗)
- Current data management and analysis:
  - Handling of flat files produced by Affymetrix software
  - Data analysis using Affymetrix tools, MS Excel
  - Manual search for annotations in public sources

# Database Requirements

---

- Storage of different types of data
- Data integration
- Annotation management
- Data normalization
- Data analysis
- Tool integration

# Data Characteristics

- Various kinds of data with different characteristics and requirements

Data		Source	Type	Characteristics	Usage
Image Data		Array scan	binary	large files	Generation of expression data
Expression Data		Image analysis	number	fast growing volume	Visualization, statistical and cluster analysis
Annotation Data	Gene	External public sources	text	regularly updated	Interpreting / Relating / Inferring gene functions
	Experiment	User input		user-specified, often free text	

# Annotation Integration

---

- Various public sources with gene annotations:
  - LocusLink and RefSeq: GO annotations, homology, organism, reference sequence
  - UniGene, GeneCards, GeneLynx, Tigr, ...
  - Vendor-specific sources: NetAffx
- However, often different gene identifiers !!!
- Manual specification of experiment annotations
  - Free text to be limited/avoided for better analysis support
- Standard support necessary, e.g., MIAME, MAGE-ML, GeneOntology, ...



# Data Integration Mechanisms

---

- Virtual integration
  - Web linkage based on accession keys
    - Navigational access
    - Annotation data not queryable
    - Little integration effort
  - Federated systems (Mediator-based)
    - Schema integration
    - On-the-fly data integration: transformation, cleaning, merging
    - Performance/Availability/Rudimentary query capabilities of public sources
- Materialized integration (Data warehousing)
  - All relevant annotation data + expression data locally stored
  - Advantages for data analysis: all data directly queryable, performance
  - High integration and update effort
- Hybrid approaches, e.g. SRS

# Management of Annotation Data

- Flexible management required
  - Coping with attribute changes / fast-evolving schemas and vocabularies
- Database representation:
  - Relational vs. EAV (Entity-Attribute-Value)

## Relational Modeling

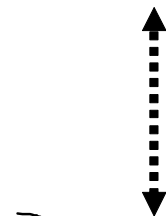
Experiment Annotation			
Experiment	Total RNA Amount	Stimulation Dosis	Time of Stimulation
1	230.46	12.5	1.02
2	225.75	25.0	10.01

## EAV Modeling

Annotation		
Item	Parent Item	Item Name
0	-	Experiment Annotation
1	0	Total RNA Amount
2	0	Stimulation Dosis
3	0	Time of Stimulation

Annotation Values		
Experiment	Item	Value
1	1	230.46
1	2	12.5
1	3	1.02
2	1	225.75
2	2	25.0
2	3	10.01

metadata as  
data instances



metadata

data instances

# Data Normalization

---

- Necessary for expression data due to
  - Fluctuations in technical experiment process
  - Comparison between multiple experiments
- Normalization for 1 experiment
  - Division by average intensity of all spots on array
  - Of control genes (housekeeping or spiked genes)
- Normalization for multiple experiments (series)
  - Normalization against a control experiment
- Storage of raw data for re-normalization

# Data Analysis

---

- Navigation/Querying/Reporting
- Online analytical processing
  - Multidimensionality of expression data
- Statistics and data mining
  - Descriptive statistics: mean, standard deviation, ...
  - Probability calculation: distributions, regression, correlation, ...
  - Inductive statistics: random sampling, estimation, tests, ...
  - Clustering: Hierarchical, K-mean, Self-Organizing Maps, ...
  - Classification: Support Vector Machines, Decision trees, ...
- Visualization
  - Display of statistical and clustering results
  - Scatter plots, dendrograms, charts, graphs, ...

# Tool Integration

---

## ■ File exchange

- Export from database, import in tool for analysis (tab-delimited ASCII format, XML etc.)
- No integration effort, but restricted / static information

## ■ API access to DBS by tools

- Use of DBS is transparent to user
- Access to current data using query language

## ■ Tight integration: Direct analysis in database systems

- Analysis / data mining approaches implemented by DBMS or as stored procedures
- Potential for high performance
- High implementation effort

# Evaluation of 8 database solutions

<b>Database</b>	<b>Organization</b>
<b>ArrayDB</b>	National Human Genome Research Institute – NHGRI <i><a href="http://genome.nhgri.nih.gov/arraydb">http://genome.nhgri.nih.gov/arraydb</a></i>
<b>ExpressDB</b>	Harvard University <i><a href="http://arep.med.harvard.edu/ExpressDB">http://arep.med.harvard.edu/ExpressDB</a></i>
<b>GeneX</b>	National Center for Genome Resources – NCGR <i><a href="http://genebox.ncgr.org/genex">http://genebox.ncgr.org/genex</a></i>
<b>GIMS</b>	University of Manchester <i><a href="http://www.cs.man.ac.uk/~norm/gims">http://www.cs.man.ac.uk/~norm/gims</a></i>
<b>M-CHIPS</b>	German Cancer Research Center <i><a href="http://www.mchips.de">http://www.mchips.de</a></i>
<b>RAD2</b>	University of Pennsylvania <i><a href="http://www.cbil.upenn.edu/RAD2">http://www.cbil.upenn.edu/RAD2</a></i>
<b>SMD</b>	Stanford University <i><a href="http://genome-www4.stanford.edu/MicroArray/SMD">http://genome-www4.stanford.edu/MicroArray/SMD</a></i>
<b>YMD</b>	Yale University <i><a href="http://info.med.yale.edu/microarray">http://info.med.yale.edu/microarray</a></i>

# Results: Data Management

---

- Supported types of data
  - Often no images stored
  - Expression data from different techniques (microarray-based and non-microarray)
- Gene annotations
  - not locally integrated/available in most cases
- Experiment annotations
  - Different content and varying degree of detail between the databases
  - Mostly free-text fields, no controlled vocabularies
- Data exchange
  - Tab-delimited used in many cases
  - XML not yet supported

# Results: Data and Analysis Integration

---

## ■ Data integration

- Web-link integration in most cases, but not sufficient for analysis
- Federated and materialized integration not yet fully exploited

## ■ Data analysis

- Canned queries widely used
- OLAP not yet applied despite multidimensionality
- Large variety of data mining approaches

## ■ Tool integration

- Advanced analysis mostly outside of database by means of stand-alone tools



# Project GeWare

---

- Specific local requirements
- Central data management and analysis platform for local users
- Data Warehouse approach
  - Data import from Affymetrix system
  - Fact tables to store both raw and derived data
  - Uniform specification of experiment annotations
  - Integration of gene annotations from public sources
  - Integration of analysis and data mining algorithms/tools

# System Architecture

## Source systems

### Experimental data

- Raw chip intensities
- Expression matrix

### Experiment annotations

- experiment, sample, ...
- MIAME

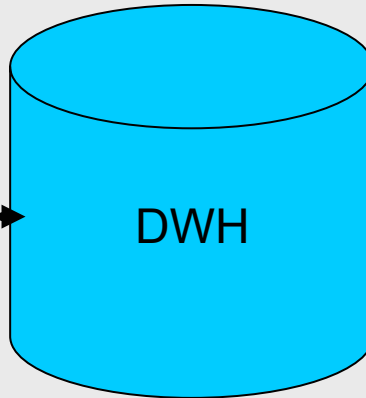
### External annotations

- Netaffx data
- Gene ontology (GO)
- LocusLink

## Data warehouse

### Core data warehouse

- multidimensional data model (star schema)



## Analysis

### Loose integration

- Export
- Download

### Transparent integration

- Use of API's
- Insightful ArrayAnalyzer
- OLAP Tools

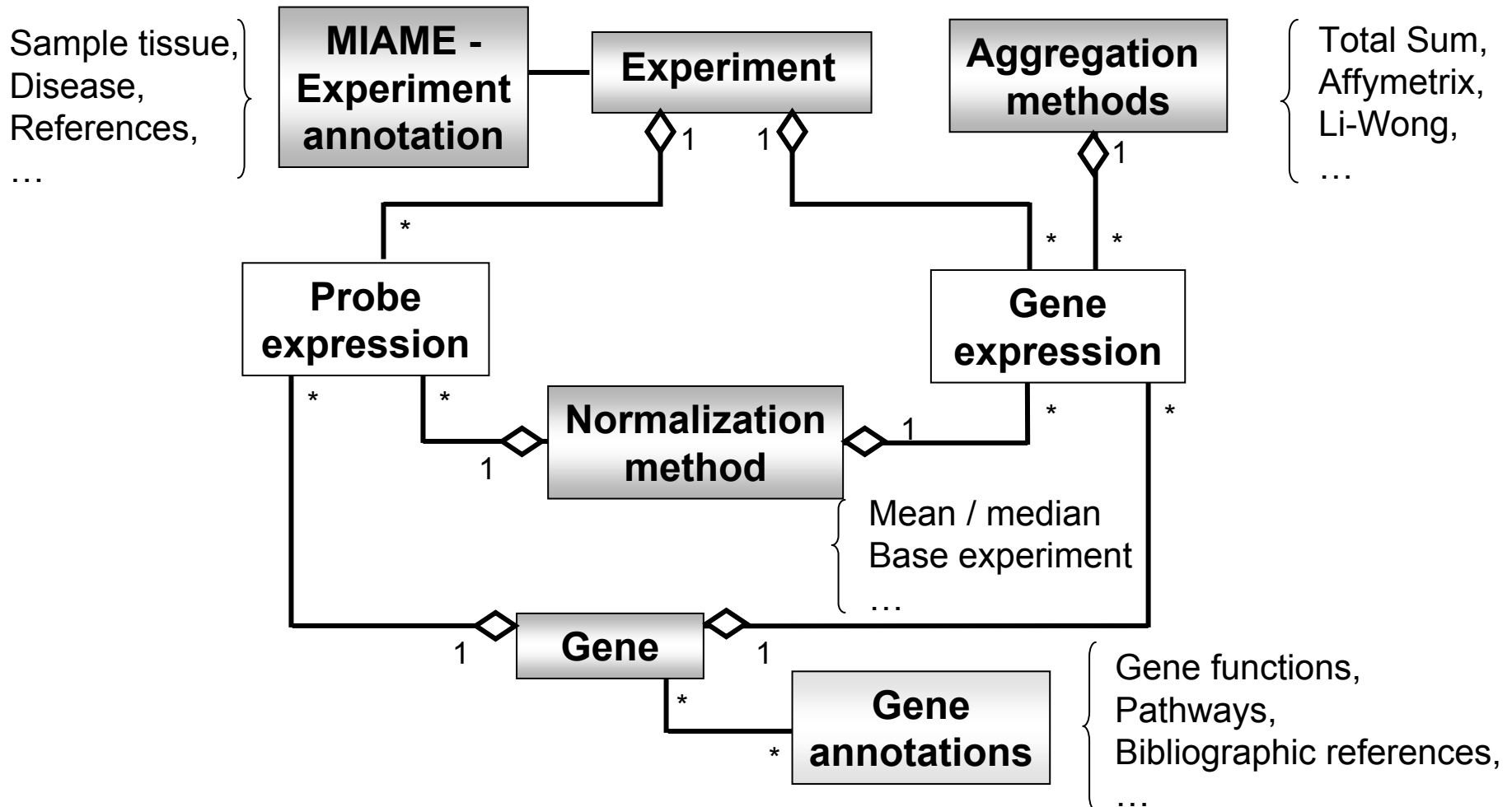
### Tight integration

- Special UDF's
- DB procedures

uniform web-based interface

# Data Warehouse Model

## ■ Multidimensional data model (star schema)



# Conclusions

---

- Microarray-based gene expression analysis
  - Promising technique for a variety of biological problems
  - High requirements for data management
- State of the art: insufficient database integration of
  - Annotations
  - Analysis approaches
- GeWare