

ROSETTA BIOSOFTWARE

www.rosettabio.com

Computational Biology: An Industrial Perspective

Lee Weng, Ph.D.
ACM Symposium on Applied Computing
March 10, 2003

Overview

Molecular Profiling

Gene Expression and Microarray Technology

Computational Challenges

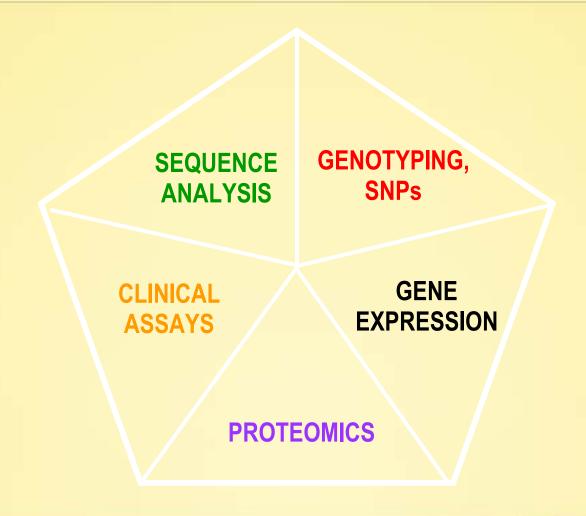
The Rosetta Resolver System

Beyond Textbooks

Data-Driven Knowledge Discoveries

Bring Data back to Biology

Molecular Profiling - A System Biology Approach



Industrialization of Molecular Biology

New high-throughput data-acquisition technologies have fundamentally changed the way to study biology

- » Genome-wide DNA sequencing
- » Microarrays for gene expression
- » LC/MS-based protein expression analysis
- » Large-scale molecular profiling becomes possible

Numerical computing is the corner stone of molecular profiling

- » More powerful data storage, retrieval and management tools
- » New computing methods and environments for data analysis and knowledge discovery
- » New standards for information exchanges

Knowledge discovery ≠ Data integration



Opportunities in Molecular Profiling

Study Biology as a System

- » Genes are not isolated. Groups of genes work together to serve a biological function. A stimulus will affect many genes.
- » Gene expression profiles provide insights to the function of the biological system.

Study Biology as a Collaborative Team Effort

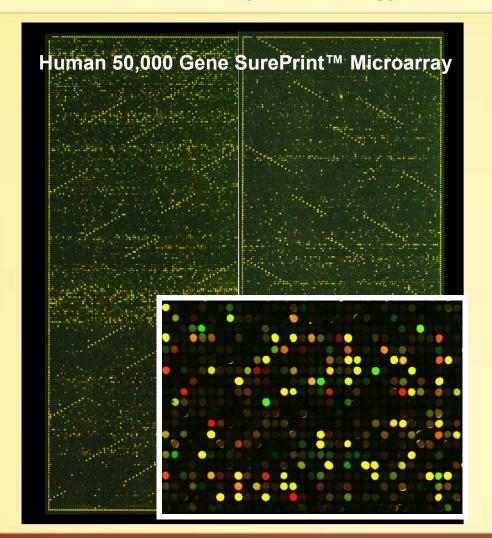
- » It has become less and less practical to study a complex biological system by one or two biologists.
- » Systematic data collection and information sharing have become inevitable.

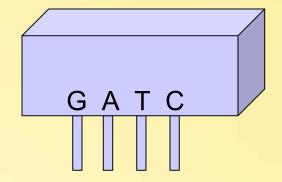
Study Biology as an Industry

- » Drug discovery becomes more difficult and costly. It is an enterprise-level effort.
- » Drug development requires enterprise solutions in data management, data analysis, and knowledge discovery.
- » These solutions should meet various regulatory requirements, such as FDA 21 CFR Part 11.



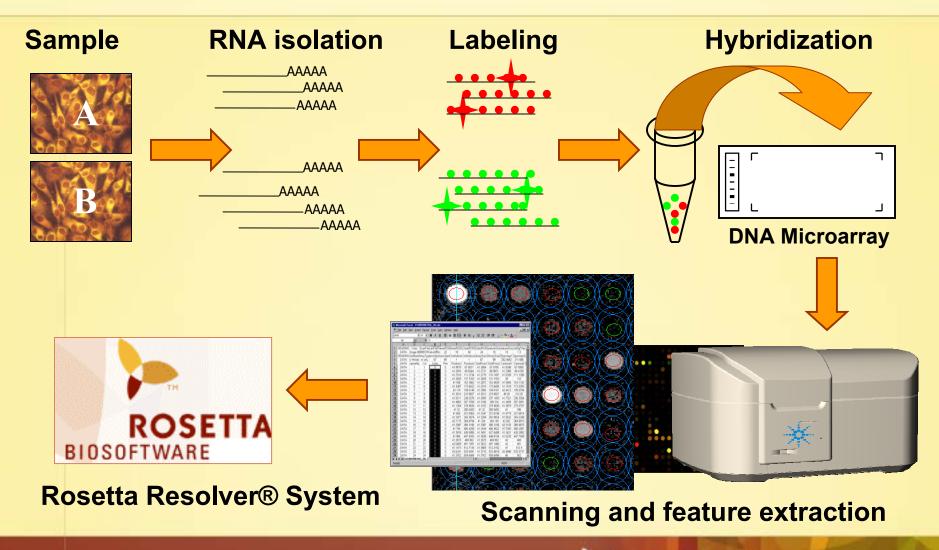
Example: Microarray Technology for Gene Expression Study



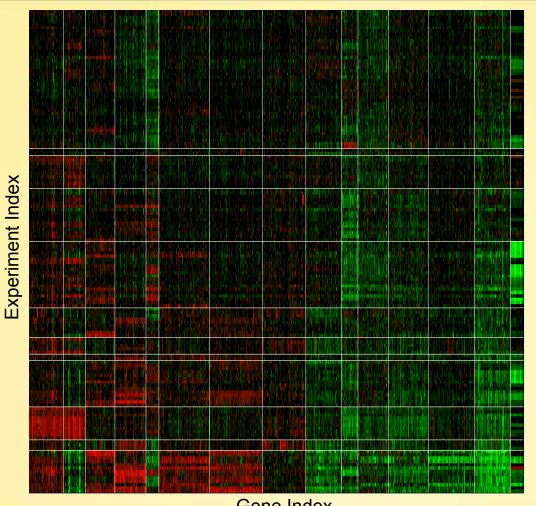


G A G T C

Processes in Analyzing Microarray Gene-Expression Data



Example: Compendium Approach to Hepatotoxicity



50+ Compounds act as Reference Standards

146 experiments ~1,800 genes

Down-regulated gene **Up-regulated gene**

Gene Index

Challenges in Molecular Profiling

Large Volume of Data

- » Microarray technology generates terabytes of gene expression data.
- » How to effectively manage the data? How to make the data available and the information searchable in an enterprise environment? How to efficiently leverage available analysis tools on the large data source?

Small Number of Replications

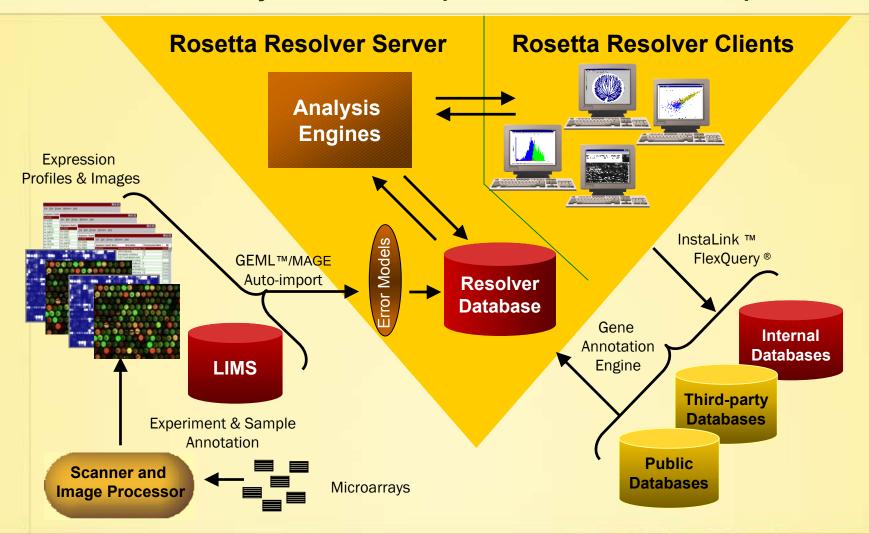
» Limited by available materials and high experiment costs, the number of replicates in molecular profiling is always very small, if any. Traditional statistical methods do not work well in this case.

Limited Knowledge

- » Need to associate profiles of cellular constituents to help interpret the biological functions.
- » Need to leverage existing biological knowledge in discovering new knowledge from the data.



Rosetta Resolver® System: An Enterprise Solution for Gene Expressions



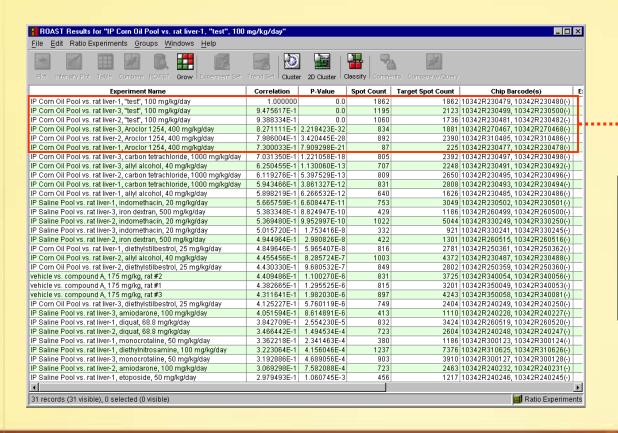
Rosetta Resolver® System v3.2

Enterprise solution to:

- » Compile gene expression information from a wide variety of technologies in a central repository.
- » Share information and analysis results collaboratively throughout the organization.
- » Perform large-scale intensity- and ratio-based analyses that leverage entire gene expression databases.
- » Publish and exchange data with collaborators in GEML™ and MAGE formats for visualization and analysis in the Rosetta Resolver system and other software applications.

Example: Rosetta Array Search Tool - ROAST

ROAST is analogous to BLAST for expression profile similarity searching. Use ROAST to search for co-regulated sequences, reporters, exons or UniGenes, as well as similar experiments.



Correlation between gene expression experiments may indicate similar mechanisms of toxicity between known and unknown compounds.

A Million-Dollar Question: Which Genes are Differentially Expressed?

Limitations of Fold-Change Method

- » Biologists often use 2-fold change as the threshold to detect differential expressions.
- » It suffers in low sensitivity and low specificity because it does not consider the error of the measurement

Limitations of Traditional Statistical Tests

» Textbook t-test or ANOVA test do not work well when the number of replicated experiments is small.

Solutions: Rosetta Resolver Error Models

Beyond Textbooks: Error Models in the Rosetta Resolver® System

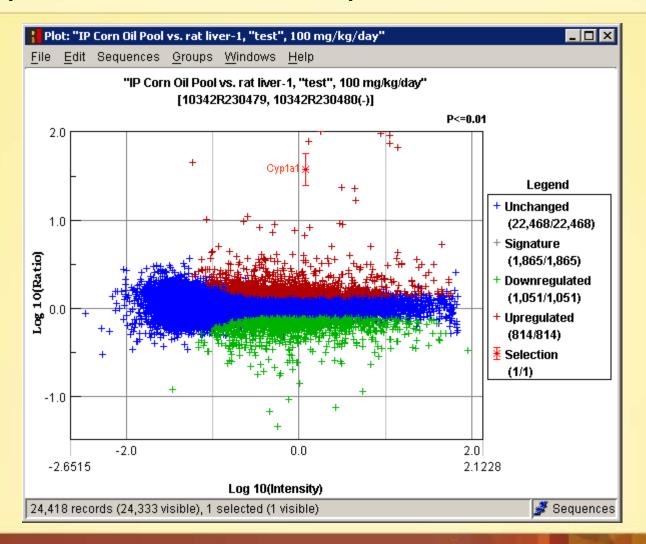
Error Models

- » Technology-specific error models are used to estimate microarray measurement errors.
- » The Rosetta Resolver system leverages the error models together with experimental replicates to compute accurate P-values and error bars for every gene expression measurement.

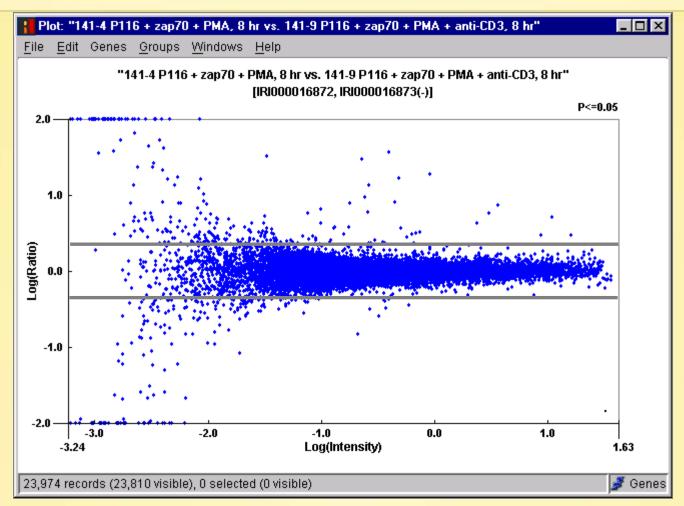
Benefits

- » Error models help to reliably estimate variance when the number of replicates is small.
- » Error models provide continued quality control in data analysis.
- » P-values and error bars are propagated and leveraged throughout the analysis, adding additional predictive power to cluster analysis, similarity searching, trend analysis, etc.

Example: Error Model Estimates Expression Measurement Error

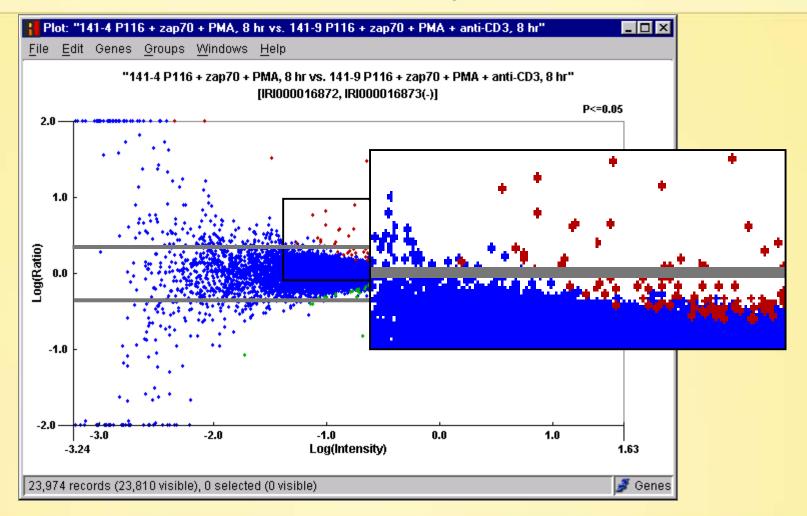


Example: Error Model Benefits in Ratio Analysis



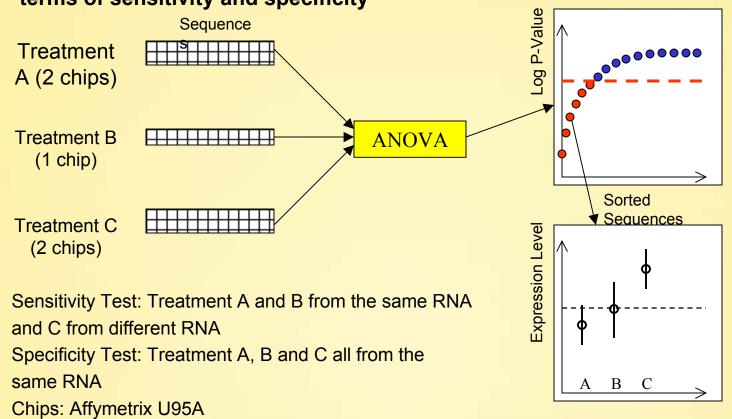
Which up/downregulations are statistically significant?

Example: Error Model Benefits in Ratio Analysis



Example: Error Model Benefits in ANOVA Tests

Compare the power of the textbook ANOVA and the improved ANOVA in terms of sensitivity and specificity



Example: Error Model Benefits in ANOVA Tests

For given threshold P-value<0.01, the improved ANOVA in the Rosetta Resolver System provides much lower false positive rate (better specificity) and much higher detection rate (better sensitivity) than the textbook ANOVA.

Analysis methods	False Positive rate	Detection rate
Textbook ANOVA	0.0093	0.17
Improved ANOVA	0.00048	0.30

Data-Driven Knowledge Discoveries

Unsupervised Learning Methods

- » Discover new knowledge by data associations based on similarities.
- » Clustering tools in the Rosetta Resolver system: agglomerative, divisive, k-mean, SOM, pattern grow, et al.

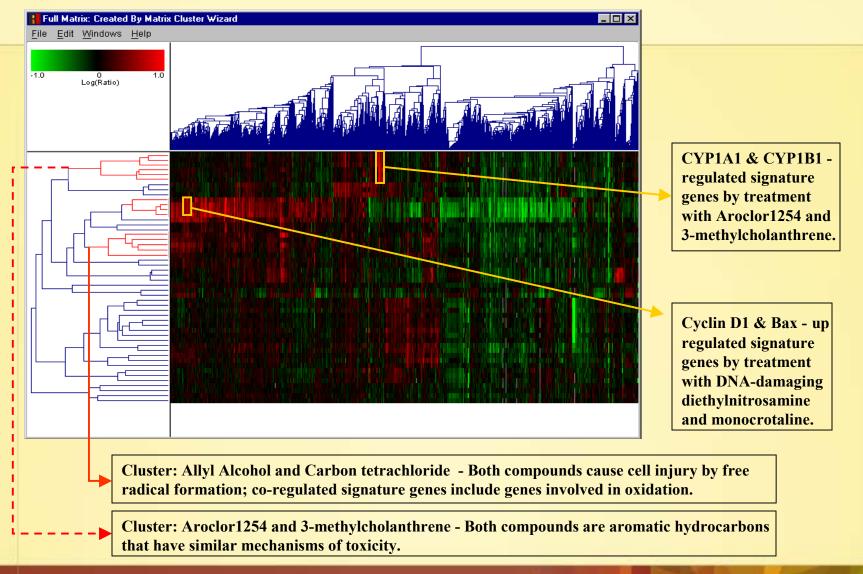
Supervised Learning Methods

- » Discover new knowledge by training.
- » Classification tools in the Rosetta Resolver system: Bayesian classifier, et al.

Example: An Unsupervised Learning Case Study

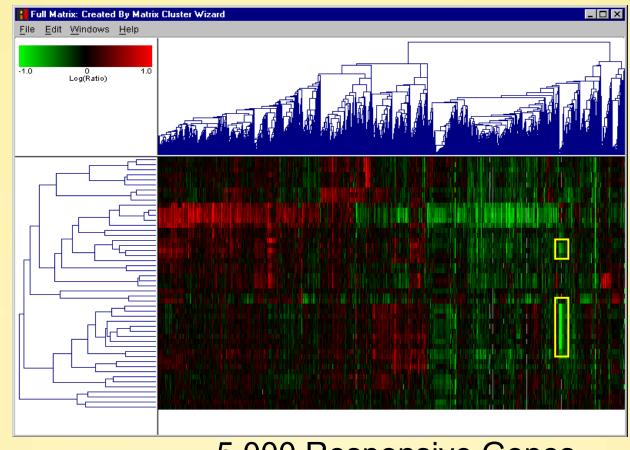
» Rosetta and Abbott Laboratories collaboration: Toxicogenomics Waring, et al. *Toxicology & Applied Pharmacology* 175, 28-42 (2001)

Toxicogenomics Signatures



Global vs. Local Similarities

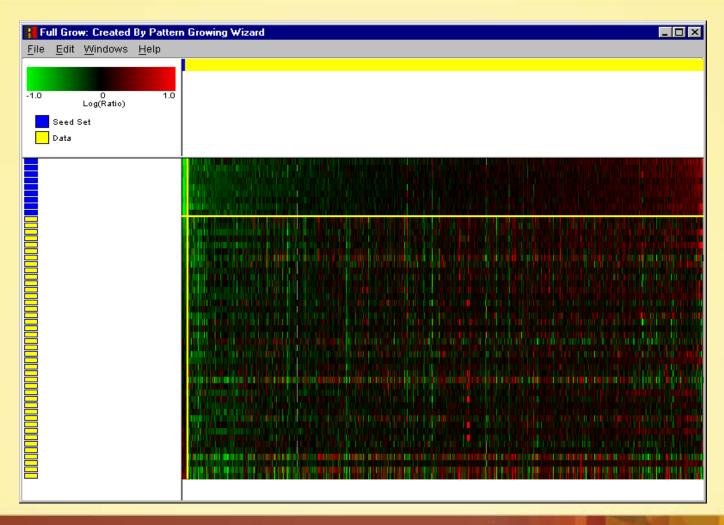
These experiments share a certain geneset response



5,000 Responsive Genes

47 Experiments

Global vs. Local Similarities: The GROW Algorithm



Example: A Supervised Learning Case Study

» Rosetta and Netherlands Cancer Institute (NKI) collaboration: Breast Cancer

van 't Veer, et al., *Nature* 415, 530-536 (2002)

Can Gene Expression Profiling be Used to Predict Clinical Outcome in Breast Cancer Patients?

Profile RNA samples from 98 breast tumors

- » 44 with >5 years disease-free survival
- » 34 with <5 years disease-free survival</p>
- » 20 with BRCA1 germ-line mutations

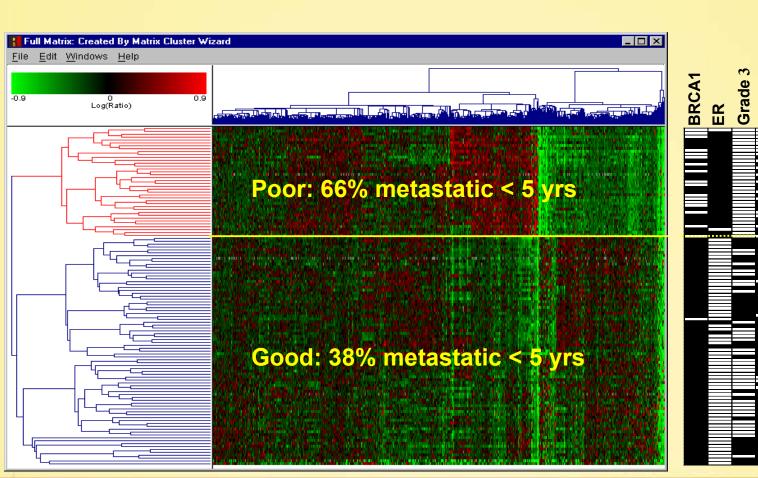
All samples carefully selected:

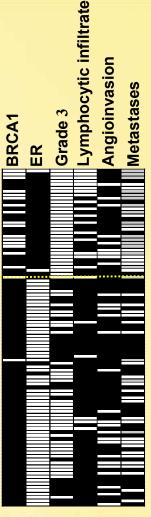
- » Patient age < 55 years</p>
- » Tumor size < 5 cm; no lymph node involvement</p>
- » Samples well annotated

Objectives:

- » Identify patterns of gene expression that correlate with prognosis
- » Confirm results with a new set of 19 similar breast tumors

Unsupervised Clustering Divides Tumors into "Good Prognosis" and "Poor Prognosis" Tumors



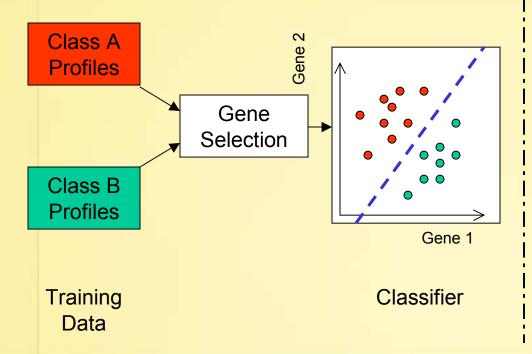


Classifier Concept

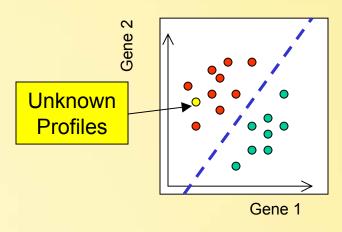
"Is this unknown profile most likely a member of the metastasis class (A) or the

non-metastasis class (B)?"

Classifier Training and Testing

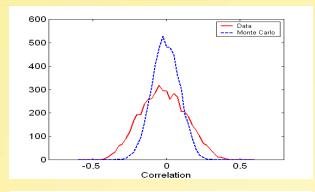


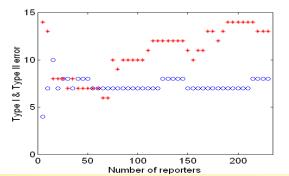
Classification/Prediction

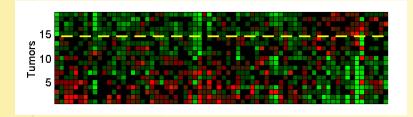


Classifier

Steps for Developing a Classifier







» Select and rank features

Define reporter genes that predict distant metastases

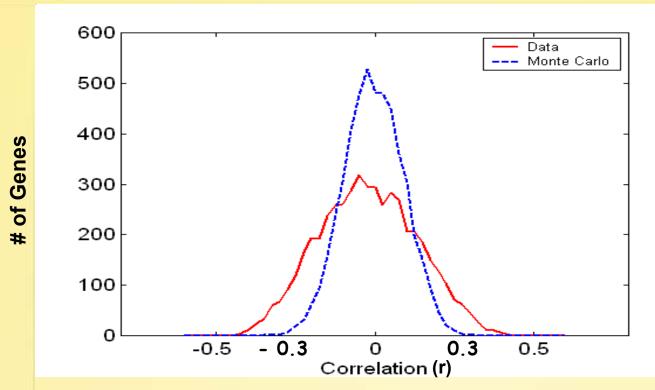
» Optimize the number of reporter genes

Leave-one-out cross-validation for defining optimal Classifier

» Evaluate the power of classifier

Independent validation set of 19 tumors

Define Reporter Genes that Predict Distant Metastases



231 genes with prognostic category red curve to r <- 0.3 & red curve to r > 0.3

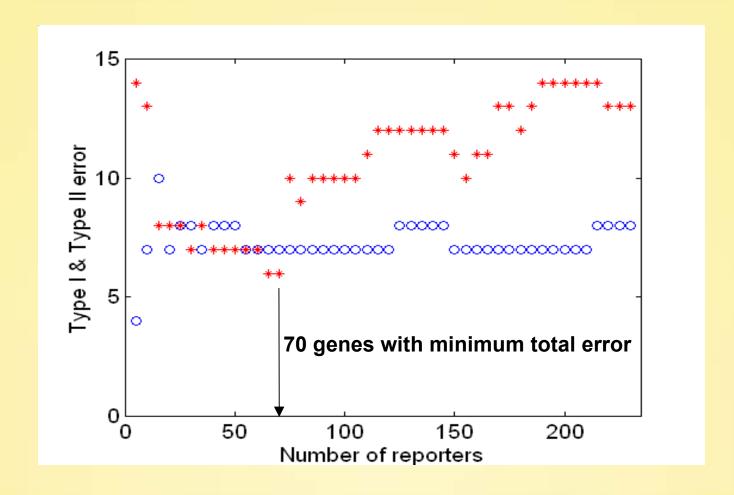
Red curve:

Histogram of correlation coefficients of genes with prognostic category (metastasis group versus nometastasis group)

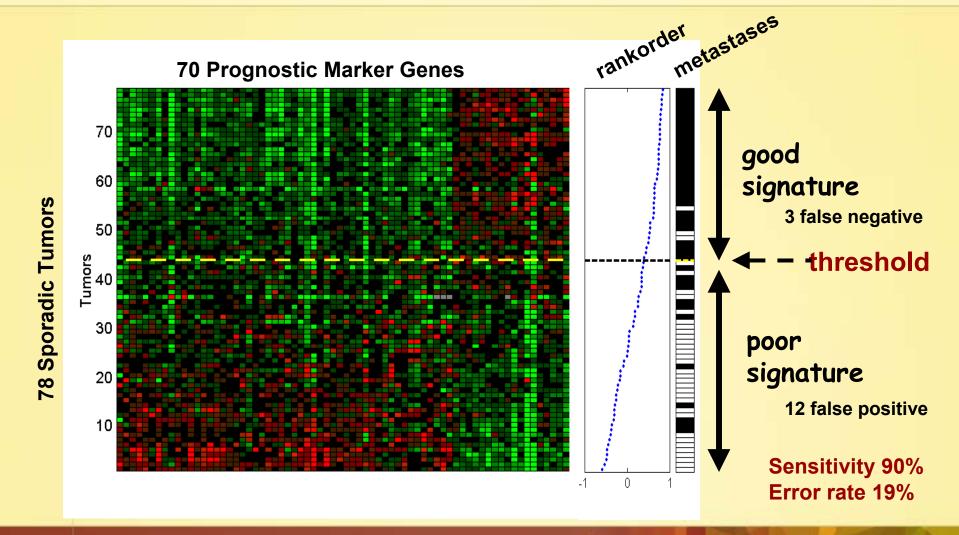
Blue curve:

Histogram of correlation coefficients of genes with prognosis from Monte Carlo analysis, where the metastasis designations have been randomized

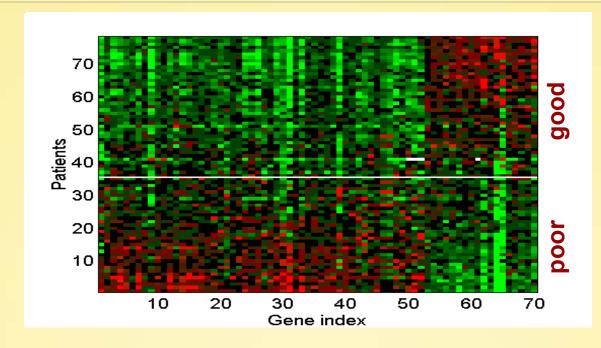
Optimize the Number of Reporter Genes Using Leave-One-Out Cross Validations



Supervised Classification Prognosis: Leave-One-Out Method for Defining Optimal Classifier



Functional Classes of Prognosis Reporter Genes



•Cell cycle: Cyclins, DNA replication

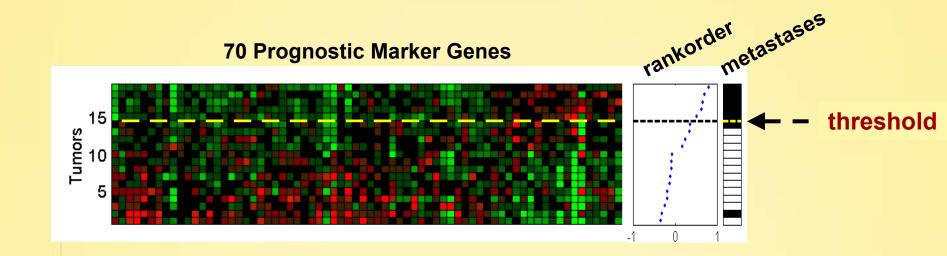
•Invasion, metastasis: (Metallo)proteinases

•Angiogenesis:

•Signal transduction:

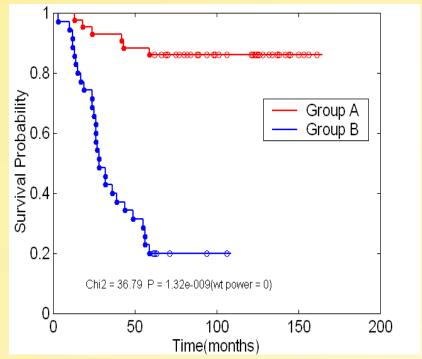
VEGF Receptor IGFBPs, Kinases

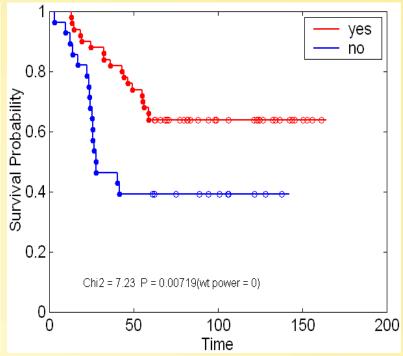
Classification of Set of 19 Tumors Using Prognosis Reporter Genes



only 2/19 tumors were misclassified: 90% accuracy, Fisher's exact p=0.0018

Microarray Reporters Predict Prognosis Better than Clinical Markers





Patients grouped by microarray-defined prognosis reporters

Patients grouped by Estrogen Receptor status

Profiling in Clinical Practice: Selection for Adjuvant Chemotherapy

Selection criteria:	Total patient group (n=78)	Metastatic disease < 5 years (n=34)	Disease free at >5 years (n=44)
St. Gallen consensus	64/78 (82%)	33/34 (97%)	31/44 (70%)
NIH-consensus	72/78 (92%)	32/34 (94%)	40/44 (91%)
Prognosis profile	43/78 (55%)	31/34 (91%)	12-18/44
			(27%-41%)

Breast cancer patients eligible for adjuvant systemic therapy



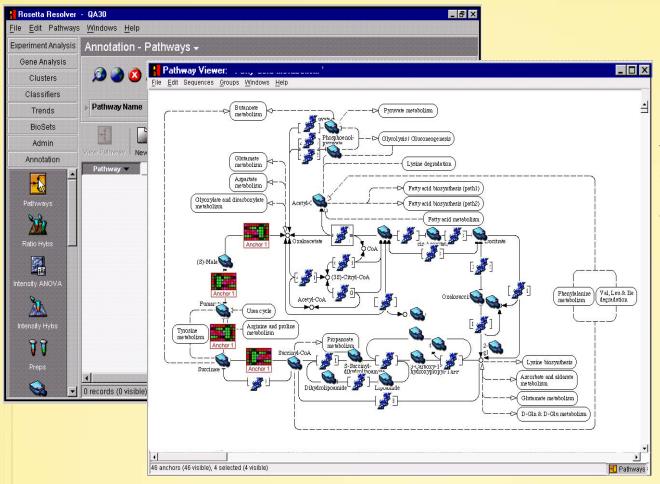
Profiling in Clinical Practice: Conclusions

Even small tumors are already programmed to metastatic phenotype.

Expression profiling:

- » selects patients that should receive adjuvant therapy similarly to conventional criteria.
- » can significantly reduce the number of patients who would receive adjuvant therapy unnecessarily.

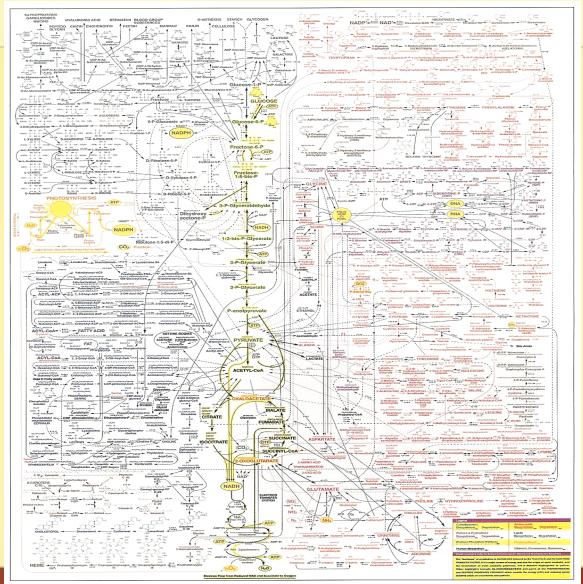
A Final Remark: Never Forget Bringing Data back to Biology



Visualize gene expression data in the context of your favorite pathways.

Download publicly available pathway maps or create your own.

The Future of Computational Biology: Can We Find It on This Map?



Never underestimate the complexity of a simple biological system.

Example:
A metabolic pathway map

Conclusions

New data acquisition technologies, such as DNA microarrays, have made molecular profiling possible. They also have created strong demands for better computational tools and systems.

To support the industrialization in molecular biology, enterprise solutions, such as the Rosetta Resolver system, have been engineered to meet the strong demand.

With these solutions, biologists can answer many challenging questions during molecular profiling. Many data-driven computational methods help biologists gain new knowledge in pharmaceutical and other bio-technology research areas.