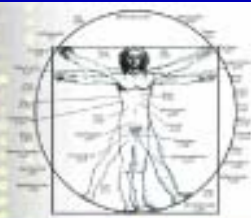


Expression to Biology Extracting Meaning from Microarrays

Montreal Microarray Symposium
18 March 2004



Microarray Expression Profiling
of Rodent Models of Human Disease

Acknowledgments

[<johnq@tigr.org>](mailto:johnq@tigr.org)

The TIGR Gene Index Team

Foo Cheung
Svetlana Karamycheva
Yudan Lee
Babak Parvizi
Geo Pertea
John Quackenbush
Razvan Sultana
Jennifer Tsai
Joseph White

Emeritus

Jennifer Cho (TGI)
Emily Chen (μ A)
Ingeborg Holt (TGI)
Feng Liang (TGI)
Kristie Abernathy (μ A)
Sonia Dharap (μ A)
Julie Earle-Hughes (μ A)
Cheryl Gay (μ A)
Jeremy Hasseman (μ A)
Priti Hegde (μ A)
Heenam Kim (μ A)
Lara Linford (μ A)
Rong Qi (μ A)
Erik Snestrud (μ A)
Shuibang Want (μ A)
Ivana Yang (μ A)
Yan Yu (μ A)

H. Lee Moffitt Center/USF

Timothy J. Yeatman
Greg Bloom

PGA Collaborators

Gary Churchill (TJL)
Greg Evans (NHLBI)
Harry Gavras (BU)
Howard Jacob (MCW)
Anne Kwitek (MCW)
Allan Pack (Penn)
Beverly Paigen (TJL)
Luanne Peters (TJL)
David Schwartz (Duke)

TIGR PGA Collaborators

Norman Lee
Renaë Malek
Hong-Ying Wang
Truong Luu
Bobby Behbahani

Funding provided by the Department of Energy
and the National Science Foundation

Funding provided by the National Cancer Institute,
the National Heart, Lung, Blood Institute,
and the National Science Foundation

TIGR Faculty, IT Group, and Staff

TIGR Human/Mouse/Arabidopsis

Expression Team

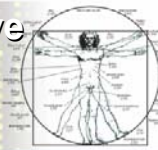
Adriana Ahumada
Tove Andersson
Joanne Emerson
Bryan Frank
Molly Freeman
Renee Gaspard
Nadeeza Ishmael
Ka Yin (Simon) Kwong
Jennie Larkin
Fenglong Liu
John Quackenbush
Yonghong Wang
Yan Yu

Array Software Hit Team

Nirmal Bhagabati
John Braisted
Tracey Currier
Jerry Li
Wei Liang
John Quackenbush
Alexander I. Saeed
Vasily Sharov
Mathangi Thiagarajan
Joseph White

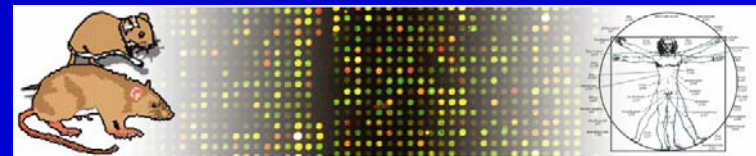
Assistant

Aseye Aboagye



Science is built with facts as a house is with stones – but a collection of facts is no more a science than a heap of stones is a house.

– Jules Henri Poincare



Levels of Biological Information

'omics

DNA

Genomics

mRNA

Functional Genomics

Proteins

Proteomics

Informational Pathways

Metabolomics

Informational Networks

Systems Biology

Cells

Molecular Biology

Organs

Medicine

Individuals

Medicine

Populations

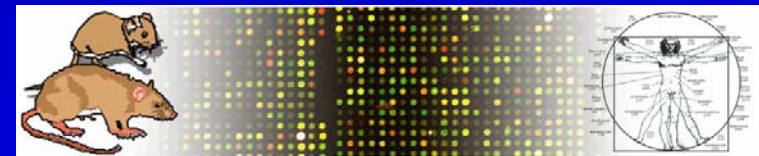
Genetics

Ecologies

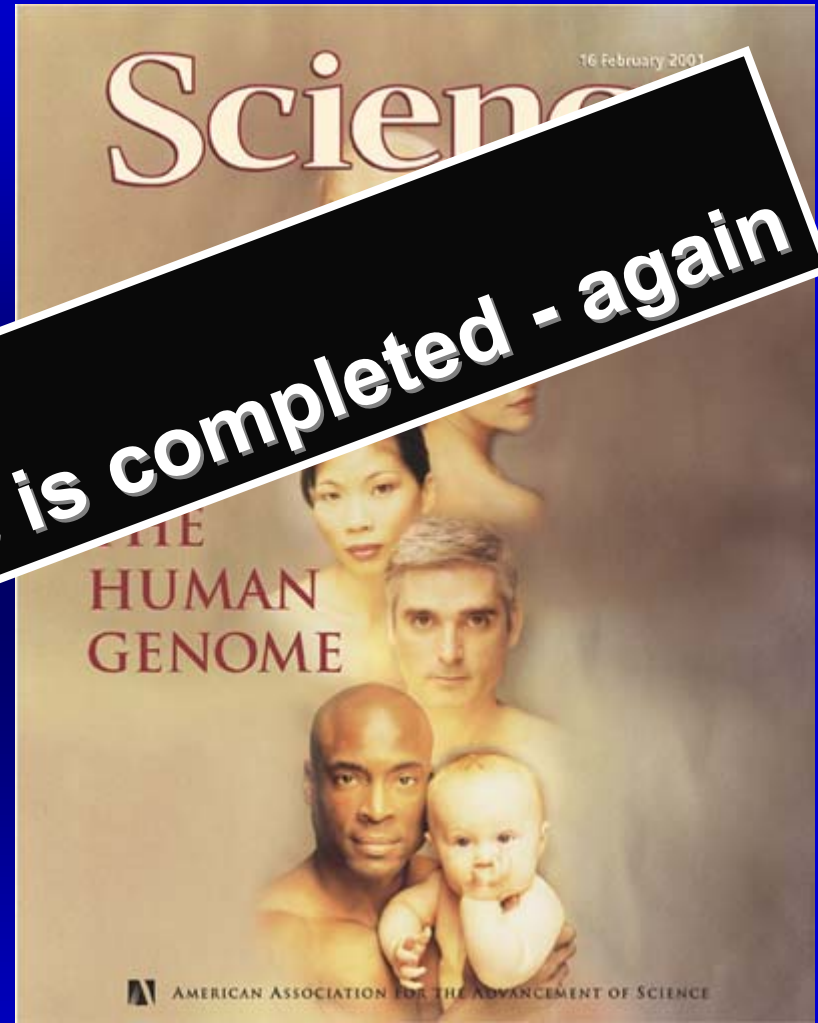
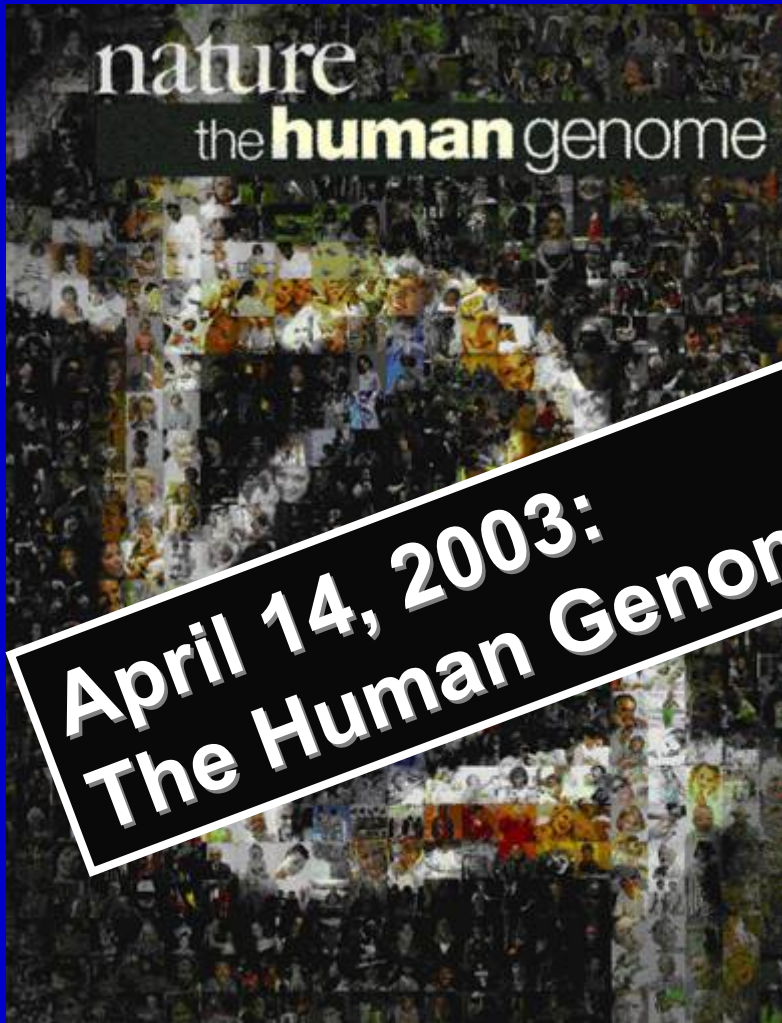
Ecology

The Future!

Traditional
Biology



February 2001: Completion of the Draft Human Genome



**April 14, 2003:
The Human Genome is completed - again**

Public HGP

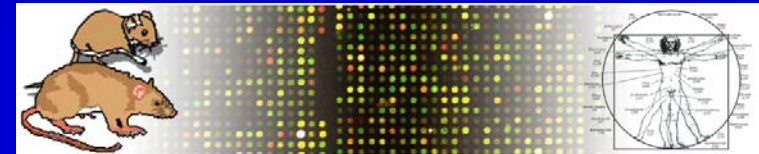
Celera Genomics

But what does *finished* mean???

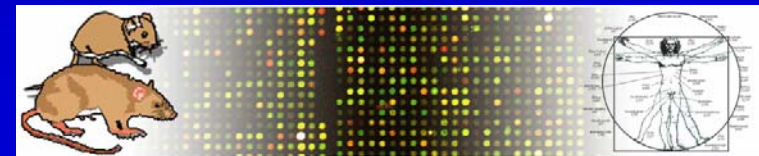


Where are the pressing questions?

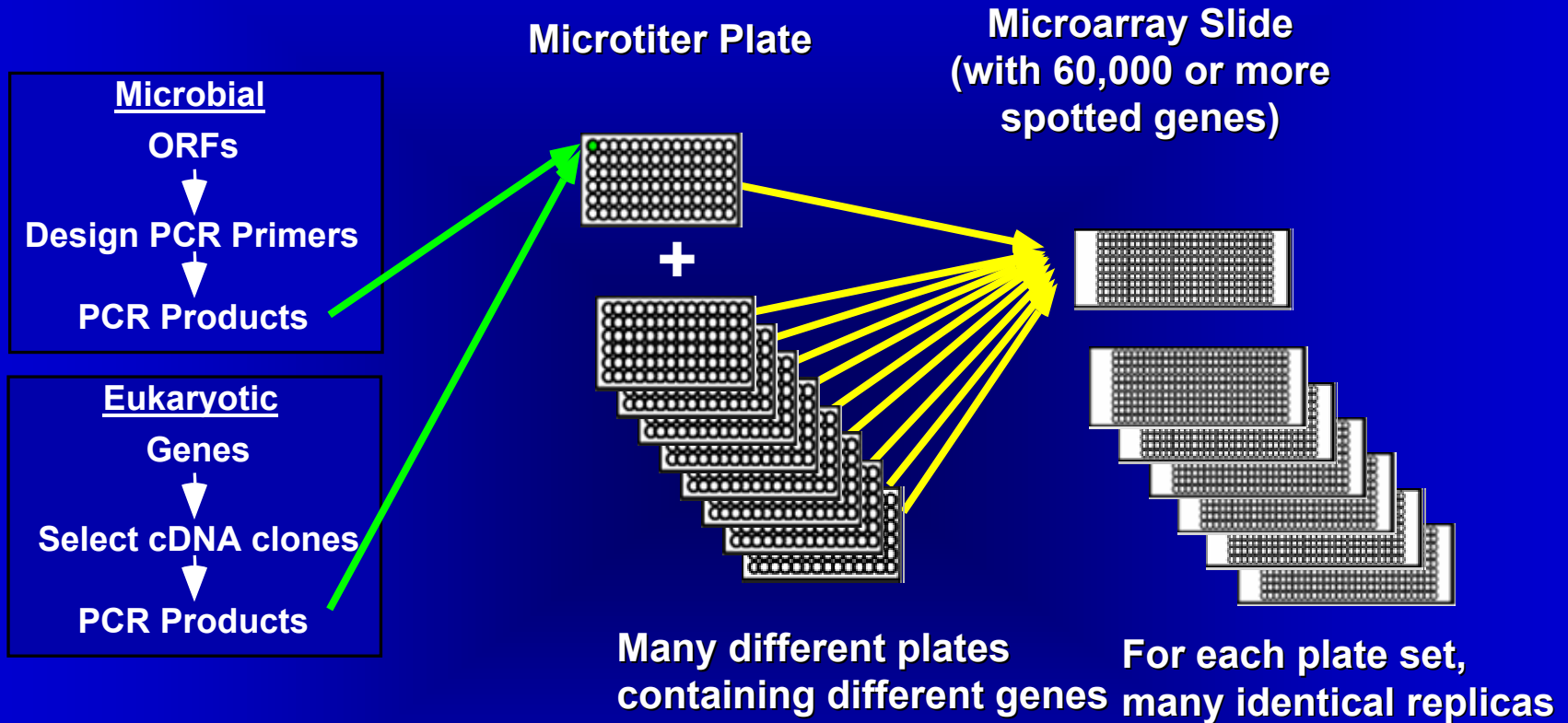
- Can we find the genes and assign them functions?
- Can we predict protein structures and functions?
- Can we reconstruct metabolic, signaling, and other pathways?
- Can we reconstruct informational networks?
- Can we link genotype to phenotype?
- Can we use genotype/phenotype to predict clinically-relevant outcome?
- Can we use cross-species comparisons to learn something?



Microarray Analysis

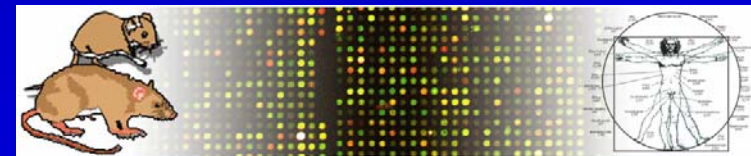


Microarray Overview I



The Beast: Microarray Robot from Intelligent Automation

<<http://www.ias.com>>

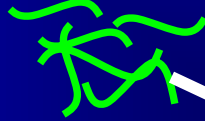


Microarray Overview II

Measure
Fluorescence
in 2 channels
red/green



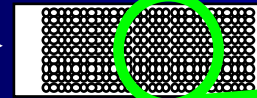
Control



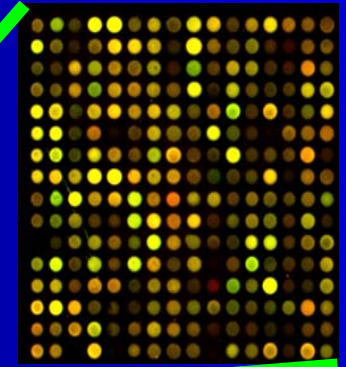
Test



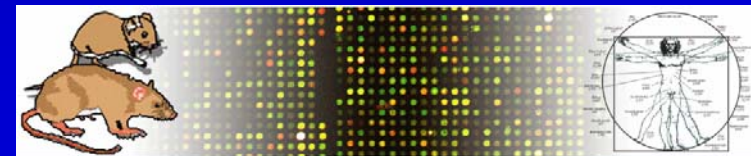
Prepare Fluorescently
Labeled Probes



Hybridize,
Wash

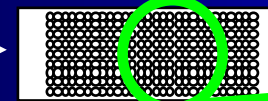
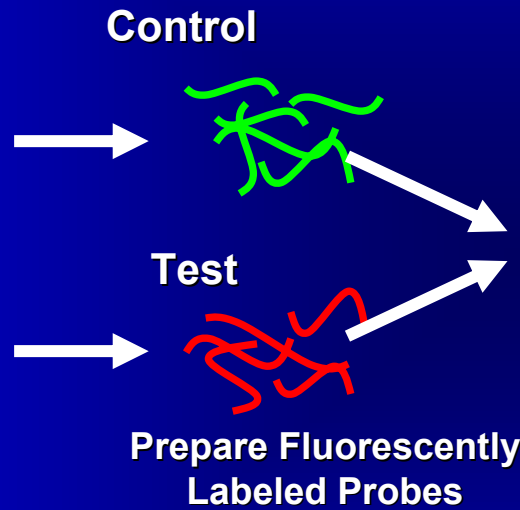


Analyze the data
to identify
patterns of
gene expression

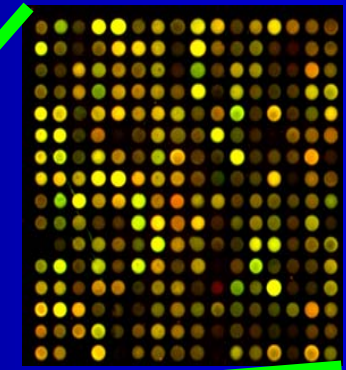


Microarray Overview II

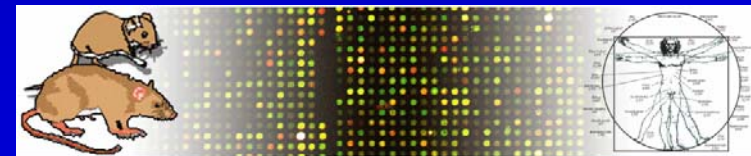
Measure
Fluorescence
in 2 channels
red/green



**Hybridize,
Wash**

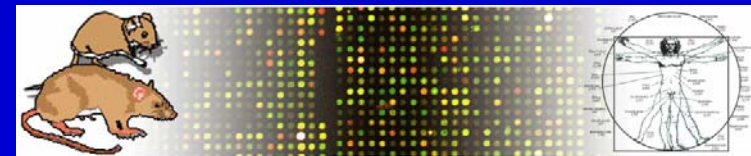


**Analyze the data
to identify
patterns of
gene expression**

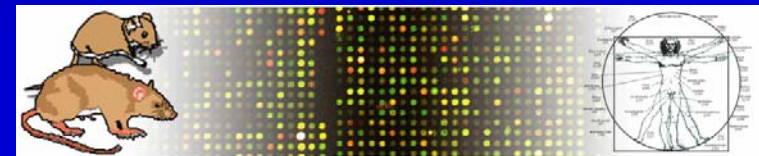


General Microarray Strategy

- Choose an experimentally interesting and tractable model system
- Design an experiment with comparisons between related variants
- Include sufficient biological replication to make good estimates
- Hybridize and collect data
- Normalize and filter
- Mine data for biological patterns of expression
- Integrate expression data with other ancillary data such, including genotype, phenotype, the genome, and its annotation



Annotating and Comparing Arrays




TIGR Gene Indices home page

www.tigr.org/tdb/tgi

~70 species

>19,000,000 sequences

 Cattle 7.0 <---Most recent version number
6-1-02 <---Date of most recent update



The TIGR Gene Index Project is supported in part by funding from the US Department of Energy, Grant #DE-FG02-99ER62852, and the US National Science Foundation, Grant #DBI-9983070. A additional funds are provided by the US National Science Foundation through grants #DBI-9813392 and #DBI-9975866.



- The TIGR Gene Indices are built using:
- **megablast** [Zheng Zhang, Scott Schwartz, Lukas Wagner, and Webb Miller (2000), "A greedy algorithm for aligning DNA sequences", J Comput Biol 2000; 7(1-2):203-14]
 - **CAP3**, developed by Dr. Xiaojin Huang
 - **Paracel TranscriptAssembler[tm]**, from Paracel Inc.
 - **DNA-Protein Search program (dps)** developed by Dr. Xiaojin Huang

- The ORF annotation of TCS is done using:
- **ESTScan** [Iseli, C. Jongeneel, C.V., and Bucher, P. (1999) ESTScan: A program for detecting, evaluating, and reconstructing potential coding regions in EST sequences. ISMB 7: 133-143.]
 - **DIANA-EST** Hatziargiou AG, Fizev P, Reczko M. Related Articles DIANA-EST: a statistical analysis. Bioinformatics. 2001 Oct;17(10):913-9.
 - **framefinder** Expressed Sequence Tag Analysis Tools Etc (c) Guy St.C. Slater 1996-1999. Human Genome Mapping Project RC, Hinxton, Cambridge, UK.

- The Genome mappings are done using:
- **blast(C)** Webb Miller - see reference
 - **Scout[tm]**, from Paracel Inc.
 - **gap2** developed by Dr. Xiaojin Huang

- The expression profiles of the ESTs are scored using:
- **'R' statistic** (Stekel, Gt, and Falciani (2000) The comparison of gene expression from multiple cDNA libraries. Genome Research 10:2055-2061).



TGICL Tools are available – with more coming

Genome Database
 Genomics
 Individual Genomes
 Research Projects
 Core Tools
 TGICL
 Software
 News
 Conferences
 TGIR International Site
 TGIR Publications
 Faculty
 Education & Training
 Genome News Review
 Related Links

TGIR Gene Indices Software Tools

TGIR Gene Indices program is making some of its software tools freely available to the scientific community. The software provided on this page represents a partial revision of the original software, adjusted to accommodate the new licensing terms. Information on the status of the software and other tools is available on the TGIR website. If you are interested in using the software, please contact us for more information. If you are interested in contributing to the development of the software, please contact us for more information. For the way of the package below. Please send the TGICL software to tgicl@tgir.org.

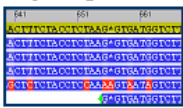
The TGIR Gene Indices Project is supported in part by funding from the US Department of Energy, Office of Biological and Environmental Research, and the National Science Foundation, Office of Biological Resources. A National Science Foundation grant is provided by the US National Science Foundation through grant #0080493 (NSF) and #0080493 (NSF).

**Geo Perteau
 Razvan Sultana
 Valentin Antonescu**



TGI Clustering tools (TGICL): a software system for fast clustering of large EST datasets

This package automates clustering and assembly of a large EST/mRNA dataset. The clustering is performed by a slightly modified version of NCBI's *megablast*, and the resulting clusters are then assembled using CAP3 assembly program. TGICL starts with a large multi-FASTA file (and an optional peer quality values file) and outputs the assembly files as produced by CAP3. Both clustering and assembly phases can be parallelized by distributing the searches and the assembly jobs across multiple CPUs, as TGICL can take advantage of either SMP machines or PVM (Parallel Virtual Machine) clusters. Here is a link to the [README](#) file which comes with the package.



clview : an assembly file viewer.

This is a graphical, interactive tool for inspecting the ACE format assembly files generated by CAP3 or phrap. Besides the ACE files, the program also supports a custom cluster layout format for the overview of a possible multiple alignments generated just from pairwise alignments, where no detailed nucleotide level alignment is needed and provided. The "containment clustering" program (nrc) mentioned in the TGI Clustering tools(TGICL) above can generate such a "cluster layout" file (*.lyt). Here is a precompiled linux version with the required dynamic FOX library included: [clview linux tar.gz](#) The program was built using the [FOX toolkit](#) by Jeroen van der Zijp, a portable and feature-rich C++ framework for developing graphical user interfaces under Unix and Windows. In order to compile the [source code](#) of this viewer, you need to download the [FOX library](#) and the [TGI C++ class library](#).



SeqClean : a script for automated trimming and validation of ESTs or other DNA sequences by screening for various contaminants, low quality and low-complexity sequences.

A precompiled Linux version is here: [seqclean tar.gz](#) Please see the [README](#) file first. Please note there is no contaminant database included in the package - you need to provide your own screening files or download and format a generic vector database like NCBI's [UniVec](#). The package also doesn't include NCBI's blastall and megablast utilities which you should obtain from the NCBI site (their full source is included in the NCBI C Toolkit). The C/C++ source for the other programs in the package is provided here:



cdbfasta/cdbyank: fast indexing/retrieval of fasta records from flat file databases. These two utilities are based on the "cdb" (Constant DataBase) concept and the file-based hashing algorithm developed by D.J. Bernstein (<http://cr.yp.to/djb.html>)

The source code is the C++ port of the original [cdb](#) library, modified to create a compact separate index file keeping the original flat file in its original format. Multi-FASTA files of up to 4GB can be indexed with cdbfasta and then any one or more records can be quickly retrieved using cdbyank. [Here](#) is a brief usage description.

TGICL Clustering tools (TGICL) a software system for fast clustering of large EST datasets

This package automates clustering and assembly of a large EST/mRNA dataset. The clustering is performed by a slightly modified version of NCBI's *megablast*, and the resulting clusters are then assembled using CAP3 assembly program. TGICL starts with a large multi-FASTA file (and an optional peer quality values file) and outputs the assembly files as produced by CAP3. Both clustering and assembly phases can be parallelized by distributing the searches and the assembly jobs across multiple CPUs, as TGICL can take advantage of either SMP machines or PVM (Parallel Virtual Machine) clusters. Here is a link to the [README](#) file which comes with the package.

The two precompiled packages below were built on Linux and Solaris, respectively. They include CAP3, megablast and all the other binaries for this platform (of course, except the Java Class unless the "jdk_*" etc.). Please note that only the Linux version was thoroughly tested at TGIR.

- [linux.tar.gz](#) Linux x86 (glibc-2.1, requires perl=5.6)
- [solaris.tar.gz](#) built on Solaris 5.8 sparc; STWFW, Ultra-250

The platform independent perl scripts can also be downloaded separately from [linux_scripts.tar.gz](#). These scripts are likely to be updated more often, as opposed to the binaries which are rather stable. In such cases there is no need to download again the full large packages provided above, although they also include these scripts.

Last update of this page was on 01/15/2002

Address/Email:

- CAP3** Huang, X. and Madan, A. (1999) CAP3: A DNA Sequence Assembly Program. *Genome Research*, 9: 107-132
- megablast** Zhang Zhang, Scott Schwartz, Lukas Wagner, and Webb Miller (2000), "A greedy algorithm for aligning DNA sequences", *J Comput Biol* 2000, 7(1): 20-32
- NCBI Toolkit** a great resource of bioinformatics source code and tools, provided by US National Center for Biotechnology Information

If you wish to build the package on other platforms or customize it, a portable C++ source code for most of the tools included in the package is given below. We do not include here the source code for the format utility and the rest of the [NCBI Toolkit](#) source and libraries which you should download from NCBI site (the C++ version, not the C++ source) in order to compile megablast. Also, we cannot distribute the source of the CAP3 assembly program - you might want to contact the author [Dr. Xiaohu Huang](#) for a precompiled version for your specific platform. Some packages:

- [cdbyank.tar.gz](#) a fast record indexing/retrieval utility (described separately)
- [index.tar.gz](#) Standalone low complexity "db" filter
- [nrc.tar.gz](#) a parallel multi-FASTA file processing tool
- [tblastn.tar.gz](#) a parallel multi-FASTA file processing tool using PVM
- [tblastx.tar.gz](#) a customized version of megablast (requires NCBI C Toolkit for compilation)

The following utilities depend on the [TGI C++ class library](#) which should be downloaded first in order to compile them. They are designed to process the custom tab-delimited output generated by the megablast program.

- [blastn.tar.gz](#) a merge-sort utility for compressed files, with multi-db output option
- [tblastn.tar.gz](#) a standalone clustering tool with merge filtering option
- [tblastx.tar.gz](#) a needed clustering tool by processing pairwise alignments
- [tblastx.tar.gz](#) a constant clustering and layout utility by processing pairwise alignments

clview - an assembly file viewer

This is a graphical, interactive tool for inspecting the ACE format assembly files generated by CAP3 or phrap. Besides the ACE files, the program also supports a custom cluster layout format for the overview of a possible multiple alignments generated just from pairwise alignments, where no detailed nucleotide level alignment is needed and provided. The "containment clustering" program (nrc) mentioned in the TGI Clustering tools(TGICL) above can generate such a "cluster layout" file (*.lyt). Here is a precompiled linux version with the required dynamic FOX library included: [clview linux tar.gz](#) The program was built using the [FOX toolkit](#) by Jeroen van der Zijp, a portable and feature-rich C++ framework for developing graphical user interfaces under Unix and Windows. In order to compile the [source code](#) of this viewer, you need to download the [FOX library](#) and the [TGI C++ class library](#).

SeqClean a script for automated trimming and validation of ESTs or other DNA sequences by screening for various contaminants, low quality and low-complexity sequences.

A precompiled Linux version is here: [seqclean tar.gz](#) Please see the [README](#) file first. Please note there is no contaminant database included in the package - you need to provide your own screening files or download and format a generic vector database like NCBI's [UniVec](#). The package also doesn't include NCBI's blastall and megablast utilities which you should obtain from the NCBI site (their full source is included in the NCBI C Toolkit). The C/C++ source for the other programs in the package is provided here:

- [index.tar.gz](#) Standalone low complexity "db" filter
- [tblastn.tar.gz](#) a multi-FASTA file processing utility (described separately)
- [tblastx.tar.gz](#) A parallel multi-FASTA file processing tool on a multi-CPU machine
- [tblastx.tar.gz](#) A parallel multi-FASTA file processing tool on a PVM cluster

cdbyank/cdbfasta, fast indexing/retrieval of fasta records from flat file databases. These two utilities are based on the "cdb" (Constant DataBase) concept and the file-based hashing algorithm developed by D.J. Bernstein (<http://cr.yp.to/djb.html>)

The source code is the C++ port of the original [cdb](#) library, modified to create a compact separate index file keeping the original flat file in its original format. Multi-FASTA files of up to 4GB can be indexed with cdbfasta and then any one or more records can be quickly retrieved using cdbyank. [Here](#) is a brief usage description.

The source code is the C++ port of the original [cdb](#) library, modified to create a compact separate index file keeping the original flat file in its original format. Multi-FASTA files of up to 4GB can be indexed with cdbfasta and then any one or more records can be quickly retrieved using cdbyank. [Here](#) is a brief usage description.

For TGI Software Tools Comments/Questions send mail to tgicl@tgir.org for specific feature requests or programming questions please write directly to the author (Geo Perteau) at gperteau@tgir.org

bioRxiv preprint doi: <https://doi.org/10.1101/000000>; this version posted January 1, 2014. 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 aCC-BY-NC-ND 4.0 International license.



Available with source

A TC Example

>TC161360 TC25195 TC29362 TC33731 TC40754 TC149101
 TGAAGCTCACAAAGAACTTTTATTCCTTTTTAAATAGACACTAAAATTATCTCCTAGTCATGAGAAATTGGTAAAGACTAAT
 TATTTGAGAATCTGACGATGACTAATGTAATAATCATTAAAGGAAATGAATTTTCAGAGAGGGGAAACTTTTCAAATGAATA
 CTGCATTTAAAACTTTTCAGCTTGACACTCCTCCTCCCACTCCCATCCTCCTCCAGGCATAGCGGTATCTTCTTTAGCT
 TAGGGTACCTTCTATGGAGAAGAATGGATATGGAGAATCGTGCTGTGGCTTGTAAAGTGGGCAGAACTTAGTAAAGACCTA
 CTGGATGAGTAACTCCTTGGGAGCATGTGTCAGATAGGTAGGAATAGCTCAATATGACTGGATGTGCCACTATTCAAAC
 ACAGGTTAGTATTATGTGGCAGAAAGCATCCCAATGTGTTAGTATGTTATGGAGAAGAAAGAAACATCCAAGGTGGAGTATCCA
 TTGCAGGCCTGCACAAAAGTTTTATTTACTTAGAGCTTGTGTTTGAAGACCACACAGGGGAAAAGGTGCTACTTCCAGTTT
 CTTTGTAAATAACAGGAAAATAACTCCCACCGGTAGCCTCTAATAAAAATAGAAAATATCCAAGGAGTGAACTTAAAGCTGTT
 CATATACCCATAATGCCTAGAAGCAGACTTGTCAATGATATGCTGATGATAGGCTATGGTGAGATCTTTTTAGGCTAACAG
 TGTCTTAGGTCAGGTGCTAGCATCCCTGCTCAGGAACAGGGGTGGGAAAGTATGGTGGCTGAGATTTAGGATTTTAAAC
 TGTGTTGTTTTTAAAGCATGATCTTTGTGTGGTAAATTTATAGTGCATATAAGATGTGTTTTGTGGTGCATCTTATAACTTTC
 CAGCTAATTGCATATTAATGTCACGACTAGTTTTCCAAATGATGTAAGATTCTGGGTGCTTTTATTCATATGGTGTCAA
 TCCAATTCGACTCTTTTGTGATTGACACATTGCTCACAAAAGTATATACTTTGATATAGCTTATACAGGCATATGGGCA
 TAGATAATTTGGTTATTTCAAACATTTCTAGAATTTGAAGAGCTGGGTTTACTCAAGTCCCCACAACATATTATTGAATA
 TCTTGGCTGCTTCTTACATTTCAACTAAATACAAAATAGTGAATAGGTTTTTGTGTTTTTGTGTTTTTGTGTTTTTCTG
 TGAATCATGAAAGACAAAATAAACTACTCTGCCATTTGGAAAGAAATCTCAAAGCACAGCTGTTGCTTAGGATTTGAGATC
 TGGGAACCCATTACATTTCTGTCTCACCACTCTTTTTCTCTTGACTAAAAGAACAAAATACATAAACGATGTGACCAGA
 AGCCAAGAACTGGAGATGGGCAAAGTTAAGACATGAACTTGCCTATGTGTAAGCTATGCTTTTTGTACAGAGACAGAC
 TTAGGACTAGGCCTTTACAGCTTCCCAAGTGGGACACAGAAGCTTGAAGAACATGCCATATTTTGTGCTTCCCCACAG
 GCATATAGGTGCTCATTTTCTCATTTATTAACAATTTCTCATTAAACAATTTCTCATATTTCTTAAAGACAATCAT
 GCAAAAAGGACTCCACAAAACATGAAGAAGGAGTCAGAGAGTTTCAGAGATGAACAGAACAGTCAATGTGCTTTTTCCCTCT
 GCTGTAACGTGCTGATGGAGGAAATGCGGAAAAACCATTACCCTACTGAGACACACATCCAAGGAAAAGGTGAATGAGAA
 TTGGATTTTCAGCGTTTCTGCCAGTGAAGACTCCCTAAGCAAGGCAAAAGACAATCTTTTATAATCACTGCATTCTTC
 AAATGTGAAGGAGAGCTAGATGTGGCTTCTCTATGCAAAAGTTAAGCTGTGAAAATATGGAGAATAGATTGTGGAAAGGCCA
 CAAGAGATGAGGGTAACTATGTGCTTGGAAAGCTCTCTAAACAATAATGCTTAGAAAAGAAAAGAGTTTAAAGGAGCGATC
 TCATTTTCCACTGGCAAGGATTCCTGCTCTGGGGAGGTGCTGGAAAAATTTATTCTGCTTCTTCTTCTTCTTCTCCTCT
 CACTGGAATAATTTCAATTTCTCTTCTTCTTTGGTTGGGAAGATACTCCACTTAGCTCAATGATTTCTTCTTCTTCTTCT
 TAATCGTCTGCTCATTTTTTCTCCAGCTGACAGCAGAACACATTAGATTTGGATCTGGTACACATTTCTTCTTCTTCTC
 TTCTCCACAATACAGCTGTCAAAAGTTTCTGGAAGAAGGCAGAGATCAGCATCGCCGACAGAAATGCCCAAGAACACAGA
 CTGAATGCTGTTTACAGTACTGTGTAGATGGGGAGTTTTCTCTGAGGAATTAGATGGTTTACAGTGCATGATATCAACAT
 ACGGCTTGGAAAGTTTGAATAAGCTCCAGTCTTCTCATTTTTAAAAAATGAGAAAAGTGCATGTTAAGTATGTCCATGATT
 GAAAGAATTATCCAGAAATGGCAGCAAAGAGGCTTAGAGAGCTCATTATCACTTTTGTCTTGAACAACTCTTCTTCTGGA
 GGTTTTTCTGCTGCAGCTGCCAGGAGTATCTGAAATAATGTACATAATGCCTCCCAGAGAGGGTACCATACTCA
 AACAGATGGGTGCGAAGACCCCTGTGGGGATCATCAGCAGTCCCCCAGGGTAAATATGGAAGAGGCCATTATGATTTGG
 ACAGCCCCAAAGCCTTTGATCCCTCATGAAGAAGCTTTGTGTGGGGCCACCAGTGAAGATGTCTTTTGGAGGTTAC
 TTTTGGAGCAGGTTGCAATGGCGAGGGGACCTTTTGTAGCTCTGCTGGGAAAAGTCCACTCATTGAGTTTCAAACGCCT
 GGGTCTCCAAGGCTCCAAGGCTCCAAGGCTTAAAGGGCCAAATAGCCACTCACAGCTGAGATCAGTTGGCCTCTTCTATC
 TCAAGTACTTGAAGATAGCTACAGTGTACTTACATATAATAATAAATAAATCTTTAAAAAAGGAAATTC

-cell differentiation antigen
 (B-LYMPHOCYTE
 (Homo
 1|X07203 CD20 receptor

Mouse Genome Database

The screenshot displays a genomic database interface with the following components:

- Header:** "The ENCODE Project Data Portal" and "ENCODE Project Data Portal: ENCODE Project Data Portal".
- Sequence Alignment:** A large block of text showing sequence alignment with various annotations and colors (blue, red, green).
- Table 1:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 2:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 3:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 4:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 5:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 6:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 7:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 8:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 9:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 10:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 11:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 12:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 13:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 14:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 15:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 16:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 17:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 18:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 19:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".
- Table 20:** A table with columns for genomic coordinates and associated data. Headers include "chr", "start", "end", "strand", "feature", "score", "start", "end", "strand", "feature", "score".



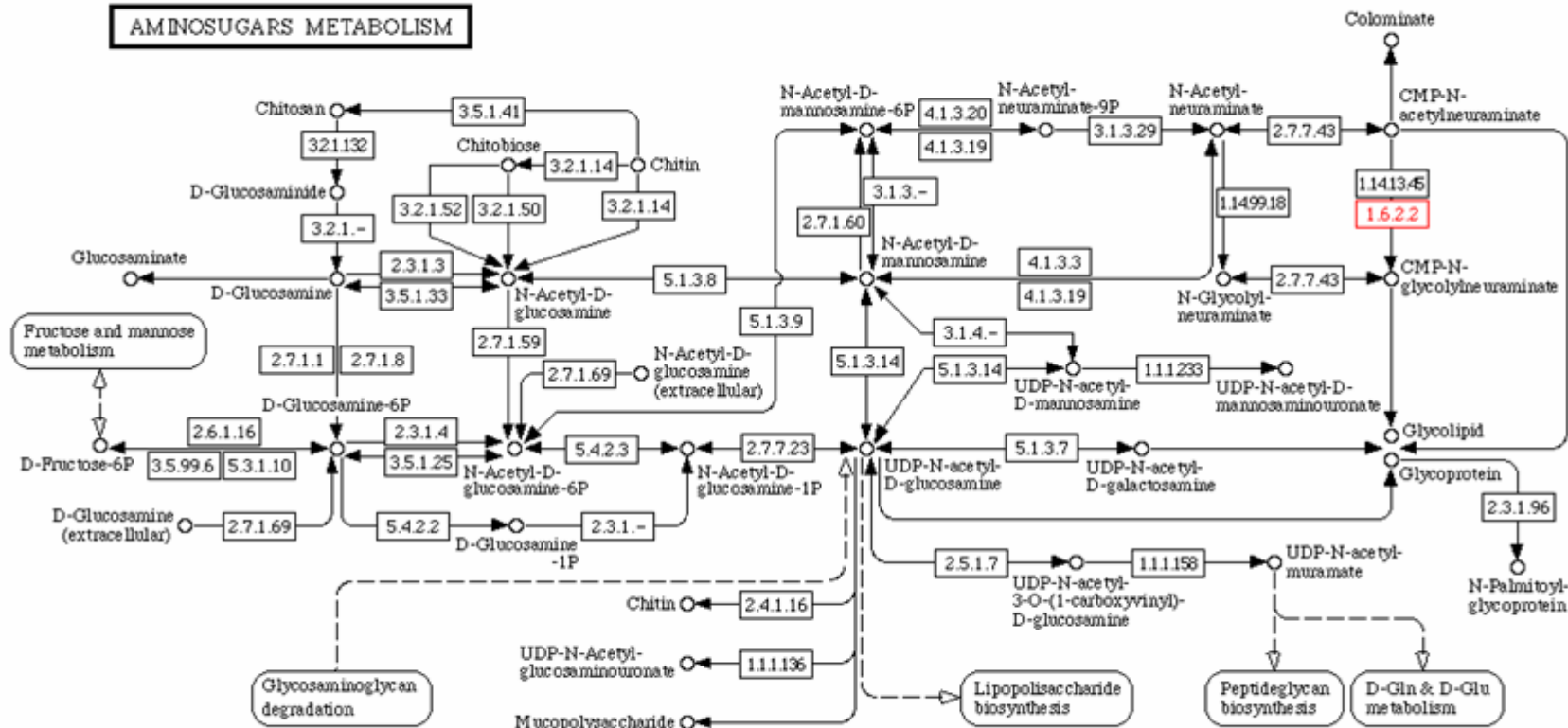
GO Terms and EC Numbers



Position of term GO:0004128

Function	GO Assignments	% of 12819 total TC/NP with GO Assignments
	31.21%	2.52%
	24.25%	1.92%

AMNOSUGARS METABOLISM



00530 10/5/01

- (I)enzyme (GO:0003624) +
- (I)oxidoreductase (GO:0016491) +
- (I)oxidoreductase), acting on NADH or NADPH (GO:0016651) +
- (I)oxidoreductase), acting on NADH or NADPH), NAD or NADP as acceptor (GO:0016652)
- (I)cytochrome b5 reductase (GO:0004128) ○

last modified on: September 05, 2001



The TIGR Gene Indices <<http://www.tigr.org.tdb/tdb/tgi>>

TIGR THE INSTITUTE FOR GENOMIC RESEARCH

Home > Databases > Gene Indices

TIGR Gene Indices

What's New | BLAST Search | TGI Software | FAQ

Integrating data from international EST sequencing, genome sequencing and gene research projects, Gene indices are an analysis of transcribed sequences represented in the world's public EST data.

EGO - linking orthologous genes across eukaryotes
Genome Maps - mapping TC sequences to eukaryotic genomes
RESOURCEER - cross referencing mammalian sequencing resources
DAS - providing distributed annotations for completed genomes

Animal Gene Indices

Arabidopsis thaliana 10 8/10/02	Bos taurus 23 2/10/02	Canis lupus 11 8/10/02	Cattle 70 8/10/02
C. elegans 18 8/10/02	Chickens 18 8/10/02	C. parvulus 18 8/10/02	Drosophila 70 8/10/02
Homo sapiens 28 8/10/02	Homo 60 8/10/02	Mus musculus 42 8/10/02	Mouse 80 8/10/02
O. latipes 18 8/10/02	Oryzias latipes 28 2/10/02	Pig 40 2/10/02	Rat 80 8/10/02
Schistosoma mansoni 48 8/10/02	Xenopus laevis 11 8/10/02	Zebrafish 80 8/10/02	

Plant Gene Indices

Arabidopsis 48 8/10/02	Eidry 48 8/10/02	Chlamydomonas reinhardtii 18 2/10/02	Cotton 10 8/10/02
Ses. plant 18 8/10/02	L. japonica 18 8/10/02	Maize 11 8/10/02	Medicago truncatula 48 8/10/02
Tritic 11 7/10/02	Rice 80 8/10/02	Rye 18 8/10/02	Sorghum bicolor 48 8/10/02
Soybean 18 8/10/02	Tomato 18 8/10/02	Wheat 48 8/10/02	

Fungal Gene Indices

Cryptosporidium parvum 18 8/10/02	Dicrocoelium viverrini 18 8/10/02	Emericella nidulans 18 8/10/02	Leishmania 18 8/10/02
Microspora 18 8/10/02	Plasmodium falciparum 18 8/10/02	Plasmodium berghei 18 8/10/02	Plasmodium vivax 18 8/10/02
Sarcocystis parva 18 8/10/02	Theileria parva 18 8/10/02	Trypanosoma brucei 18 8/10/02	Trypanosoma cruzi 18 8/10/02
Trichoplax 18 8/10/02	Trichomonas axosporus 18 8/10/02		

Fungal Gene Indices

Aspergillus nidulans 18 8/10/02	Coccidioides immitis 11 8/10/02	Cryptosporidium 18 8/10/02	Microspora 18 8/10/02
Microspora 18 8/10/02	Sarcocystis parva 18 8/10/02	Sarcocystis parva 18 8/10/02	

Cattle 70 — Most recent version available
 Cattle 70 — Date of most recent update

The TIGR Gene Indices Project is supported in part by funding from the US Department of Energy, Grant #CE-000-00000000, and the US National Science Foundation, Grant #DCB-00000000. Additional funds are provided by the US National Science Foundation through grants #DCB-9973366 and #DCB-9973366.

The TIGR Gene Indices are built using:

- **megalign** (Sheng Zhang, David Sidman, Luke Wagner, and Walsh Miller (2000), "A greedy algorithm for aligning DNA sequences", J. Comput. Biol. 9(8): 573-581)
- **PARCEL**, developed by Dr. Steven Wang
- **DNA Feature Search program (dfe)** developed by Dr. Xiang Wang

The ORF identification of TCs is done using:

- **REVERSE** (Liu, C., Ingendahl, C.T., and Bracken, P. (1995) ESTScan: A program for detecting, evaluating, and reconstructing potential coding regions in EST sequences. *EMBO J.* 14:1417-1423)
- **DIANA-EST** (Suzanne Aujay, Peter P. Bork, M. Balasubramanian, DIANA-EST: a statistical analysis. *Bioinformatics* 2001 17(12):1502-1510)
- **FrameIndex** (Suzanne Aujay, Peter P. Bork, M. Balasubramanian, FrameIndex: a statistical analysis. *Bioinformatics* 2001 17(12):1502-1510)

The database mappings are done using:

- **MEMO** (Walsh Miller, see reference)
- **GeneMap** from Paracel Inc.
- **EST** developed by Dr. Steven Wang

The expression profiles of the ESTs is done using:

- **W** (Walden, O., and Faloutsos (2003) The comparison of gene expression from multiple cDNA libraries. *Genome Research* 13(2):203-209)

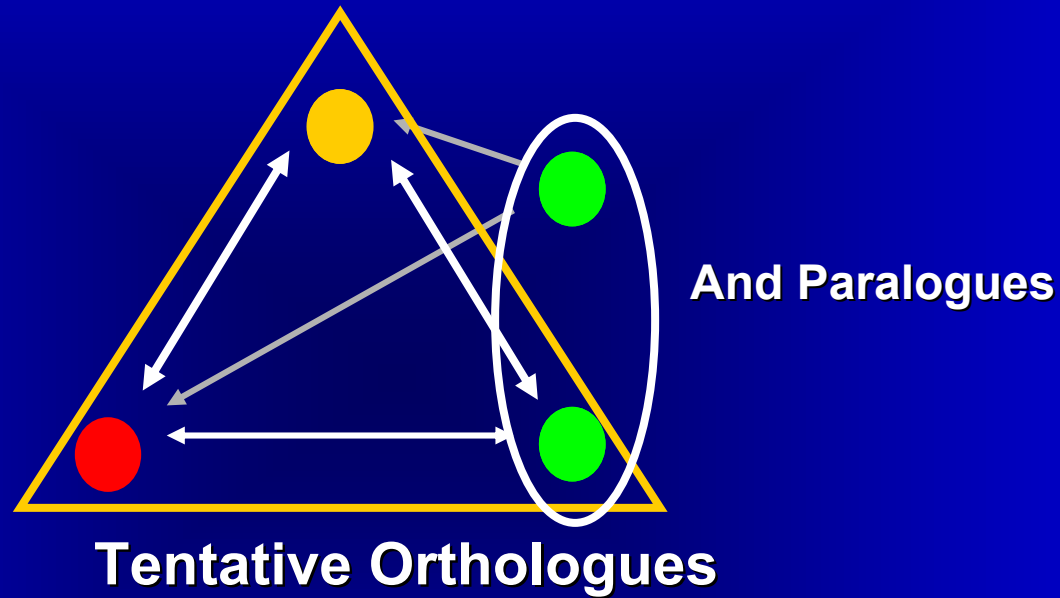
Eukaryotic Gene Orthologs

The **Eukaryotic Gene Orthologs (EGO)**, is a database for orthologous genes in eukaryotes. EGO is generated by pair-wise comparison between the Tentative Consensus (TC) sequences that comprise the TIGR Gene Indices from individual organisms. The reciprocal pairs of the best match were clustered into individual groups and multiple sequence alignments were displayed for each group. The EGO database can be accessed through the **SEARCH** function. The release notes for the current EGO can also be referenced.



Dan Lee, Ingeborg Holt

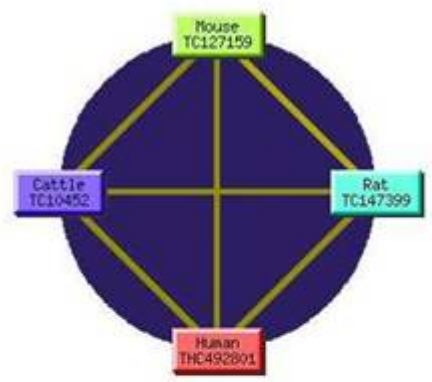
Building TOGs: Reflexive, Transitive Closure



TOGA: An Sample Alignment: bithoraxoid-like protein

Tentative Ortholog 3220

 [Cattle|TC10452](#)
 [Rat|TC147399](#)
 [Mouse|TC127159](#)
 [Human|THC492801](#)



Sequence1	Sequence2	PID	Match length
Rat TC147399	Cattle TC10452	89.49	408
Mouse TC127159	Cattle TC10452	88.73	407
Human THC492801	Cattle TC10452	89.93	406
Mouse TC127159	Rat TC147399	92.93	646
Human THC492801	Rat TC147399	89.63	375
Human THC492801	Mouse TC127159	89.43	387

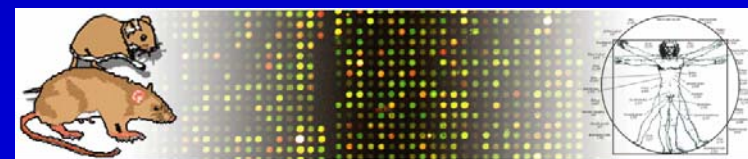
CLUSTAL W (1.8) multiple sequence alignment

```

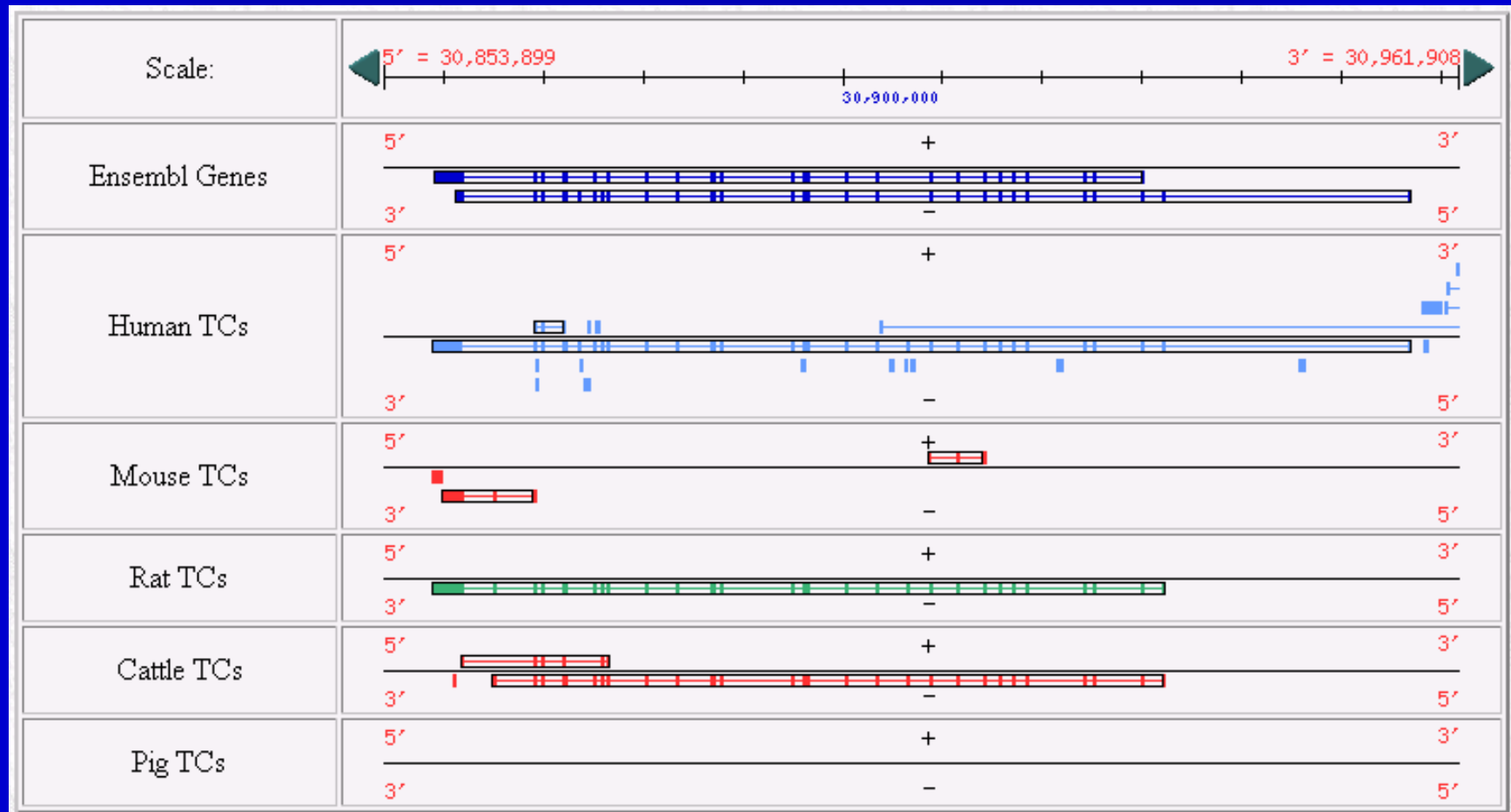
mouse| TC127159      TGGTCTACACAGGCTC-AG-GTGGCCACCACGTGC----CCACTGACATGATTAGCACTA
human| THC492801    TGGTGTGAGTGGGTTCCAA-GCGACTGCCATGTGCTAGTCC&CTGACATGATTGACATTA
cattle| TC10452     CAGCCTGGGAGGGGCTCCAACGTGCCTTCCACGTGCCCGTCAATGGACATGATTAACGCTA
rat| TC147399      A&GGTCTGCATGGCTCCAGGCAGCC--ACATGTGCC---ACTGACATGATTAACGCTA
      *          ** * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *

mouse| TC127159      TTATTCCTGGGGGACATTA&ATTAAGGAATGACACAGGAAGCCAGACAGTGGCCTTATTC
human| THC492801    ACATTCCTGGGGGGCATTAA&ATTAAGGAATGACACAGGGAGCC&AGAGTGGCCTTATTC
cattle| TC10452     TTATTCCTG&GGGGCATTAA&ATTAAGGAATGACCGCAGGGAGCC&AGAAAGCAGCTTATTC
rat| TC147399      TTATTCCTGGGGG&CATTAA&ATTAAGGAATGACACAGGAAGCC&AGACAGTGCCTTATTC
      ***** * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *

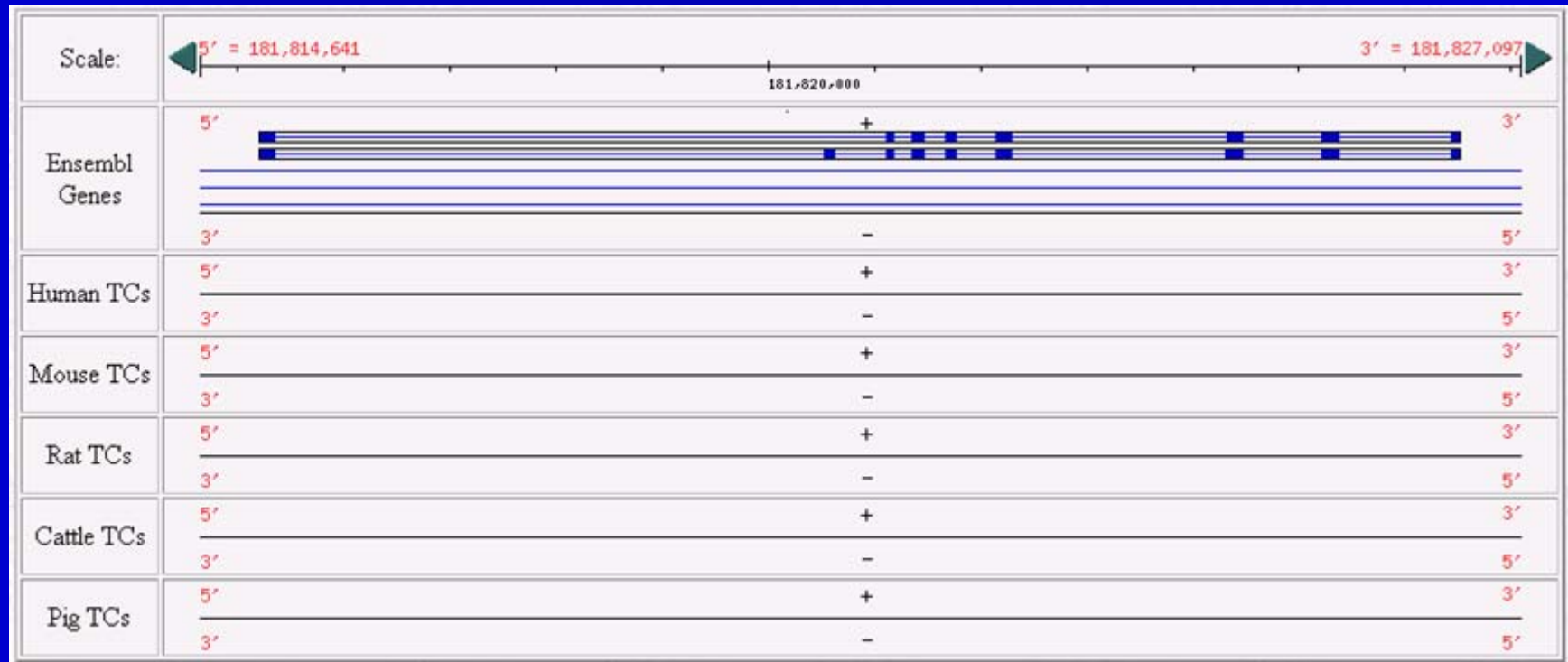
mouse| TC127159      AGTTGGATTCTGGATCACAATCAGGAA&ATAGTCTTTATCTGGTGCCACCATA&TTTCATT
human| THC492801    G&TTGGATTCTG&ATCACAATCAGGAA&ATAGTCTTTATCTGGTGCA&CCATA&TTTCATT
cattle| TC10452     AGTTGGATTCTG&ATCACAATCAGGAA&ATAGTCTTTATCTGGTGCA&CCATA&TTTCATT
rat| TC147399      AGTTGGATTCTG&ATCACAATCAGGAA&ATAGTCTTTATCTGGTGCCACCATA&TTTCATT
      * * * * * * * * * * * * * * * * * * * * * * * * * * * * * *
    
```



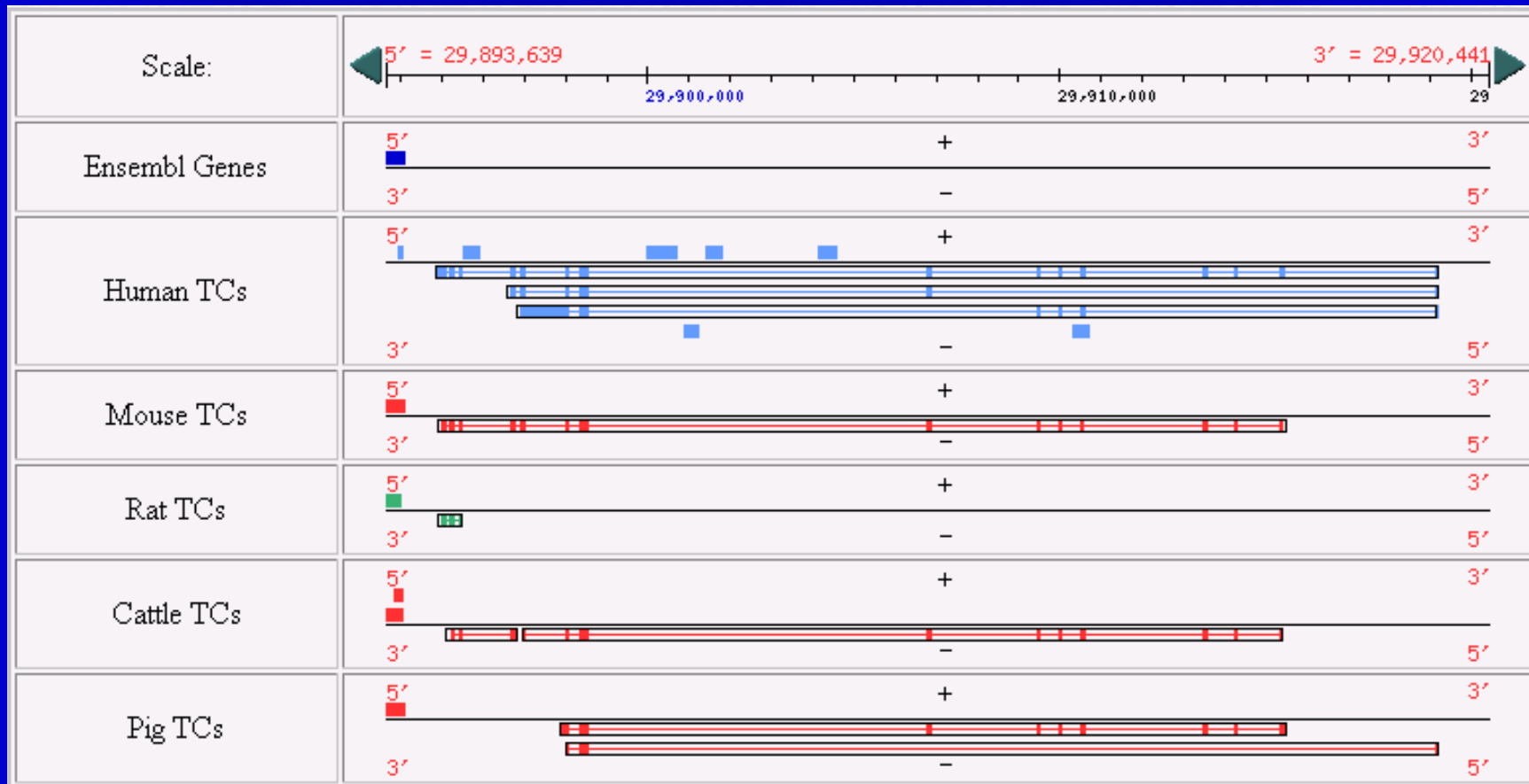
Gene Finding in Humans is easy!



Gene Finding in Humans is easy?



Gene Finding in Humans is difficult?



Gene Finding in Humans is difficult?

Scale:			
Ensembl Genes	5' + 3'		
Human TCs	5' + 3'		
Mouse TCs	5' + 3'		
Rat TCs	5' + 3'		
Cattle TCs	5' + 3'		
Pig TCs	5' + 3'		

A genome and its annotation is *only* a hypothesis that must be tested.



RESOURCERER


Jennifer Tsai

<http://pga.tigr.org> The Institute for Genomic Research's Program for Genomic Application



Microarray Expression Profiling
of Rodent Models of Human Disease

The Institute for Genomic Research | Duke University | Boston University | The Jackson Laboratory | Medical College of Wisconsin | University of Pennsylvania

- What's New
 - Description
 - Contacts
 - Data
 - Analysis Tools
 - Targets
 - Protocols
 - Outreach
 - NHLBI PGA Participants
- email the
PGA_webmaster
- 
- ©1999-2003 TIGR

RESOURCERER 6.0

- [READ ME](#)
- [GB Search](#)
- [Gen Marker Search](#)
- [BLAST Search](#)



RESOURCERER(Genome Biology 2001 PDF) provides annotation based on the TIGR Gene Indices (TGI) for commonly available microarray resources, including widely used clone sets and Affymetrix GeneChip Arrays. RESOURCERER also allows comparisons between resources from the same species using either the TGI or UniGene and between species using the EGO database.

- Physical map for Rat Sets **NEW**
- Rat Genetic Marker Search **NEW**
- Array search based on physical chromosome coordinates (under Gen Marker Search). **NEW**
- Addition -- Human: Illumina_Human, Agilent_Human1A. **NEW**
- Addition -- Mouse: NIA_7.4K, FANTOM1, FANTOM2. **NEW**

Comments are welcome

Select a single resource in "Data Set A" (while leaving "Data Set B" as "None") to generate hyperlinked annotation.	Data Set A: <input type="text" value="Human: affy_HG-U95Av2"/>
To Compare Two Resources: Select both, choose the basis for comparison (EGO or UniGene), and the type of comparison to perform (Intersection, A_unique, or B_unique).	Data Set B: <input type="text" value="None"/>
	<input checked="" type="radio"/> EGO <input type="radio"/> UniGene
	<input checked="" type="radio"/> Intersection <input type="radio"/> A_unique <input type="radio"/> B_unique
	<input type="button" value="GetTable"/>

- [What's New](#)
- [Description](#)
- [Contacts](#)
- [Data](#)
- [Analysis Tools](#)
- [Targets](#)
- [Protocols](#)
- [Outreach](#)
- [NHLBI PGA Participants](#)



<http://pga.tigr.org/tools.shtml>

RESOURCERER: An Example

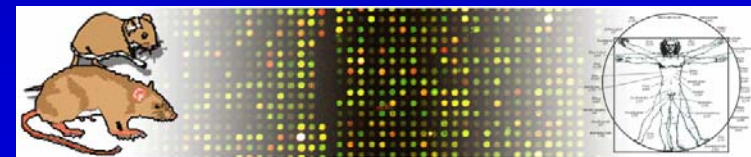


BMP & affy_HG-U95Av2

Based On: EGO

There are 956 rows in this table. [Download](#) [Jump to page](#)
 Page 1100 is currently displayed. [Next](#)

Dataset A	Rearray ID	Clone Name	GenB Acc	BMAP TC	Dataset B	Probe ID	Clone Name	GenB Acc	affy_HG-U95Av2 TC
BMAP	1-a-4	UI-M-AQ0-aaa-k-0-UI	AI835193	TC654714	affy_HG-U95Av2	38871 at		AJ006288	THC1358632
BMAP	1-a-7	UI-M-AQ0-aaa-i-0-UI	AI835201	TC662425	affy_HG-U95Av2	41449 at		AJ000534	THC1264723
BMAP	1-a-9	UI-M-AQ0-aaa-l-0-UI	AI835206	TC774982	affy_HG-U95Av2	40425 at		M57730	THC1251303
BMAP	1-b-1	UI-M-AQ0-aaa-c-0-UI	AI835214	TC754936	affy_HG-U95Av2	36542 at		AF030409	THC1252352
BMAP	1-b-3	UI-M-AQ0-aaa-e-0-UI	AI835221	TC773710	affy_HG-U95Av2	33577 at		AC004079	THC1357542 THC1228447 THC1359621 THC1395026
BMAP	1-b-3	UI-M-AQ0-aaa-e-0-UI	AI835221	TC773710	affy_HG-U95Av2	41321 s at	IMAGE-965756	AA528077	THC1357542
BMAP	1-b-3	UI-M-AQ0-aaa-e-0-UI	AI835221	TC773710	affy_HG-U95Av2	41322 s at	IMAGE-2517632	AI816034	THC1357542
							GO:0008656 (caspase activator) GO:0016064 (humoral defense mechanism (sensu Vertebrata)) GO:0016066 (cellular defense response (sensu Vertebrata))		leukemia/lymphoma 10 [Mus musculus]



RESOURCERER: Using Genetic Markers



Mapped to Mouse Genome: Chr16

Range: 23628083-28355680

There are 14 rows in this table.

[Download](#)

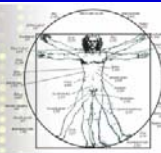
Jump to page

1

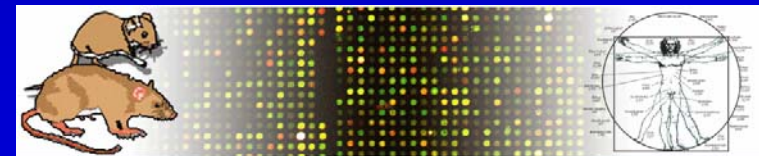
Page 1 of 1 is currently displayed.

Marker Name	UniSTS ID	Genetic Map	Data Set	GenBank Acc	TIGR TC	Chr Left	Chr Right	TGI Annotation
.MMHAP12FLA1.seq	122645					23628083	23628267	
			TIGR_25K_Mouse_Set	AI848192	TC689579	23651614	23652949	preprosomatostatin, somatostatin [Mus musculus]
			TIGR_25K_Mouse_Set	AI853264	TC673044	23687727	23688118	
			TIGR_25K_Mouse_Set	AU017440	TC763498	24433266	24434248	
			TIGR_25K_Mouse_Set	AU024130	TC721291	24754238	24758505	
			TIGR_25K_Mouse_Set	AW553532	TC749004	25732192	25760355	homologue to GP 7022921 dbj BAA91769 unnamed protein product (Homo sapiens), partial (77%
			TIGR_25K_Mouse_Set	AU014937	TC717150	25804670	25813571	
			TIGR_25K_Mouse_Set	AI854418	TC676007	27152177	27210459	RIKEN cDNA 2610529H0 gene [Mus musculus]
			TIGR_25K_Mouse_Set	AW544111	TC669191	27212167	27214345	
			TIGR_25K_Mouse_Set	AU017969	TC671593	27728646	27729800	
			TIGR_25K_Mouse_Set	AI837800	TC683108	27916927	28202210	fibroblast growth factor-related protein FGF-12A; fibroblast growth factor homologous factor 1; fibroblast growth factor 12 [Mus musculus]
			TIGR_25K_Mouse_Set	AI847382	TC683108	27916927	28202210	fibroblast growth factor-related protein FGF-12A; fibroblast growth factor homologous factor 1; fibroblast growth factor 12 [Mus musculus]
01.MMHAP38FLA1.seq	122654					28355515	28355680	

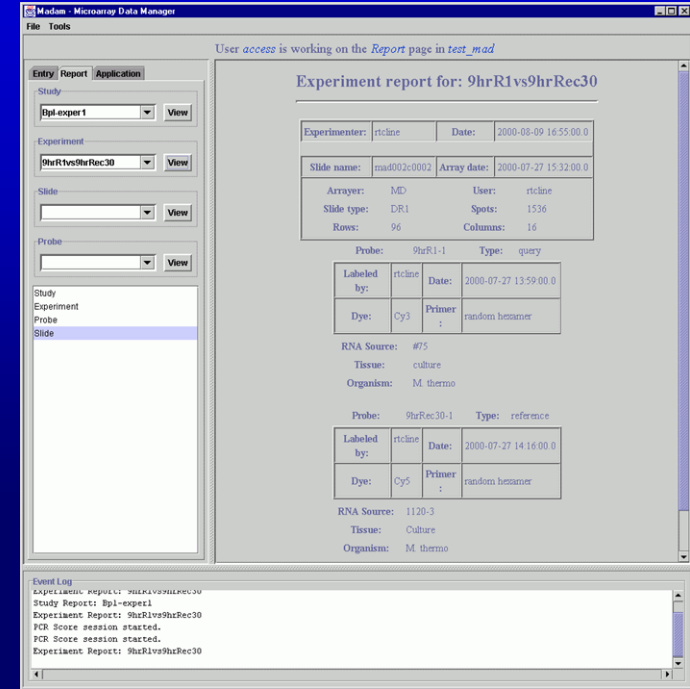
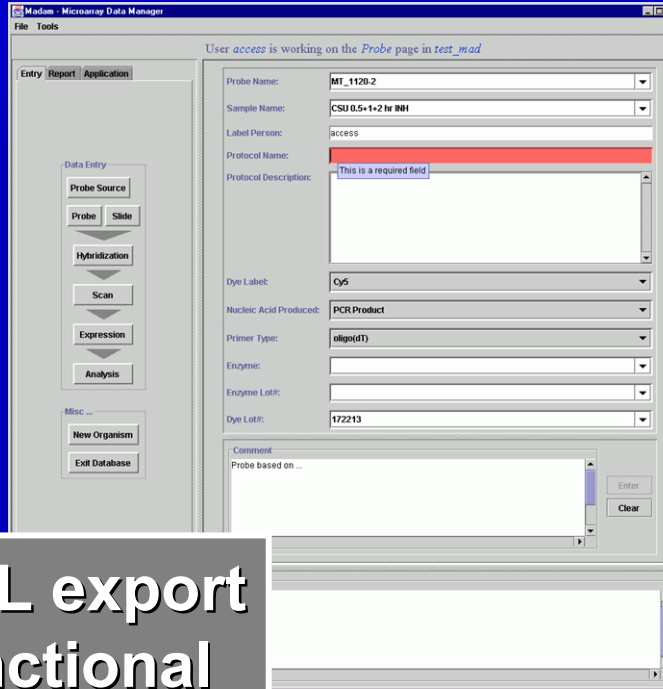
Recently added: QTL-based searches



Tools for Array Analysis

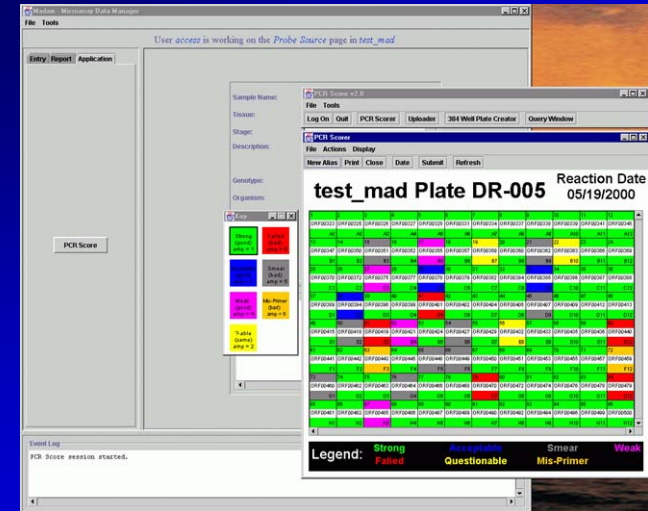
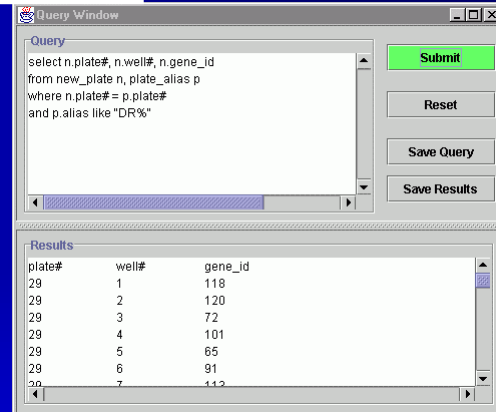


MADAM: Microarray Data Manager



**MAGE-ML export
now functional**

**Joseph White
Jerry Li
Alexander Saeed
Vasily Sharov
Syntek Inc.**



Available with OSI source and MySQL



MIDAS: Data Analysis

Wei Liang

Work Flow

Parameters

Parameter	Value
Slice Data Population	500
Data Keep Range >	+/-2.00 SD
Data Keep Range <	-----

Tools

View all parameters

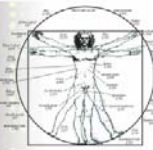
Processing Status

- Performing Intensity Filter ... Done!
- Performing Locfit(Lowess) normalization ... Done!
- Performing Flip-Dye Consistency Checking ... Done!
- Performing Slice Analysis ... Done!
- Writing output file(s) ... Done!

Process finished.

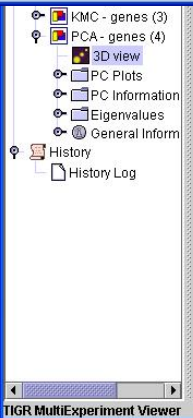
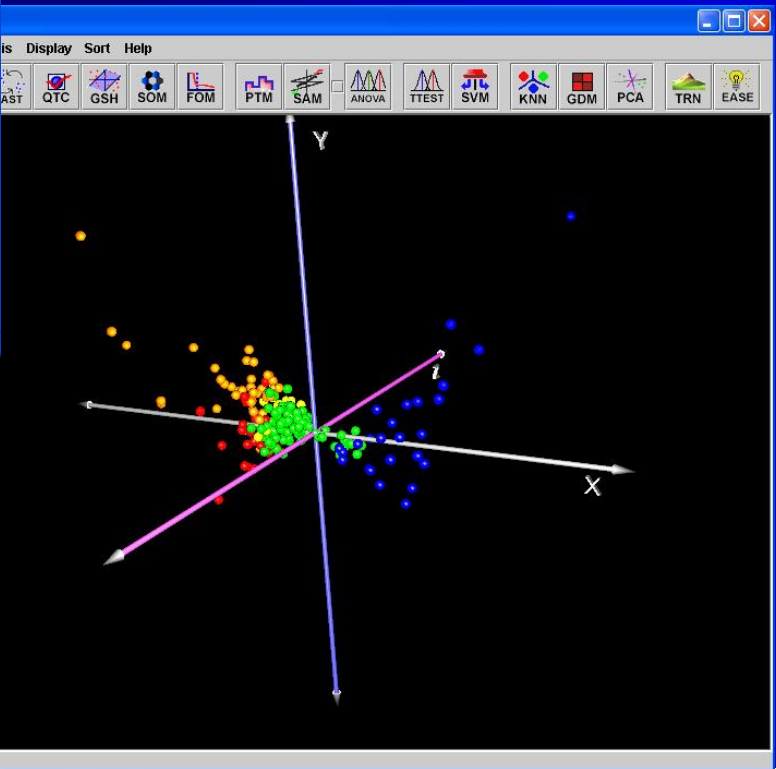
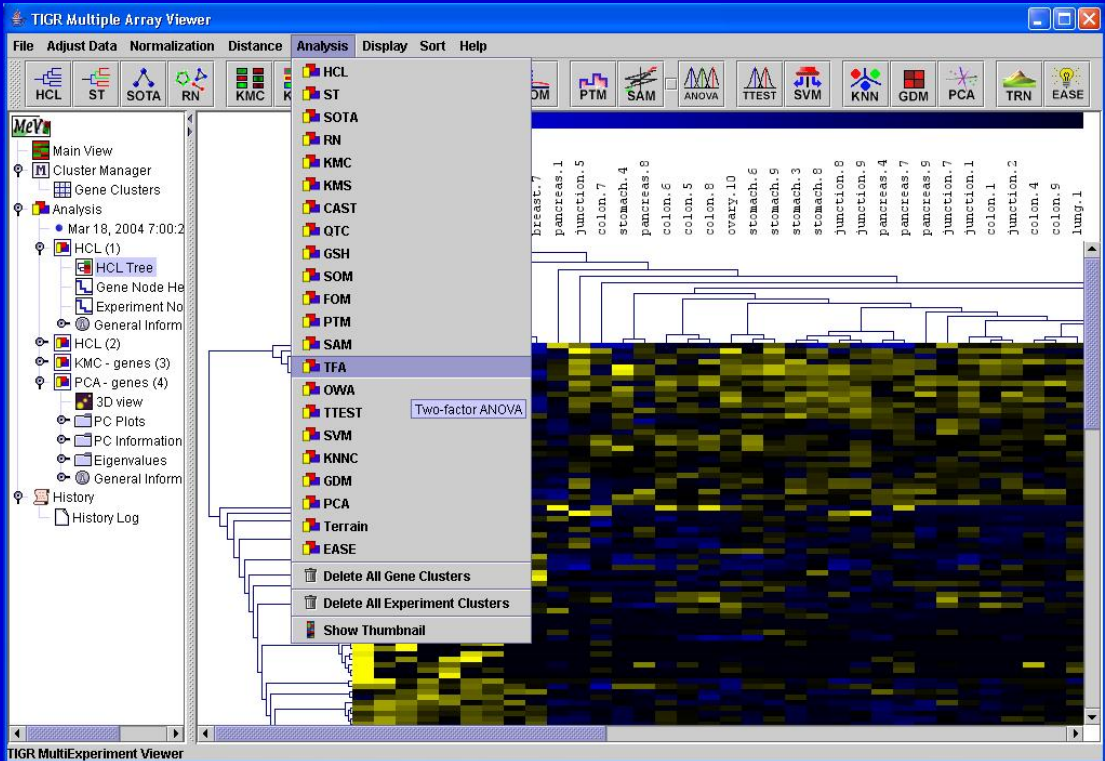
TIGR MIDAS

Variance Stabilization,
Adding Error Models,
MAANOVA,
Automated Reporting

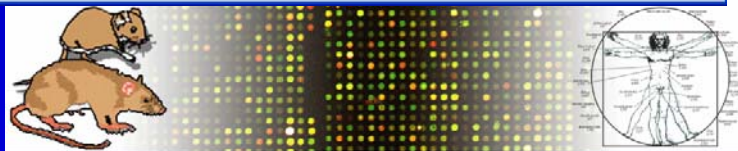


MeV: Data Mining Tools

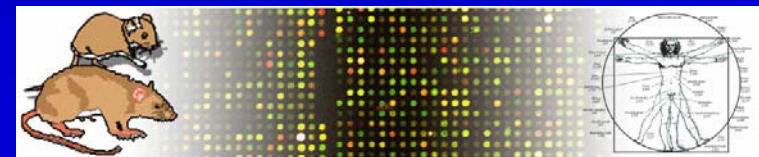
Alexander Saeed
Alexander Sturn
Nirmal Bhagabati
John Braisted
Syntek Inc.
Datanaut, Inc.



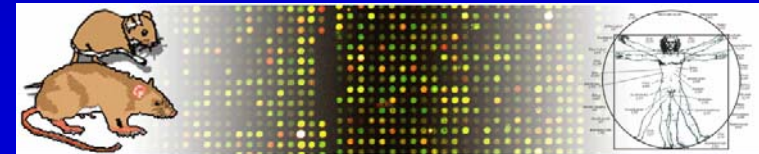
Available with OSI source



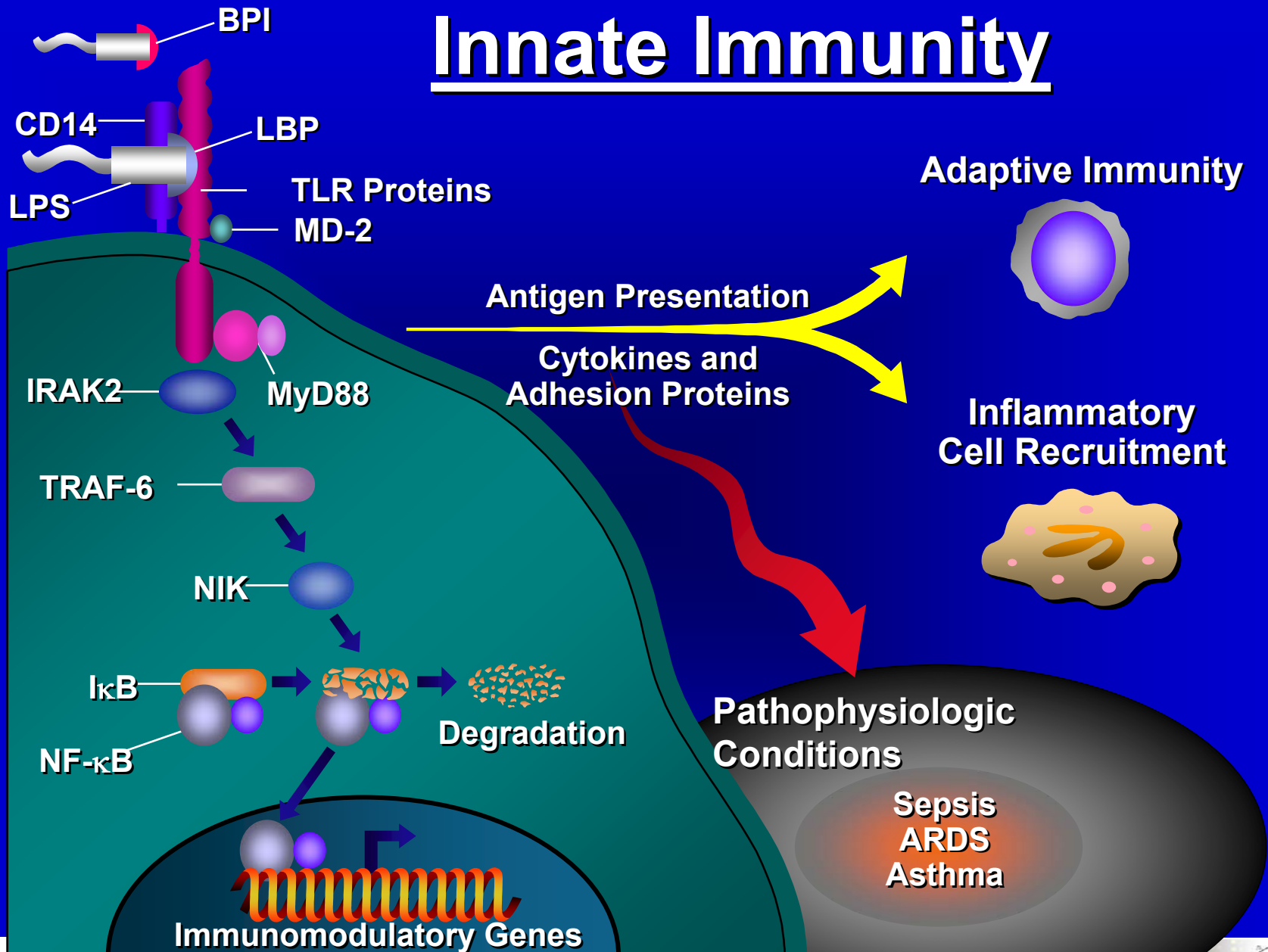
Science



Integrating Expression with other data



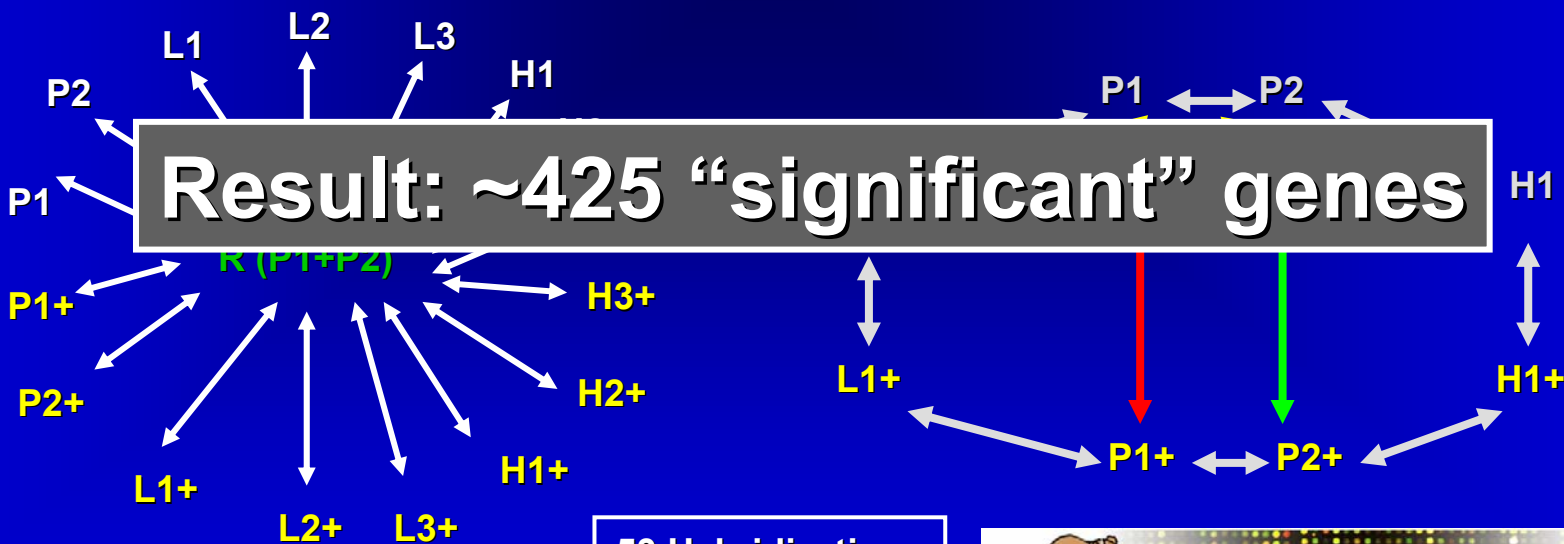
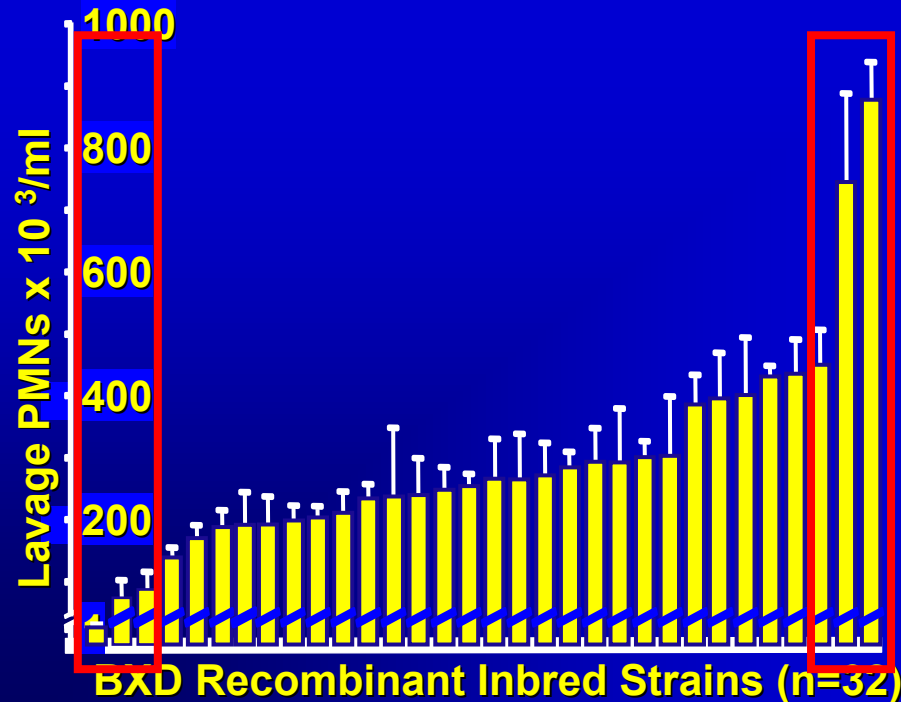
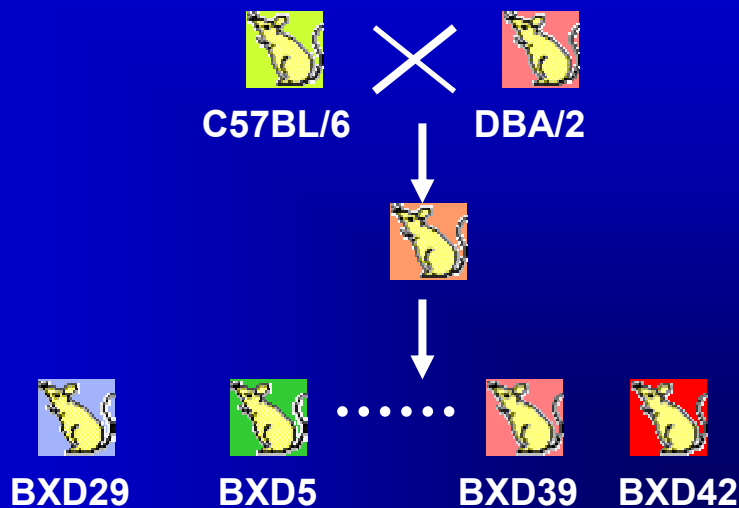
Innate Immunity



David Schwartz

Adapted from Godowski. *NEJM* 1999; 340:1835

Examples

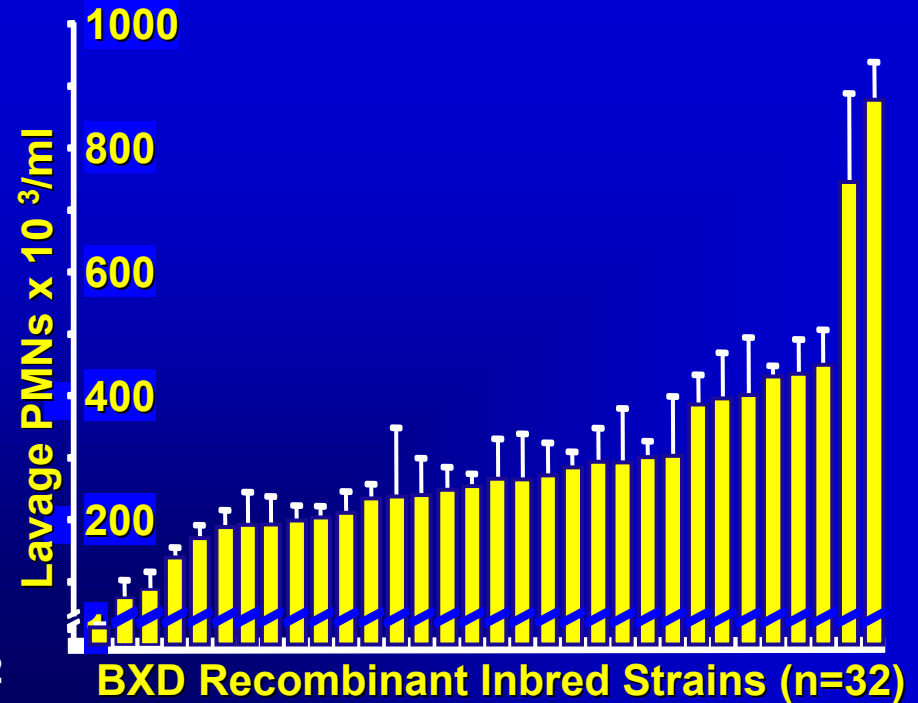
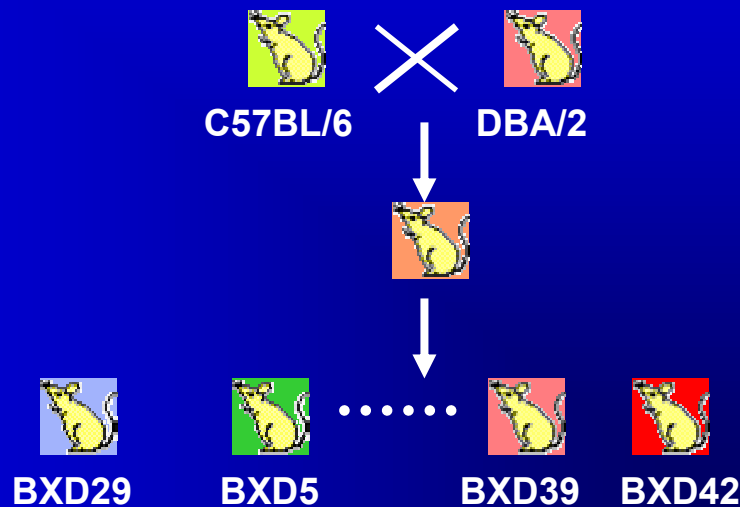


53 Hybridizations



Don Cook, Shuibang Wang, David Schwartz

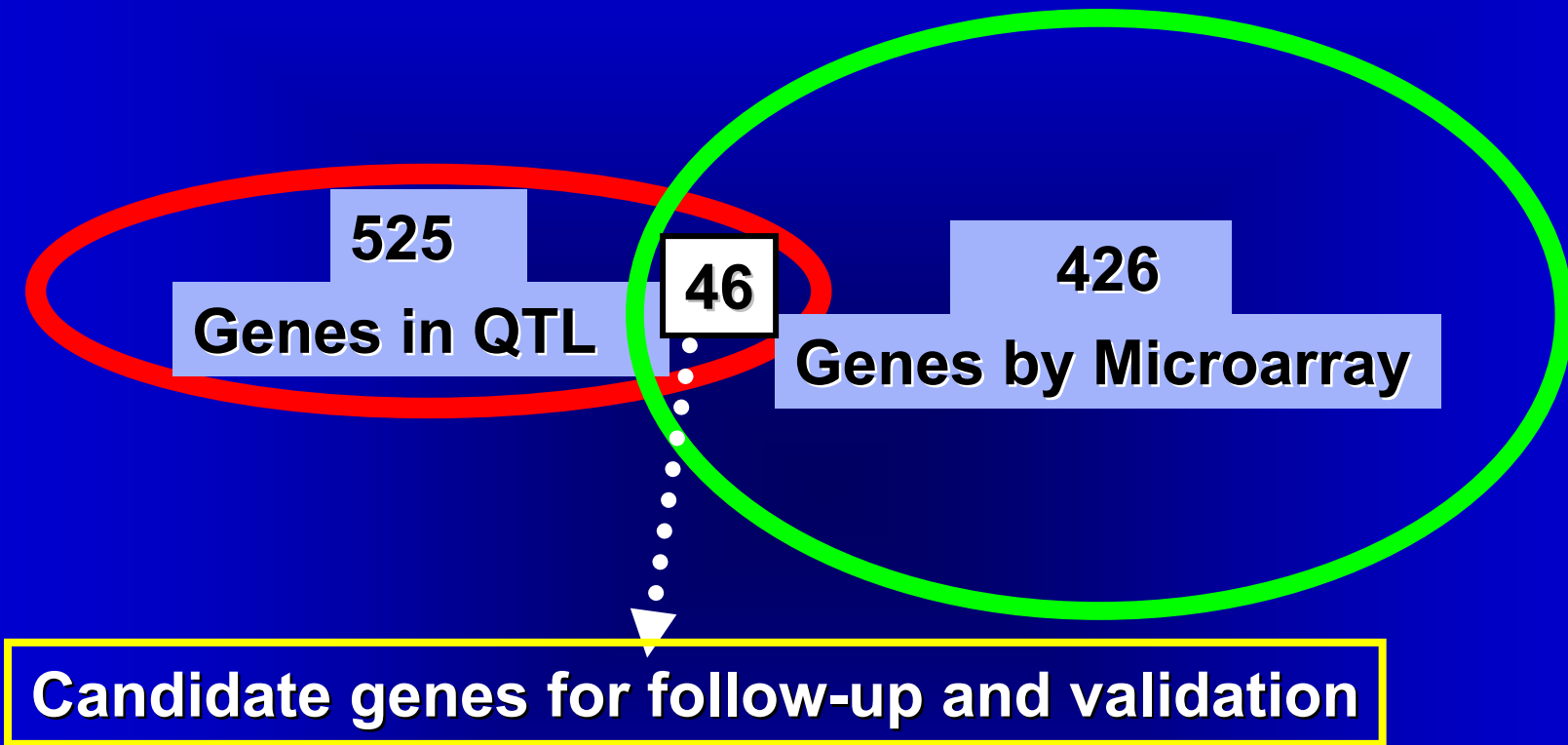
Examples



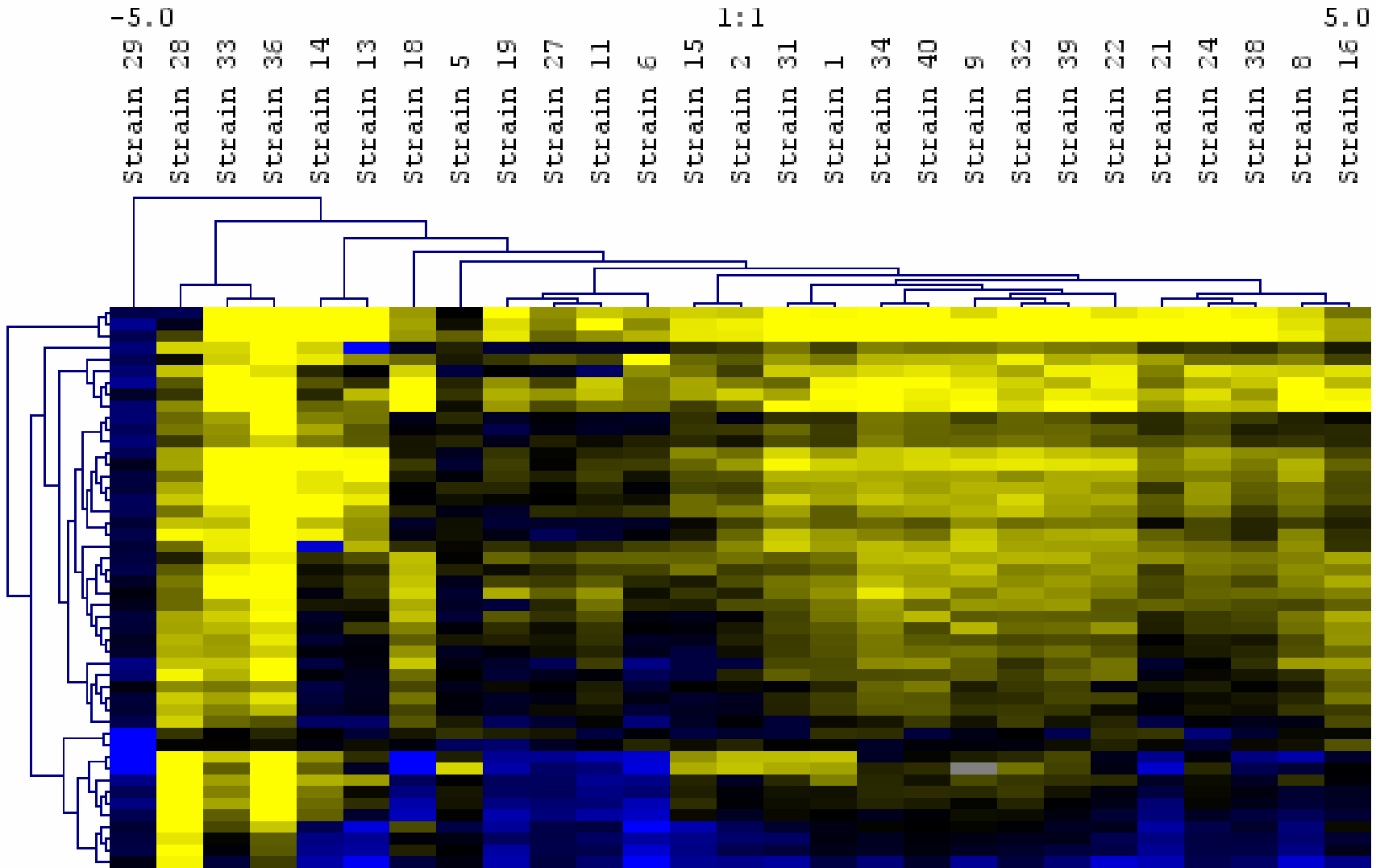
IDEA: Build QTL Maps and use those to filter expression data

Goal: Find differentially expressed genes genetically linked to response

Microarray Expression-QTL Consensus Candidate Genes



Candidate Gene Set for LPS response



Don Cook, Bryan Frank, David Schwartz

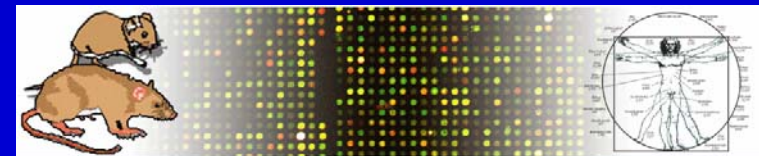
What have we learned?

■ Well...

- **Genetics and Expression can be powerful when used together**
- **This yields genes that are differentially expressed but also genetically linked to traits**
- **The expression fingerprint itself can also be used as a quantitative trait - eQTLs**

■ But...

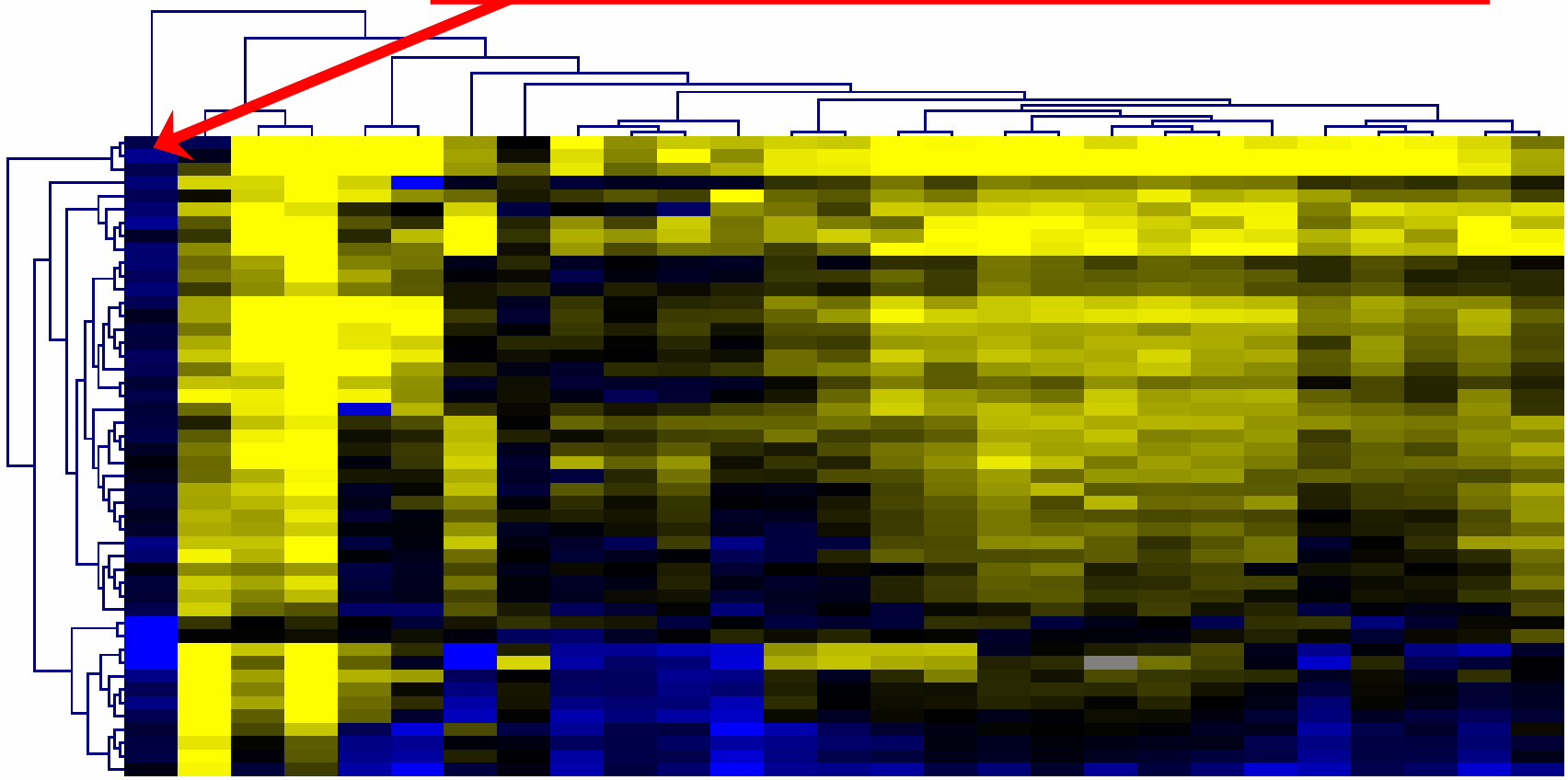
- **QTL+Expression may miss regulatory and other genes whose expression does not change but which contain mutations that are causative**



Candidate Gene Set for LPS response

Strain 29 5.0
Strain 28
Strain 33
Strain 36
Strain 14
Strain 13
Strain 18
5
19
27
11
6
15
2
31
1
34
40
9
32
39
22
21
24
38
8
Strain 16 5.0

BXD29 has spontaneous *tlr 4* mutation



Sleep Deprivation Studies in Mouse

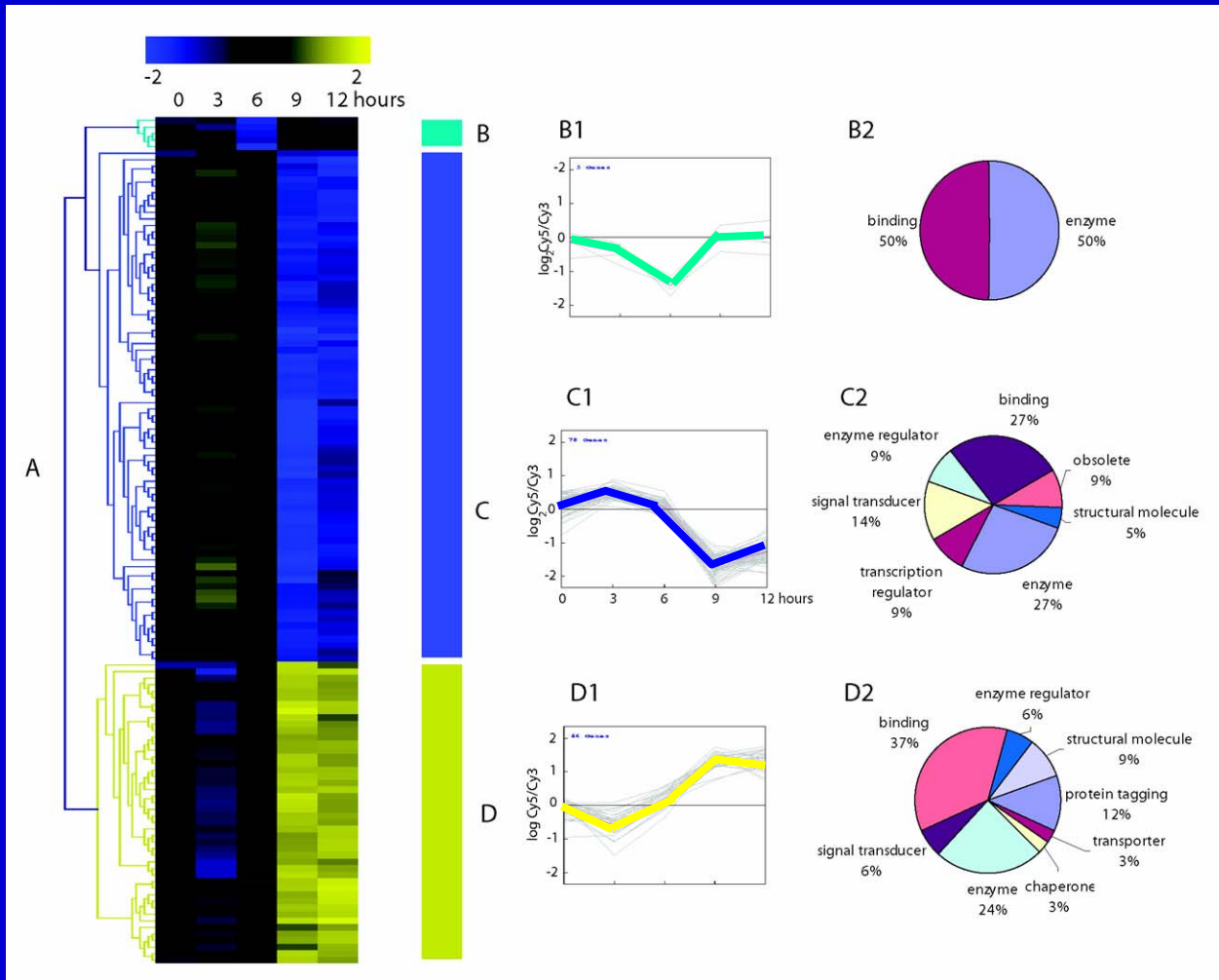


Experimental Paradigm

- Compare gene expression between sleeping and sleep-deprived mice in cortex and hypothalamus
- Perform 3 biological replicates
- Normalize and filter data and use data mining techniques to select distinct patterns of gene expression
- Use Gene Ontology (GO) assignments to classify genes by cellular localization, molecular function, biological process
- Use GO analysis to develop an understanding of response



Differential Expression in Hypothalamus



Sleep signaling

EASE Analysis of GO terms

Cortex – Up-regulated Genes

<i>GO Class</i>	<i>GO Category</i>	<i>p-value</i>
GO Cellular Component	endoplasmic reticulum	6.06×10^{-03}
GO Molecular Function	heat shock protein activity	8.78×10^{-04}
	pyruvate dehydrogenase (lipoamide) phosphatase activity	3.17×10^{-03}
	chaperone activity	7.38×10^{-03}

Themes:

General biological trends based on representation of functional roles on the array

Problem:

Requirement of functional class assignment limits utility for discovery of new functional networks

	mitochondrial inner membrane	3.70×10^{-03}
GO Molecular Function	structural constituent of ribosome	6.46×10^{-39}
	RNA binding activity	4.83×10^{-21}
	cytochrome c oxidase activity	9.79×10^{-04}
	hydrogen ion transporter activity	1.88×10^{-03}

Thanks to Doug Hosack and Glynn Dennis, NIAID

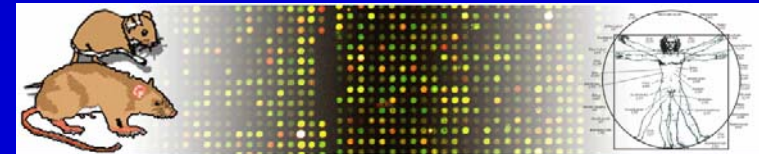
What have we learned?

■ Well...

- **Functional assignments provide a powerful filter on the measured expression**

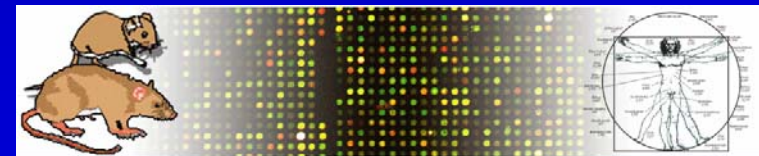
■ But...

- **Functional classes are not functions**
- **Significant additional work is necessary to translate these classes to pathways and responses**

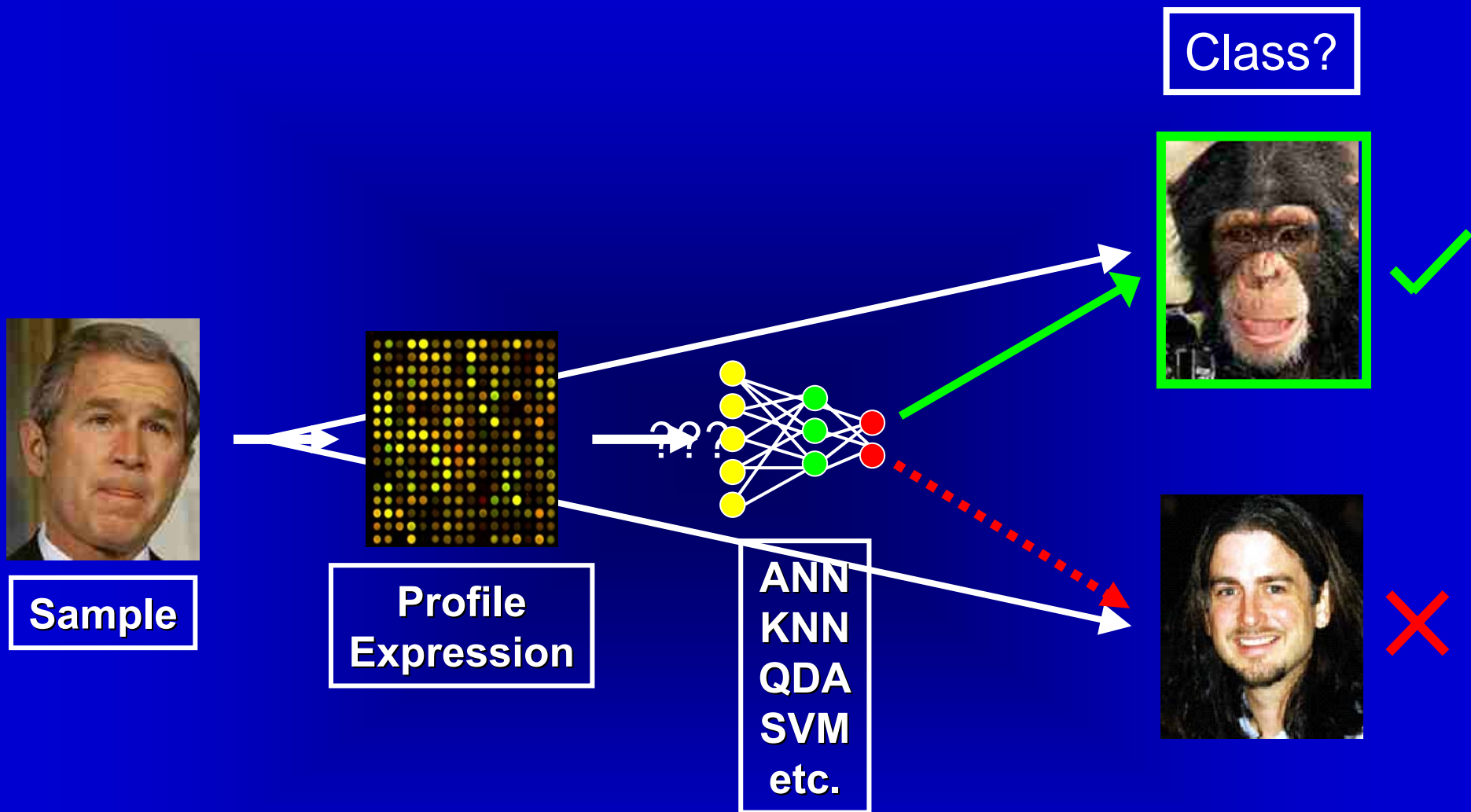


Predicting Outcome

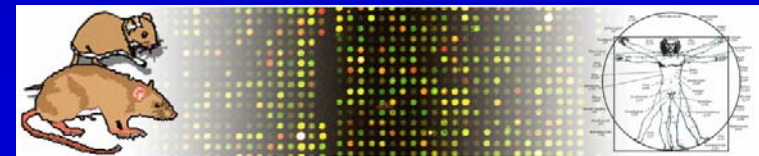
work in collaboration with
Timothy J. Yeatman
H. Lee Moffitt Cancer Center



The Classification Problem

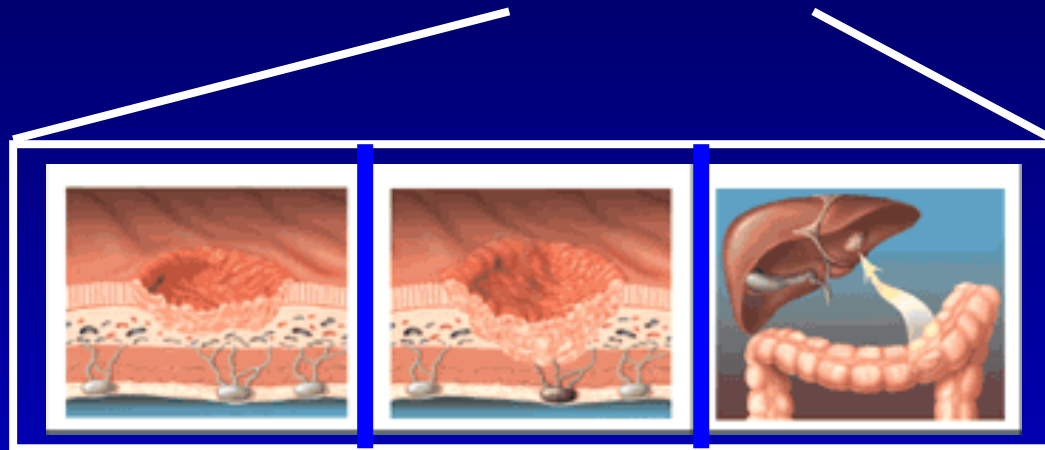


RNA and Protein Correlations



Colon Cancer Progression

Normal → Adenoma → Primary Tumor → Metastasis



Dukes' B

Dukes' C

Dukes' D

5 year survival

70-85%

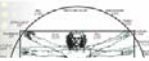
25-60%

5%



Experimental Design

- 6 groups: normal, adenoma, Dukes B, Dukes C, Dukes D, metastases
- 10 samples per group
- **Microarray:**
 - 2 color 32k array
 - Samples hybridized to common reference
 - Dye-reversal replicates performed
 - Total RNA, no amplification
- **Microarray:**
 - 2D-PAGE followed by mass-spectrometry



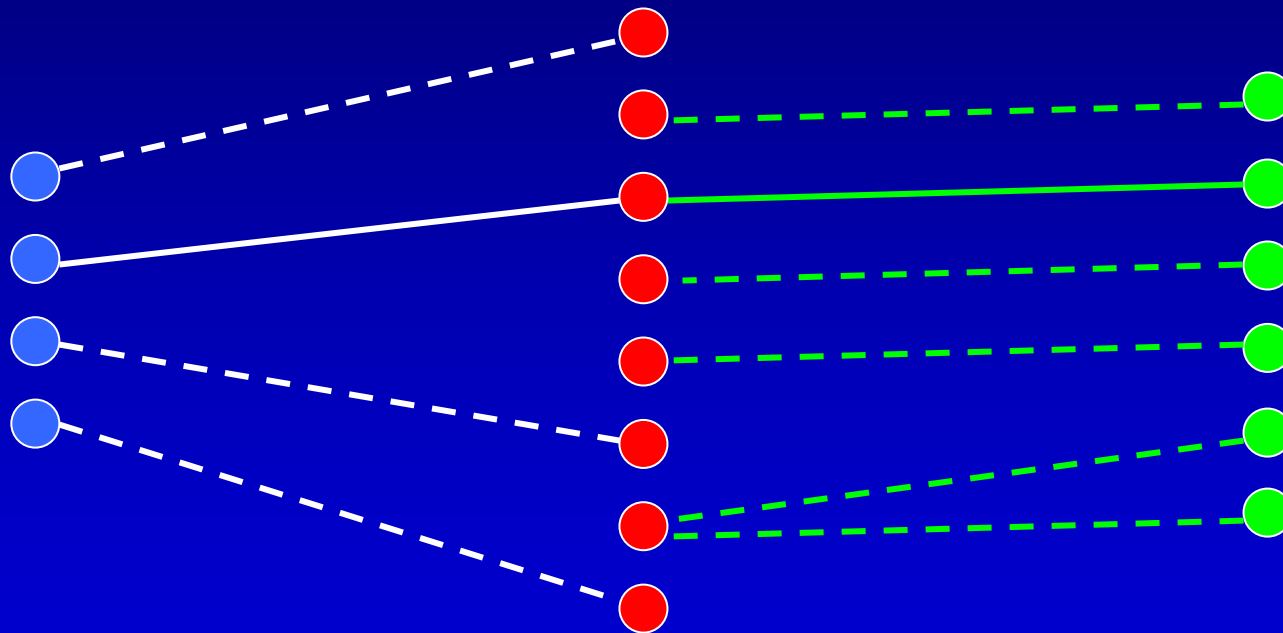
Linking genes to proteins

Array Probes are cDNAs, not genes

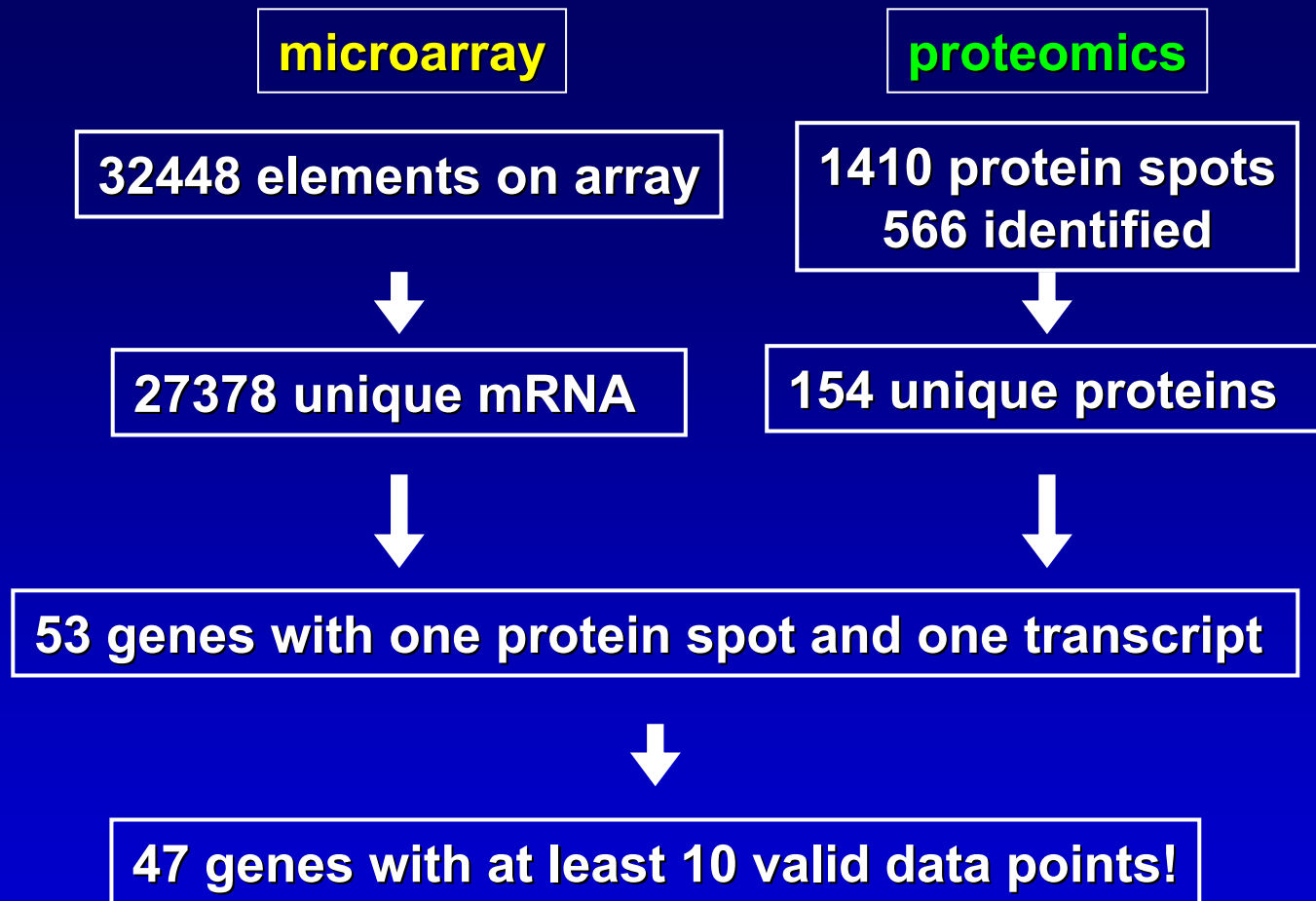
Protein Sequence
SwissProt

TIGR
Tentative Consensus (TC)

cDNA Sequence
Genbank



Linking genes to proteins

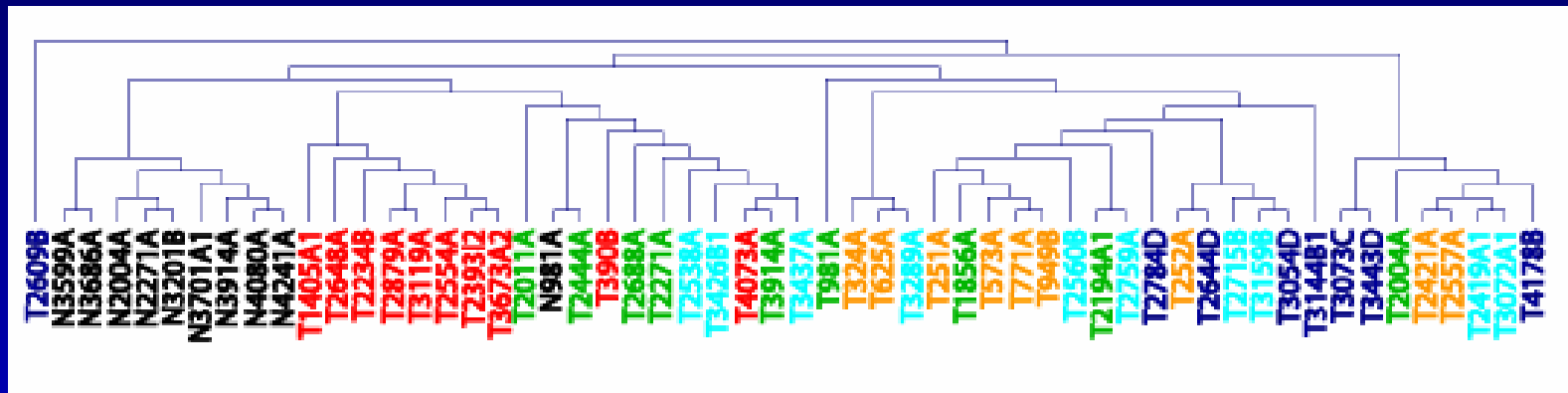


“Highly” correlated genes (12/47)

Description	r	p values
Proteasome activator complex subunit 2	0.682	2.53E-08
Metastasis inhibition factor nm23	0.579	8.14E-06
Proliferating cell nuclear antigen (PCNA)	0.567	1.08E-05
Proliferation-associated protein 2G4	0.558	1.88E-05
Fatty acid-binding protein, liver (L-FABP)	0.797	3.79E-05
Peroxiredoxin 4	0.517	7.18E-05
Creatine kinase, B chain	0.545	9.40E-05
Mitochondrial 3-hydroxyisobutyrate dehydrogenase	0.808	2.35E-04
Cadherin-17 precursor	0.516	2.09E-03
Argininosuccinate lyase	0.408	3.29E-03
F-actin capping protein alpha-2 subunit	0.336	0.01
Mitochondrial acyl-CoA dehydrogenase	0.487	0.03



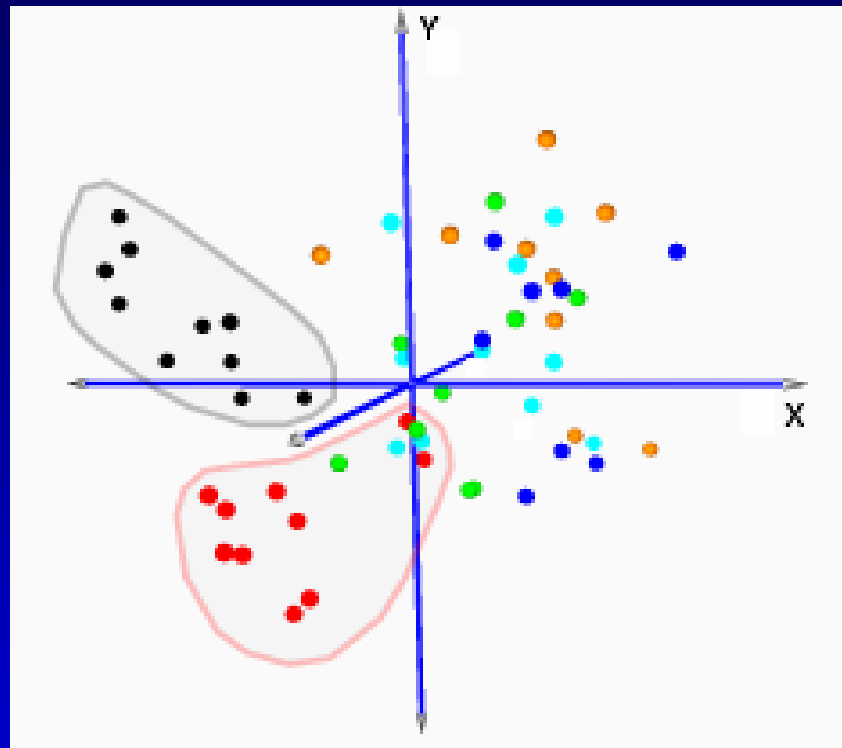
Hierarchical Clustering by mRNA expression



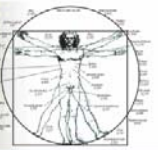
— normal — adenoma — Duke B — Duke C — Duke D — metastasis



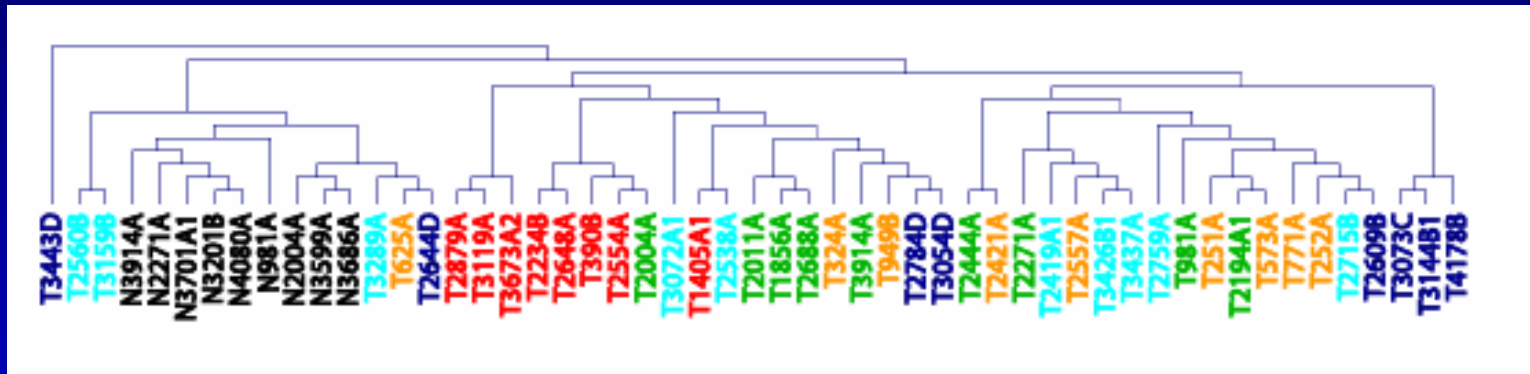
Microarray PCA



— normal — adenoma — Duke B — Duke C — Duke D — metastasis



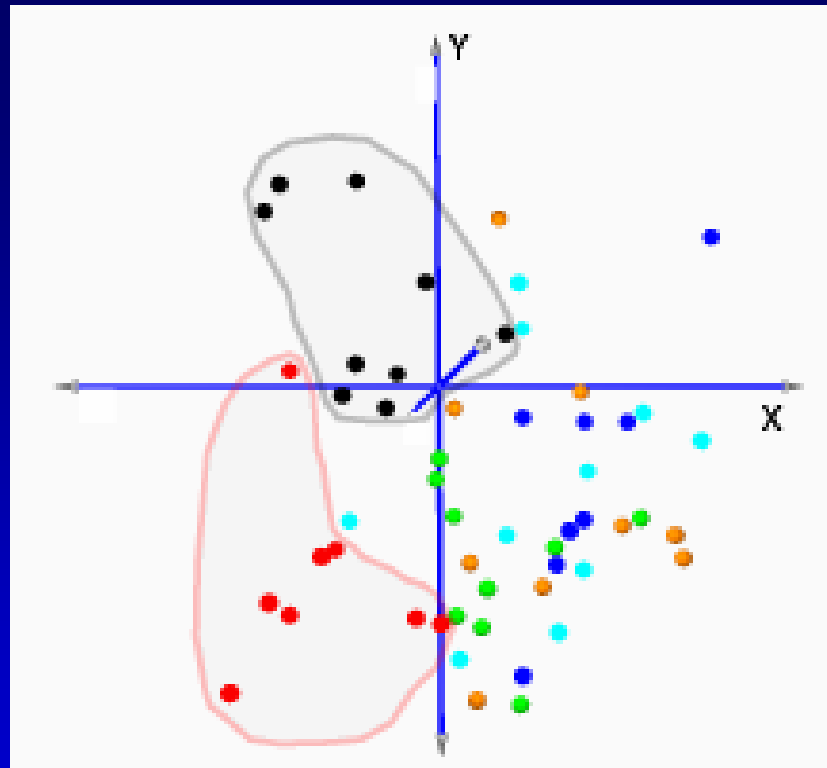
Hierarchical Clustering by protein expression



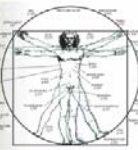
— normal — adenoma — Duke B — Duke C — Duke D — metastasis



Proteomics PCA



— normal — adenoma — Duke B — Duke C — Duke D — metastasis



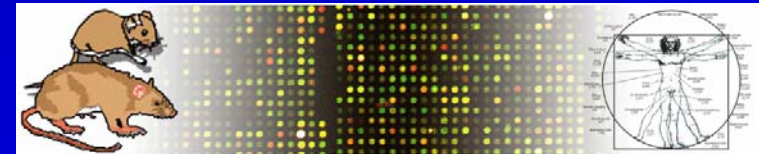
What have we learned?

- RNA and proteins do not correlate well
- There are numerous possible technical explanations for this – requires further analysis
- But RNA and protein paint the same fundamental picture of the disease – primary tumor fingerprints do not respect clinical staging classes



Predicting Survival

work in collaboration with
Timothy J. Yeatman
H. Lee Moffitt Cancer Center



Survival Classification Set

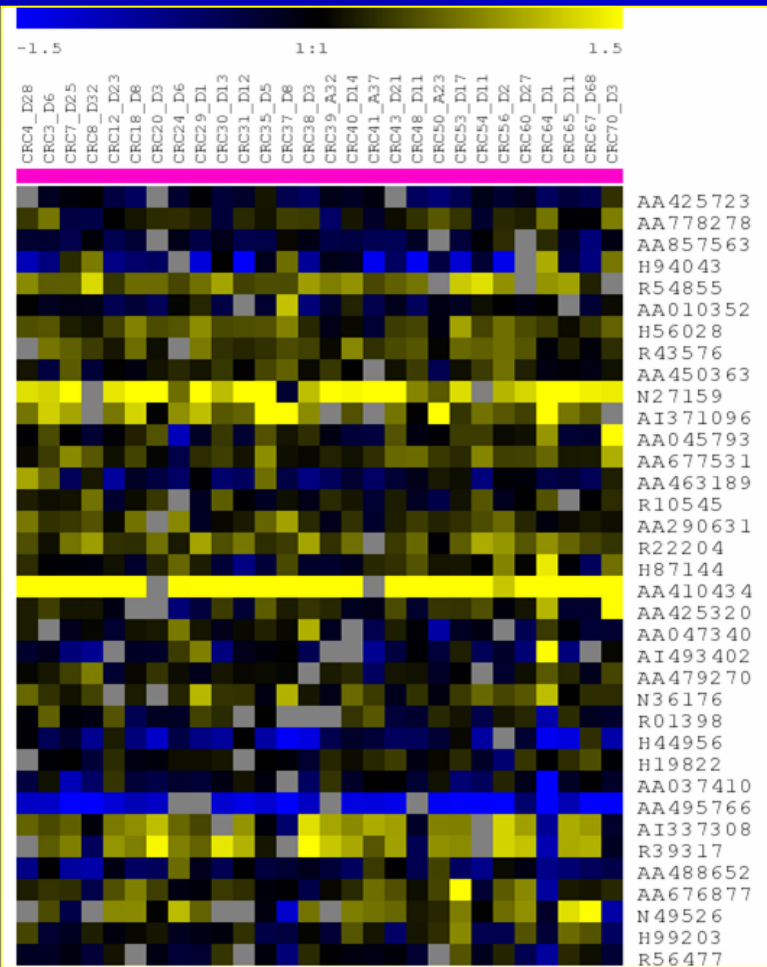
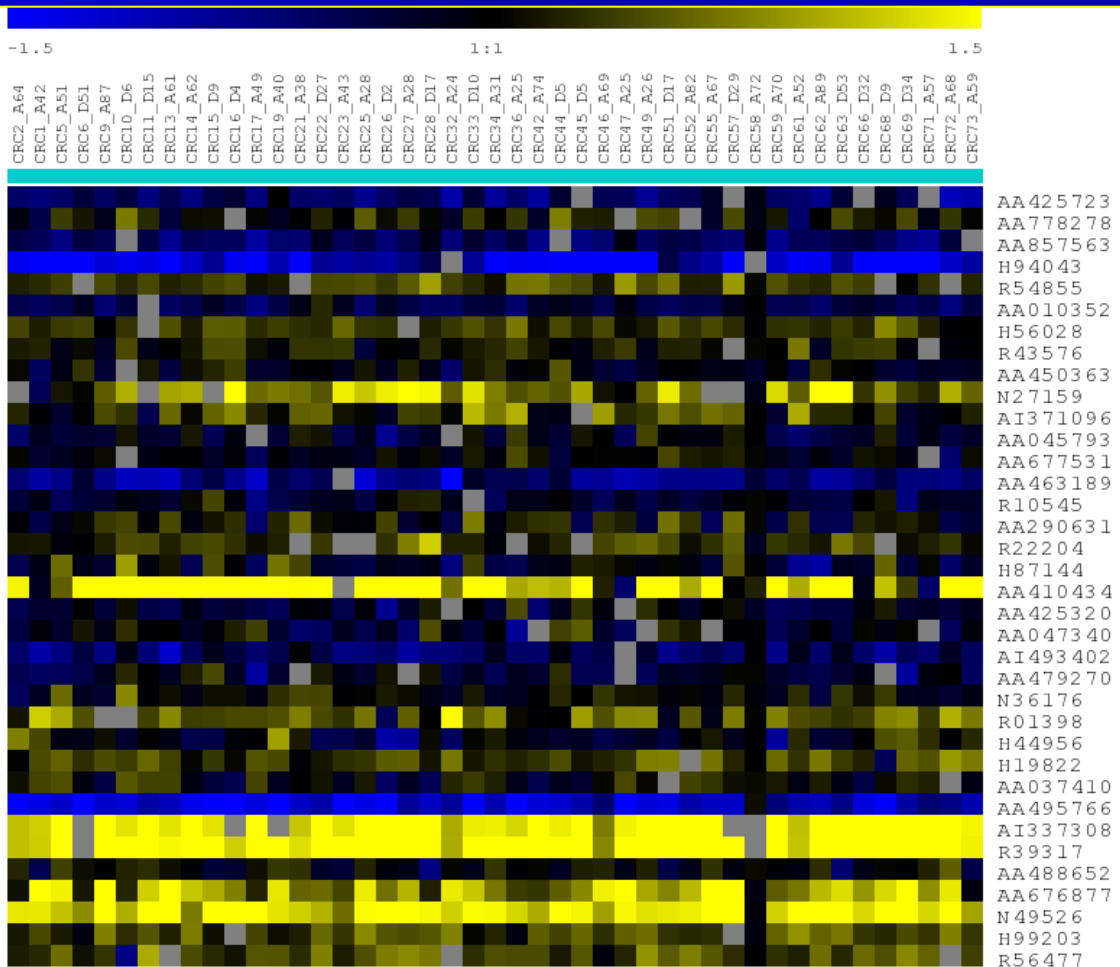
- Start with 73 Dukes' B and C tumors with survival data
- Profile expression on arrays
- Use SAM with Censored Survival Analysis to select significant genes
- Use unsupervised k-means support to separate tumors into 2 classes

SAM SIGNIFICANT GENES

Original row	Clone Name	GB#	TC#	Putative Role
23667	Image:488033	AA045793	THC915844	DnaJ homolog subfamily B member 9 (Microvascular endothelial differentiation gene-1 protein) (Mdc)
4697	Image:460395	AA677531	THC925057	Similar to hypothetical protein FLJ22625 {Homo sapiens}
11029	Image:128457	R10545		
17380	Image:130854	R22204	THC987381	
5152	Image:242797	H94043	THC1013186	yes-associated protein homolog DKFZp586I1419.1 - human (fragment)
28052	Image:2043415	AI371096	THC1022305	Death-associated protein kinase 1 (EC 2.7.1.-) (DAP kinase 1). [Human] {Homo sapiens}
9637	Image:700461	AA290631	THC918268	unnamed protein product {Homo sapiens}
25310	Image:796878	AA463189	THC942191	brain-specific GTP-binding protein {Homo sapiens}
26499	Image:269815	N27159	THC899426	Inhibin beta A chain precursor (Activin beta-A chain) (Erythroid differentiation protein) (EDF).
6925	Image:773278	AA425320	THC915844	DnaJ homolog subfamily B member 9 (Microvascular endothelial differentiation gene-1 protein) (Mdc)
14102	Image:252453	H87144	THC958545	CGI-86 protein {Homo sapiens}
6146	Image:509516	AA047340	THC908165	unnamed protein product {Homo sapiens}
29198	Image:2115602	AI493402	THC1014756	transcription factor SL1 - human
28743	Image:754250	AA479270	THC862179	KIAA1253 protein {Homo sapiens}
30718	Image:272694	N36176	THC960848	hypothetical protein {Homo sapiens}
10497	Image:753428	AA410434	THC986960	Similar to RIKEN cDNA 1110014B07 gene {Homo sapiens}
1922	Image:123742	R01398	THC1024446	unknown {Homo sapiens}
18930	Image:183200	H44956	THC863066	Fumarylacetoacetase (EC 3.7.1.2) (Fumarylacetoacetate hydrolase) (Beta-diketonase) (FAA). [Human]
20064	Image:24067	R39317	THC987467	putative {Mus musculus}
24682	Image:172495	H19822	THC986991	hypothetical protein {Homo sapiens}
29976	Image:2062345	AI337308	THC987466	protein-tyrosine kinase EPHB2v {Homo sapiens}
27256	Image:897107	AA676877	THC1021656	mitochondrial citrate transport protein {Homo sapiens}
28497	Image:768316	AA495766	THC925935	RCC1-like G exchanging factor RLG [imported] - human
30633	Image:321271	AA037410	THC864588	rho GTPase activating protein 8 isoform 1 {Homo sapiens}
23089	Image:243549	N49526	THC1031960	Myb proto-oncogene protein (C-myb). [Human] {Homo sapiens}
21529	Image:741474	AA401111	THC897309	Glucose-6-phosphate isomerase (EC 5.3.1.9) (GPI) (Phosphoglucose isomerase) (PGI)
9651	Image:416280			
20264	Image:843263	AA488652	THC862216	ribosomal protein L2 {Homo sapiens}
3184	Image:509458	AA056375	THC899187	
17742	Image:261829	H99203	THC968032	Ubiquitin carboxyl-terminal hydrolase 7 (EC 3.1.2.15) (Ubiquitin thiolesterase 7)
30132	Image:41170	R56477	THC911668	
19840	Image:178818	H49455	THC882357	Apical-like protein (APXL protein). [Human] {Homo sapiens}
26573	Image:179276	H50323	THC862010	fatty-acid synthase (EC 2.3.1.85) (version 2) - human
26161	Image:810133	AA464251	THC950559	

Red are positive and green are negative genes. A positive score means that higher expression is associated with higher risk, ie, shorter survival!!!

Clustering based on *k*-means support

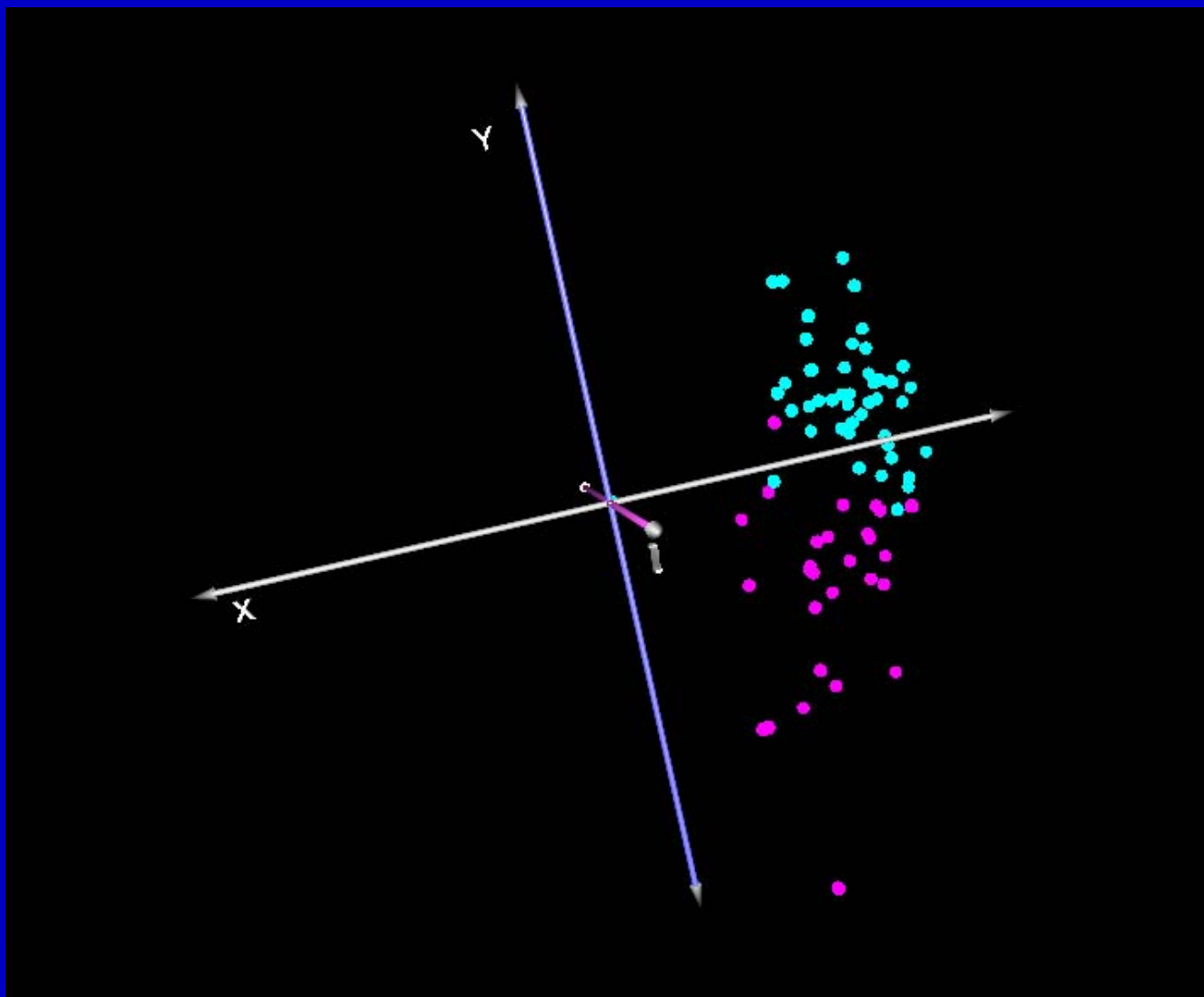


Good Prognosis

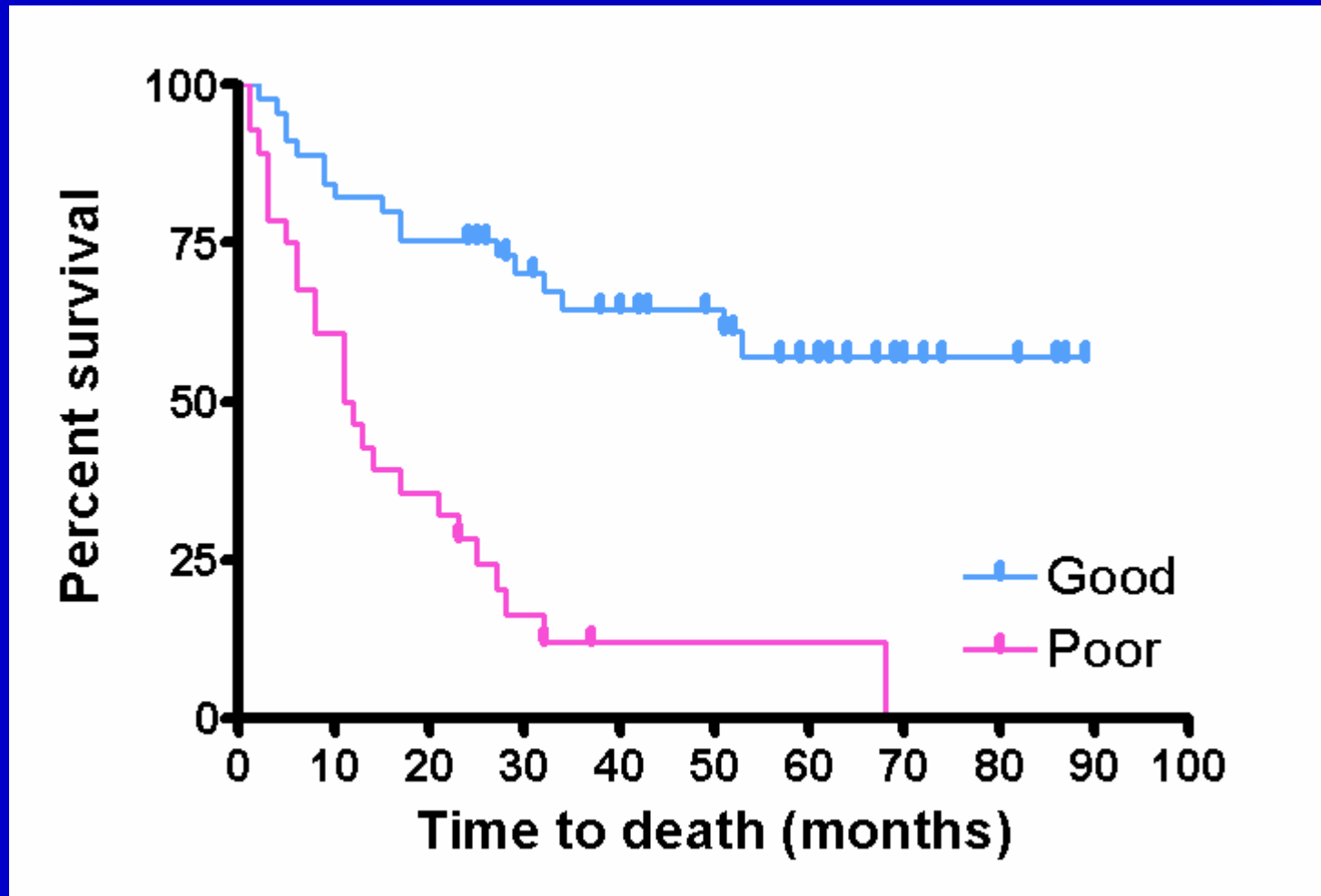
Poor Prognosis



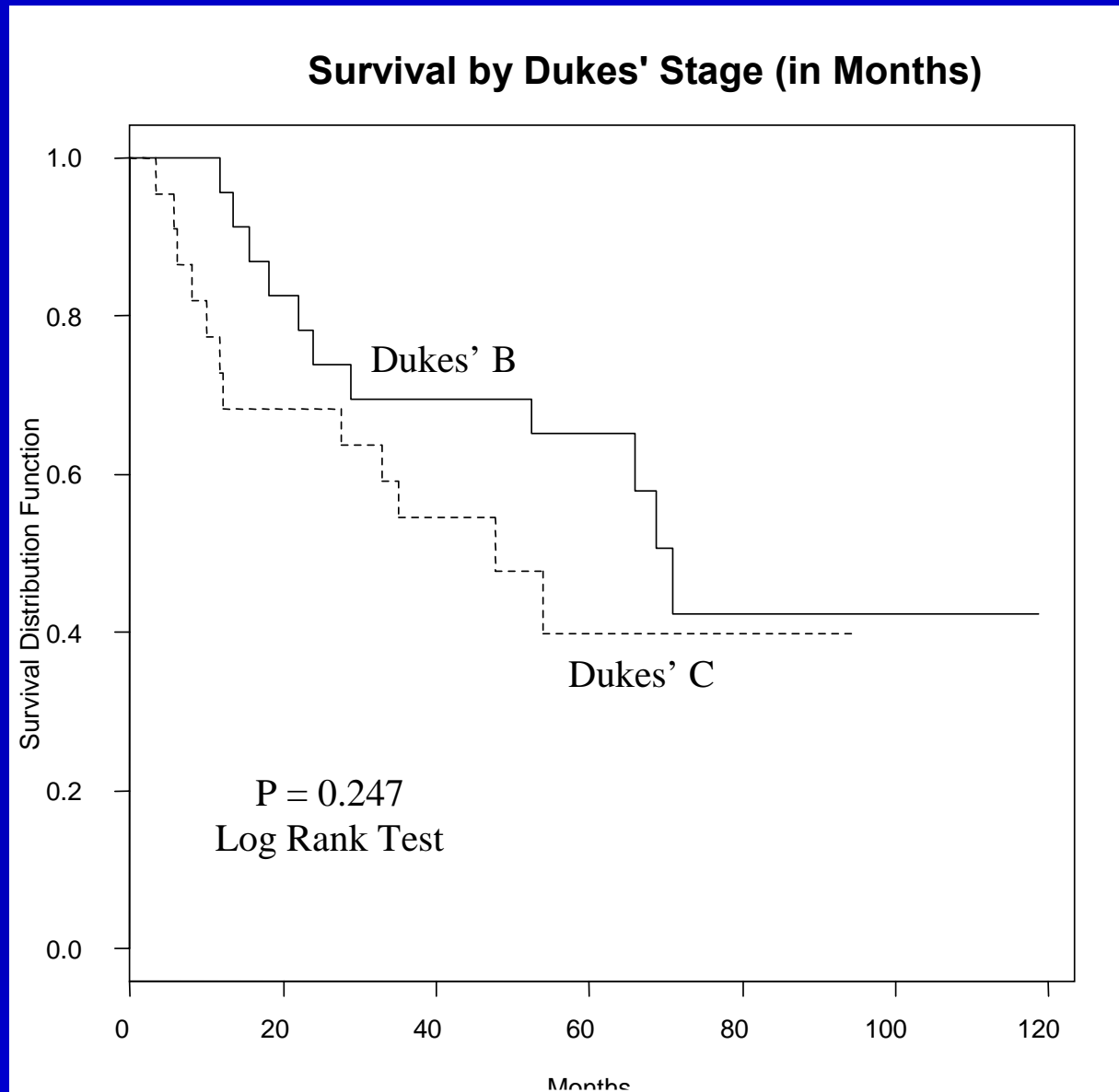
PCA analysis of survival samples



Survival Classes based on *k*-means clustering



Survival based on Clinical Stage

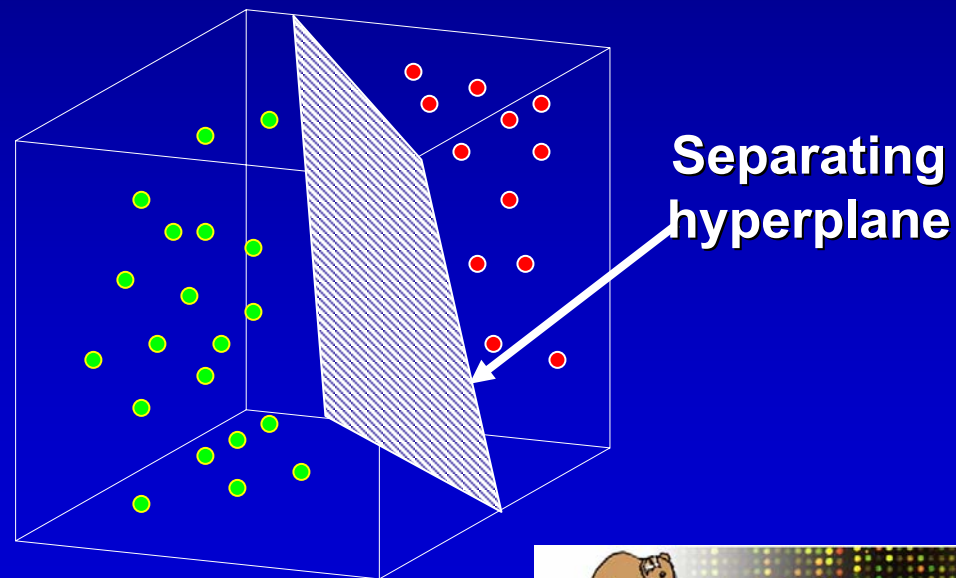


The Challenge

- Supervised gene selection followed by unsupervised clustering works quite well
- We anticipate that supervised approaches would perform better
- Use *complete* leave one out cross-validation to classify samples and to identify a core set of genes for classification
 - Iteratively leave out one sample and redo gene selection, algorithm training, and testing

SVM Classification

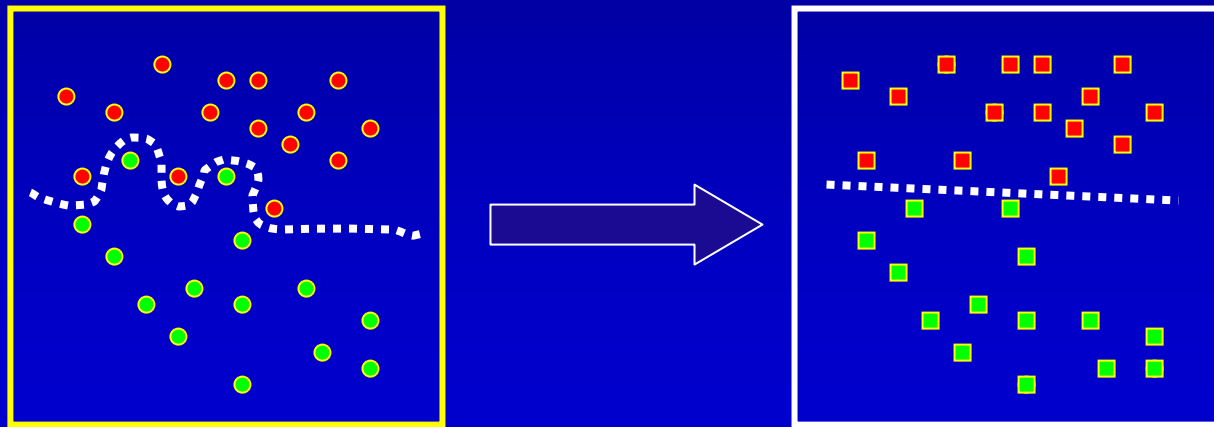
- SVM attempts to find an optimal separating hyperplane between members of the two initial classifications.



SVM Kernel Construction

The expression data can be transformed to a higher dimensional space (feature space) by applying a kernel function.

This transformation can have the effect of allowing a “separating hyperplane” to be found.

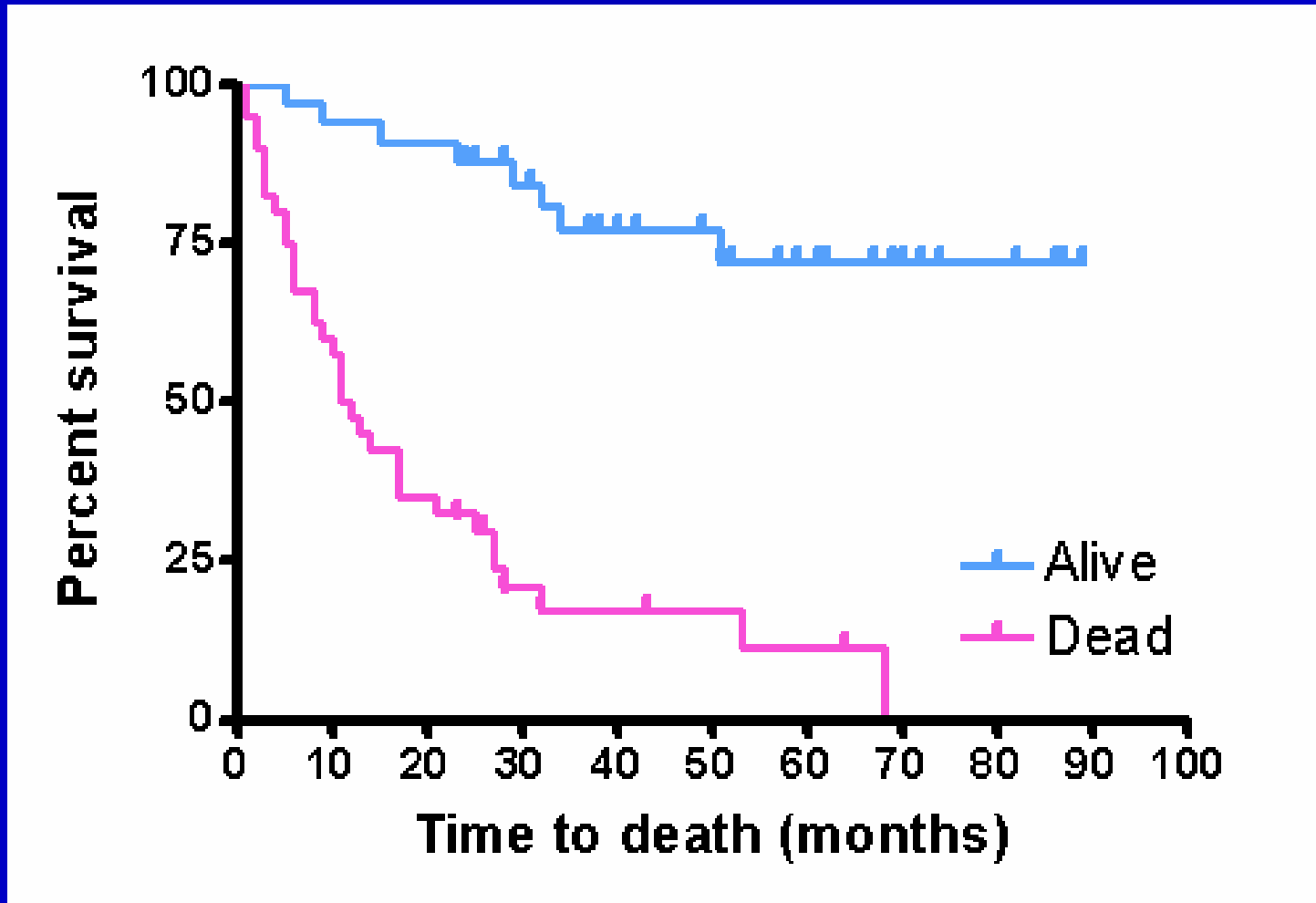


Consensus genes

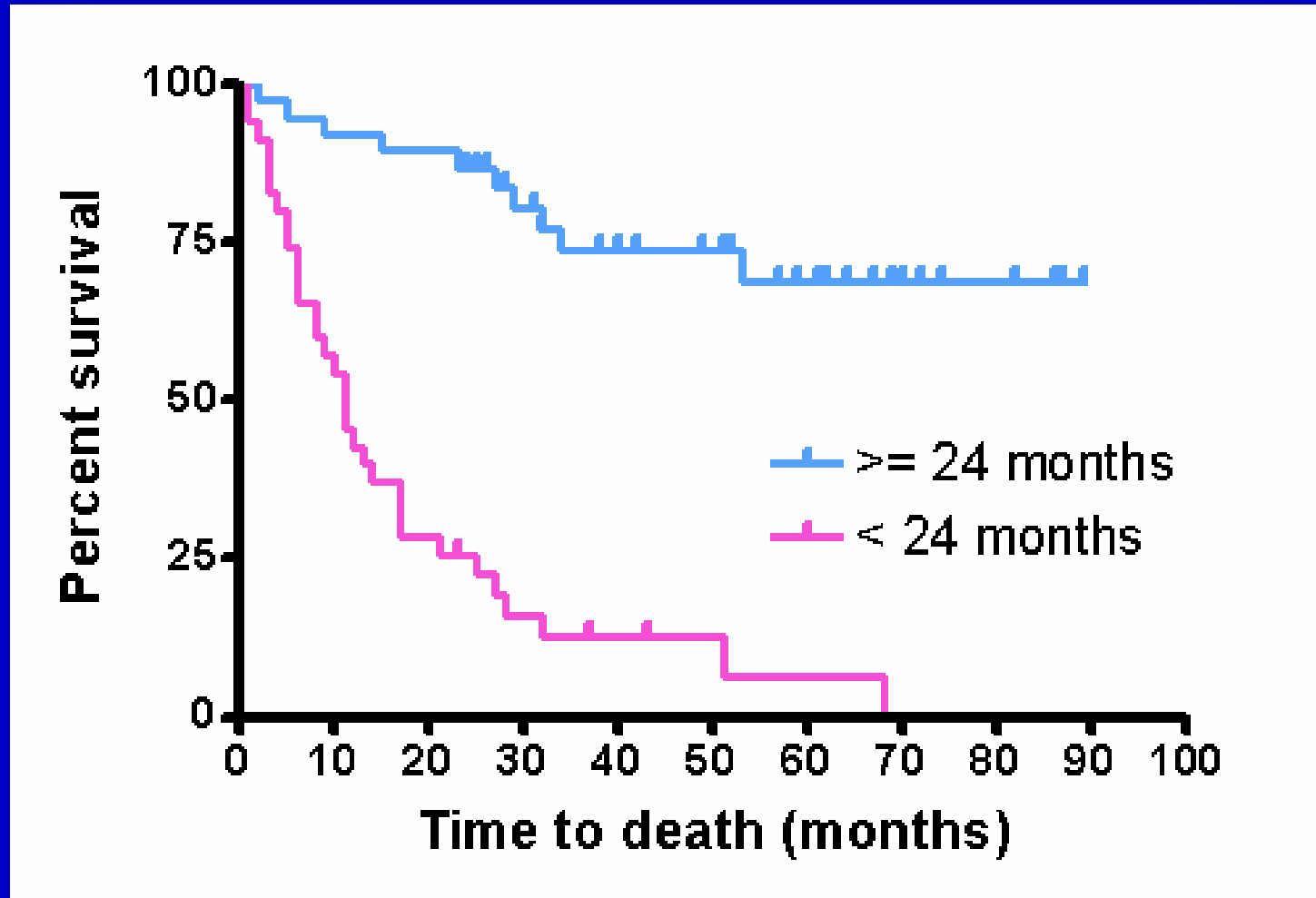
<u>GB#</u>	<u>TIGR TC#</u>	<u>Putative Role</u>
R01398	THC1024446	unknown {Homo sapiens}
H44956	THC863066	Fumarylacetoacetase (EC 3.7.1.2) (Fumarylacetoacetate hydrolase) (Beta-diketonase) (FAA).
H19822	THC986991	hypothetical protein {Homo sapiens}
AA037410	THC864588	rho GTPase activating protein 8 isoform 1 {Homo sapiens}
AA495766	THC925935	RCC1-like G exchanging factor RLG [imported] - human
AI337308	THC987466	protein-tyrosine kinase EPHB2v {Homo sapiens}
R39317	THC987467	putative {Mus musculus}
AA488652	THC862216	ribosomal protein L2 {Homo sapiens}
R56477	THC911668	
N49526	THC1031960	Myb proto-oncogene protein (C-myb). [Human] {Homo sapiens}
H99203	THC968032	Ubiquitin carboxyl-terminal hydrolase 7 (EC 3.1.2.15) (Ubiquitin thiolesterase 7)
AA401111	THC897309	Glucose-6-phosphate isomerase (EC 5.3.1.9) (GPI) (Phosphoglucose isomerase) (PGI)
AA448641	THC889589	transcription factor {Homo sapiens}



Survival Classes based on SVM – alive/dead



Survival Classes based on SVM – Two year cutoff

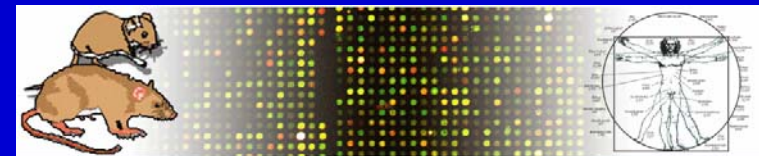


What have we learned?

- **Microarray fingerprints can provide clinically important clues about disease progression**
- **Many of these require additional validation**
- **Many key genes are not obviously related to the underlying biology**
- **More data are needed to validate these findings**
 - **This approach and the classification gene set have since been validated using an *independent* test set and an Affymetrix GeneChip™ data set**

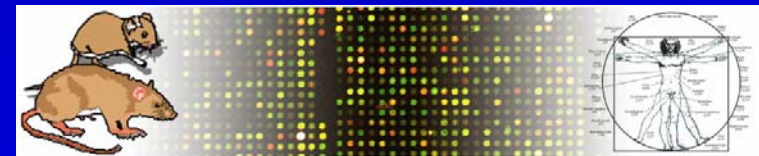
Where are we going?

- **There is still a role for biology!**
- **We are approaching a time in which we can begin to look at cells and organisms holistically.**
- **We also need to begin to think about integrating diverse data types in an intelligent way.**
- **This must include cross-species comparisons and inclusion of environmental effects.**
- **We may soon be in a position to begin development of a theoretical biology.**
- **Theoretical biology will require a transition from a Deterministic to a Stochastic approach.**



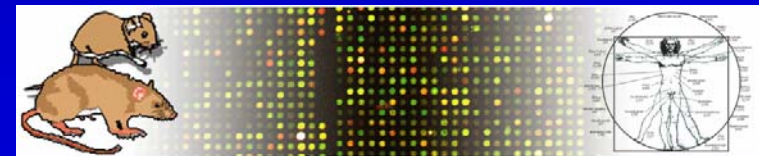
A theory has only the possibility of being right or wrong. A model has a third possibility; it may be right but irrelevant.

– Manfred Eigen



**Nobody in the game of football
should be called a genius.
A genius is somebody like Norman Einstein.**

**- Joe Theisman,
Former Washington Redskins quarterback**



Acknowledgments

[<johnq@tigr.org>](mailto:johnq@tigr.org)

The TIGR Gene Index Team

Foo Cheung
Svetlana Karamycheva
Yudan Lee
Babak Parvizi
Geo Pertea
John Quackenbush
Razvan Sultana
Jennifer Tsai
Joseph White

Emeritus

Jennifer Cho (TGI)
Emily Chen (μ A)
Ingeborg Holt (TGI)
Feng Liang (TGI)
Kristie Abernathy (μ A)
Sonia Dharap (μ A)
Julie Earle-Hughes (μ A)
Cheryl Gay (μ A)
Jeremy Hasseman (μ A)
Priti Hegde (μ A)
Heenam Kim (μ A)
Lara Linford (μ A)
Rong Qi (μ A)
Erik Snestrud (μ A)
Shuibang Want (μ A)
Ivana Yang (μ A)
Yan Yu (μ A)

H. Lee Moffitt Center/USF

Timothy J. Yeatman
Greg Bloom

PGA Collaborators

Gary Churchill (TJL)
Greg Evans (NHLBI)
Harry Gavras (BU)
Howard Jacob (MCW)
Anne Kwitek (MCW)
Allan Pack (Penn)
Beverly Paigen (TJL)
Luanne Peters (TJL)
David Schwartz (Duke)

TIGR PGA Collaborators

Norman Lee
Renaë Malek
Hong-Ying Wang
Truong Luu
Bobby Behbahani

Funding provided by the Department of Energy
and the National Science Foundation

Funding provided by the National Cancer Institute,
the National Heart, Lung, Blood Institute,
and the National Science Foundation

TIGR Faculty, IT Group, and Staff

TIGR Human/Mouse/Arabidopsis

Expression Team

Adriana Ahumada
Tove Andersson
Joanne Emerson
Bryan Frank
Molly Freeman
Renee Gaspard
Nadeeza Ishmael
Ka Yin (Simon) Kwong
Jennie Larkin
Fenglong Liu
John Quackenbush
Yonghong Wang
Yan Yu

Array Software Hit Team

Nirmal Bhagabati
John Braisted
Tracey Currier
Jerry Li
Wei Liang
John Quackenbush
Alexander I. Saeed
Vasily Sharov
Mathangi Thiagarajan
Joseph White

Assistant

Aseye Aboagye

