

Microarray Data Analysis - II

FIOCRUZ Bioinformatics Workshop
6 June, 2001

Challenges in Microarray Data Analysis

- **Spot Identification and Quantitation.**
- **Normalization of data from each experiment.**
- **Identification of Differentially Expressed Genes**
- **Identified of genes with correlated patterns of expression.**
- **Interpretation of data with respect to pathways.**
- **Literature filtered analysis.**

Image Processing Issues

- **Spot Finding**
- **Background Subtraction**
- **Reproducibility**
- **Measure - median vs. mean (integrated intensity)**
- **Quality measures**

TIGR Spotfinder Loading Image Data

General | Channel A | Channel B | Overlay | Selection | Mask | Converter | Batch File

Reset

● ●

Bkg

Image Info

Image Width=2200,
Image Length=6000,
Document Name:
ScanArray, Software:
GenScan, DateTime:
1999:03:03 10:17:30,
Artist: priti

OpenA F:\Priti\CH-c-12C-L4a\CH-c19-ch1.Tif

OpenB F:\Priti\CH-c-12C-L4a\CH-c19-ch2.Tif

Image Info

Image Width=2200,
Image Length=6000,
Document Name:
ScanArray, Software:
GenScan, DateTime:
1999:03:03 10:17:32,
Artist: priti

A	B	Overlay

Spacing
50

SpotSize
20

Use BKG

Draw Grid Manual Grid

Process

Excel

TIGR Spotfinder Zooming In

The screenshot displays the TIGR Spotfinder software interface. At the top, a menu bar includes 'General', 'Channel A', 'Channel B', 'Overlay', 'Selection', 'Mask', 'Converter', and 'Batch File'. The main window shows a large grid of spots, with a zoomed-in view of a specific region. On the left side, there is a control panel with the following elements:

- Reset** button and two radio buttons (one selected).
- Two vertical sliders for adjusting parameters.
- Bkg** label and a white rectangular input field.
- A small thumbnail image of the spot array.
- Three small preview windows labeled **A**, **B**, and **Overlay**.
- Spacing** dropdown menu set to **50**.
- SpotSize** dropdown menu set to **20**.
- Use BKG** checkbox.
- Draw Grid** and **Manual Grid** buttons.
- Process** button.
- Excel** button.

Buttons for **Zoom In** and **Zoom Out** are located to the left of the main image area. The main image area shows a grid of spots with a zoomed-in view of a specific region, and a small inset window in the top left corner of the main image area shows a zoomed-in view of a specific region of the spot array.

TIGR Spotfinder Image Overlay

General | Channel A | Channel B | **Overlay** | Selection | Mask | Converter | Batch File

Reset

A-Green

B-Red

A	B	Overlay

Spacing: 35

SpotSize: 20

Use BKG

Draw Grid Manual Grid

Process

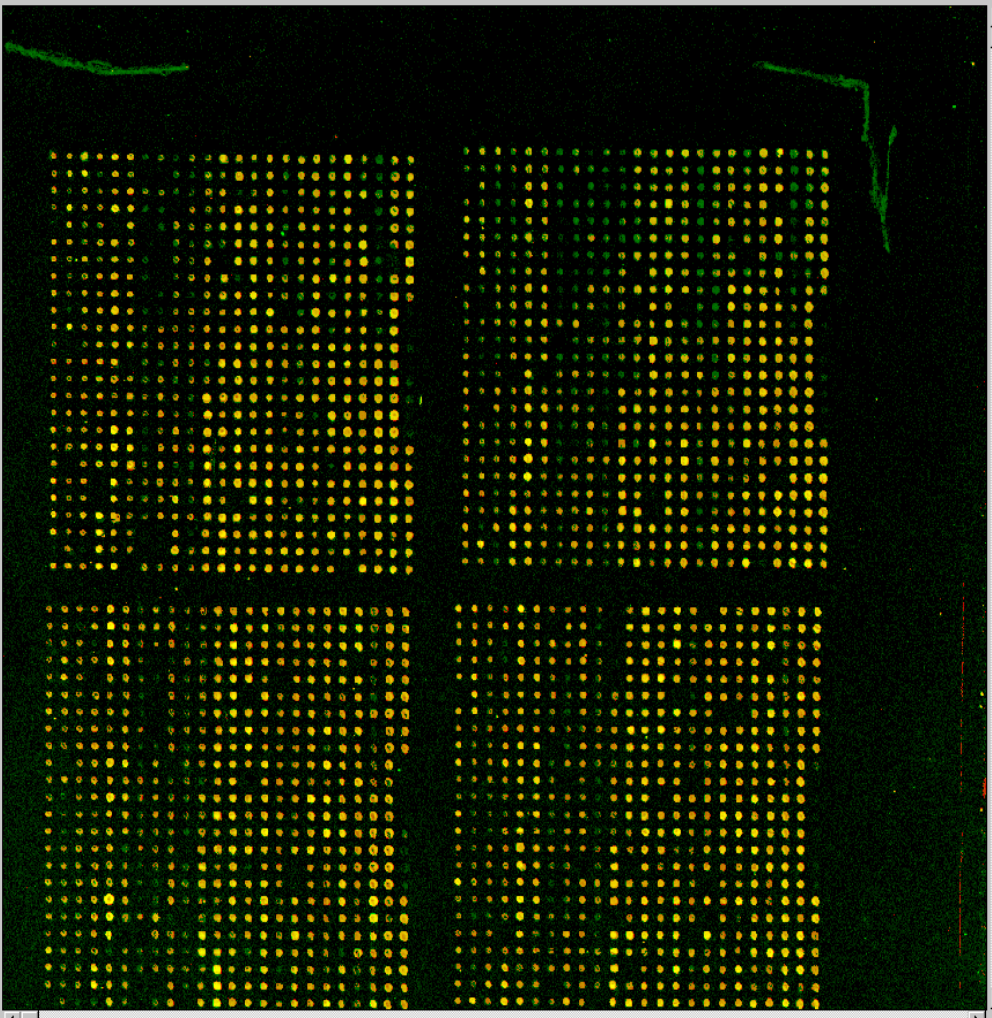
Excel

Zoom In

Zoom Out

Reset

Copy Image



TIGR Spotfinder Region Selection

General | Channel A | Channel B | Overlay | Selection | Mask | Converter | Batch File

Reset

Bkg

Zoom In

Zoom Out

Show contour

Fix Grid

Fix Row

Fix Column

Delete Row

Delete Col

Delete Spot

Copy Image

A B Overlay

Spacing
50

SpotSize
20

Use BKG

Draw Grid Manual Grid

Process


Excel

TIGR Spotfinder Grid Determination

General | Channel A | Channel B | Overlay | Selection | Mask | Converter | Batch File

Reset

Bkg



Zoom In

Zoom Out

Show contour

Fix Grid

Fix Row

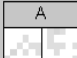
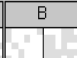
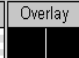



Fix Column

Delete Row

Delete Col

Delete Spot

Copy Image

A	B	Overlay
		
		

Spacing
35

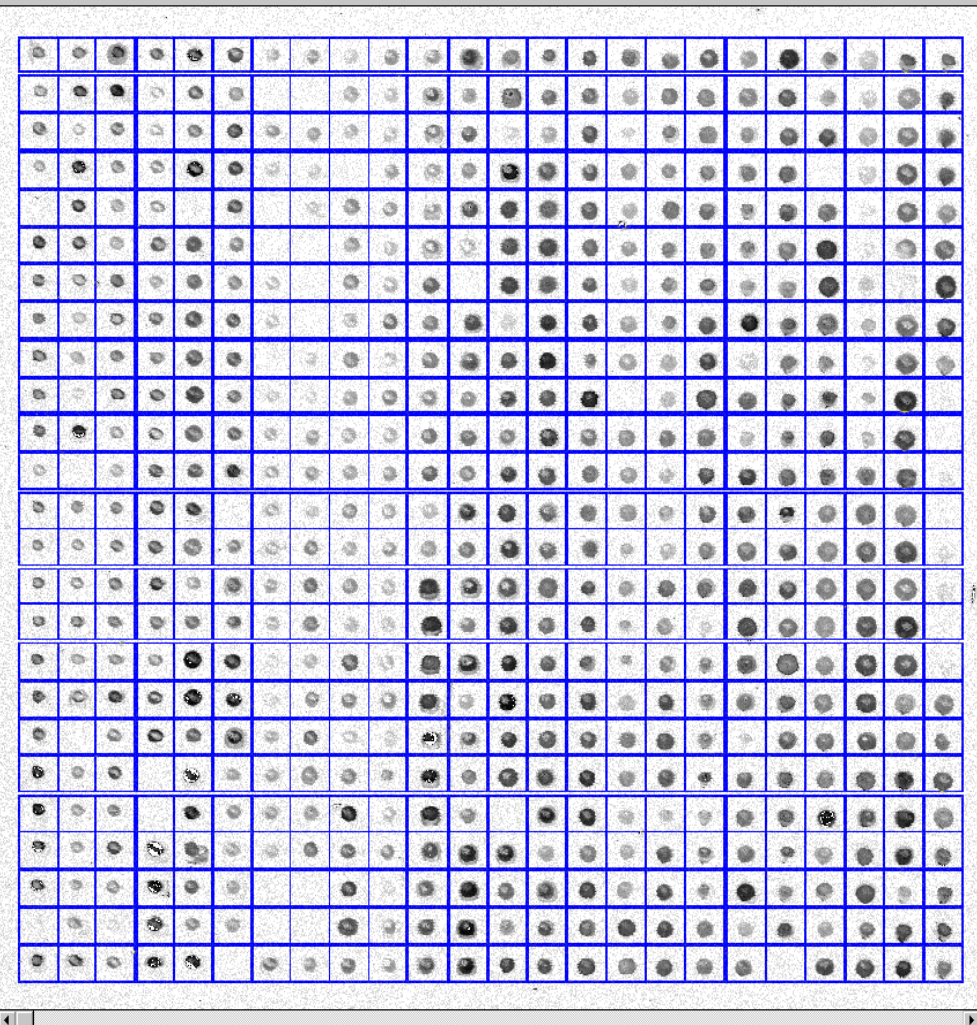
SpotSize
20

Use BKG

Draw Grid Manual Grid

Process

Excel



TIGR Spotfinder Grid Adjustment

General | Channel A | Channel B | Overlay | Selection | Mask | Converter | Batch File

Reset

Bkg

Zoom In

Zoom Out

Show contour

Fix Grid

Fix Row

Fix Column

Delete Row

Delete Col

Delete Spot

Copy Image

A	B	Overlay

Spacing
35

SpotSize
20

Use BKG

Draw Grid Manual Grid

Process

Excel

TIGR Spotfinder Spot Determination

General | Channel A | Channel B | Overlay | Selection | Mask | Converter | Batch File

Reset

Bkg

Zoom In

Zoom Out

A	B	Overlay

Spacing
35

SpotSize
20

Use BKG

Draw Grid Manual Grid

Process

Excel

TIGR Spotfinder Batch Mode

Reset

Bkg

General Channel A Channel B Overlay Selection Mask Converter Batch File

Save Selection MetaColumn MetaRow Clear All Selections

Save File 2 6 Use Grid

Load File

Run All

Use Fix Grid

Run OneByOne

Run Selection

A B Overlay

Spacing 35

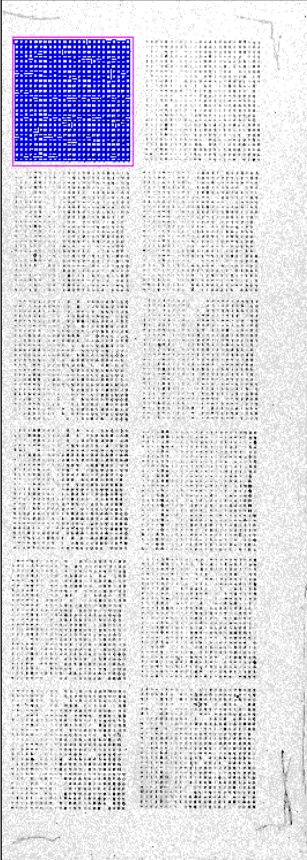
SpotSize 20

Use BKG

Draw Grid Manual Grid

Process

Excel



Comparison of Mean, Median, and Mode Ratios

Mean ratio

	1.012	0	0.966	0.987	0.897
	1.135	0	1.037	1.034	1.015
	1.008	0	1.058	1.008	1.058
	1.079	0	1.059	1.061	1.026
	1.022	0	1.069	1.031	1.019
	1.070	0	1.032	1.024	1.139
	0.986	0	1.058	1.064	1.047
	1.057	0	0.990	1.063	1.022
	0.935	0	1.105	1.069	1.079
	1.094	0	1.024	1.057	0.892
	1.014	0	1.040	0.997	1.019
	0.985	0	1.005	1.067	1.035
	1.011	0	1.033	1.035	1.143
	1.232	0	0.996	1.169	1.077
	0.819	0	1.085	1.118	1.039
	0.942	0	0.999	1.129	1.061
average	1.025		1.035	1.057	1.036
stdev	0.092		0.037	0.049	0.067

Median ratio

	1.000	0	0.930	1.053	0.898
	1.067	0	1.053	1.056	1.015
	1.008	0	1.042	1.056	1.003
	1.098	0	1.047	1.034	1.026
	0.980	0	0.998	1.004	0.955
	1.041	0	1.040	1.045	1.074
	1.030	0	1.081	1.029	1.134
	1.013	0	0.987	1.085	1.014
	0.896	0	1.025	1.007	1.079
	1.026	0	1.028	1.072	0.874
	1.001	0	1.067	1.020	1.014
	0.977	0	0.928	1.094	0.979
	1.061	0	1.105	1.027	1.096
	1.136	0	0.963	1.067	1.020
	0.929	0	1.033	1.083	1.062
	0.877	0	0.974	1.133	1.110
average	1.009		1.019	1.054	1.022
stdev	0.068		0.051	0.034	0.072

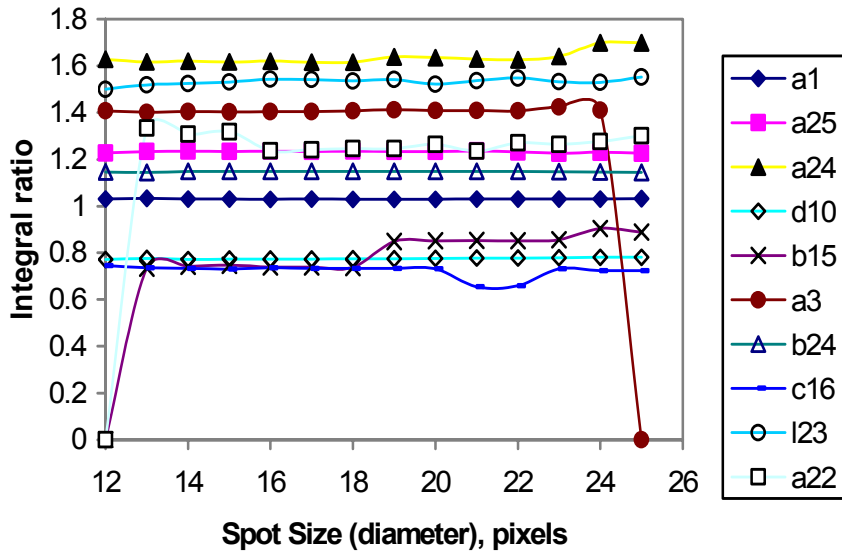
Mode ratio

	1.398	0	1.113	0.830	0.518
	1.984	0	0.554	0.906	5.721
	1.536	0	1.113	1.152	0.570
	1.051	0	2.490	1.684	1.437
	2.794	0	0.976	1.544	1.651
	1.095	0	1.564	1.203	1.516
	1.332	0	1.253	0.614	94.797
	1.697	0	16.921	2.039	0.788
	0.550	0	1.065	0.873	0.916
	1.022	0	0.742	0.377	0.681
	1.707	0	1.714	1.411	2.167
	1.784	0	0.392	0.976	1.080
	2.269	0	0.318	1.708	0.615
	1.932	0	0.604	1.331	0.579
	0.052	0	1.618	0.380	1.254
	0.641	0	0.543	1.373	1.022
average	1.428		2.061	1.150	7.207
stdev	0.690		4.004	0.479	23.391

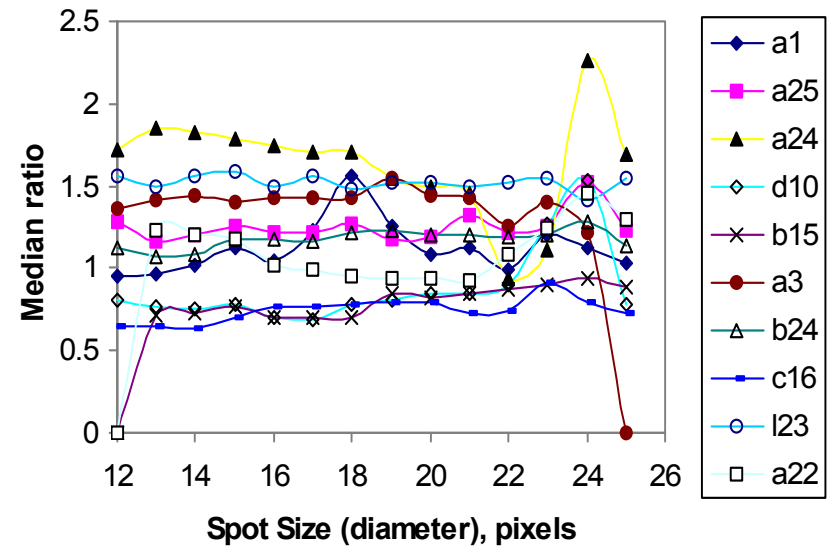
A comparison of Cy3/Cy5 ratios calculated using the mean, median, and mode ratios for control spots that should have a measured ratio of 1 for the 1st, 3rd, 4th, 5th columns.

Integral (Mean) Ratio vs. Median Ratio

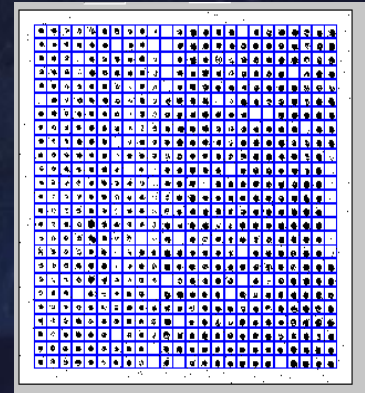
Integral Ratio vs. used Spot Size



Median ratio vs. used Spot Size



A comparison of Cy3/Cy5 ratios for various spot sizes using either the integrated intensity or the pixel median. In this case, the actual spot size is approximately 15 pixels in diameter.



TIGR

THE INSTITUTE FOR GENOMIC RESEARCH

Microarray Expression Analysis

Species Selection

Differential Growth Conditions

RNA Preparation and Labeling

Competitive Hybridization

Gene

Spot on a Slide

Fluorescence Intensity

Expression Measurement

Data Analysis Issues

- **Presentation**
- **Multiple Views**
- **Normalization**
- **Identification of Differentially Expressed Genes**
- **Multiple Experiments**

Why Normalize Data?

- Goal is to measure ratios of gene expression levels
 $(\text{ratio})_i = R_i/G_i$
where R_i/G_i are, respectively, the measured intensities for the i th spot.
- In a self-self hybridization, we would expect all ratios to be equal to one:
 $R_i/G_i = 1$ for all i . But they may not be.
- Why not?
 - Unequal labeling efficiencies for Cy3/Cy5
 - Noise in the system
 - Differential expression
- Normalization brings (appropriate) ratios back to one.

Normalization Approaches

- **Total Intensity**
- **Linear Regression**
- **Ratio statistics described by Chen, Dougherty, & Bittner**
J. Biomed. Optics (1997) 2(4) 364-374
- **Iterative log(ratio) mean centering**

Any of these using:

- **Entire Data Set**
- **User-defined Data Set/Controls**

Normalization Approaches

Entire Data Set

- **Probe Quantification less important**
- **No assumption on which genes constitute “housekeeping” set**
- **Uses all the data**
- **No independent confirmation**

User-defined Data Set/Controls

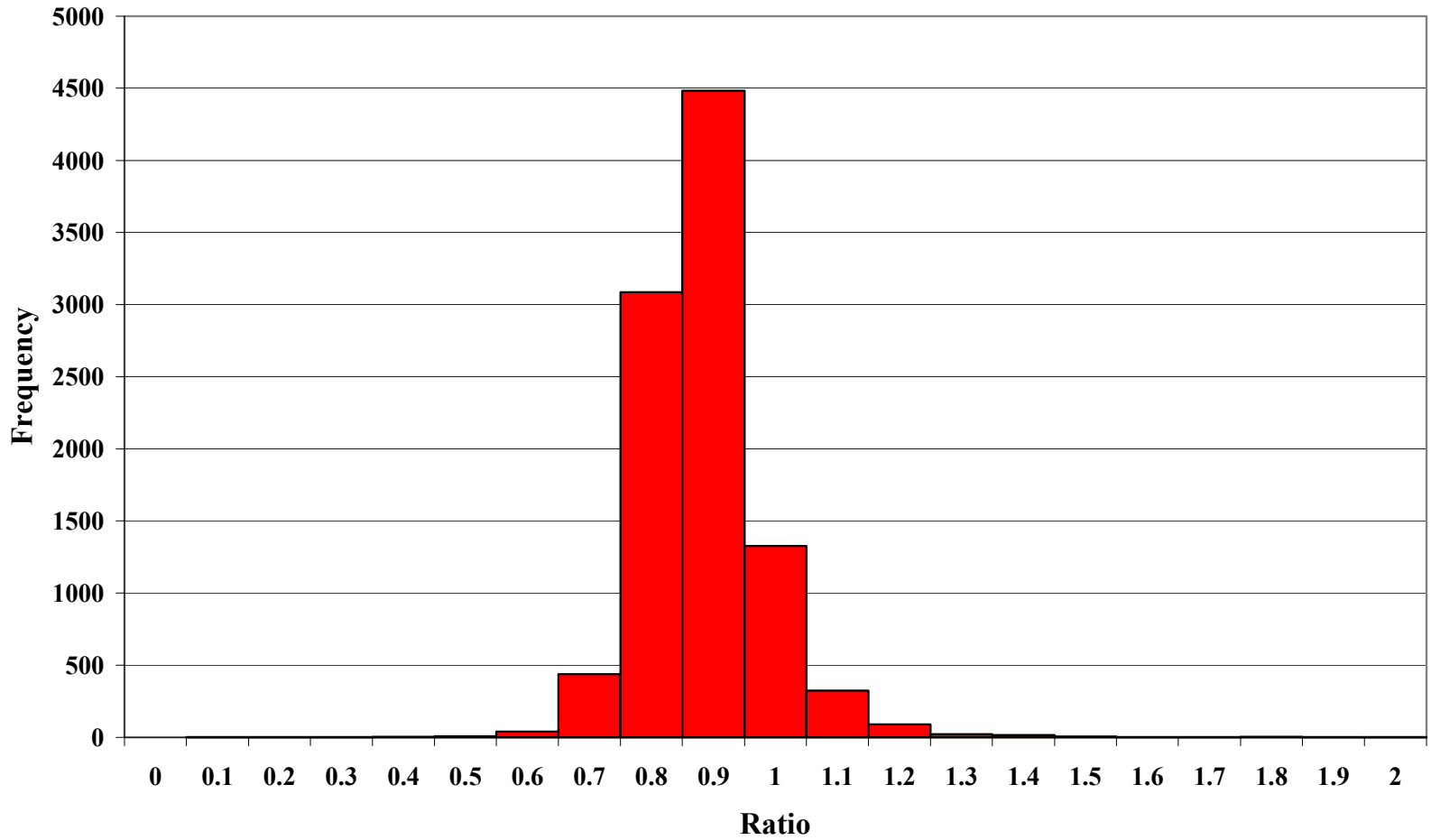
- **Requires definition of “housekeeping” set
or good added controls**
- **Requires good RNA quantitation**
- **Ignores much data**

Normalization Approaches

Solution(?)

- **Experiment dependent**
- **Use a combination of techniques**
- **SMART Experimental design**

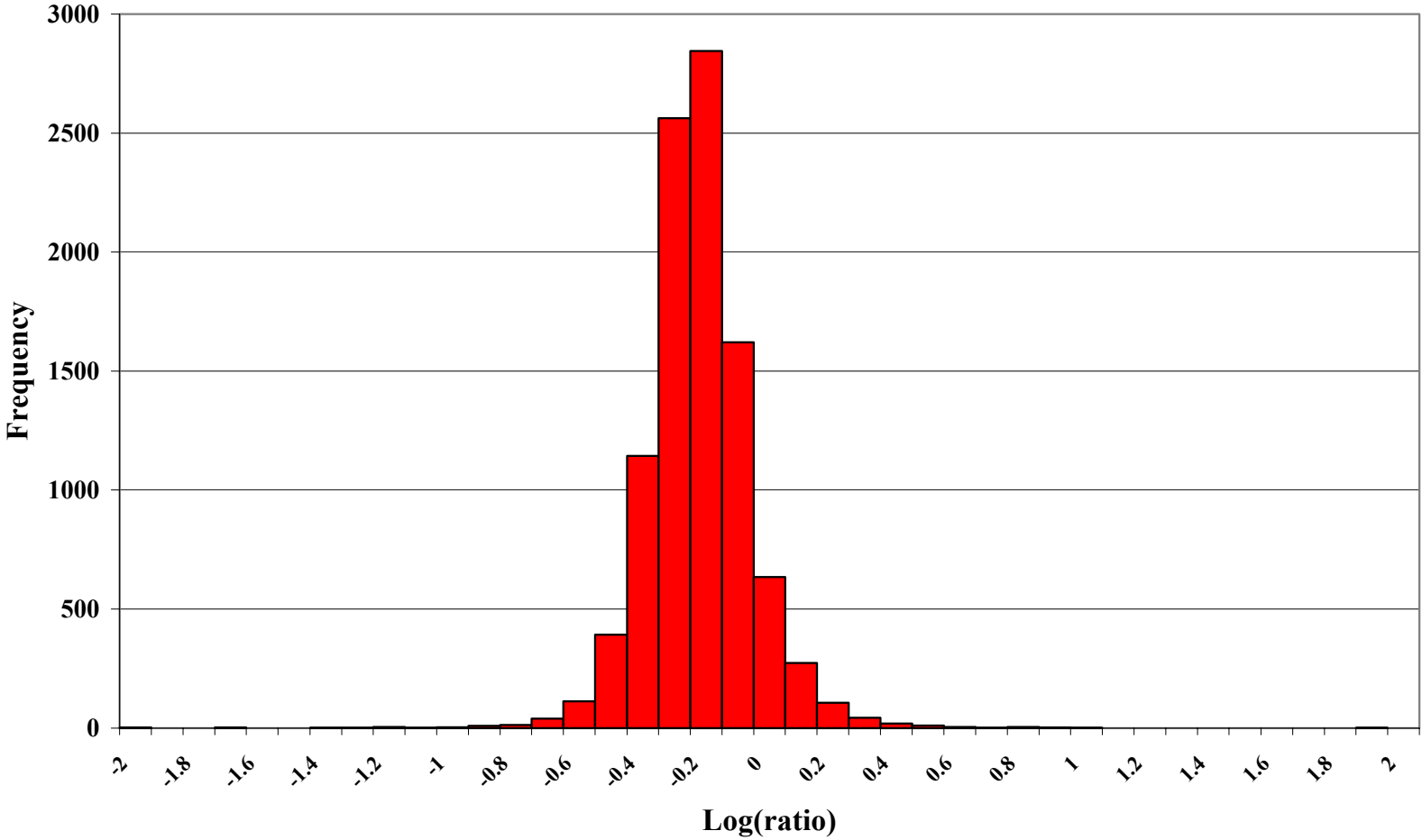
Ratio Histogram



TIGR

THE INSTITUTE FOR GENOMIC RESEARCH

Log(ratio) Histogram



Normalization Approaches: Total Intensity

- **Assumption: Total RNA (mass) used is same for both samples.**
- **So, averaged across thousands of genes, total hybridization should be the same**

Normalization Factor:
$$N = \frac{\sum_{k=1}^{N_{array}} R_k}{\sum_{k=1}^{N_{array}} G_k}$$

Normalization: $G'_k = NG_k$ and $R'_k = R_k$.

Normalization Approaches: Linear Regression

Assumption: Total RNA used is constant, some genes expressed with ratio of 1, slope of best fit line normalized to 1

$$R_k = \beta_0 + \beta_1 G_k + u_k$$

Normalization Factor:

$$S(\beta_0, \beta_1) = \sum_{k=1}^n u_k^2 = \sum_{k=1}^n (R_k - \beta_0 - \beta_1 G_k)^2$$

The values of β_0 and β_1 that minimize $S(\beta_0, \beta_1)$, b_0 and b_1 , are given by

$$b_1 = \frac{\sum_{k=1}^n (R_k - \bar{R})(G_k - \bar{G})}{\sum_{k=1}^n (G_k - \bar{G})^2} \quad \text{and} \quad b_0 = \bar{R} - b_1 \bar{G},$$

$$\text{where } \bar{R} = \frac{\sum R_k}{n} \quad \text{and} \quad \bar{G} = \frac{\sum G_k}{n}.$$

Normalization: $G'_k = \left[\frac{1}{b_1} \right] G_k$ and $R'_k = R_k$.

Normalization Approaches: Ratio Statistics (1)

Assumption: Total RNA used is constant, some genes expressed with ratio of 1, variations are functions of the common mean

$$\sigma_{G_k} = c\mu_{G_k} \text{ and } \sigma_{R_k} = c\mu_{R_k}, \text{ with } \mu_{G_k} = \mu_{R_k} = \mu_k.$$

Probability Density for Ratio T_k : $f_{T_k}(t) \approx \frac{(1+t)\sqrt{1+t^2}}{c(1+t^2)\sqrt{2\pi}} \exp\left[\frac{-(t-1)^2}{2c(1+t^2)}\right]$

This density can be used to calculate the mean, standard deviation and confidence interval limits for the distribution of measured ratio values. As functions of c , these parameters can be estimated using a polynomial approximation

$$y = a_3c^3 + a_2c^2 + a_1c + a_0$$

with constants are chosen appropriately:

$$\mu: (a_3, a_2, a_1, a_0) = (0.364, 1.279, -0.0427, 1.001)$$

$$\sigma: (a_3, a_2, a_1, a_0) = (-2.805, 2.911, -2.706, 0.979)$$

$$\text{lower limit at 95\% confidence: } (a_3, a_2, a_1, a_0) = (28.644, -2.830, 3.082, 0.989)$$

$$\text{upper limit at 95\% confidence: } (a_3, a_2, a_1, a_0) = (-5.002, .4.462, -3.496, 0.9968)$$

Normalization Approaches: Ratio Statistics (2)

Assume that the population mean $\mu_0 = 1$ and let the first approximation of the normalization parameter m_1 be equal to the calculated sample mean. A first approximation of c , \hat{c}_1 , is calculated using

$$\hat{c}_i = \left[\frac{1}{n} \sum_{j=1}^n \frac{(t_j - 1)^2}{(1 + t_j^2)} \right]^{1/2}$$

where the sum is over the n elements taken initially between the one-half and twice the sample mean.

Upper and lower limits at the 95% confidence level, θ_1 and θ_2 , are then calculated using \hat{c}_1 and the previous approximation.

A normalization factor \hat{m}_1 is calculated using

$$\hat{m}_i = \frac{1}{\hat{\mu}_{i-1}} \left(\frac{1}{n} \sum_{j=1}^n t_j \right),$$

where, again, we take $\hat{\mu}_0 = 1$, the sum is over the n array elements used to estimate \hat{c}_1 , and i is an index used to count the number of iterations.

The individual ratios are then rescaled using

$$t'_k = \frac{t_k}{\hat{m}_i} = \frac{R_k}{(\hat{m}_i G_k)} = \frac{R'_k}{G'_k}.$$

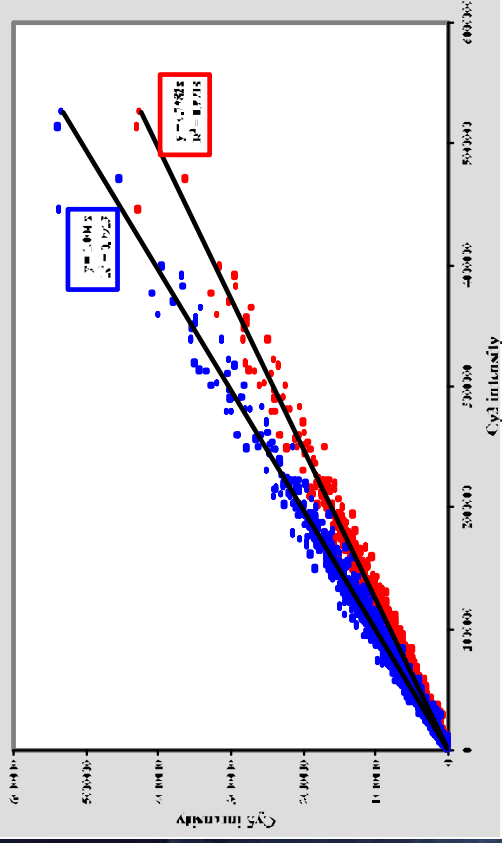
This process is then iterated until the calculated value of the mean estimator converges to a fixed value.

The upper and lower confidence limits for the normalized experimental data are then calculated as

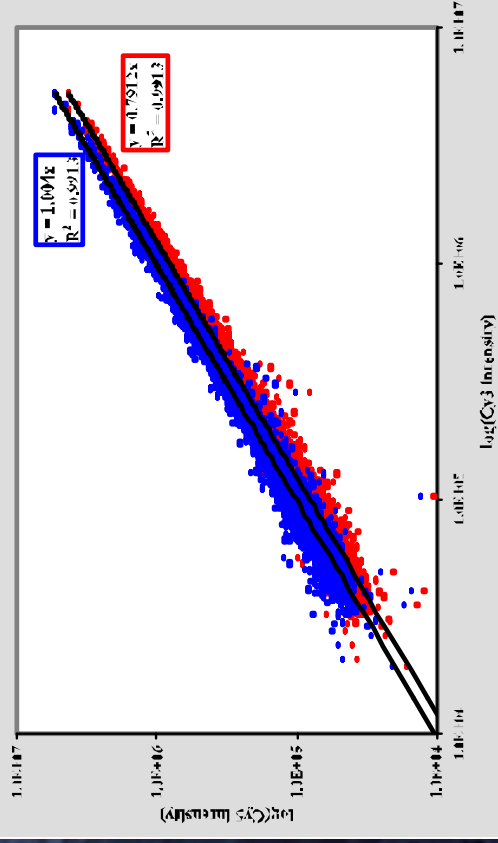
$$\theta'_1 = \hat{m} \theta_1 \text{ and } \theta'_2 = \hat{m} \theta_2$$

and (θ'_1, θ'_2) are used to define the limits for identification of differentially expressed genes

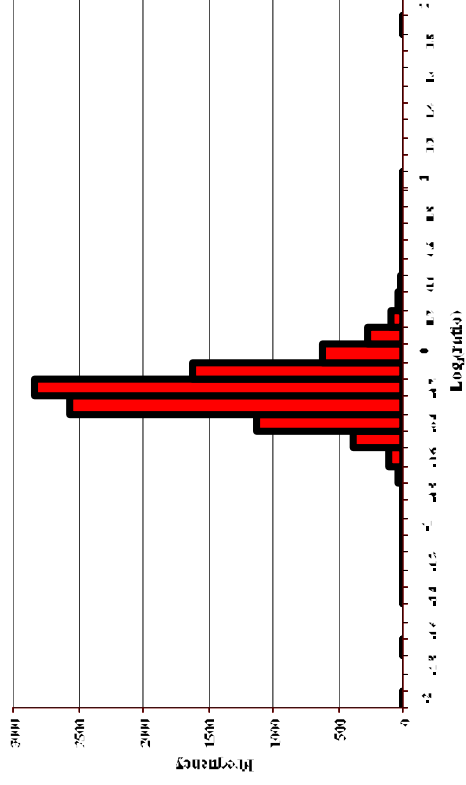
Scatterplot of Intensities for a self vs. self hybridization



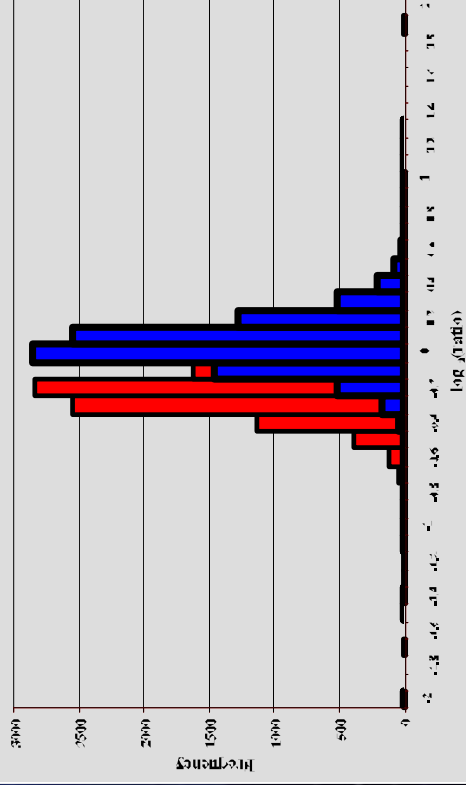
Scatterplot of log(intensities) for a self vs. self hybridization



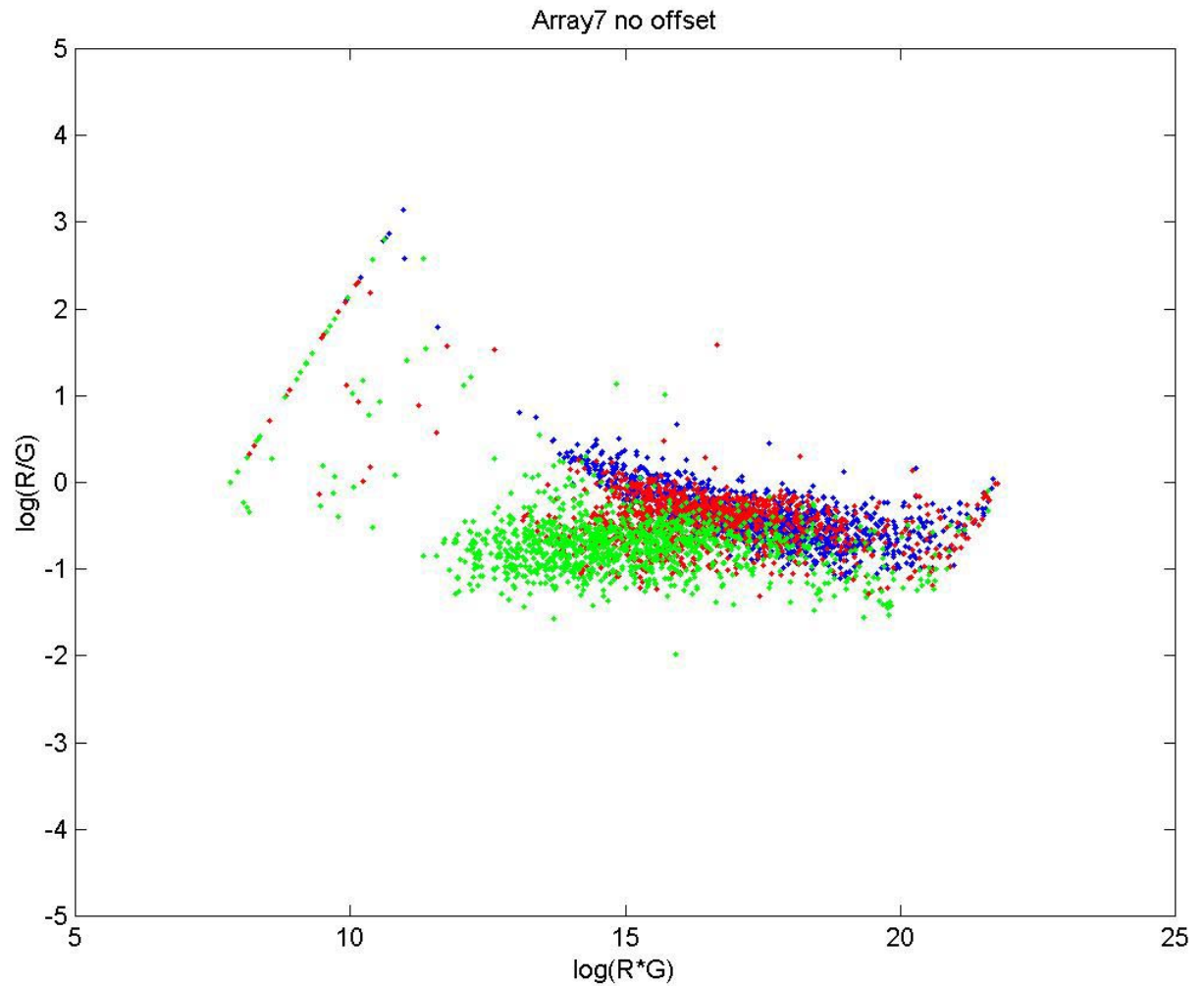
Logratio Histogram



logratio Histogram

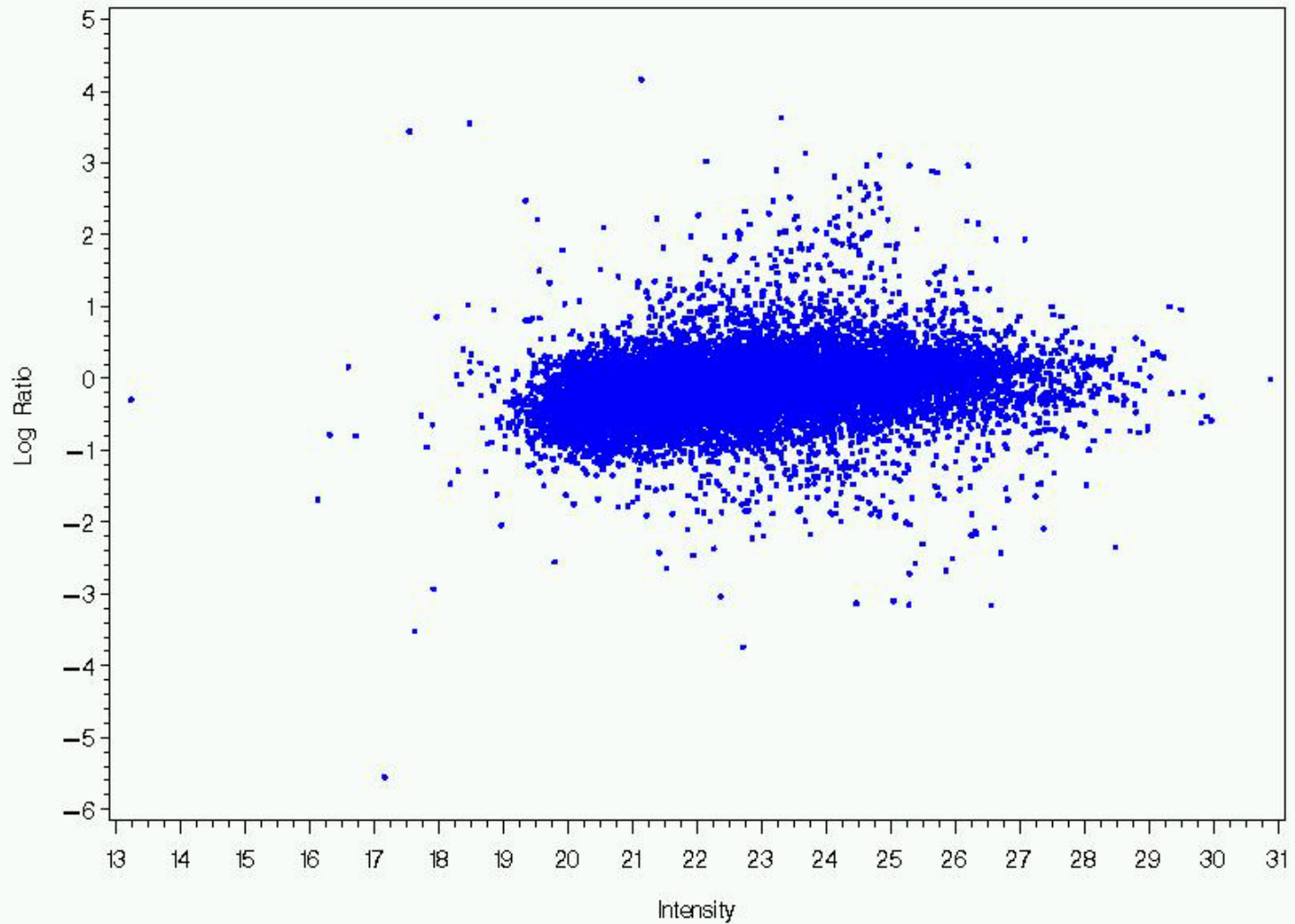


Bad Data from Parts Unknown



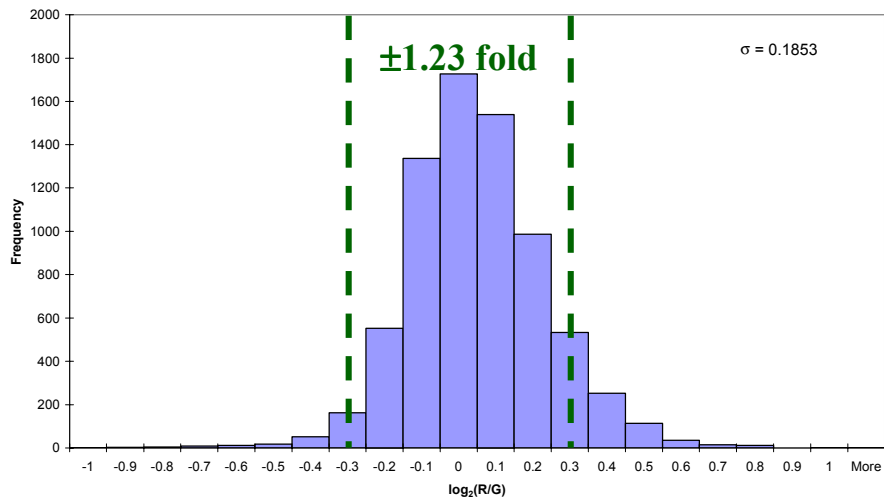
Good Data from TREG

Log Ratio vs Intensity Plot of Slide 10 (A5 to C3) with All Values Greater than One

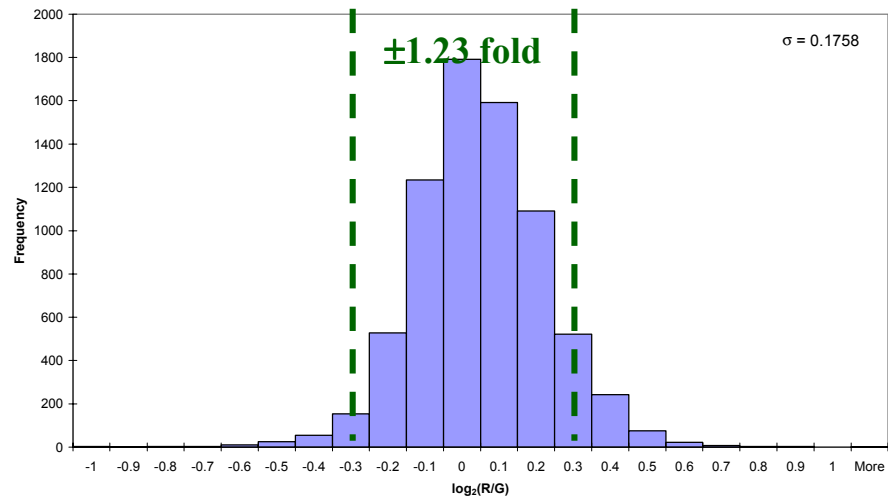


Normalization using local linear regression

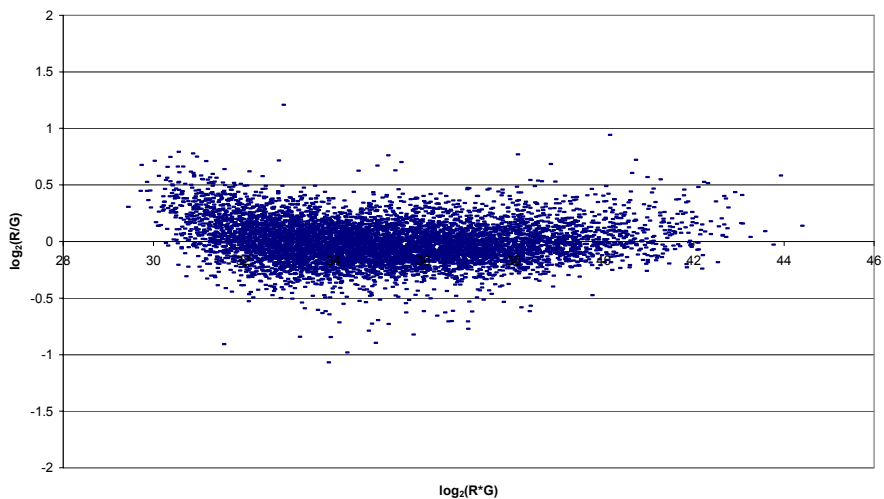
L4A
01-04-01-41
Normalized



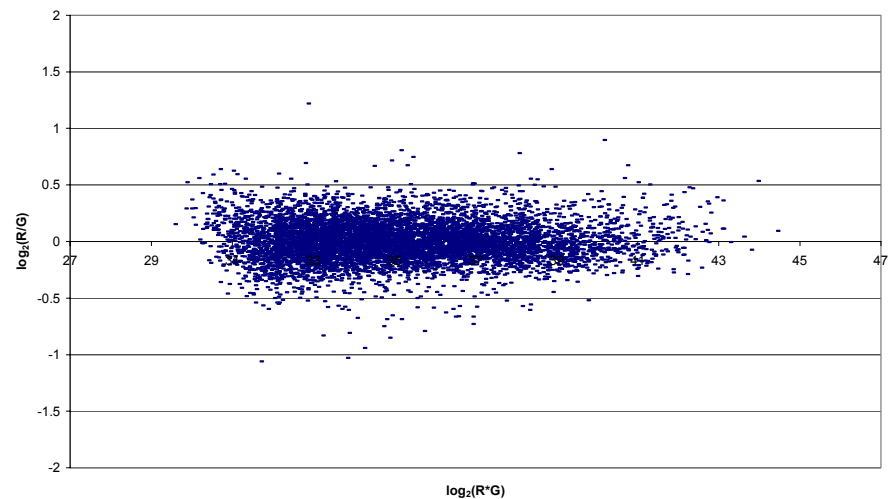
L4A
01-04-01-41
Lowess correction



L4A
01-04-01-41
Normalized

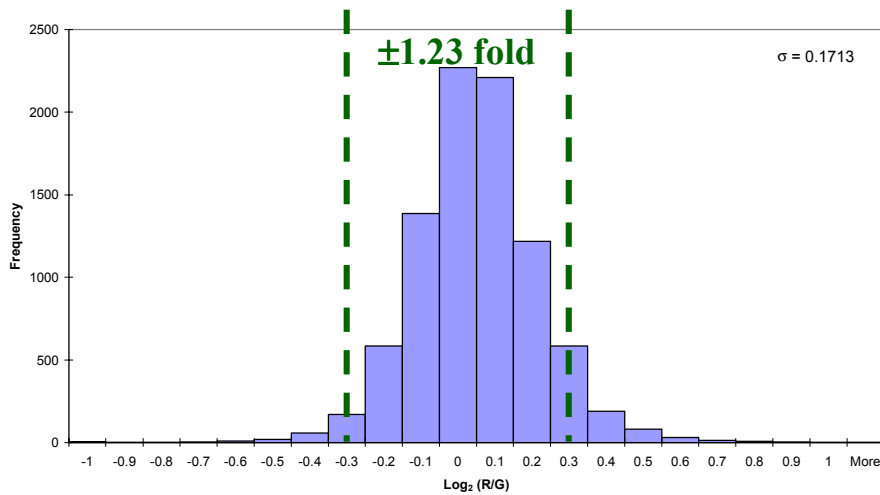


L4A
01-04-01-41
Lowess correction

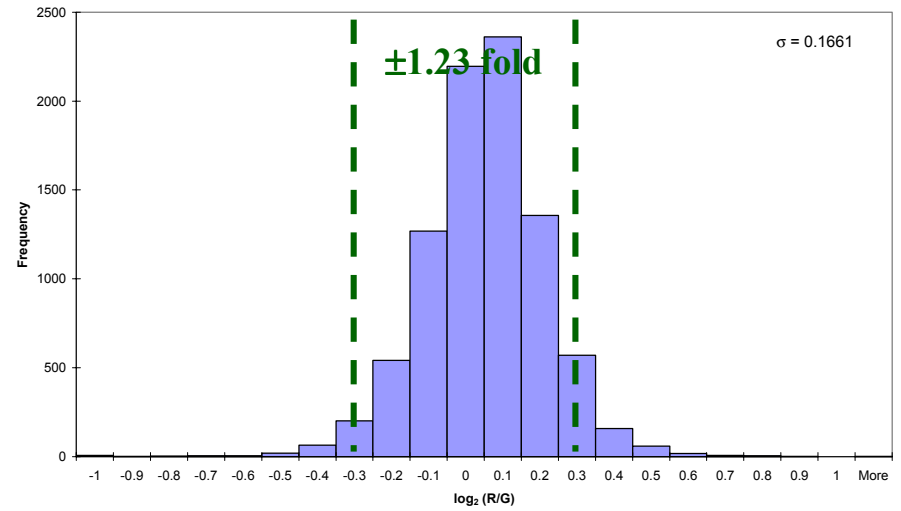


Normalization using local linear regression

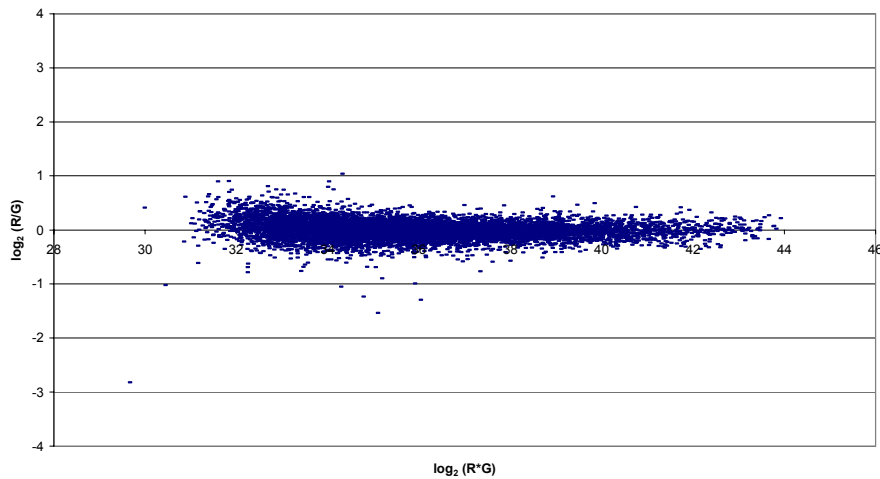
OVCAR3
01-04-01-16
Normalized



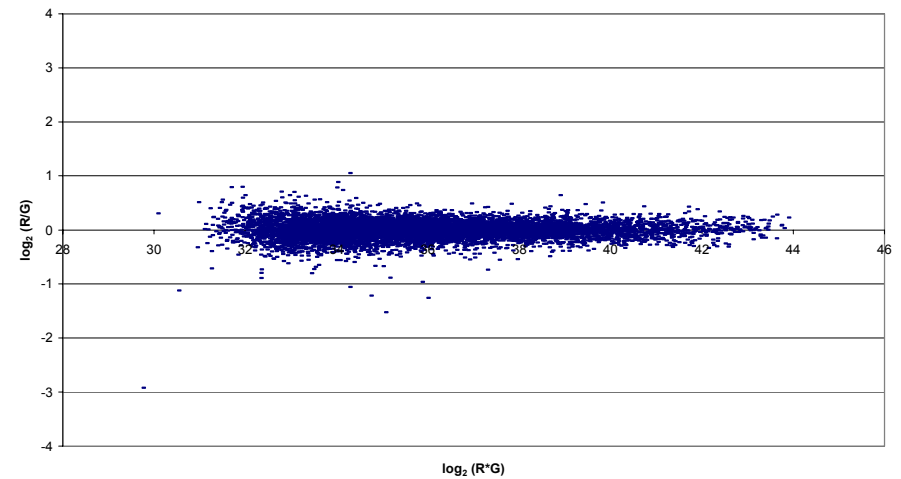
OVCAR3
01-04-01-16
Lowess correction



OVCAR3
01-04-01-16
Normalized

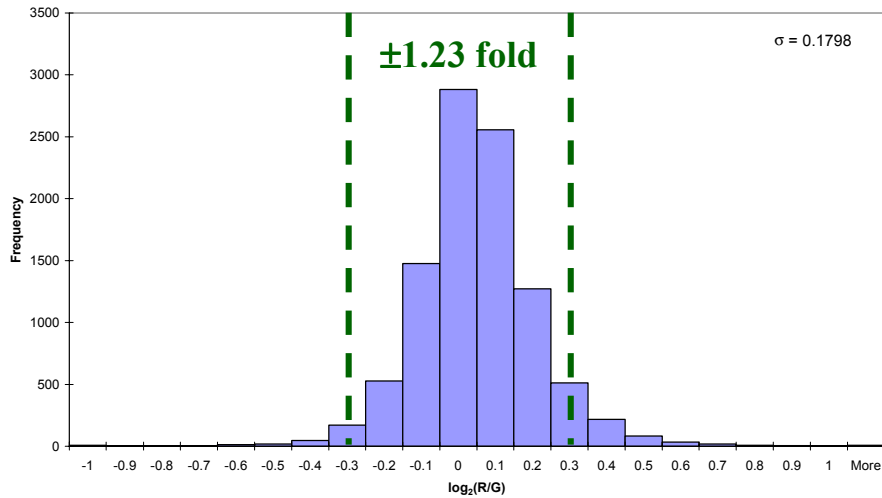


OVCAR3
01-04-01-16
Lowess correction

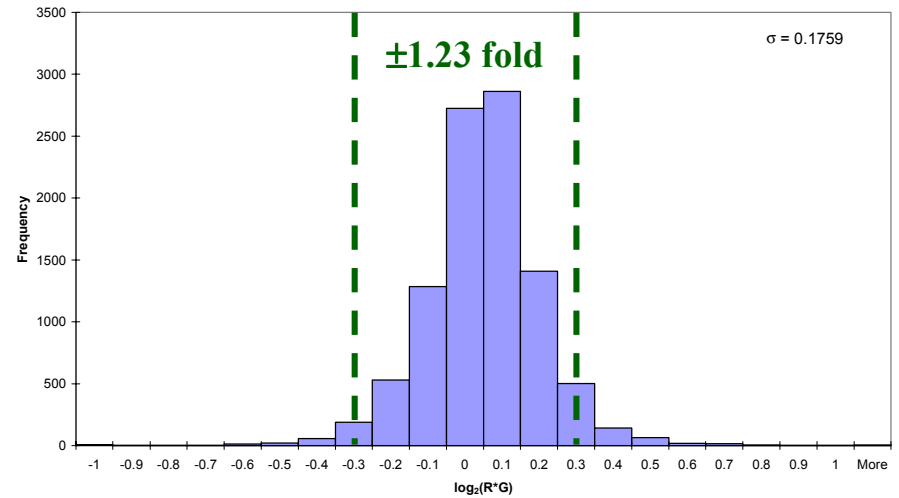


Normalization using local linear regression

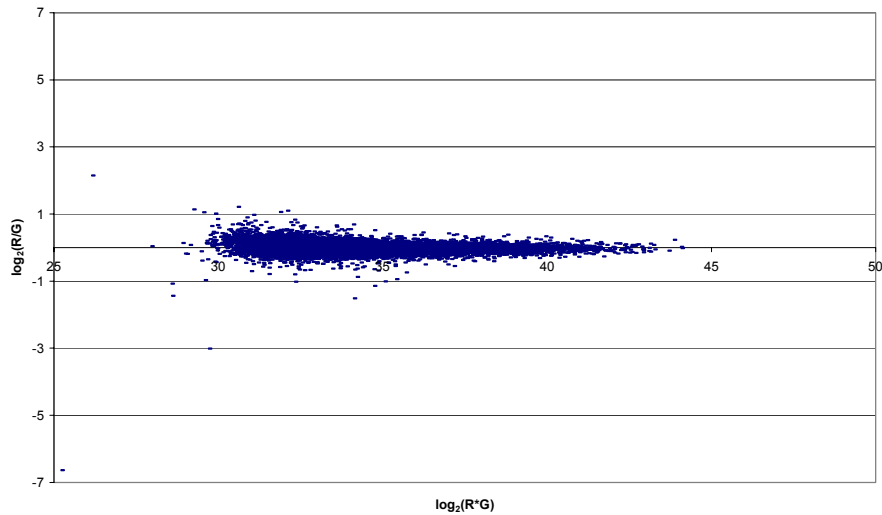
SW480
01-04-01-7
Normalized



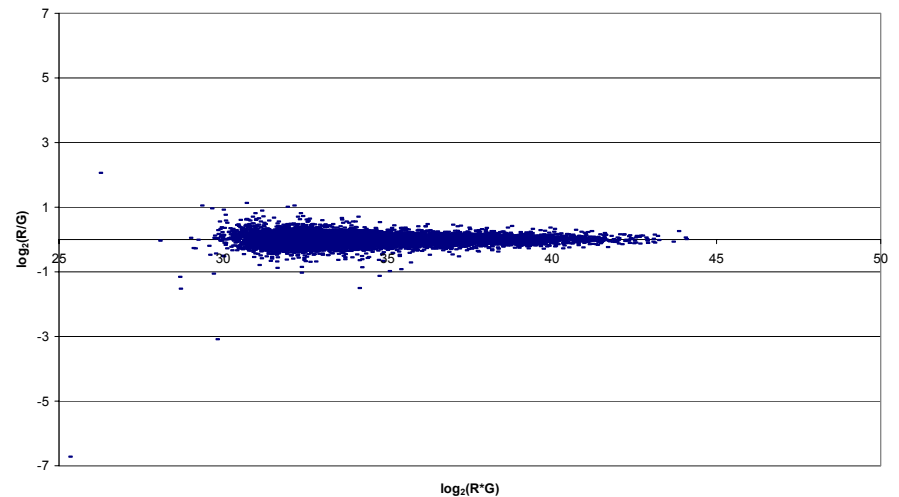
SW480
01-04-01-7
Lowess correction



SW480
01-04-01-7
Normalized



SW480
01-04-01-7
Lowess correction



Multiple Experiments?

- Goal is identify genes (or experiments) which have “similar” patterns of expression
- This is a problem in data mining
- “Clustering Algorithms” are most widely used
- Types
 - Agglomerative: Hierarchical
 - Divisive: k -means, SOMs
 - Others: Principal Component Analysis (PCA)
- All depend on how one measures distance

Expression Vectors

- Crucial concept for understanding clustering
- Each gene is represented by a vector where coordinates are its values $\log(\text{ratio})$ in each experiment
 - $x = \log(\text{ratio})_{\text{expt1}}$
 - $y = \log(\text{ratio})_{\text{expt2}}$
 - $z = \log(\text{ratio})_{\text{expt3}}$
 - etc.
- For example, if we do six experiments,
 - $\text{Gene}_1 = (-1.2, -0.5, 0, 0.25, 0.75, 1.4)$
 - $\text{Gene}_2 = (0.2, -0.5, 1.2, -0.25, -1.0, 1.5)$
 - $\text{Gene}_3 = (1.2, 0.5, 0, -0.25, -0.75, -1.4)$
 - etc.

Expression Matrix

- These gene expression vectors of $\log(\text{ratio})$ values can be used to construct an expression matrix

	Expt 1	Expt 2	Expt 3	Expt 4	Expt 5	Expt 6
Gene ₁	-1.2	-0.5	0	0.25	0.75	1.4
Gene ₂	0.2	-0.5	1.2	-0.25	-1.0	1.5
Gene ₃	1.2	0.5	0	-0.25	-0.75	-1.4
etc.						

- This is often represented as a red/green colored matrix

Distance metrics

- Distances are measured “between” expression vectors
- Distance metrics define the way we measure distances
- Many different ways to measure distance:
 - Euclidean distance
 - Pearson correlation coefficient(s)
 - Manhattan distance
 - Mutual information
 - Kendall’s Tau
 - etc.
- Each has different properties and can reveal different features of the data

Distance Matrix

- Once a distance metric has been selected, the starting point for all clustering methods is a “distance matrix”

	Gene ₁	Gene ₂	Gene ₃	Gene ₄	Gene ₅	Gene ₆
Gene ₁	0	1.5	1.2	0.25	0.75	1.4
Gene ₂	1.5	0	1.3	0.55	2.0	1.5
Gene ₃	1.2	1.3	0	1.3	0.75	0.3
Gene ₄	0.25	0.55	1.3	0	0.25	0.4
Gene ₅	0.75	2.0	0.75	0.25	0	1.2
Gene ₆	1.4	1.5	0.3	0.4	1.2	0

- The elements of this matrix are the pair-wise distances. Note that the matrix is symmetric about the diagonal.

Hierarchical clustering

- **Select the data you want to cluster**
- **“Filter” (normalize) the data appropriately and select distance**
- **Apply method:**
 - 1. Search through the distance matrix and find the two most similar clusters. This is the first true stage in the “clustering” process. If several pairs share the same similarity, use a predetermined rule to decide between alternatives.**
 - 2. Fuse the two selected clusters to produce a new cluster that now contains at least two objects.**
 - 3. Calculate the distances between this new cluster and all other clusters. There is no need to calculate *all* distances since only those involving the new cluster have changed.**
- **Repeat steps 1-3 until all objects are in one cluster.**

k-means clustering

- Select the data you want to cluster and filter; select distance
- Apply method:
 1. All initial objects are randomly assigned to one of k clusters (where k is an input parameter to the algorithm).
 2. An average expression vector is then calculated for each cluster and this is used to compute the distances between clusters.
 3. Objects are moved between clusters and intra- and inter-cluster distances are measured with each move. Objects are allowed to remain in the new cluster only if they are closer to it than to their previous cluster.
 4. Following each move, the expression vectors for each cluster are recalculated.
 5. The shuffling proceeds until moving any more objects would make the clusters more variable.

Self Organizing Maps (SOMs)

- Select the data you want to cluster and filter; select distance
- Apply method:
 1. Random vectors are constructed and assigned to each partition. (where the number and geometry are input parameters).
 2. A gene is picked at random and using a selected distance metric, the reference vector that it is closest to the gene's is identified .
 3. The reference vector is then adjusted so that it is more similar to the randomly picked gene's. The reference vectors that are nearby on the two dimensional grid are also adjusted so that they too are more similar to the randomly selected gene .
 4. Steps 2 and 3 are iterated several thousand times, decreasing the amount by which the reference vectors are adjusted and increasing the stringency used to define closeness in each step. As the process continues, the reference vectors are converge to fixed values .
 5. Finally, the genes are mapped to the relevant partitions depending on the reference vector to which they are most similar.

Principal Component Analysis (PCA)

- Select the data you want to cluster and filter
- Apply method:
OK, this gets a bit complicated. . . .

Basically:

- 1. We find the eigenvectors of the expression matrix**
- 2. We select those with the greatest eigenvalues**
- 3. We project our data on the eigenvectors with the three greatest eigenvalues**
- 4. And make pretty pictures**

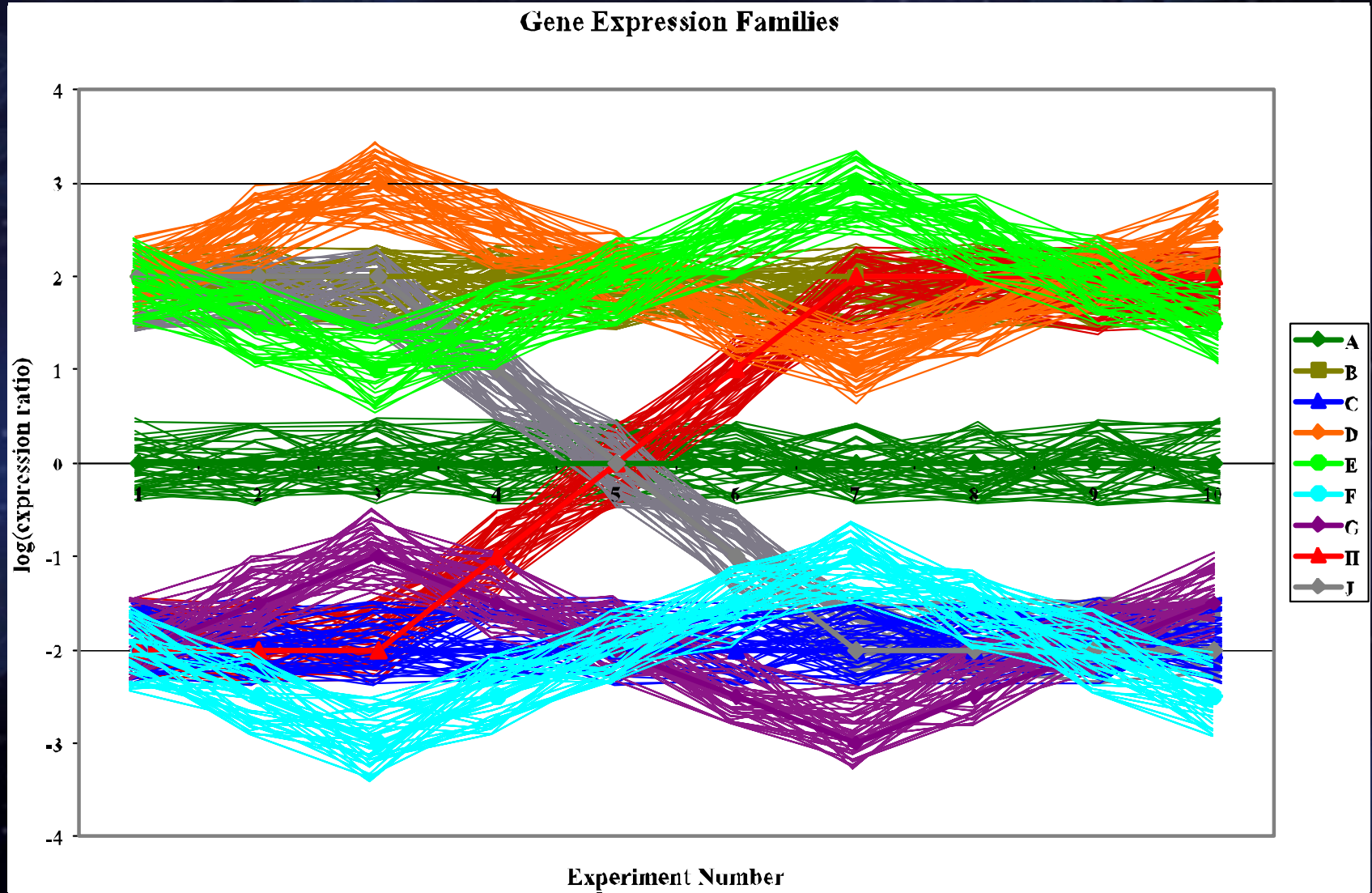
Support Vector Machines (SVM)

- **Select the data you want to cluster and filter**
- **Apply method:**
OK, this gets even more complicated. . . .

Basically this is a neural network approach to finding dividing your data into genes “like” and “unlike” a training set. . . .

- 1. Pick a set of genes you are know about (your training set)**
- 2. Train the SVM. This produces a pattern that can be recognized**
- 3. Screen the data using the SVM model**

TIGR MeV: Test Data Set



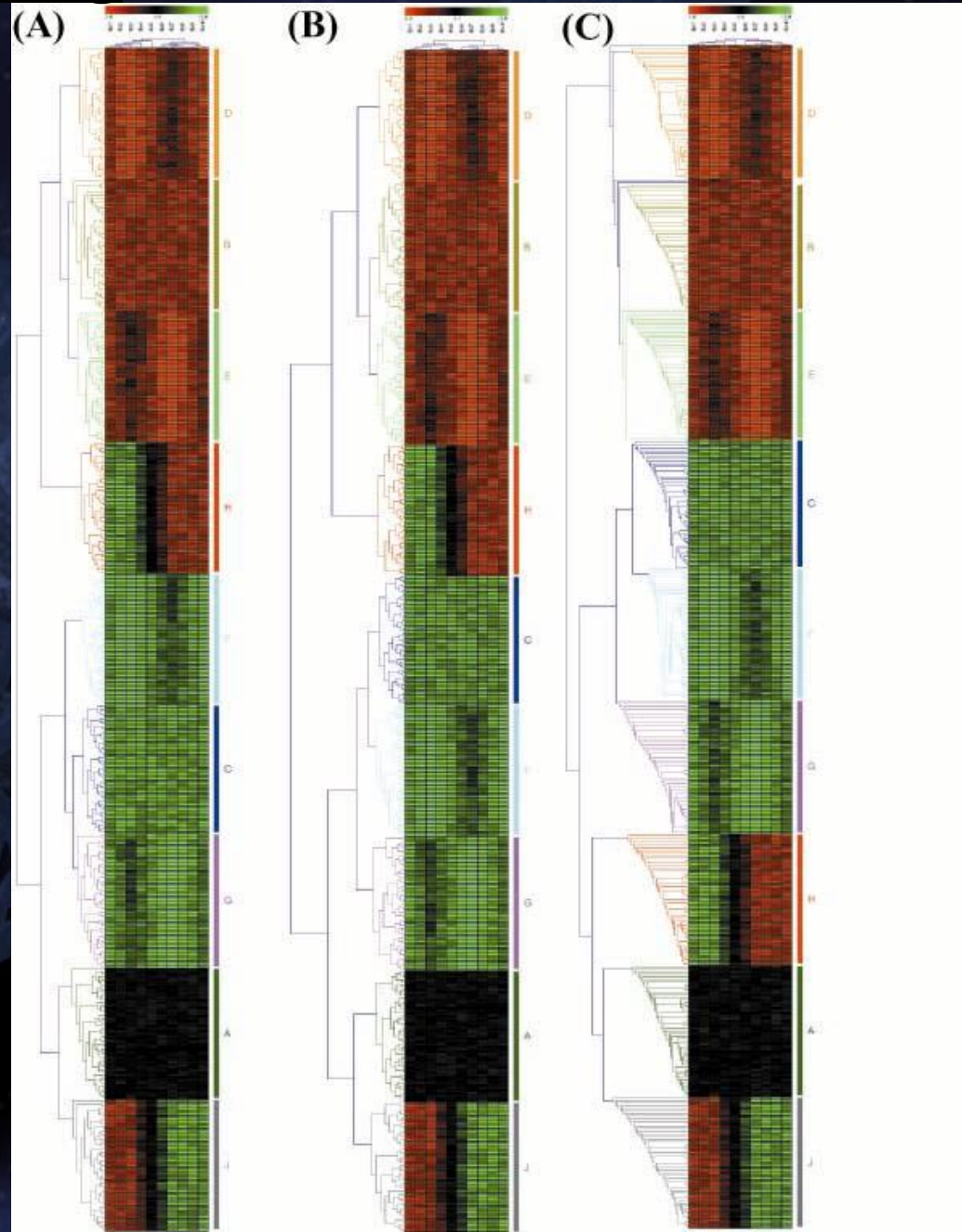
Hierarchical Clustering

(A) Average Linkage

(B) Complete Linkage

(C) Single Linkage

Even related algorithms produce slightly different views of the data.

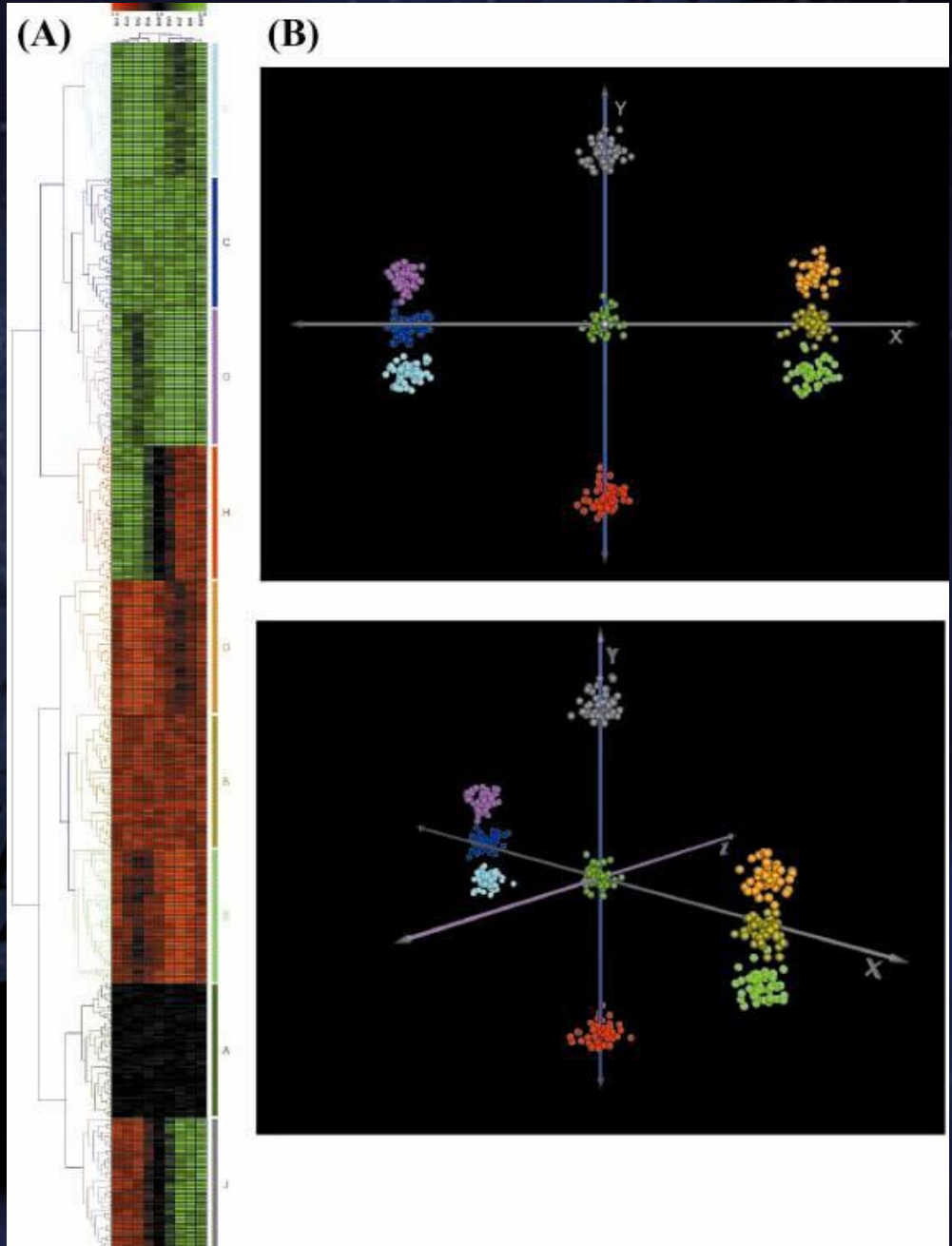


Hierarchical Clustering and PCA

(A) Average Linkage

(B) PCA

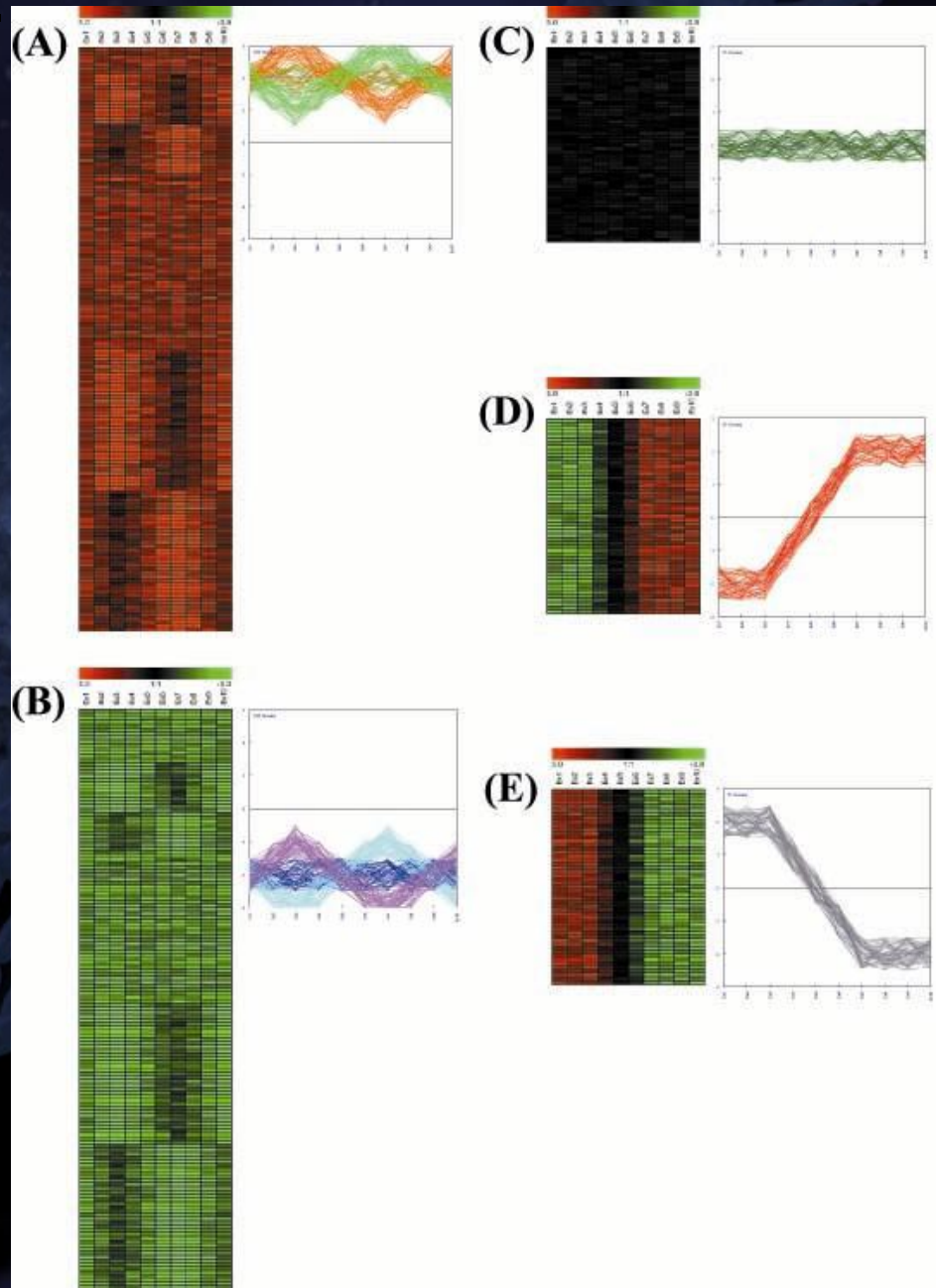
Separate clusters may have more or less support when using different algorithms.



k-means Clustering

Separate clusters may have more or less support when using different algorithms.

Note colors are based on hierarchical clustering Results.



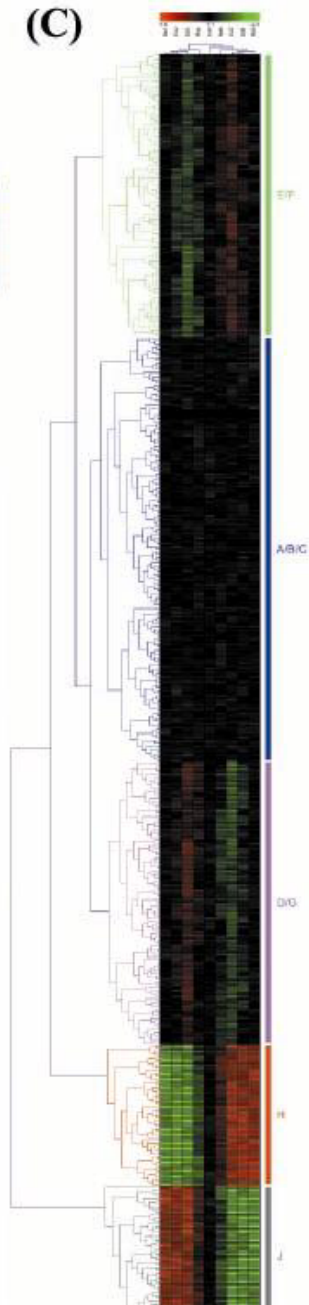
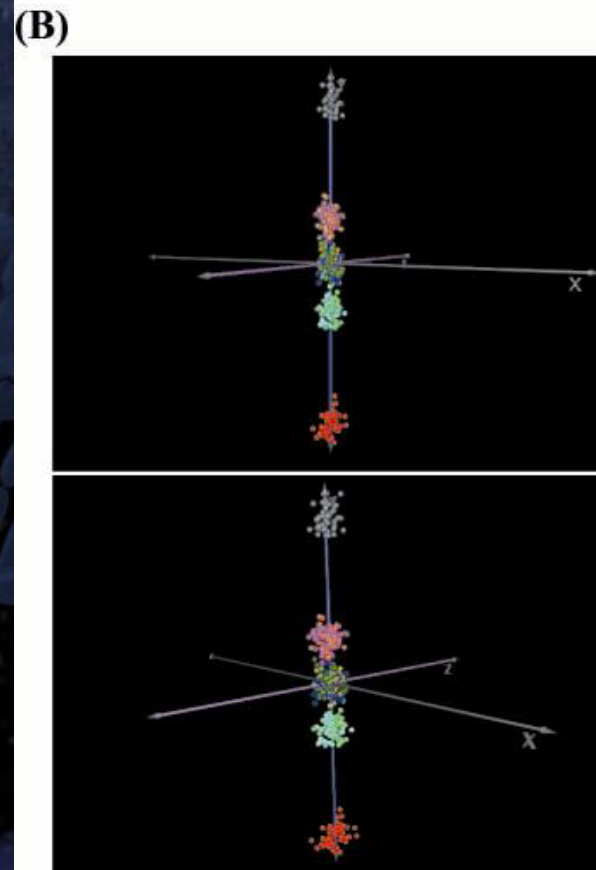
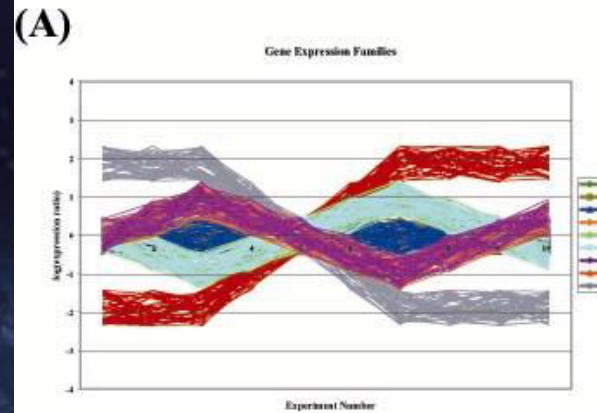
The effects on Mean Centering

(A) Expression Profiles

(B) PCA

(C) Hierarchical Clustering

Adjusting the data can have profound effects, but allow different patterns to be seen.



Very Useful Microarray URLs

Leming Shi	http://www.gene-chips.com
TIGR	http://pga.tigr.org/tools
MGED	http://www.mged.org
Wentian Li	http://linkage.rockefeller.edu/wli/microarray
EBI	http://industry.ebi.ac.uk/~alan/MicroArray
Terry Speed	http://stat-www.berkeley.edu/users/terry/zarray/Html
Joe Derisi	http://www.microarrays.org/index.html
Pat Brown	http://cmgm.stanford.edu/pbrown/mguide/
NCGR	http://www.ncgr.org/research/genex/other_tools.html
Stanford	http://www.dnachip.org
HAPI	http://array.ucsd.edu

Acknowledgments

The TIGR Gene Index Team

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

TIGR Collaborators

Norman Lee
Rena Malek
Hong-Ying Wang
Truong Luu
Nnenna U. Nwokekeh

TIGR Human/Mouse/Arabidopsis

Expression Team

Emily Chen
Renee Gaspard
Jeremy Hasseman
Heenam Kim
John Quackenbush
Erik Snesrud
Shiubang Wang
Ivana Yang
Yan Yu
Baoping Zhao

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

Array Software Hit Team

Jerry Li
John Quackenbush
Alex Saeed
Vasily Sharov
Alexander Sturn
Joseph White

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

johnq@tigr.org

Assistant
Mary Mulholland

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

TIGR Faculty, IT Group, and Staff