# Clustering and displaying microarray data

George Bell, Ph.D.

Bioinformatics and Research Computing
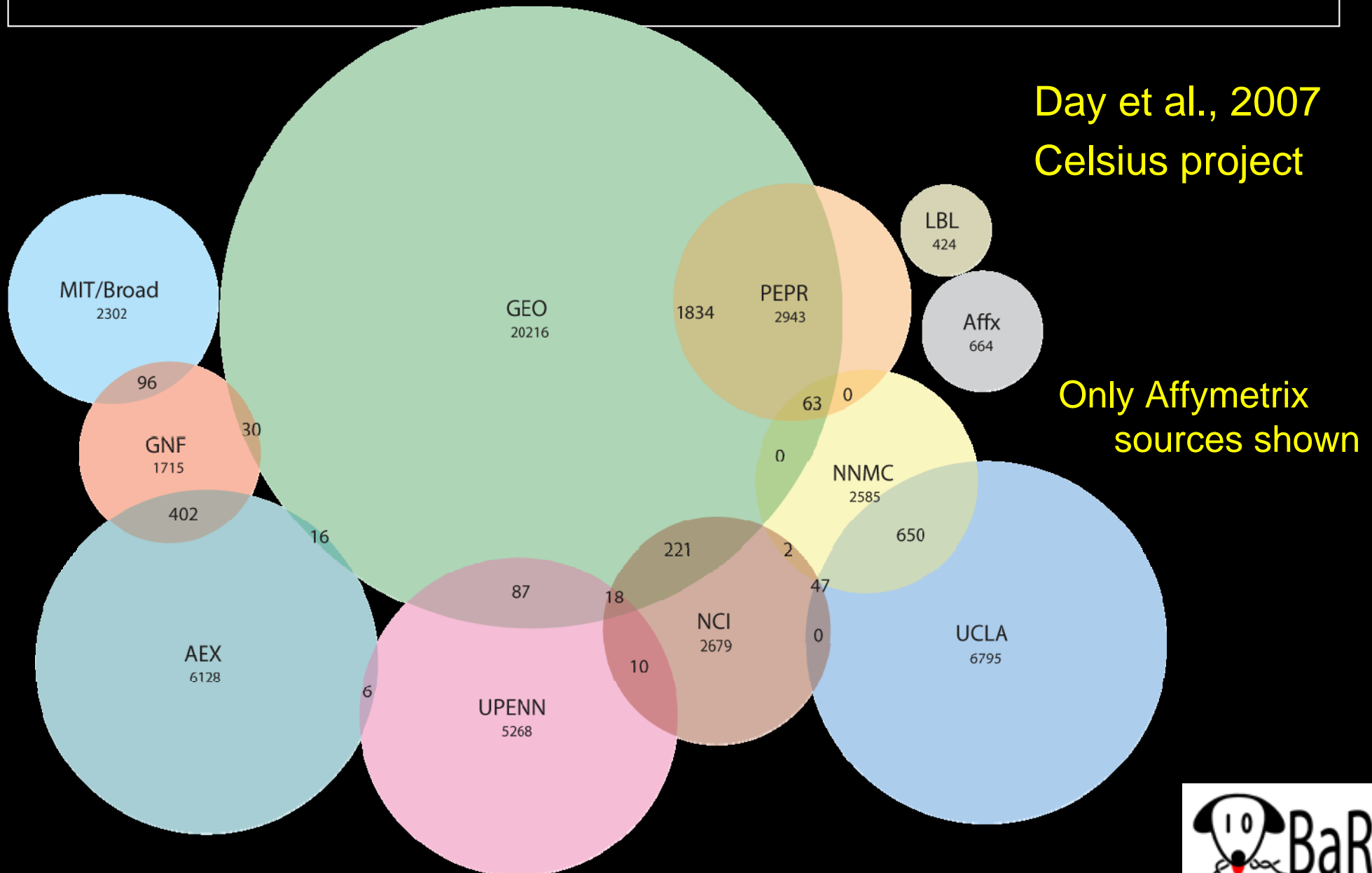
Hot Topics – March 2008

# Why?

- Explore a large amount of expression or other data

- Get experiment-wide look at interesting subset of data

- Visually identify patterns for further analysis

- Order genes and/or experiments in a sensible way

- Split genes and/or experiments into a predefined number of groups

# Why not?

- Clustering is not a substitute for rigorous statistics
- Clustering cannot identify
  – differentially expressed genes
  – profiles that are correlated with a reference profile
- Any data – even noise – can be clustered
- Clustering is not an essential step for most analyses

# Where to get the data?



Day et al., 2007
Celsius project

Only Affymetrix
sources shown

# Types of data

- **Single-color arrays (mainly Affymetrix)**
  - Data reported as expression values
  - Raw values or log2-transformed values (RMA; GCRMA)
- **Two-color arrays**
  - Data reported as expression ratios
  - Raw ratios or log2-transformed ratios

# Clustering with Cluster 3.0

- Based on original clustering program by Michael Eisen

- Code updated by Michiel de Hoon

- Runs on Windows, Mac, and Linux

- Free from
  http://bonsai.ims.u-tokyo.ac.jp/~mdehoon/software/cluster/software.htm

- Hierarchical, k-means, SOMs

- Other option for large datasets:
  – XCluster, a command-line tool by Gavin Sherlock

BaRC
Bioinformatics and Research Computing

# Getting Cluster 3.0

# Cluster data import

- ## Minimal matrix  (text, not Excel format)

| Probe | Amygdala | Heart | Kidney | Liver | Lung |
|---|---|---|---|---|---|
| 1000_at | 0.85 | 0.19 | -0.92 | -0.32 | -0.27 |
| 1009_at | 0.02 | 0.44 | 0.32 | 0.53 | -0.80 |
| 1014_at | -0.25 | 0.17 | -5.83 | -5.83 | 0.93 |
| 1030_s_at | -0.25 | | 0.13 | -2.09 | 0.21 |
| 1031_at | -0.35 | -0.19 | -0.22 | -5.00 | |

- ## Matrix with annotation and cluster weights

| GeneID | NAME | GWEIGHT | Amygdala | Heart | Kidney | Liver | Lung |
|---|---|---|---|---|---|---|---|
| EWEIGHT | | | 0 | 0 | 1 | 1 | 1 |
| 1000_at | MAPK3 | 1 | 0.85 | 0.19 | -0.92 | -0.32 | -0.27 |
| 1009_at | HINT1 | 1 | 0.02 | 0.44 | 0.32 | 0.53 | -0.80 |
| 1014_at | POLG | 1 | -0.25 | 0.17 | -5.83 | -5.83 | 0.93 |
| 1030_s_at | TOP1 | 0 | -0.25 | | 0.13 | -2.09 | 0.21 |
| 1031_at | SRPK1 | 0 | -0.35 | -0.19 | -0.22 | -5.00 | |

# Data filtering

- **Why filter?**
  - Noise (unexpressed or uninteresting genes) can hide signal
  - A complete dataset is too much to visually process
- **What are you looking for?**
  - Differentially expressed genes
  - Most variable genes
  - Most interesting profile (expression pattern)
- **Select list of genes of interest**
- **Select set of genes with GO annotation of interest**
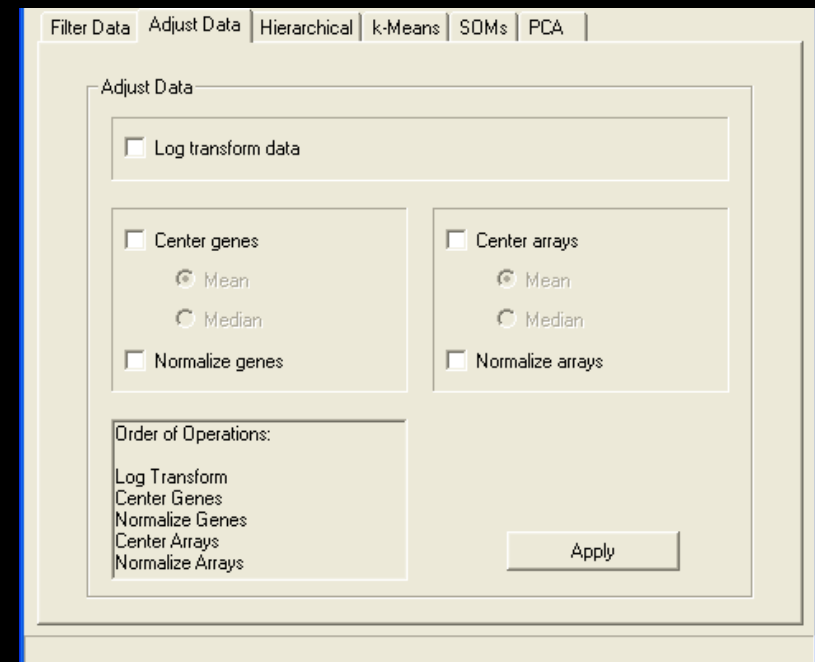- **Do in spreadsheet or Cluster ("Filter Data" tab)**

# Transforming data

- Do in spreadsheet or Cluster ("Adjust Data" tab)
- Common methods
  - Log-transformation
  - Converting values into ratios
  - Centering:
    - value – mean (row or column)
    - value – median (row or column)
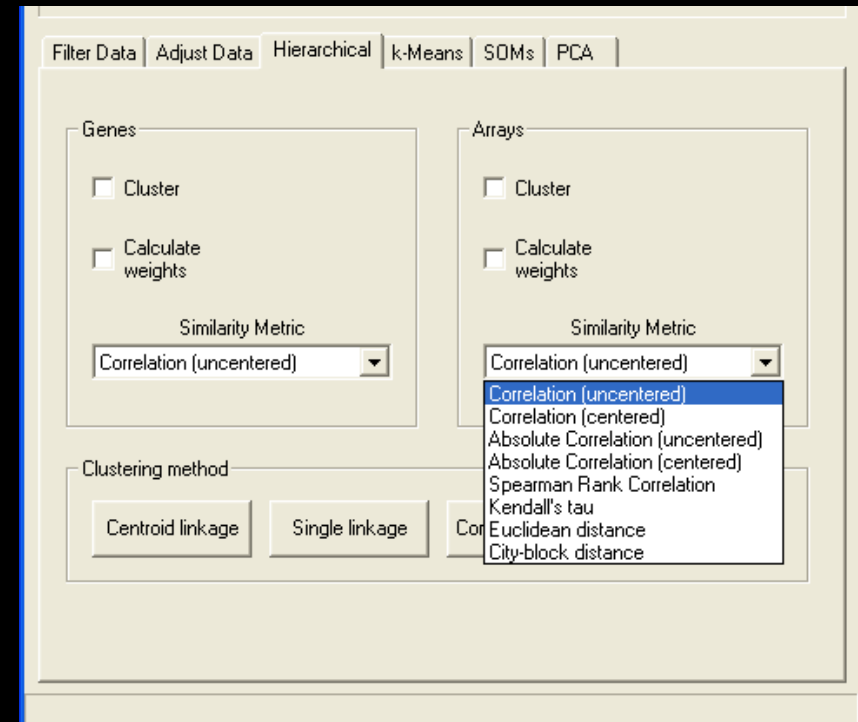  - Many normalization methods (from elsewhere)

# Clustering goals and caveats

- Potential goal: organize a set of data to show relationships between data elements
- With microarray analysis: genes and/or chips
- Most data does not inherently exist in clusters
- Most effective with optimal quantity of data
- Interpretation of data in obvious clusters: is it filtered?
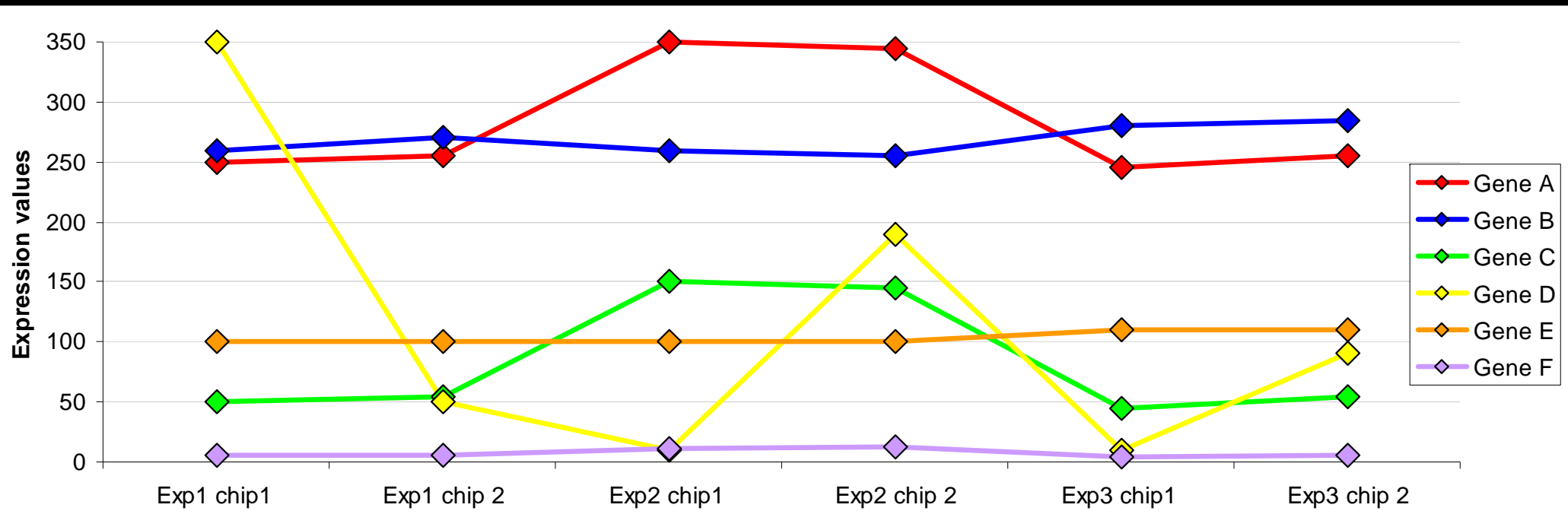- Clustering vs segmenting

# Hierarchical clustering

- Agglomerative, unsupervised analysis
- Steps
  1. Create an all vs. all distance matrix
  2. Fuse closest objects
  3. Compare fused object to all others
  4. Repeat steps 2-3 until one inclusive cluster is created
- Can be performed on genes and/or arrays
- Efficiency = $O(n^2m)$
- Need to select:
  - Similarity Metric
  - Clustering method

# Measuring similarity between profiles

- Similarity (distance) metric is an important choice when comparing genes and/or experiments
- What are you trying to group?

# Common similarity metrics

- Pearson correlation
  - Measures the difference in the shape of two curves
  - modifications:
    - uncentered correlation: for offset profiles, coefficient < 1
    - absolute correlation: opposite profiles cluster together

- Euclidean distance: multidimensional Pythagorean Theorem
  - Measures the distance between two curves

- Nonparametric or Rank Correlation
  - Similar to the Pearson correlation but data values are replaced with their ranks
  - Ex: Spearman Rank, Kendall's Tau
  - Good idea if distribution of data is not normal
  - More robust (against outliers) than other methods

# Clustering methods

How can groups of objects be represented?
How is distance measured to a cluster of objects?

- **Single linkage** (b)
  - minimum distance
- **Complete linkage** (r)
  - maximum distance
- **Centroid linkage** (p)
  - distance to "centroid" of group
- **Average linkage** (x)
  - average distance

**Weighting?**
  - GWEIGHT, EWEIGHT

b

p

r    y

2

1

$x = mean\ (b,y,r)$

BaRC
Bioinformatics and Research Computing

# Cluster data output

- For hierarchical clustering by genes and arrays, 3 output files are created:
  - .cdt ("clustered data table")
  - .gtr ("gene tree")
  - .atr ("array tree")
- All are tab-delimited text and can be opened as a spreadsheet
- Create your own 'cdt' file and bypass Cluster 3.0:
  - Tab-delimited text
  - First 2 columns are gene identifiers

| Gene ID | Symbol | Amygdala | Heart | Kidney | Liver | Lung |
|---------|--------|----------|-------|--------|-------|------|
| 1000_at | MAPK3 | 0.85 | 0.19 | -0.92 | -0.32 | -0.27 |
| 1009_at | HINT1 | 0.02 | 0.44 | 0.32 | 0.53 | -0.80 |

# Representation of clustered data

- Hierarchical clustering produces a dendrogram(s) showing relationships between objects

- Order of leaves: $2^{N-1}$ choices

- How can objects be partitioned into groups?
  - k-means clustering
  - self-organizing maps
  - How many clusters (k)?

- Are the data really hierarchical?

- Original distance matrix may be informative

# Visualizing clustered data with Java TreeView

- Based on original clustering program by Michael Eisen

- Code updated by Alok Saldanha

- Runs on Windows, Mac, and Linux

- Free from
  http://sourceforge.net/project/showfiles.php?group_id=84593

# Getting Java TreeView

- http://sourceforge.net/project/showfiles.php?group_id=84593

# Java TreeView main view

# Java TreeView: settings

# Java TreeView: exporting images

# Displaying other types of data

# Demo ?