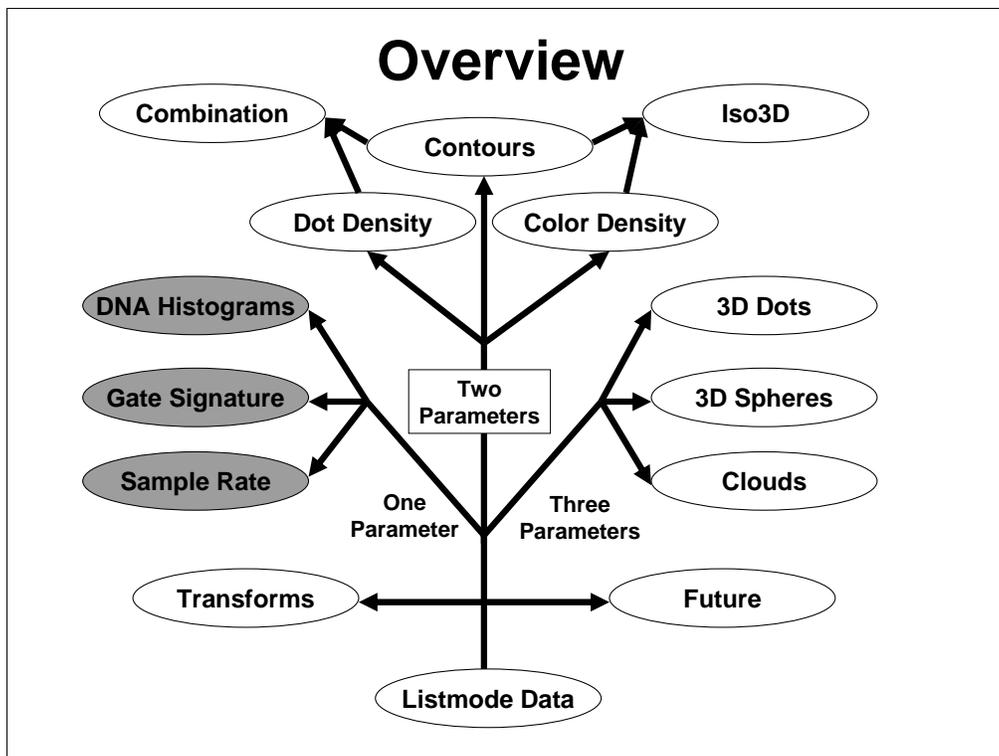# Displaying Flow Data

C. Bruce Bagwell MD, Ph.D.
Verity Software House, Inc.

Our understanding of flow cytometry data is highly influenced by how we visualize our data.  Many display techniques have advantages and disadvantages that dictate when they should be used or not used.  This lecture describes many of the more common visualization techniques in flow cytometry.

The general outline of this talk is shown above.

# Listmode Data Structure

**FCS Keywords**          **Event Index**          **FCS Listmode Data**

| Name | KWValue |
|------|---------|
| PRGNAM | ACQ8  V5.2c |
| $DATE | 16-APR-91 |
| $TOT | 30300 |
| $BTIM | 09:59:01.00 |
| $ETIM | 09:59:53.00 |
| $MODE | L |
| $PAR | 6 |
| RATE | 590.1849 |
| OVERRUNS | 0 |
| KINPTS | 0 |
| KINVEC | 0 |
| KINETICS | n |
| KINTIC | 0 |
| $P1R | 1024 |
| $P1B | 16 |
| $P2R | 1024 |
| $P2B | 16 |
| $P3R | 1024 |
| $P3B | 16 |
| $P4R | 1024 |
| $P4B | 16 |
| $P5R | 1024 |
| $P5B | 16 |
| $P6R | 1024 |
| $P6B | 16 |
| PARIN | 6 |

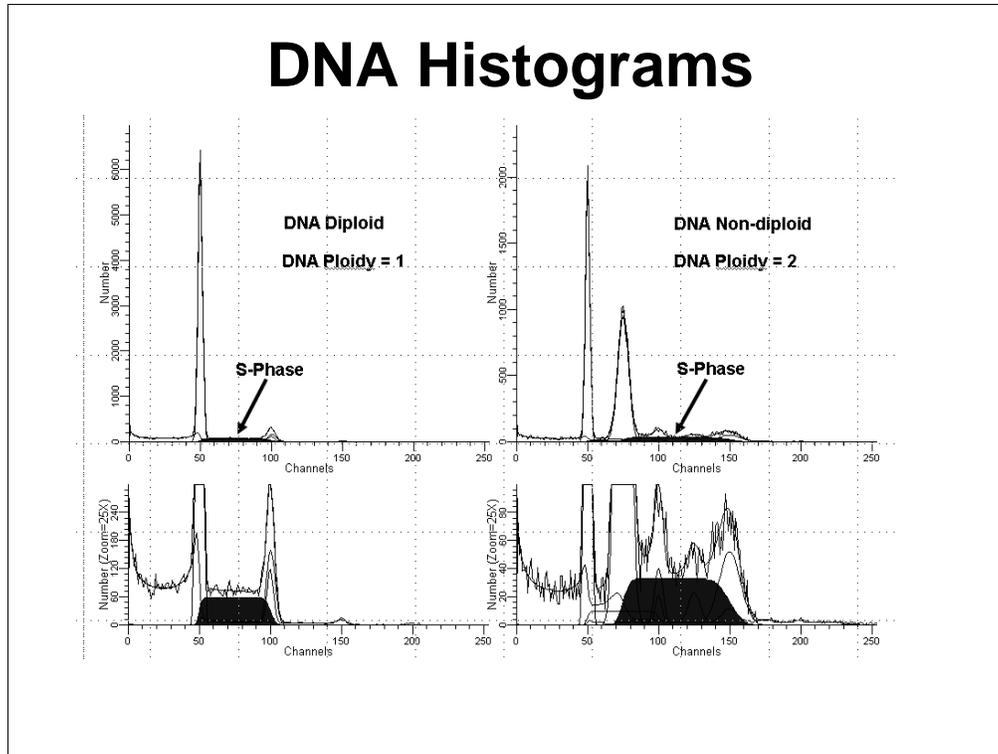|    | 1 | 2 | 3 | 4 | 5 | 6 |
|----|-----|-----|-----|-----|-----|-----|
| 1 | 232 | 392 | 149 | 68 | 253 | 67 |
| 2 | 216 | 387 | 549 | 317 | 210 | 83 |
| 3 | 231 | 346 | 561 | 321 | 215 | 126 |
| 4 | 367 | 450 | 281 | 274 | 585 | 583 |
| 5 | 236 | 356 | 121 | 138 | 154 | 100 |
| 6 | 240 | 427 | 187 | 141 | 130 | 52 |
| 7 | 235 | 376 | 190 | 503 | 105 | 102 |
| 8 | 226 | 361 | 125 | 97 | 193 | 70 |
| 9 | 221 | 391 | 192 | 102 | 160 | 125 |
| 10 | 337 | 484 | 165 | 279 | 570 | 573 |
| 11 | 216 | 403 | 138 | 91 | 133 | 136 |
| 12 | 216 | 414 | 215 | 520 | 102 | 47 |
| 13 | 227 | 474 | 82 | 70 | 642 | 301 |
| 14 | 351 | 535 | 177 | 230 | 558 | 581 |
| 15 | 270 | 443 | 89 | 72 | 480 | 219 |
| 16 | 227 | 322 | 299 | 494 | 171 | 143 |
| 17 | 341 | 591 | 174 | 285 | 587 | 588 |
| 18 | 217 | 345 | 146 | 529 | 137 | 92 |
| 19 | 261 | 357 | 79 | 76 | 431 | 134 |
| 20 | 341 | 577 | 225 | 317 | 546 | 574 |

One parameter histograms generally use one listmode data column*

■ ■ ■          **ADC Value**

Most flow cytometry data are stored in a listmode data format (see above).  A listmode data structure can be viewed as a long table with parameters as columns and events as rows.  The individual cells of the table are ADC (analog-to-digital) values.  A description of the ADC parameters' resolution and byte size, as well as other important information, is encoded in a list of keywords (see left inset).  To create single parameter histograms, only one column of the listmode data needs to be examined (e.g. highlighted area).
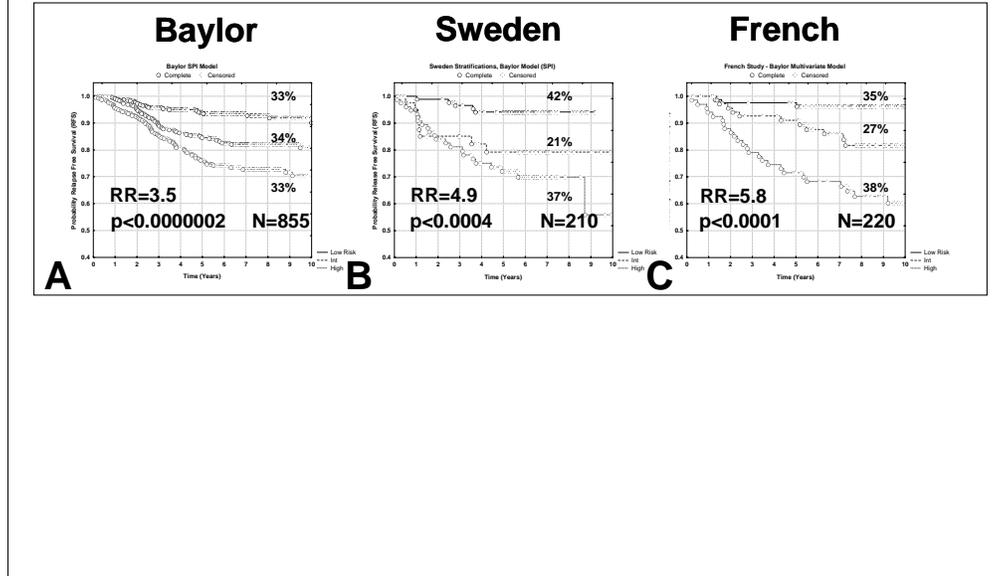
To create a one parameter or single parameter histogram, each ADC value in the parameter column is processed sequentially. Processing generally includes a reduction in scale, usually a power-of-two, to a histogram channel. It can also include a transformation to log or HyperLog/Biexponential space. Once a histogram channel is calculated, it serves as a lookup index into a special array of numbers that keeps track of the frequency of that channel. We generally refer to the plot of channel number versus the frequency as a histogram, or in this case, a single parameter histogram.
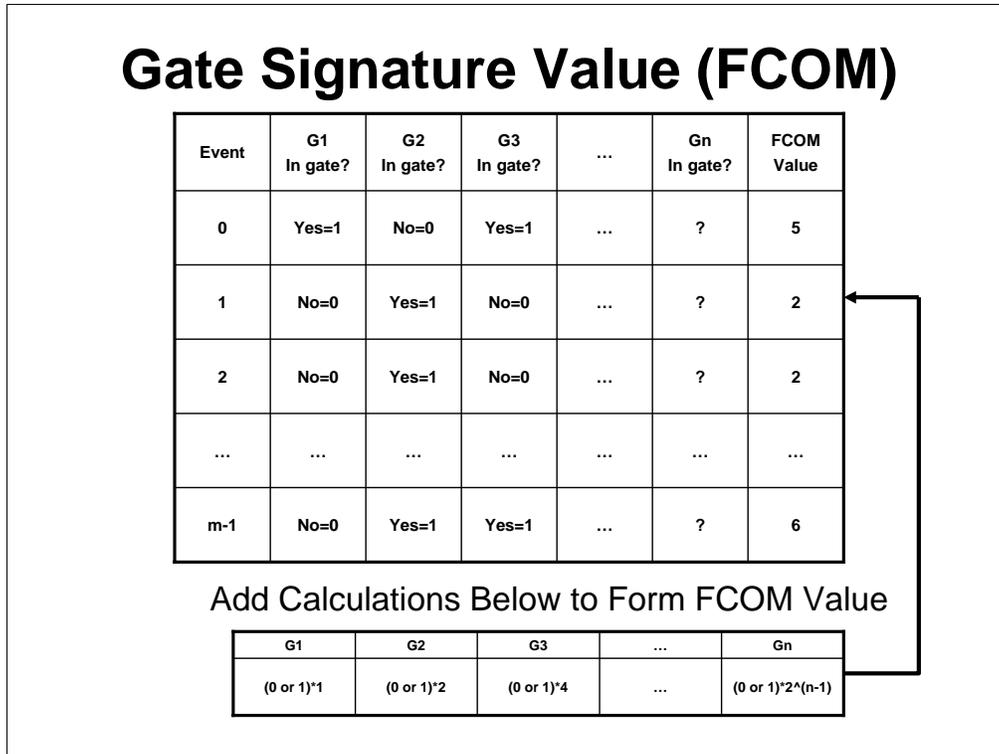
Single parameter histograms have limited use for phenotyping but there are a few applications where single parameter histograms are quite useful and, in some cases, necessary. Single parameter DNA histograms (see above) are necessary to model G0G1, S, and G2M populations, where it is very important to account for population overlap. Once DNA ploidy and S-phase are estimated, they can be incorporated into a prognostic model that provides important risk of relapse information to oncologists.
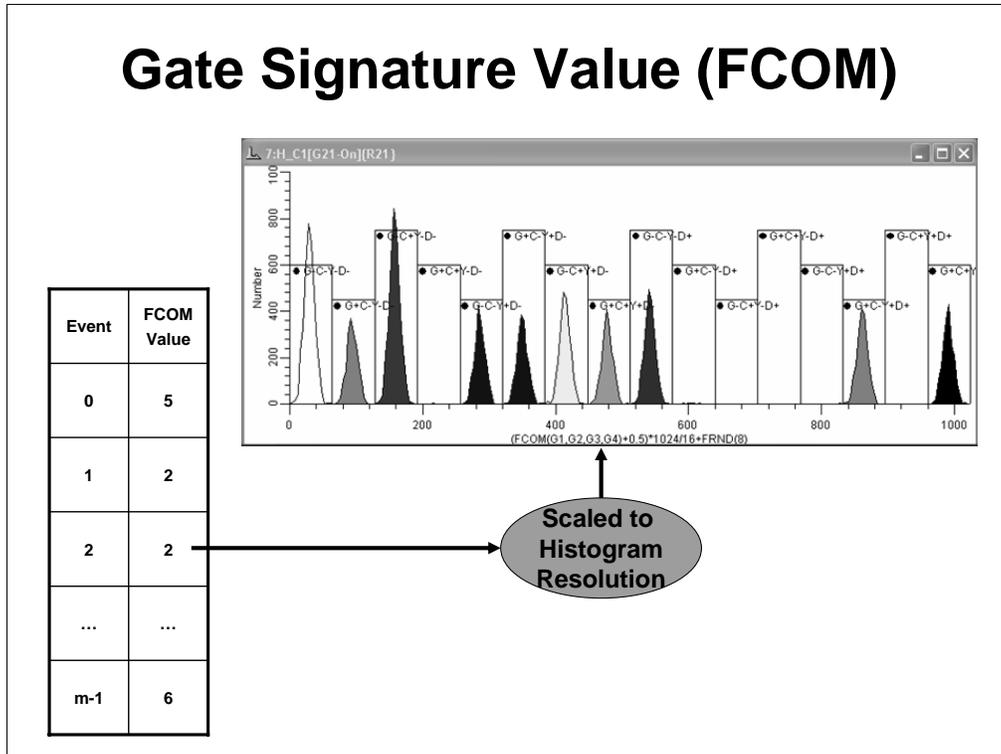
The above set of three panels shows an application of these DNA S-Phase/Ploidy prognostic models in stratifying node-negative breast cancer patients into low, intermediate, and high risk groups. Ref:  Bagwell CB, Clark GM, Spyratos F, Chassevent A, Bendahl P-O, Stål O, Killander D, Jourdan ML, Romain S, Hunsberger B, and Baldetorp B:  Optimizing Flow Cytometric DNA Ploidy and S-Phase Fraction as Independent Prognostic Markers for Node-Negative Breast Cancer Specimens, Communications in Clinical Cytometry, Vol 46:3, pps 121-135, 2001.  This paper is provided in the laboratory notes for the DNA laboratory.
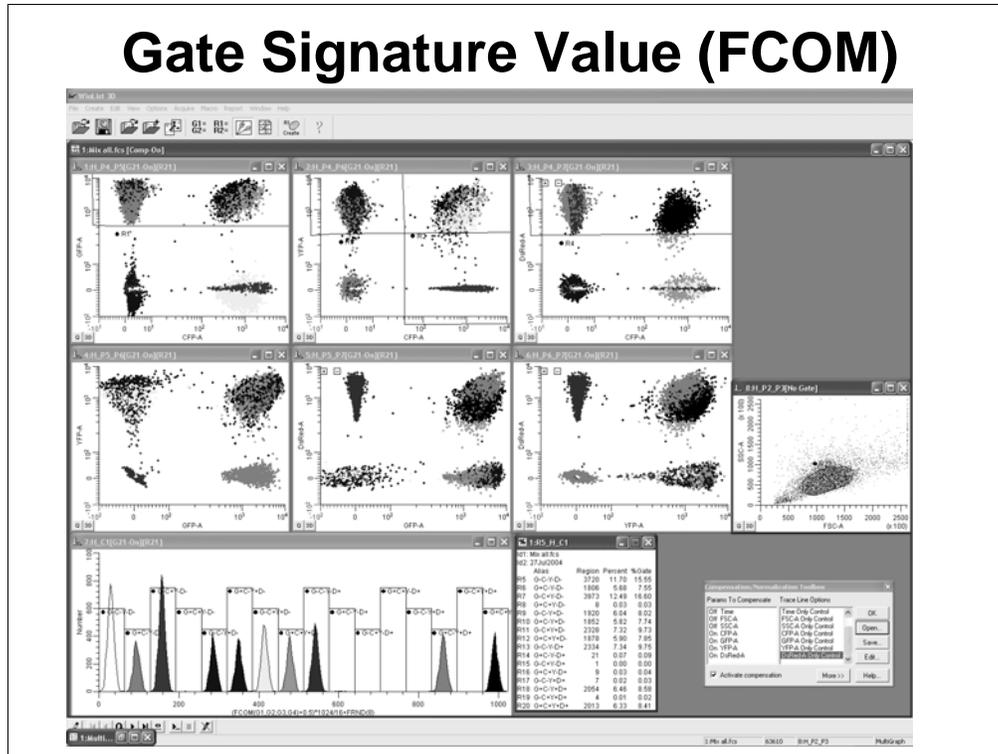
# Gate Signature Value (FCOM)

| Event | G1 In gate? | G2 In gate? | G3 In gate? | ... | Gn In gate? | FCOM Value |
|---|---|---|---|---|---|---|
| 0 | Yes=1 | No=0 | Yes=1 | ... | ? | 5 |
| 1 | No=0 | Yes=1 | No=0 | ... | ? | 2 |
| 2 | No=0 | Yes=1 | No=0 | ... | ? | 2 |
| ... | ... | ... | ... | ... | ... | ... |
| m-1 | No=0 | Yes=1 | Yes=1 | ... | ? | 6 |

### Add Calculations Below to Form FCOM Value

| G1 | G2 | G3 | ... | Gn |
|---|---|---|---|---|
| (0 or 1)*1 | (0 or 1)*2 | (0 or 1)*4 | ... | (0 or 1)*2^(n-1) |

Another interesting example of an important use of single parameter histograms is the Gate Signature histogram or FCOM. The above figure shows how it works. Each gate that will be involved in clustering events appears as a column in our listmode data structure, much like our primary parameters (see slide 3). If the event is in the gate, it gets a score of 1; otherwise, it gets a zero. The above figure shows an example of three gates; G1, G2, and G2, being involved as FCOM gates. These scores of 1 or 0's form a binary expression that ultimately produces a number. For example, the first event has bits 1 0 1 (see above), which forms the number 5 ($1*1 + 0*2 + 1*4 = 5$, see key at the bottom of the graph). These numbers represent unique signatures or combinations for each combination of gates, thus the name, Gate Signature or FCOM (function combination).
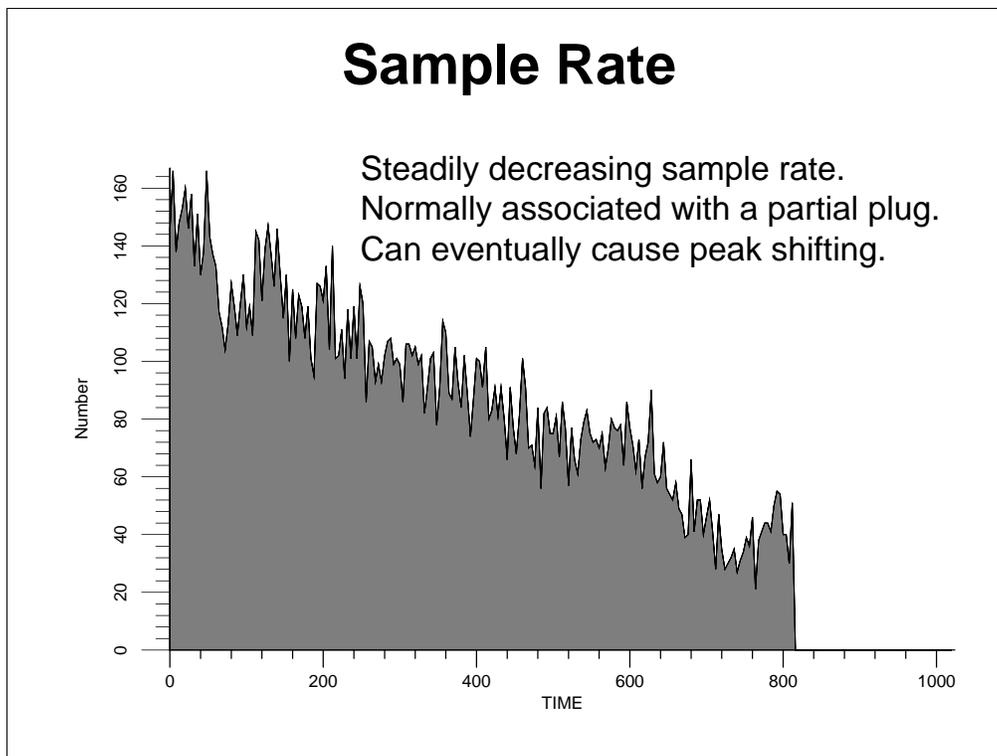
The Gate Signature or FCOM parameter was originally conceived by Jim Wood at Coulter Electronics in the 80's. The product was called PRISM and was restricted to two or three rectangular gates.
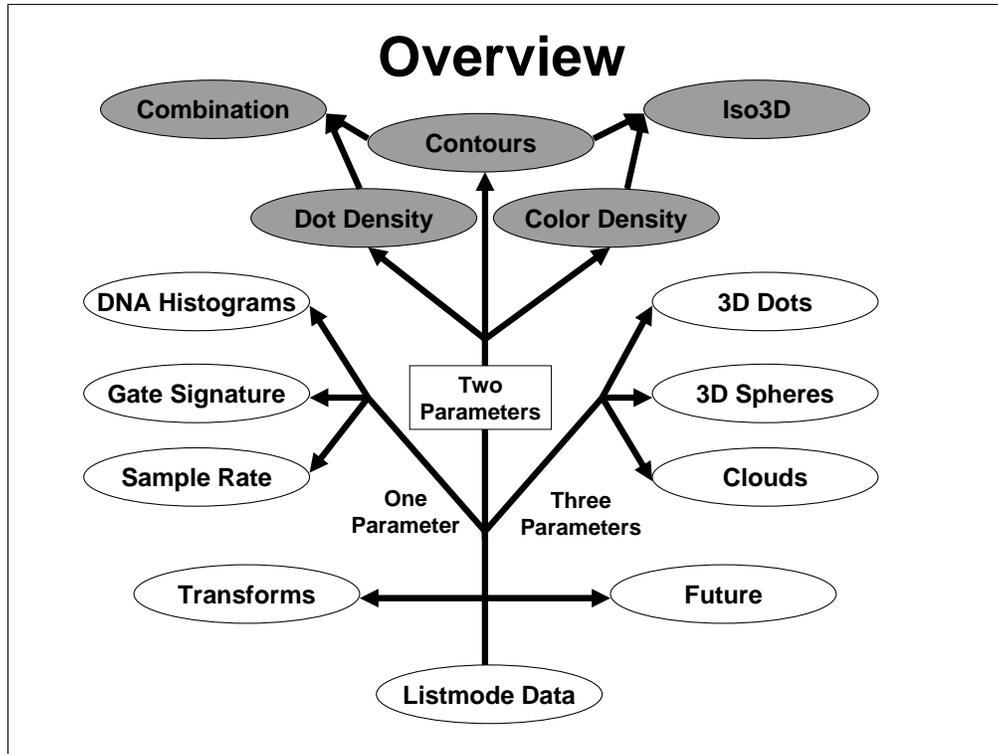
These Gate Signature or FCOM values are then appropriately scaled to histogram channels to form a special type of single parameter histogram that shows all combination of events that satisfy a specific set of gates (see above for a four gate FCOM example).
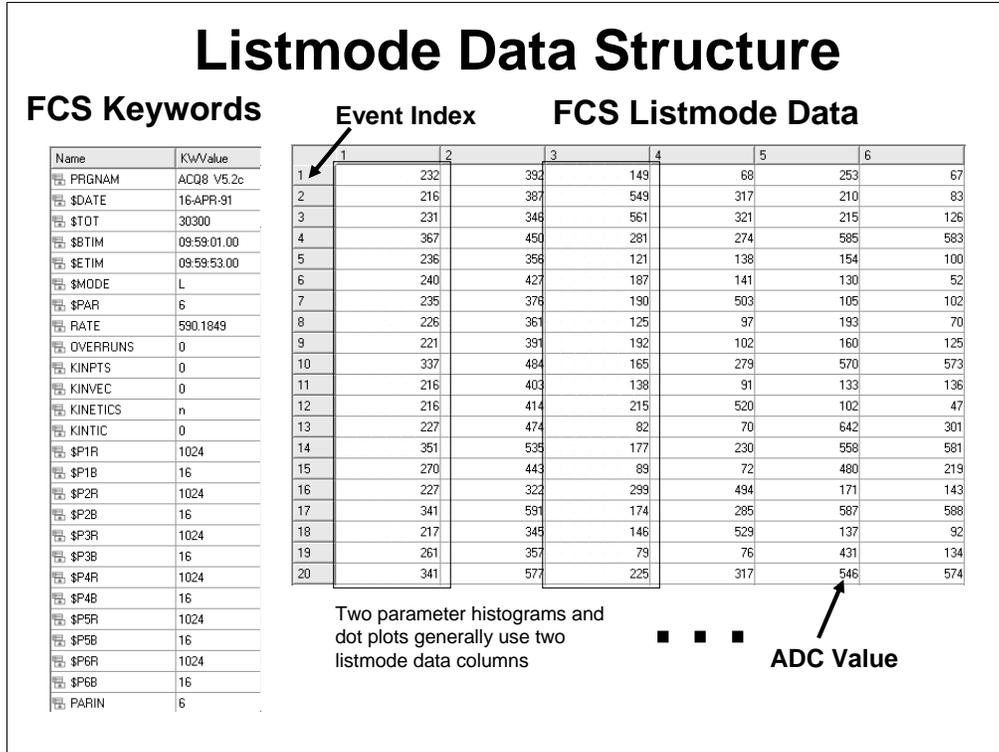
**Gate Signature Value (FCOM)**

This type of histogram allows a very compact way of color coding and enumerating clusters in other types of displays (see above). The data shown above was kindly provided by Teresa and Robert Hawley and shows four fluorescent proteins as surrogates for gene expression.

# Sample Rate

Steadily decreasing sample rate.
Normally associated with a partial plug.
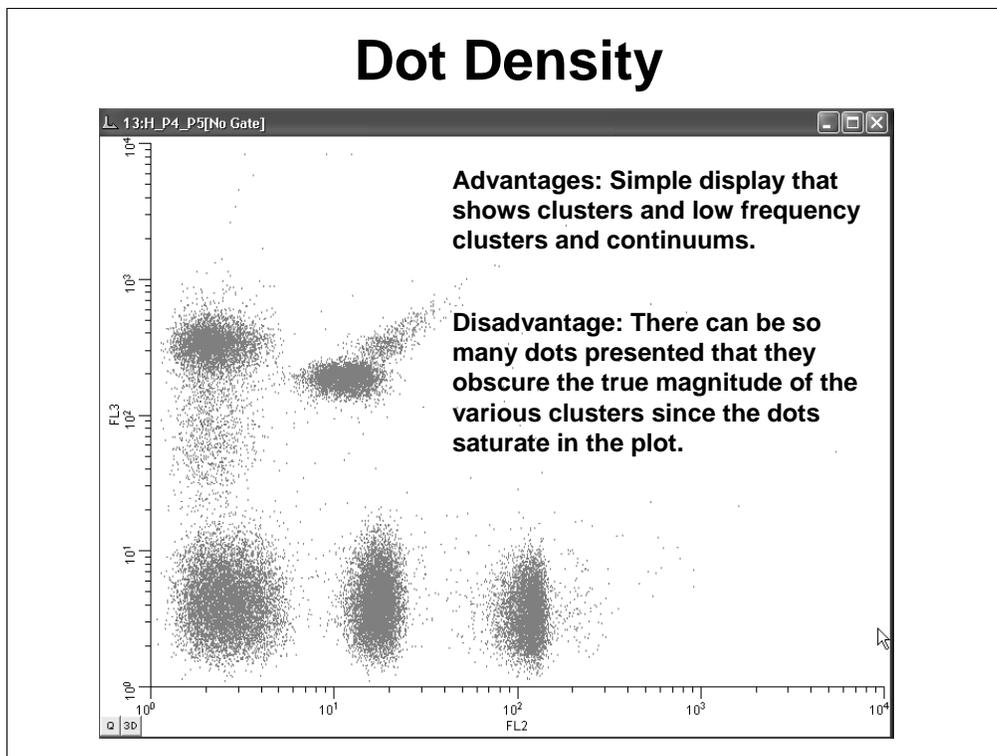Can eventually cause peak shifting.



Another interesting use of a single parameter histogram is the examination of sample rate as a function of time. Time as a parameter is an option that is available on most cytometers today and is one of the very best quality control parameters available. If you create a single parameter Time histogram, you can assess the sample rate over the duration of an acquisition. As shown above, if the cytometer is starting to form a partial plug, it can easily be detected with this type of plot.
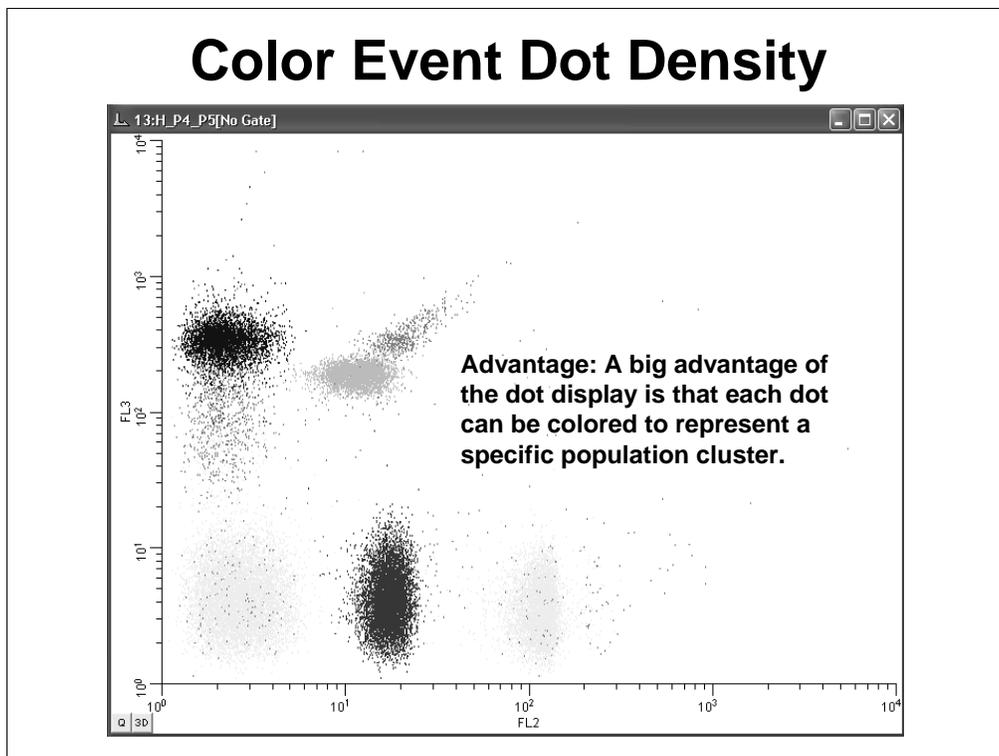
# Overview



Let's now examine some popular two parameter displays.
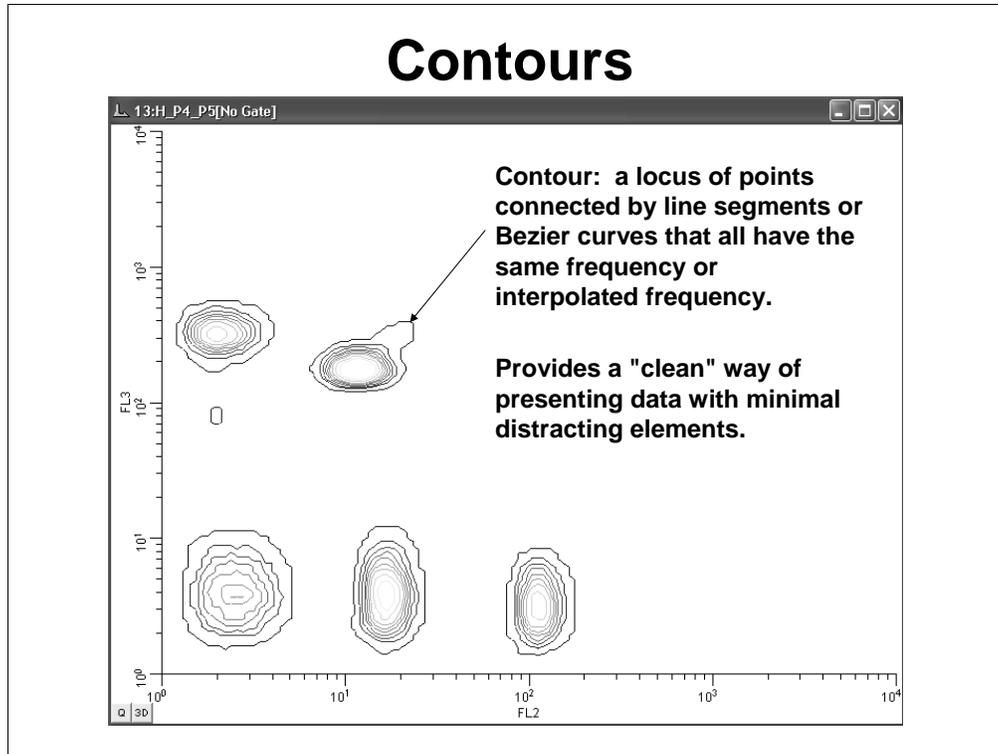
**Listmode Data Structure**

In order to create a two parameter histogram or plot, we need to use two parameters in our listmode data structure (see highlighted parameter columns above).  The real advantage of two parameter displays over single parameter displays is that you can visualize two parameter correlations.  In most cases, it is necessary to create a two parameter histogram via the same mechanism as described earlier for the single parameter histogram.  In other cases, like the dot density plot, we can plot the data points directly.

# Dot Density



The dot density plot (see above) is probably one of the most popular ways of displaying flow data.  The dots form density clusters that can be directly visualized. Small continuums and clusters are easily discernable using this type of display format.

## Color Event Dot Density

13:H_P4_P5[No Gate]

**Advantage: A big advantage of the dot display is that each dot can be colored to represent a specific population cluster.**
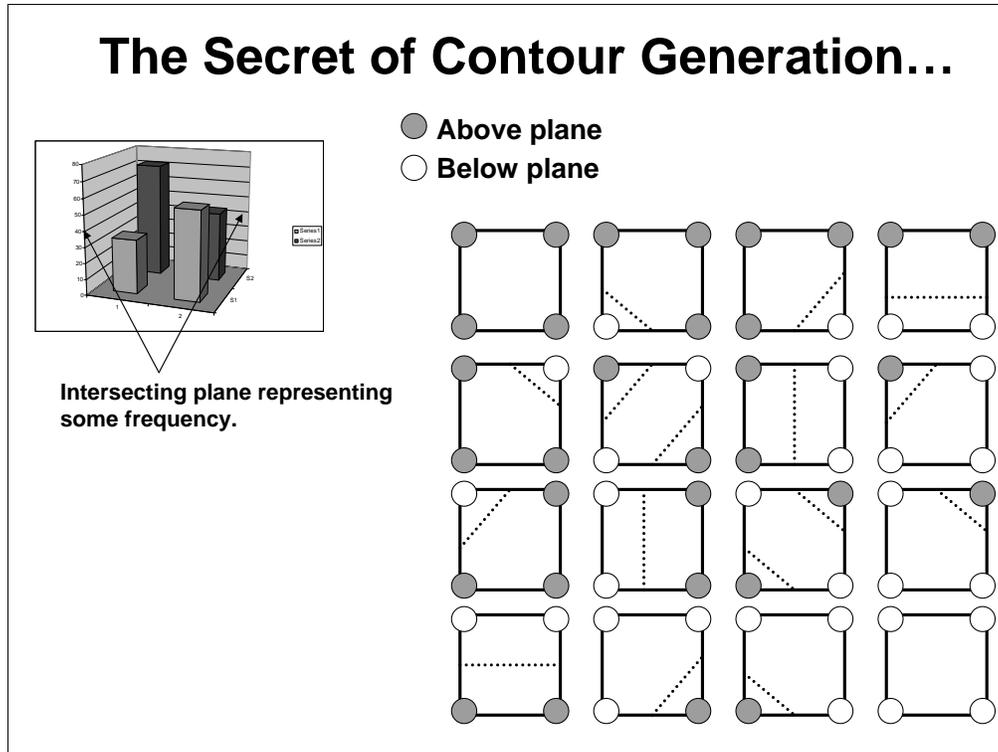
The big advantage to using a dot display is that each dot can be color coded to represent its inclusion within a specific population cluster.  These same colors appear in other displays allowing the visualization of various clusters in different histograms.  The original idea behind this form of display was Mike Logen's and appeared first in Becton Dickinson's Paint-A-Gate product.
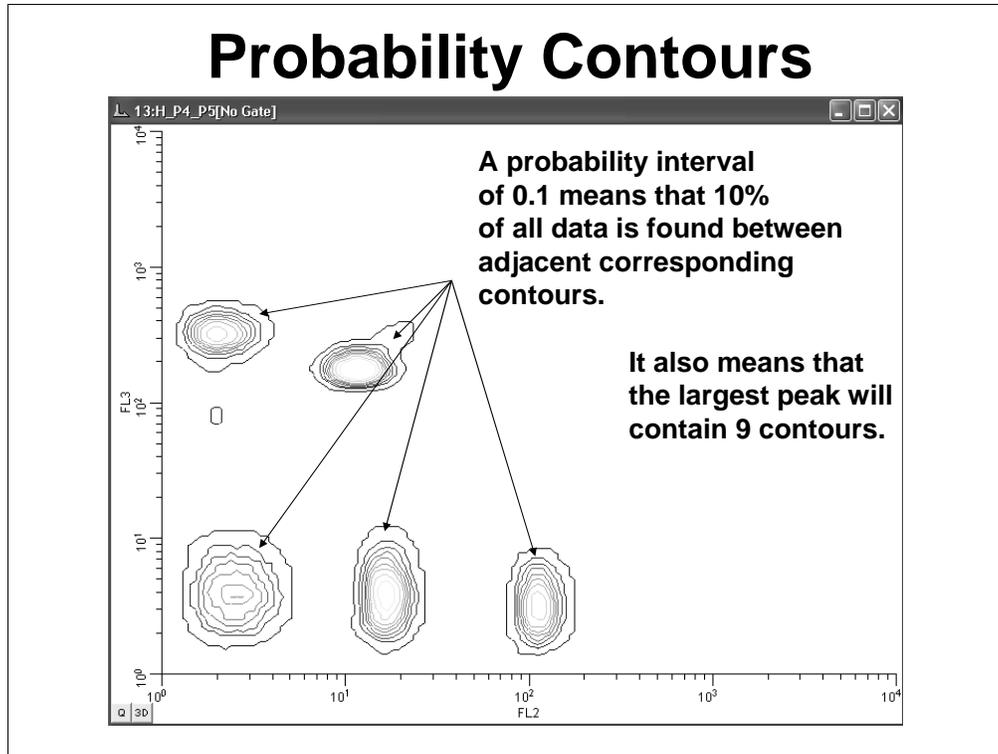
Another common method of presenting two parameter data is via contours (see above). A contour represents a locus of points connected by line segments or Bezier curves that all have the same frequency or interpolated frequency. Just like a quadrangle or weather map, the more contours a cluster has, the bigger the cluster. Contours provide a very clean way of viewing populations. There are not a lot of visual distractions in appreciating the number or the positions of the clusters present in a sample.
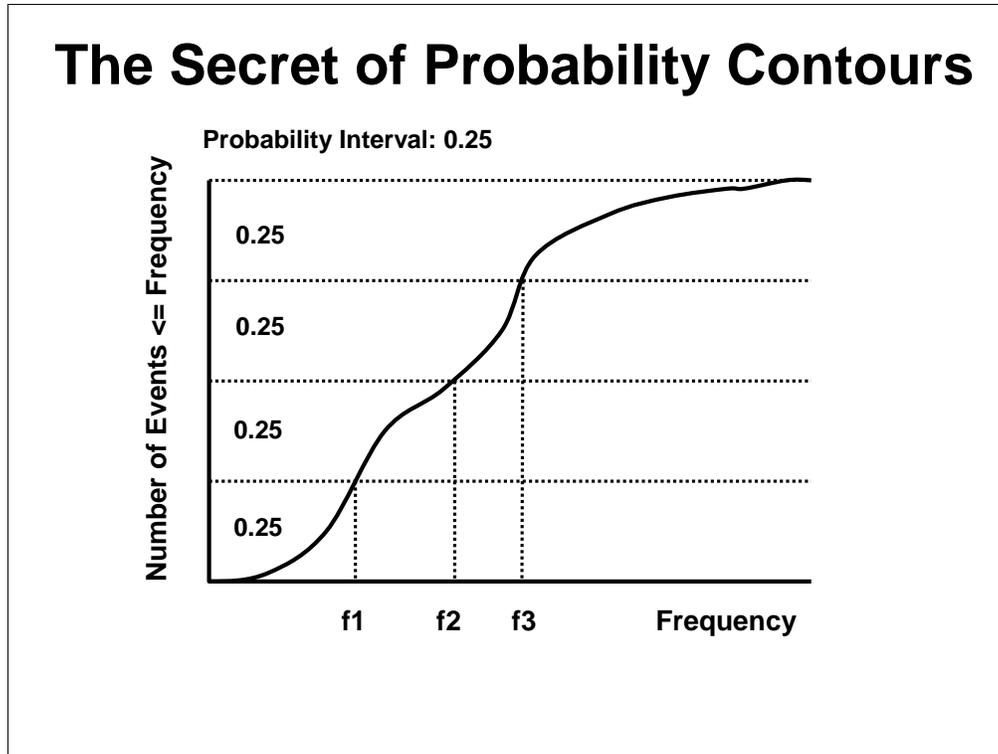
To create a contour map such as shown above, may seem like it would involve very complex calculations, but there is a secret to doing contours. The next slide shows how we do it.

**The Secret of Contour Generation…**

Intersecting plane representing some frequency.

Generating contour segments begins by considering only four neighboring channels (shown above). The initial four neighboring channels shown above are those that are nearest the origin. Imagine a plane at a particular frequency intersecting the four neighboring histogram cells. If the cell has a frequency that is greater than the frequency plane we assign it a dark color (see above key). If it is lower, we assign it a white color. We end up with 16 possible combinations (see above). If all the cells are either above or below the intersecting plane, no line segments (dotted lines) are drawn. For the other combinations, one or two line segments are possible. The system appropriately interpolates the end points of these line segments to define their exact location. The system then considers the next set of four neighboring points. By definition, the end of one line segment will be the beginning of another. By doing this state-type of analysis on all the cells in a two parameter histogram we obtain a complete set of contour segments. The above method is very fast and really not that hard to do in software.
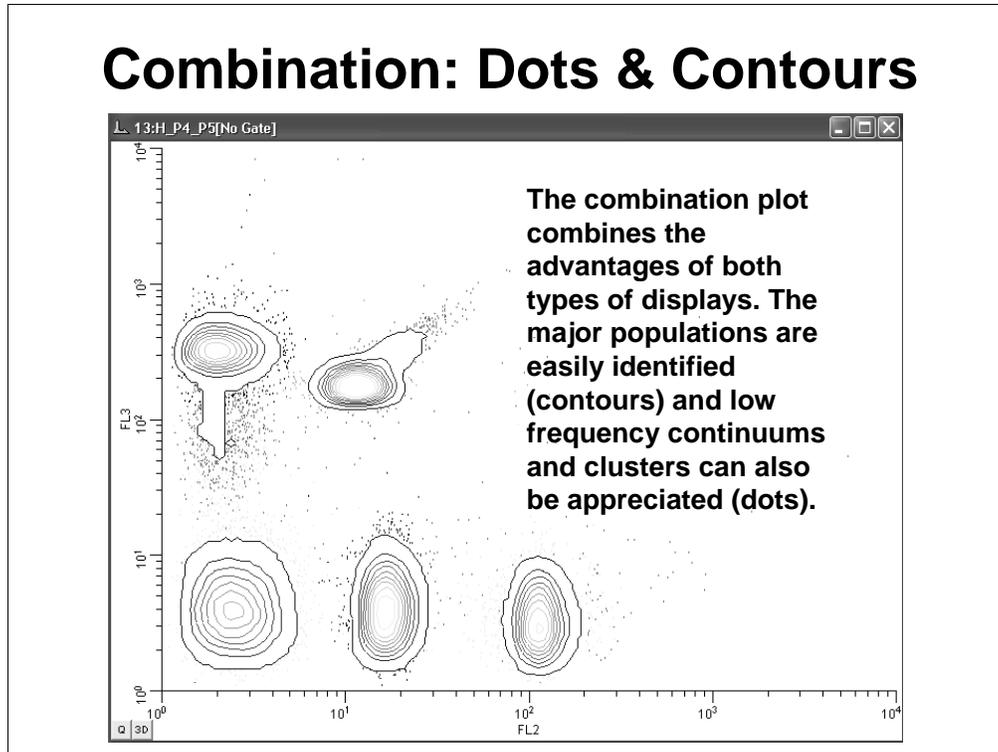
The most common way of spacing contours is by means of probability. Probability contours were initially conceived of by Wayne Moore and provide a good way of spacing the contours. The user generally selects a "probability interval", which determines the fraction of events that are found between adjacent corresponding contours. The higher the number, the less contours presented. The probability value is either a value ranging between 0 and 1 or a percent, ranging between 0 and 100%.

# The Secret of Probability Contours

**Probability Interval: 0.25**

**Number of Events <= Frequency**

0.25

0.25

0.25

0.25

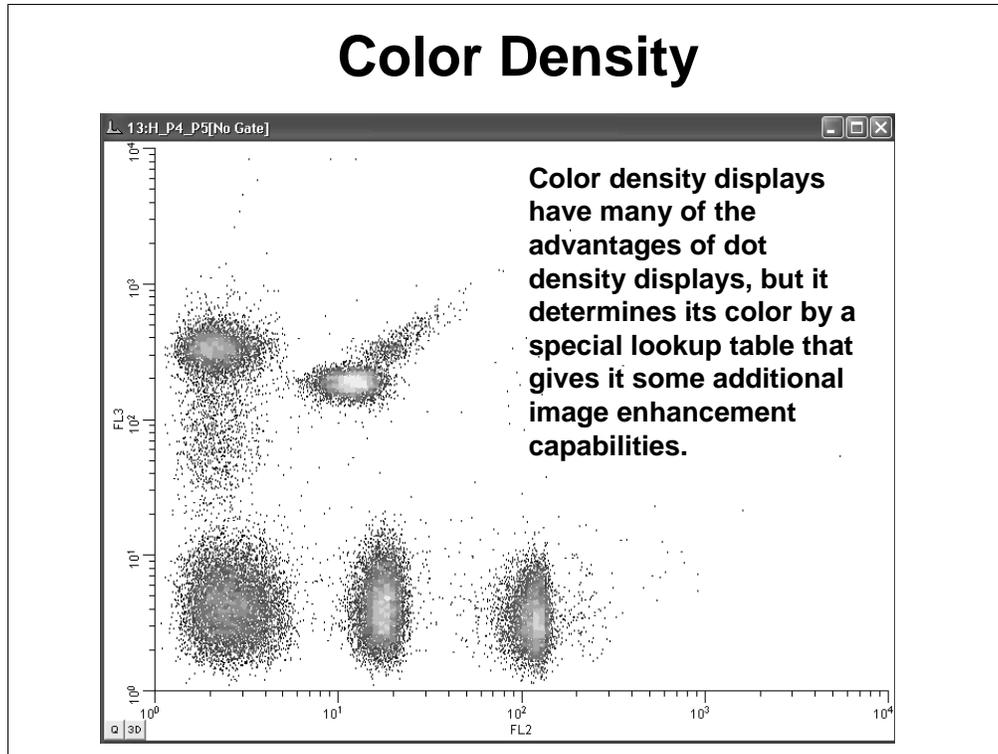f1    f2    f3          **Frequency**

The secret to creating probability contours is to first create a curve that represents the number of events equal to or below a frequency threshold. When the threshold frequency is zero, the sum is zero. As the frequency threshold increases (moving left-to-right on the x-axis, see above), the number of events under the threshold increases and approaches the total number of events in the listmode file.

If the Probability Interval is 0.25, we divide the y-axis into four equal partitions. The intersection of these partitions with the curve gives us the frequency thresholds to use for our contours. Thus, the interval between f1 and f2 comprises 0.25 or 25% of all the events. It's that simple!
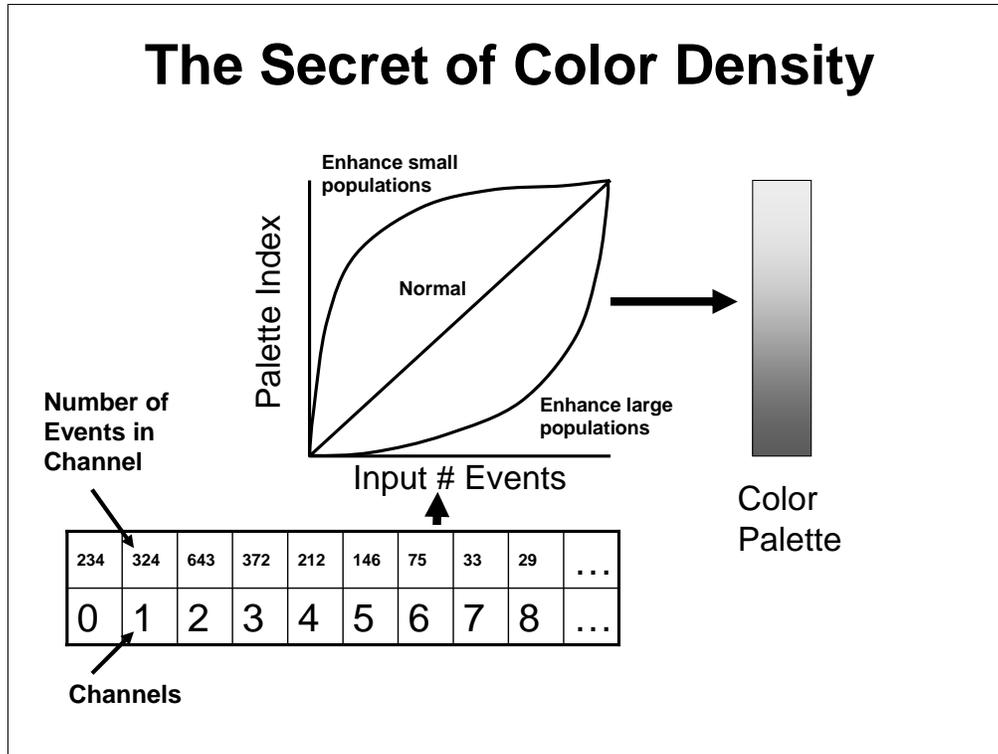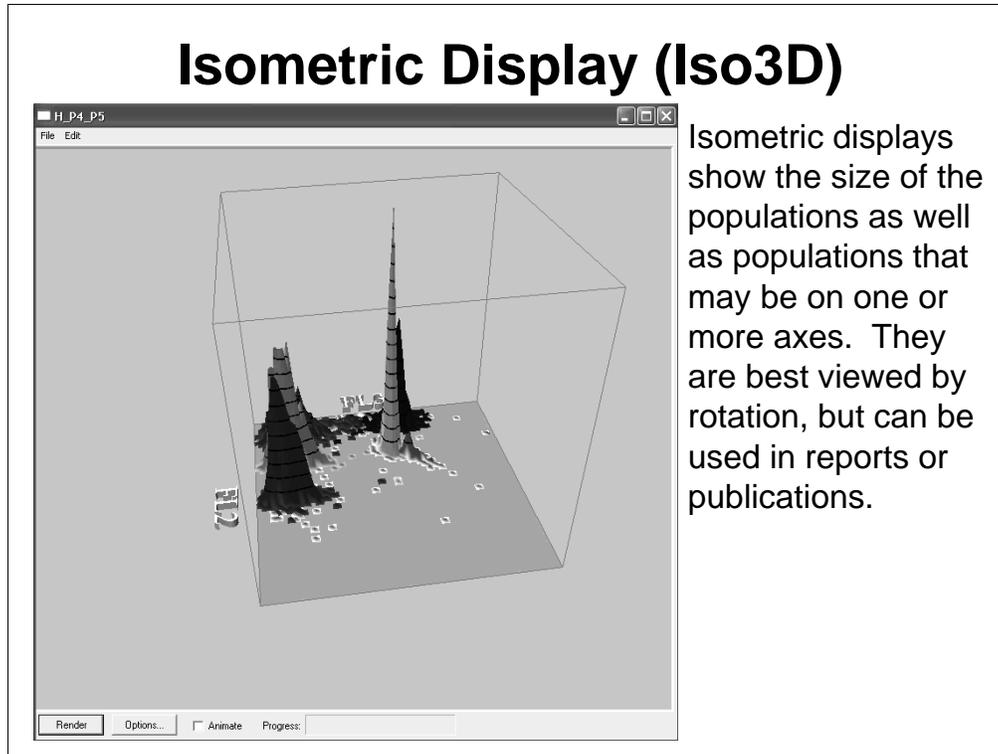
The combination of dots and contours gives the user the best of both types of display methods.  This type of display nicely shows the location of the populations (contours) but also shows subtle low frequency populations and continuums (dots).
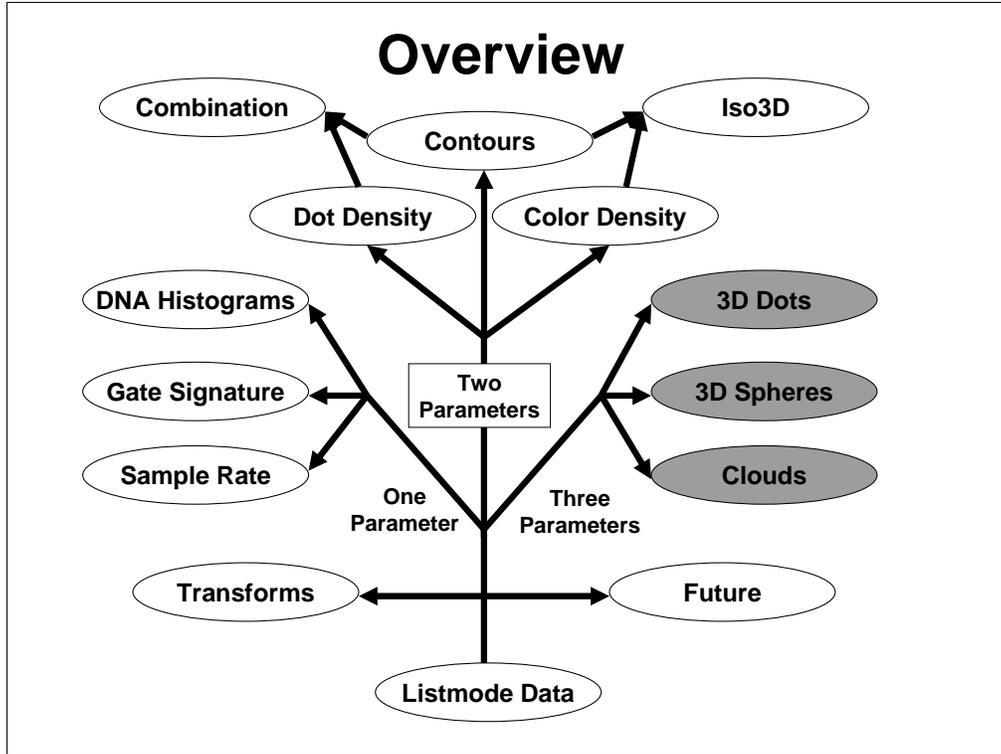
The color density display is very similar to the dot density display, except that it uses a special lookup table to find the color of the dots. This lookup table allows the user to enhance small populations or large populations. The way it works is shown on the next slide.
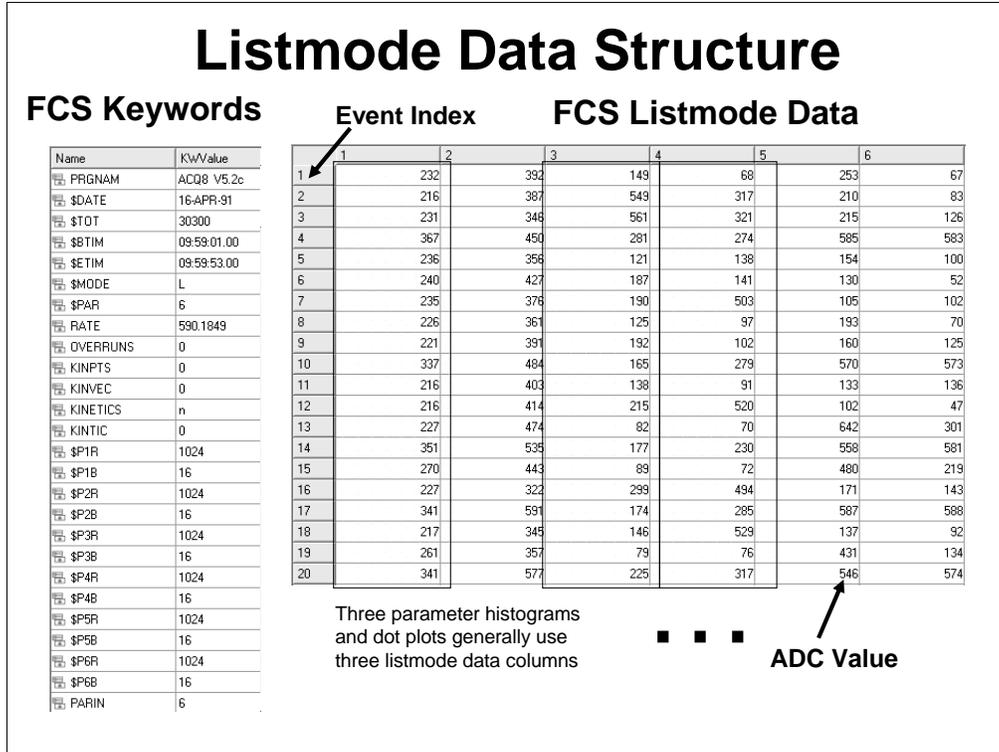
**The Secret of Color Density**

In the color density plot, the color for each dot (or rectangular square) comes from the frequency of the histogram at the dot's location, the type of color mapping used, and a color palette. Normally the color mapping is linear, but non-linear mappings (see middle graph) can enhance either small or large populations depending on the shape of the mapping function. The output of this mapping is then used to form an index into a color palette, which is usually a continuous set of colors starting at one color and ending at another. Thus, with this type of plot, you can choose to enhance different characteristics of your populations.

# Isometric Display (Iso3D)

Isometric displays show the size of the populations as well as populations that may be on one or more axes. They are best viewed by rotation, but can be used in reports or publications.
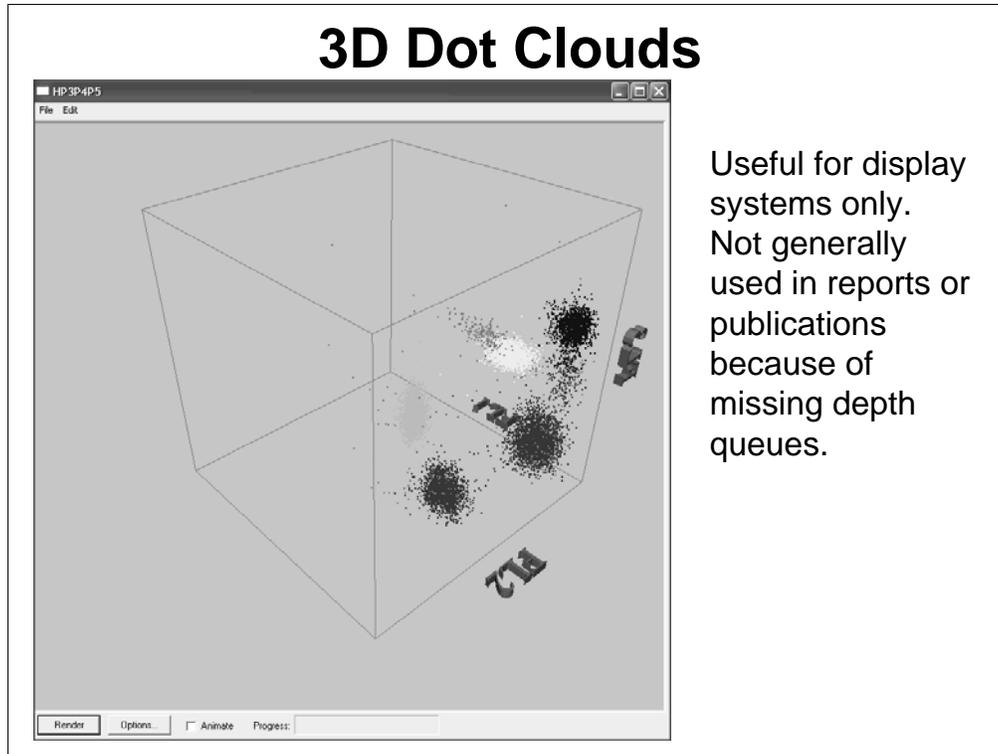
Isometric display are useful to appreciate population sizes as well as events that are piled on the axes (not shown above). They are best viewed rotating, but can be used in reports or publications.
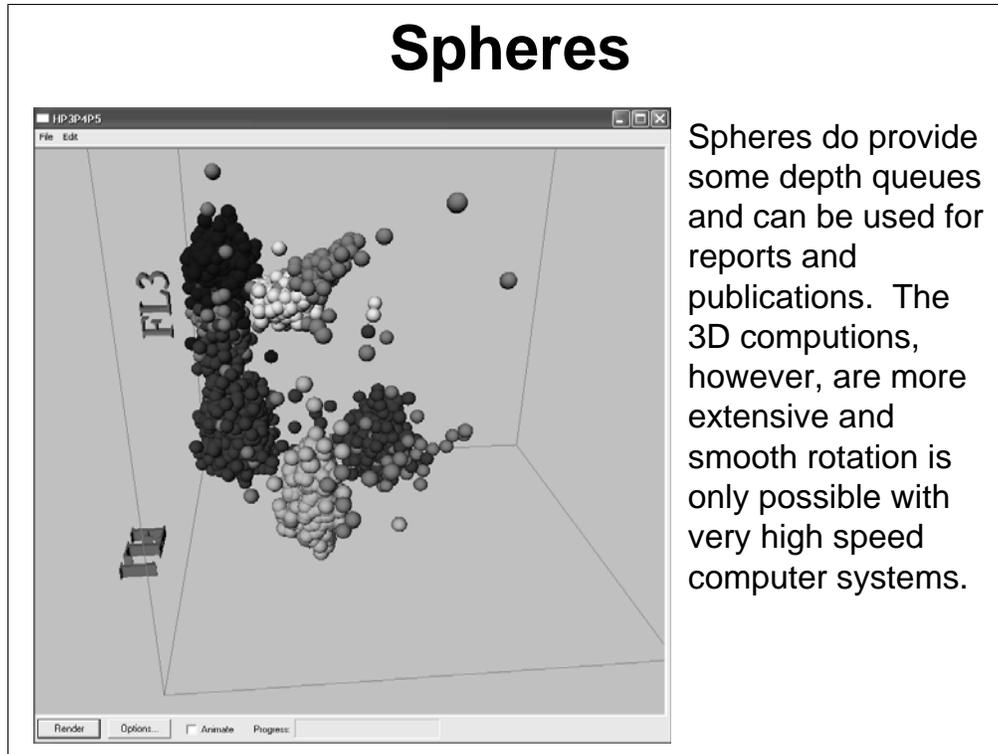
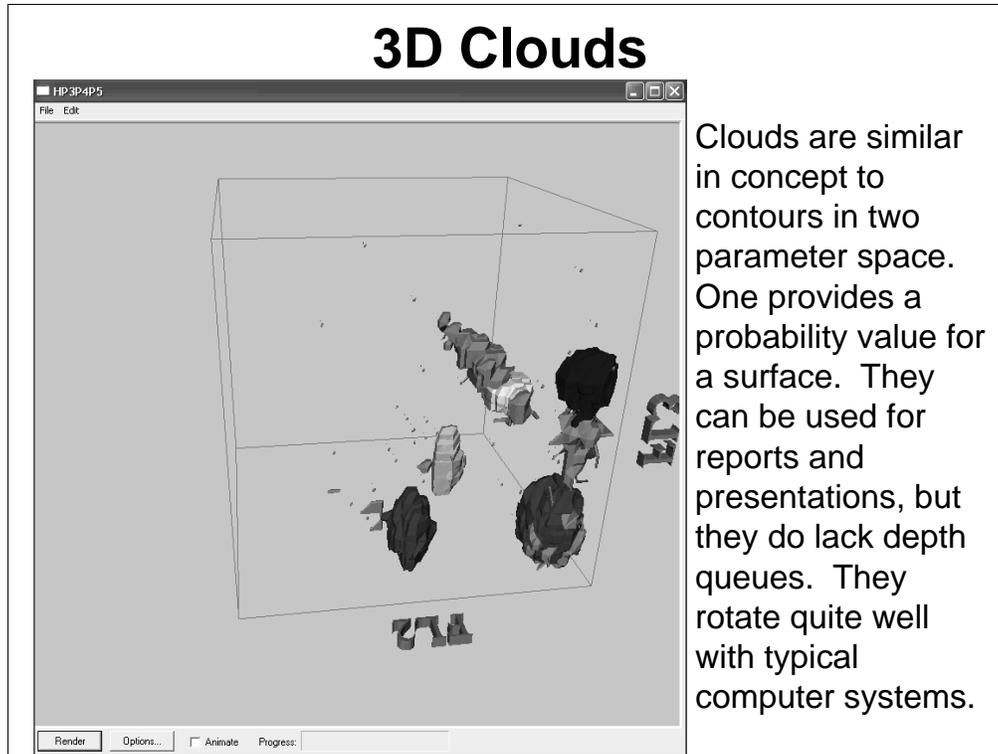We now look at some techniques for displaying three parameter data.

In order to create a three parameter histogram or plot, we need to use three parameters in our listmode data structure (see highlighted parameter columns above).  The real advantage of three parameter displays over single and dual parameter displays is that you can visualize parameter correlations in multiple dimensions.  In most cases, it is necessary to create a three parameter histogram via the same mechanism as described earlier for the single parameter histogram.  In other cases, like the dot density plot, we can plot the data points directly (3D dots and Spheres).

# 3D Dot Clouds

Useful for display systems only. Not generally used in reports or publications because of missing depth queues.
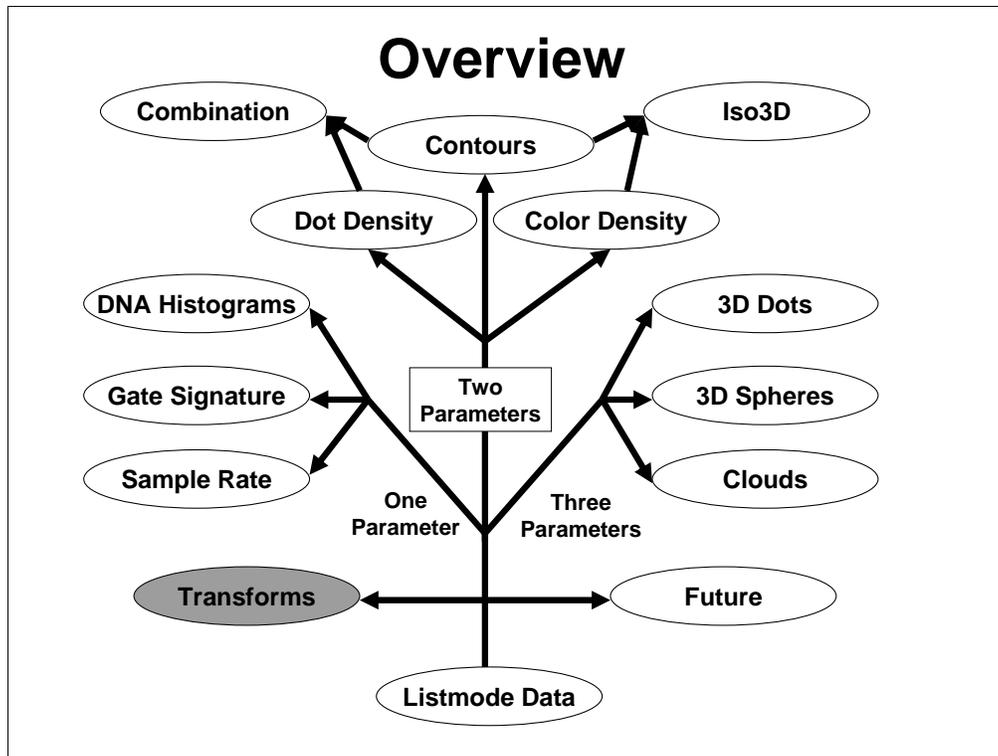
3D dots are really only useful when the display is rotating. They provide a means of appreciating the correlation between three parameters. It is usually not desirable to use these plots in reports or publications because the depth queues are missing and it is impossible to appreciate which clusters are in front and which are in back of the display.

# Spheres



Spheres do provide some depth queues and can be used for reports and publications. The 3D computions, however, are more extensive and smooth rotation is only possible with very high speed computer systems.
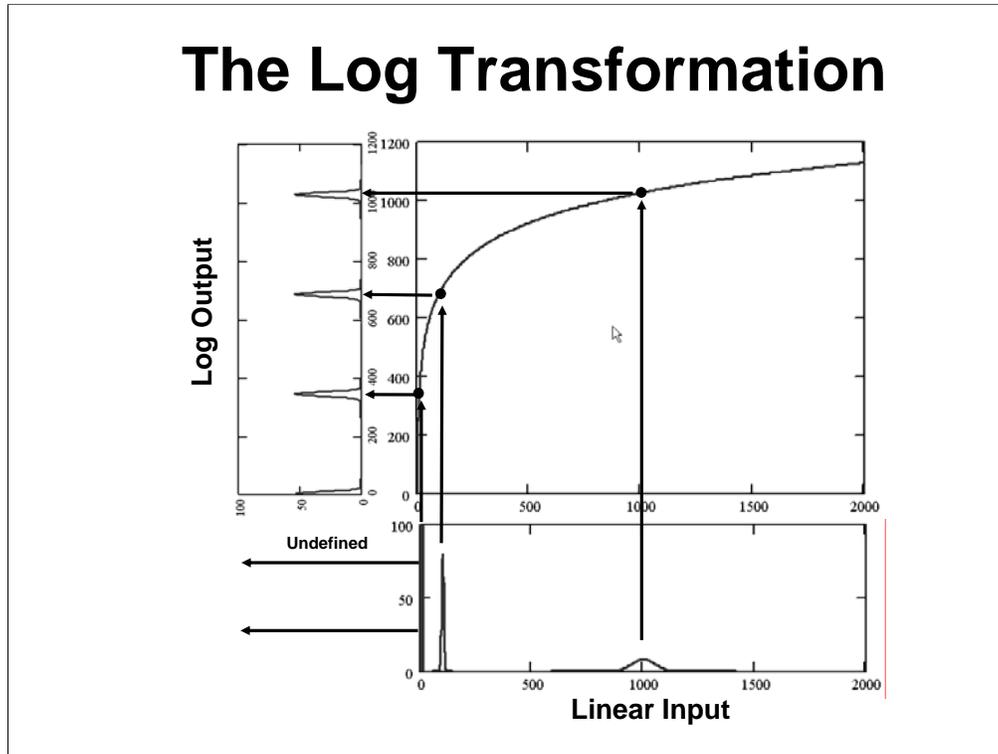
Of all the 3D options, spheres probably are the best for publication since they provide some visual depth queues with the radius of the spheres. The closer spheres are larger and the distant spheres are smaller. Their disadvantage, however, is that the number of computations is quite high for the spheres and smooth rotation of displays systems is really only possible with very high speed computer systems.

## 3D Clouds

Clouds are similar in concept to contours in two parameter space. One provides a probability value for a surface. They can be used for reports and presentations, but they do lack depth queues. They rotate quite well with typical computer systems.
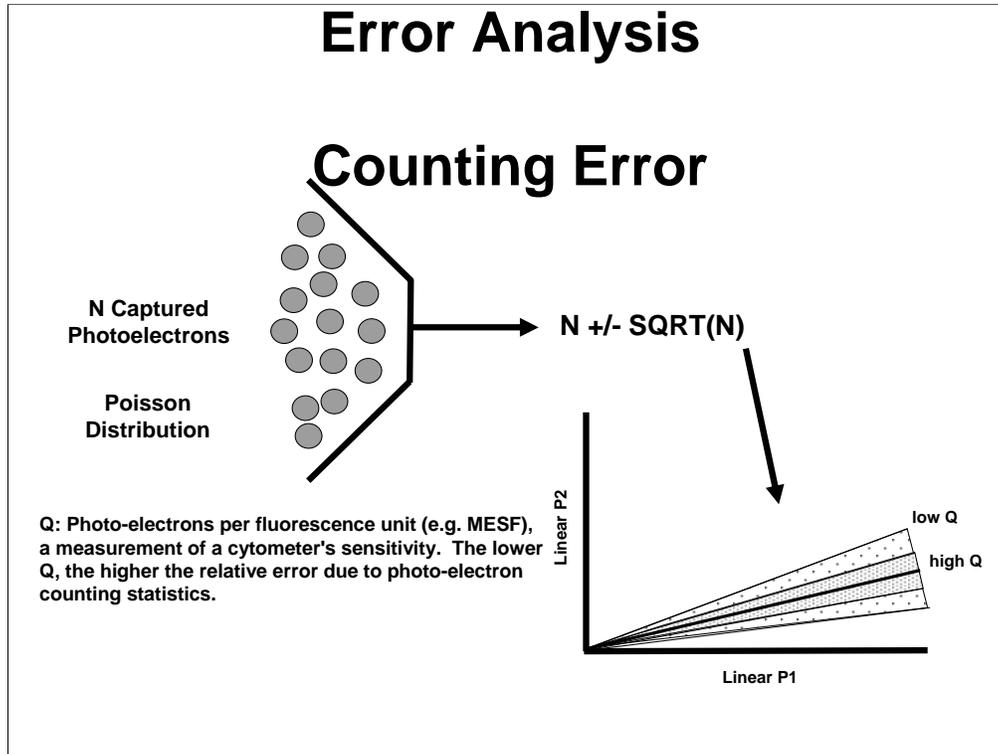
3D clouds are similar to contours in concept. A surface is defined that encloses a designated fraction of events. In the above implementation, a separate surface is defined for each color event cluster. Probability clouds can be used for reports and presentations since their complexity is normally fairly low, but they do lack the visual depth queues necessary to completely appreciate the dimensionality of the plot. They are very fast to compute and are therefore excellent to rotate in real-time.

# Overview

We now we discuss some popular transforms for flow data.
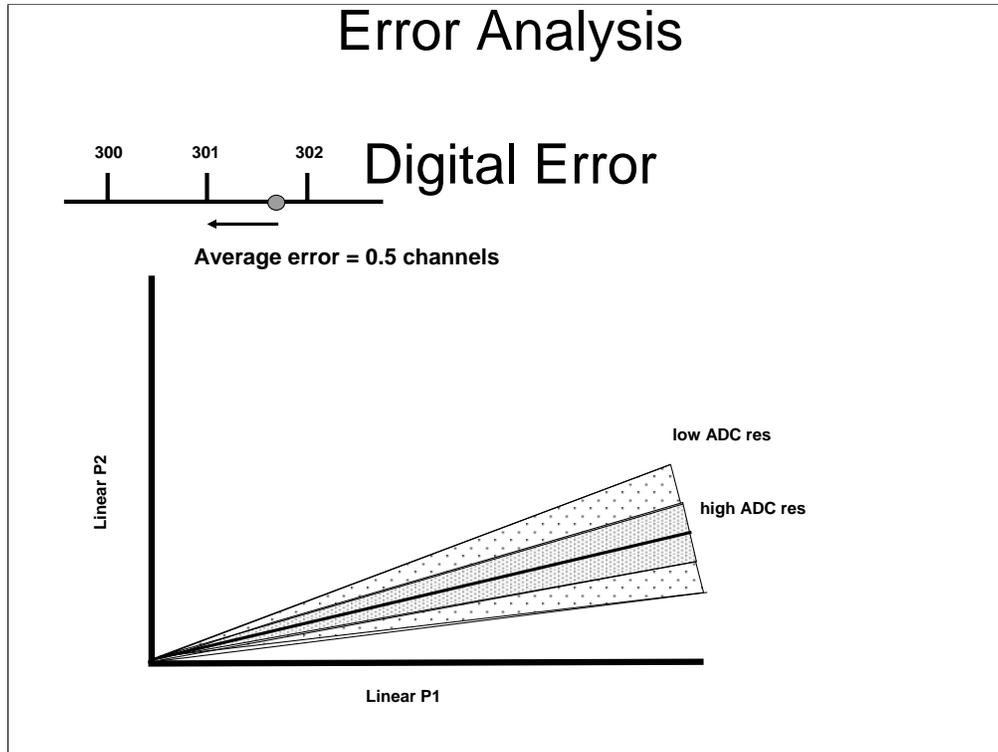
The Log Transformation

Since the very beginning of cytometry it has been recognized that the log transform is useful for displaying a wide dynamic range of signals. The other important characteristic of the log transform is that it normalizes populations of linearly increasing standard deviations to constant standard deviations. If the log transform didn't have this characteristic, it would be very difficult to appreciate populations separated by more than two decades of intensity. Unfortunately, the log transform is undefined at 0 or negative values (see above), which creates a few problems in the proper interpretation of compensated data.
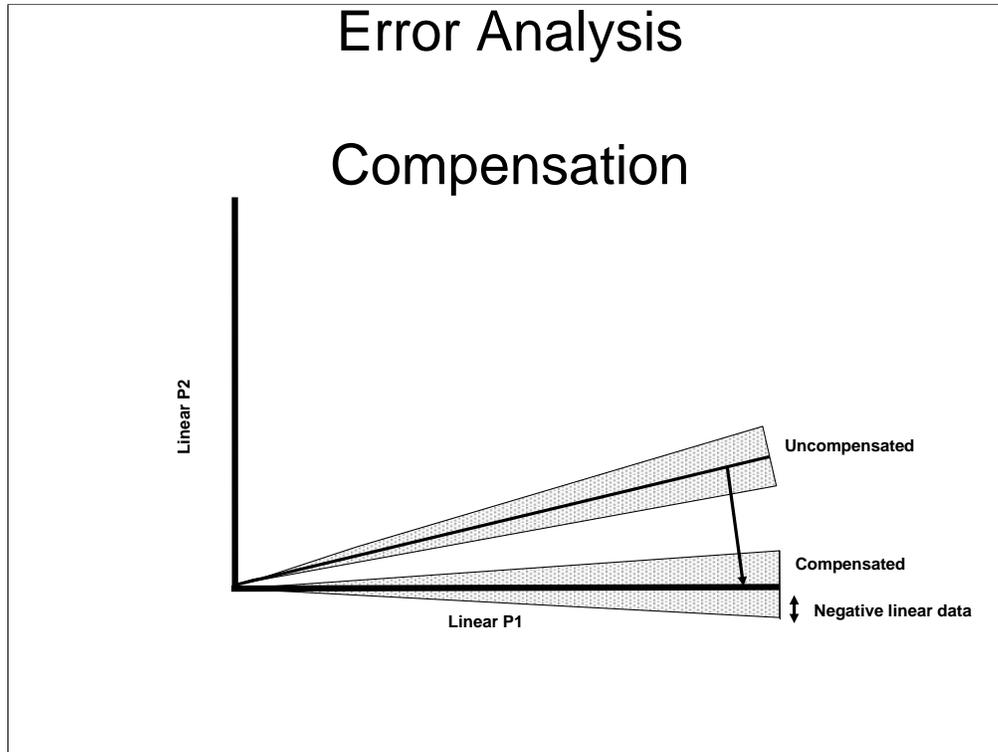
There are two major types of errors that greatly influence the display of our data, especially when its compensated. The first error is at the level of signal detection. Since the number of photo-electrons ultimately captured by the cytometer's PMT is finite, the signal quantitation is stochastic with a variance governed by the Poisson distribution. The square-root of the number of captured photo-electrons is approximately the standard deviation of the error distribution.

Q is a measurement of the cytometer's efficiency and has units of number of photo-electrons per unit fluorescence intensity (e.g. MESF). The lower the Q value, the higher the amount of relative counting error associated with the signal. The plot in the lower-right depicts this error for a single-color control primarily fluorescing in P1 and crossing-over into P2 with low and high Q values.
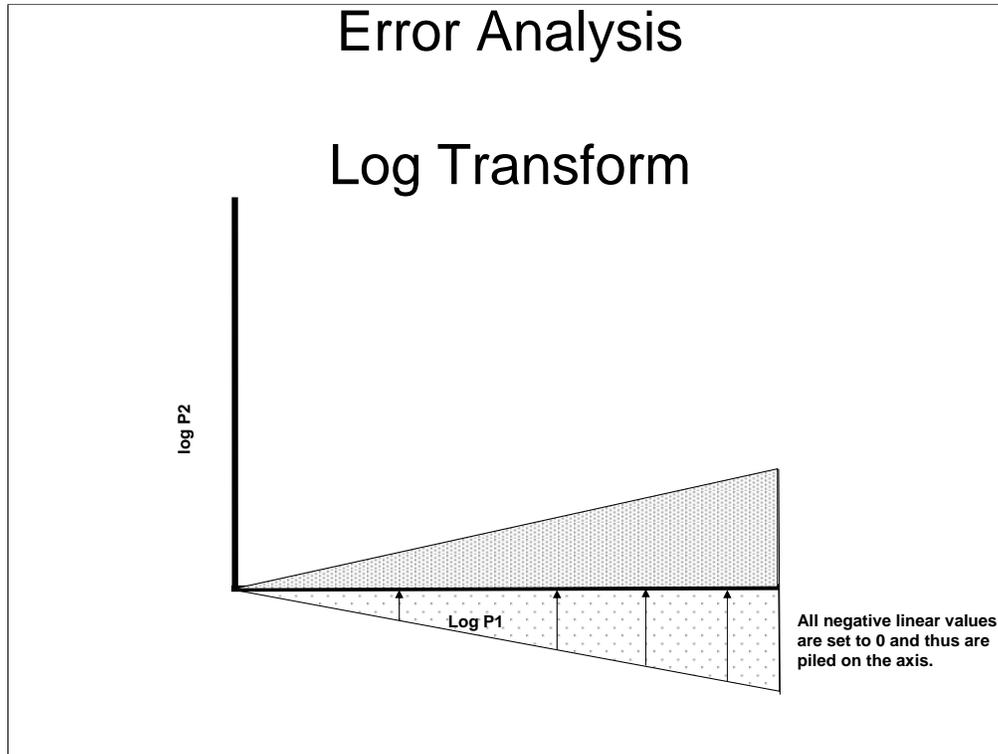
The second error we'll consider is the truncation of the analog log amplified signal to an integer value, digital error. On the average, digital error yields, on average, an error of ½ ADC channel. Unfortunately, this error is augmented by the log transform so that truncations at high ADC values can have a substantial error.

# Error Analysis

# Compensation

Linear P2

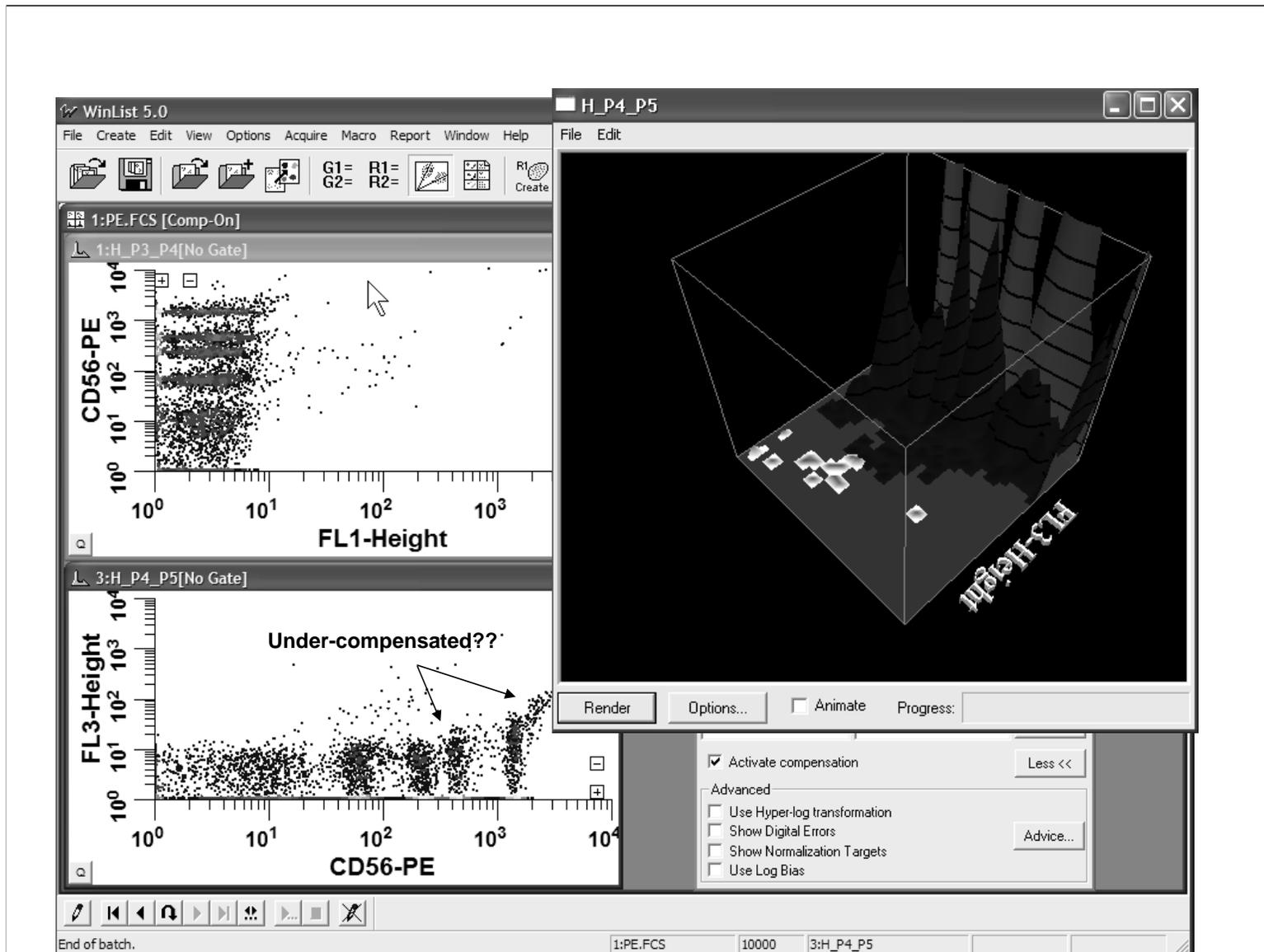Uncompensated

Compensated

Negative linear data

Linear P1

The combination of counting and digital error creates a symmetric error distribution about the single-color trace line.  When this data is properly compensated, this distribution becomes symmetric about a horizontal line with some events defined below the axis.  These negative values create several problems for us when we transform the axes to log.
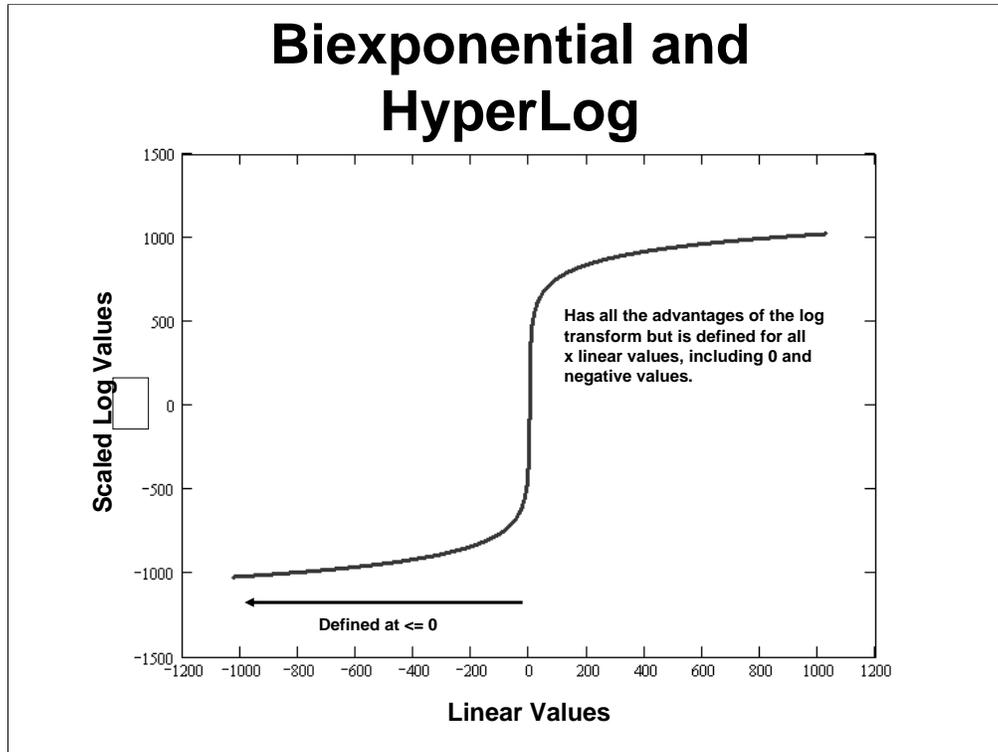
Note that the error exists before and after compensation.  Compensation does not create this error, it simply exposes it.

# Error Analysis

## Log Transform



**log P2** (vertical axis label)

**Log P1** (horizontal axis label)

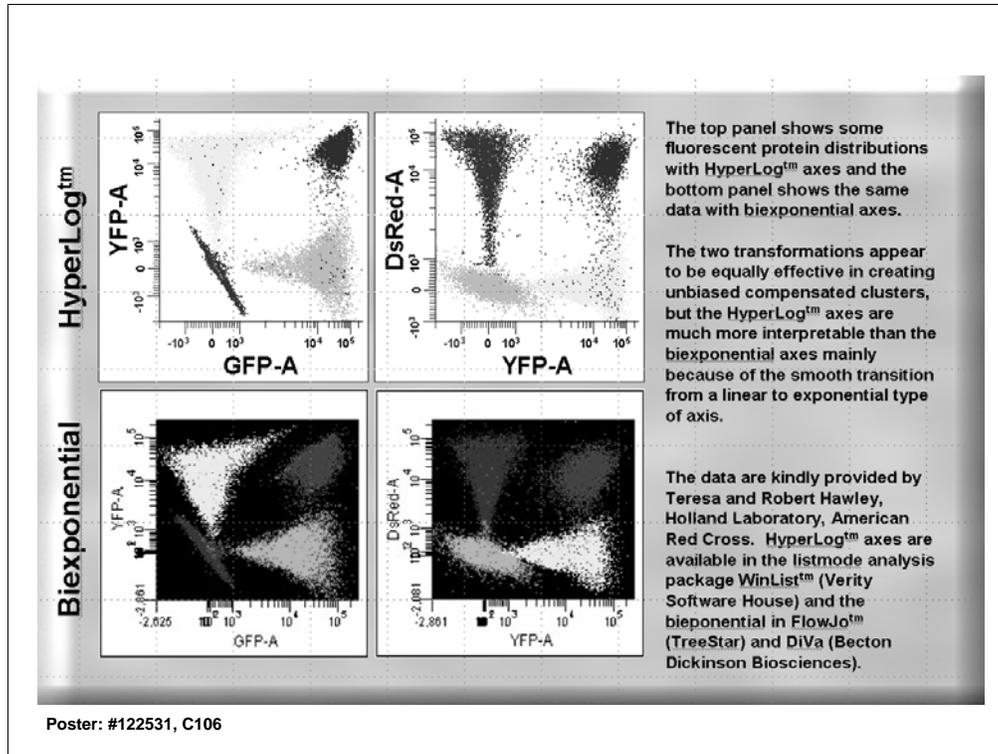**All negative linear values are set to 0 and thus are piled on the axis.**

Because the log transform is undefined at negative values, software algorithms must set these values to zero, creating a relatively large number of events directly on the axis. The loss of the negative value information creates two display problems. The first is that many events are now piled on the axes, which essentially masks their presence. A worse problem, however, is that from the cytometrists point-of-view, the data looks like it is under-compensated since the symmetric negative error envelop is now gone. This under-compensated appearance is the single most important reason why traditionally listmode data has been over-compensated.
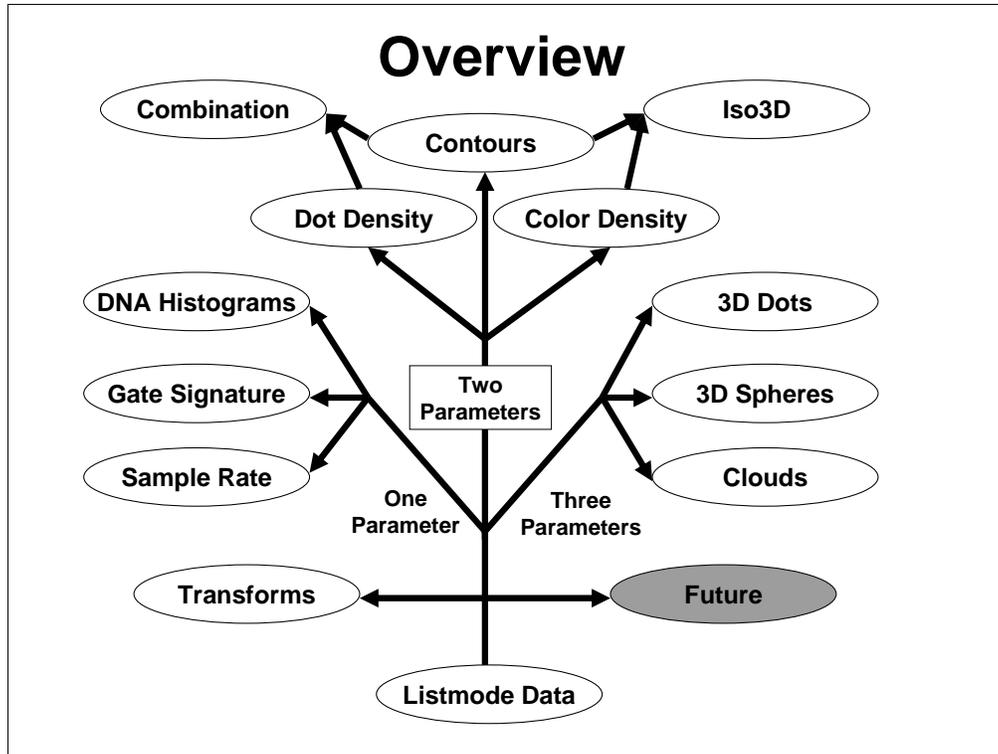
Thus, the data that appear to be under-compensated are really not. The upward trend of the upper envelope of the data is due to both counting and digital errors. Is there a better display method that avoids this under-compensated appearance?

## Biexponential and HyperLog



A solution to this display problem was first proposed by using a generalized form of the hyperbolic sin function, also known as bi-exponential and Logicle (presented at the 2002 Asimilar Conference, Moore and Parks). More recently, another type of transformation, HyperLog, has been implemented successfully in a listmode analysis software package. Both of these transforms have similar shape characteristics as shown above. The HyperLog transform was designed to be more linear near the origin and not to distort the higher decades of the transform. The advantage of these transforms is that the entire domain of linear values can be transformed, preserving the symmetry of the error distribution about our compensated data. For the first time, a compensated dot display, does not appear to be under-compensated. These transforms will play an important future role in all listmode acquisition and post-acquisition systems.

Poster: #122531, C106

The above two panels are a comparison between the HyperLog and Biexponential transforms for some four-color data kindly provided by Teresa and Robert Hawley. Both transforms nicely cluster some fluorescent protein distributions. If you examine the axes near the origin, you can appreciate the increased linearity of the HyperLog as it crosses the origin.

Overview

We now discuss the real future of these display systems.

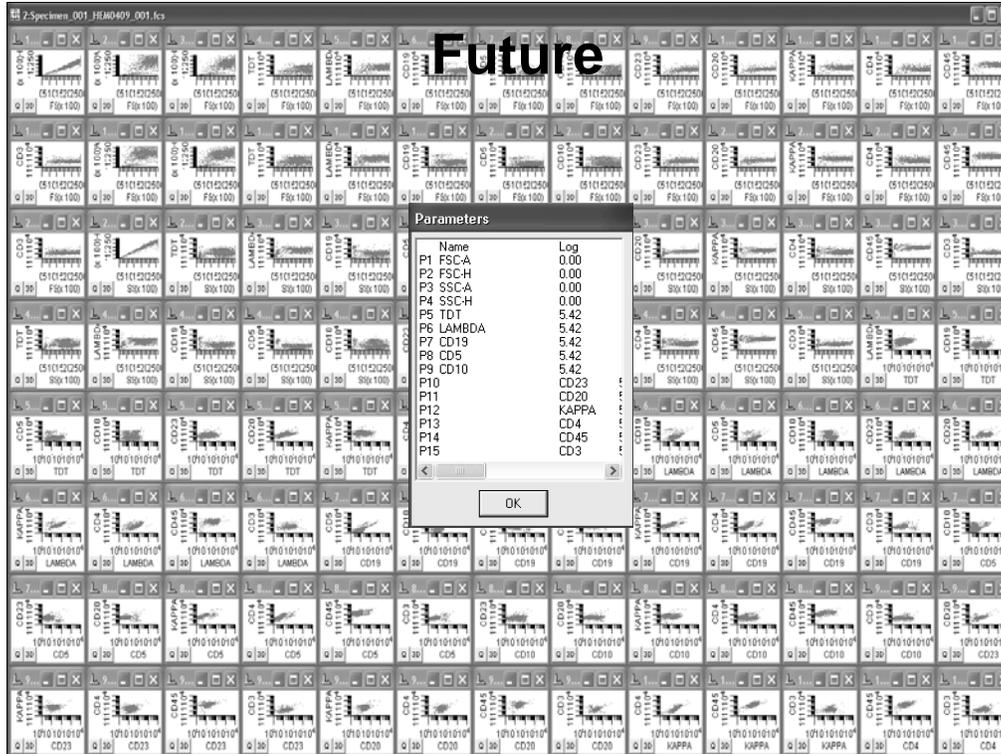# Flow Cytometry Display and Analysis Paradigm

**Create a one or two parameter histogram**

**Create one or more regions that can be used in a gate to define a population.**

**Repeat the above two steps until all populations have been defined.**

The above slide gives the modern paradigm for displaying and analyzing flow cytometry data. One normally uses one and two parameter histograms to provide a context for gating. Gates are typically some boolean combination of regions. The user typically starts defining populations at a high level and works towards defining subpopulations by a divide-and-conquer strategy.

This strategy has served us well, but with the advent of high dimensionality data (>10 parameters), we are now at a point where this strategy is becoming increasingly cumbersome.

If we look at some data provided by Fred Preffer of 11-color data (15 parameters), we can see the problem that faces cytometry today.  The number of two parameter plots necessary to form our gating regions can be in the hundreds.  It is difficult, if not impossible, with this paradigm to really appreciate all the various populations that are in the sample.

# Future

**We can still apply our display/analysis paradigm to high dimensionality data but at great expense and possible loss of important information.**

**Cytometry is still in great need of new display and analysis systems that can visualize and analyze high dimensionality clusters.**

See above.