# 23

# Visualization of High Content Screening Data

**Matthew J. Anstett**

## Summary

Visualization is essential to the understanding of complex data derived from high content screening. It is necessary to present information in a way that captures patterns and trends in the data in order to answer specific questions while also providing a way to formulate new questions and hypothesis. Specific types of visualizations can provide information on the quality of the data, temporal, and spatial patterns of cellular response, cell phenotype, and the relationship to additional data such as the chemical structure of test compounds. Interacting with this data through linked visualizations and visual filtering facilitates exploration and hypothesis generation to better understand biological systems.

**Key Words:** Cell-based assays; data visualization; drug discovery; HCS; high throughput; screen development.

## 1. Introduction

Over the past decade high-throughput technologies have provided researchers with advantages in experimental speed and throughput, but also problems of analysis and interpretation inherent with large-scale experimental results. With the emergence of these technologies, traditional small scale experiments to answer specific hypothesis have evolved into mass exploration utilizing high-throughput systems that combine robotics, high density formats, and assay miniaturization. These advances promised the ability to rapidly answer more questions, generate new hypotheses and ultimately discover new drugs. However, interpreting large datasets and complex interactions became a central problem for those involved in making decisions about which compounds to pursue from these efforts.

High content screening (HCS) is a relatively new technology that provides a wealth of information designed to better describe complex biological responses within living cells *(1,2)*. Coupled with high-throughput technology, data interpretation reaches a new level of complexity. HCS data is typically captured as an image of cell cultures whose pixel intensity is interpreted, often by specialized algorithms, to generate derived data that represent the experimental results. The data that is generated from HCS provides added dimensionality that attempts to capture the complexity of responses within living cells as a quantitative measurement. For example, HCS results often include cellular morphological measurements, temporal and spatial patterns, and membrane potentials captured over time and in the presence of test compounds. Thousands of measurements can be captured for each cell. When numerous events or states are measured within populations of living cells in a high-throughput mode even greater amounts of data are captured that require interpretation.
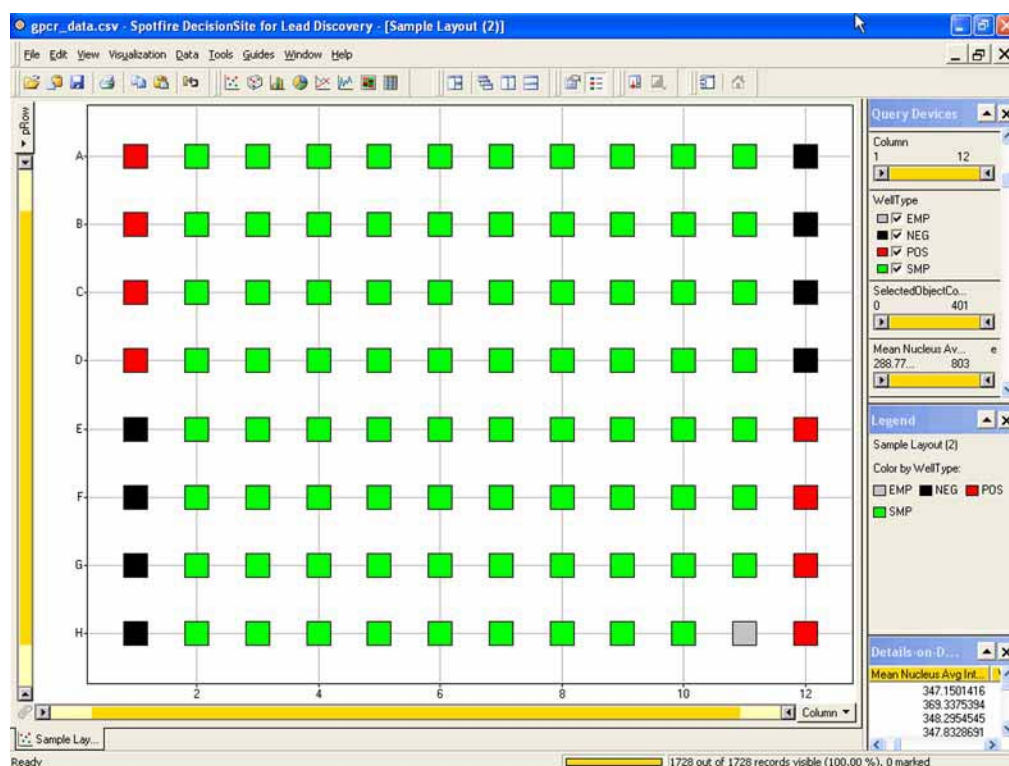
Fig. 1. Sample layout view: a single 96-well plate layout view containing negative controls (black), positive controls (red), sample (green), and empty (gray) describe metadata.

Visualization is a key factor for understanding complex information derived from HCS. HCS data presents a unique challenge for visualization and interpretation owing to its high dimensionality. When faced with such complex, large datasets it becomes necessary to provide a way in which to capture events in dimensions that can be effectively interpreted. Visualization can be used to provide different levels of understanding of events from assay quality to the response of individual test compounds. These observations can lead to a series of questions that must be addressed by incorporating additional visualizations, analyses and related data in a way that drives the analysis and provides valuable insight (*see* **Note 1**).

In this chapter we will describe:

- Data sources.
- Visualizing transformed data.
- Combining related information.
- Visualization of spatial and temporal patterns.
- Visualization of plate and well level data.
- Interacting with data to better understand biological response.

## 2. HCS Data
### *2.1. Data Sources*

Optimizing data integration is critical to all areas of scientific research as data volume and complexity increases (*3*). Challenges of working with HCS data include the large volume, varied types, storage and access needs. This data is often generated from instruments that capture the intensity of fluorescent biosensors as images (*4*). These images contain raw data from
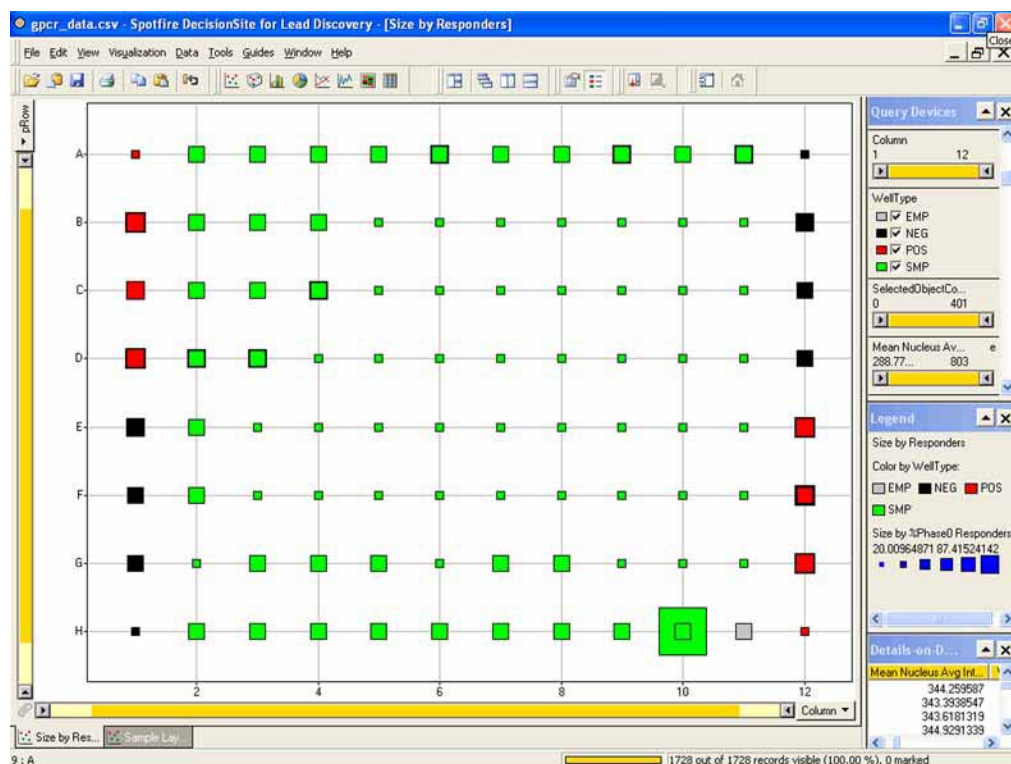
Fig. 2. Spatial view of bias: in this plot color is used as in **Fig. 1**, but the added dimension of size reflects the experimental result of "%Phase0 Responders." Here one notes that the controls and sample wells that surround the plate are larger in size, reflecting higher values for these wells and thus, potential edge effect bias in control wells.

the experiment represented as pixel intensity. Algorithms have been developed that use intensity information to generate quantitative measurements or derived data representing, for example, nuclear volume for every cell in a single test well within a 96-well plate. The primary or raw data is often stored as large image files. Derived data can be described as the data that is refined and presented at a higher level in which it is aggregated, characterized statistically, and interpreted to drive decision making *(3)*. Images and derived data might be stored in files or database schemas. As HCS technology advances this data is approaching the Terabyte size range and beyond thus, effective storage and access becomes essential. For example, access to the derived data combined with the ability to reference back to original images can be useful in confirming a particular finding and identifying problems with images that might not be evident from the derived data. Metadata is also captured which contains summary information, plate name, assay protocol, information about controls, and experimental design. We will describe how visualizations of primary data, derived data and metadata can be constructed to reveal systematic bias and improve the identification and understanding of features that drive biological events.

### 2.2. Transformation

Data transformation might be necessary in order to visualize a particular trend, normalize to a control, for example, by subtracting background noise, divide by a baseline value to view fold change, remove systematic bias, or summarize multiple measurements. A view of the metadata as a plate layout can be useful before performing normalization. Selected controls might be
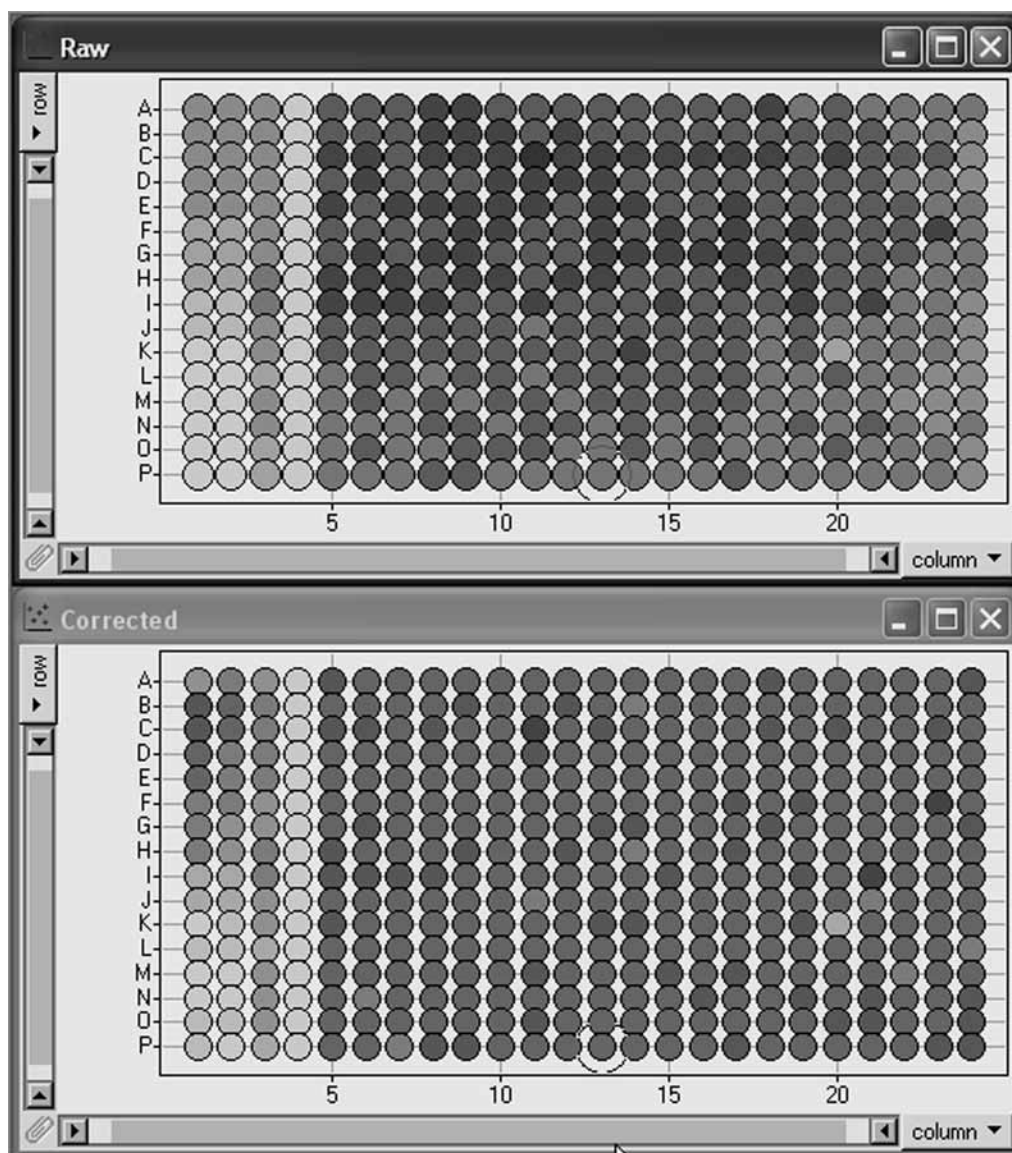
Fig. 3. Raw and corrected data: a plate view, shaded by raw and corrected experimental results removes bias that masks true differences in wells which are more clearly visible in the corrected view. (Please *see* the companion CD for the color version of this figure.)

grouped and used to normalize the data. Coloring by sample type in a virtual plate view is often the first visualization that is used in beginning an analysis (**Fig. 1**).

Visualization of the results can be added to the virtual plate map by sizing markers by an experimental result while retaining color, which denotes sample type in each well. In doing so, one might notice systematic bias, such as edge effect that might occur when stacking plates in a cell culture incubator. The response of the cell culture is reflected in a measured result, which is recognized visually as a bias to large values at the plate edge (**Fig. 2**).

Methods for reducing systematic bias in such experiments have recently been developed *(5,6)*. Applying these data correction methods can reduce bias that masks true results. For example, systematic bias that yields higher values in a pattern across a plate can mask the real differences
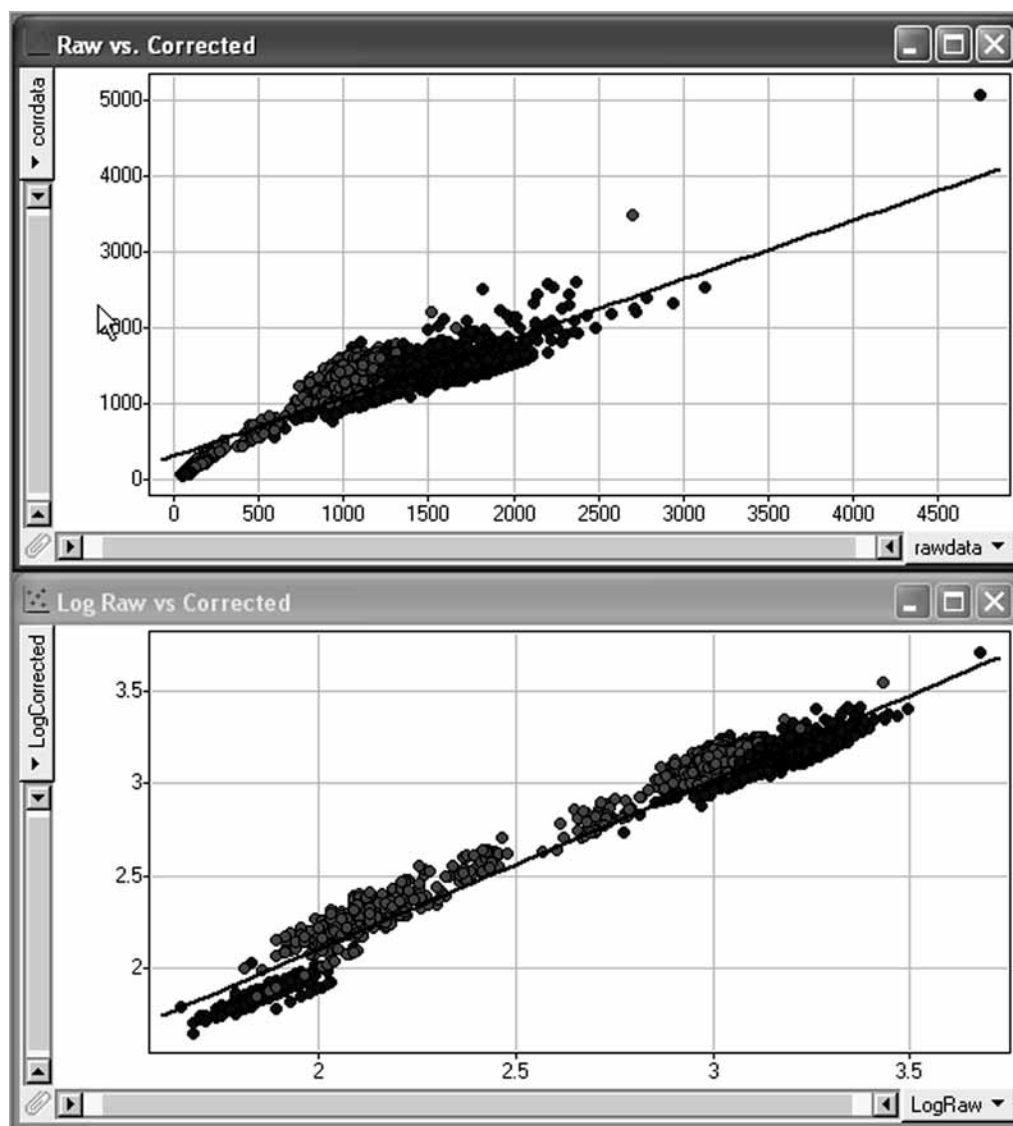
Fig. 4. Log transformation: scatter plots show the correlation of raw and corrected data on the *x*- and *y*-axis, respectively, of both original and log transformed data (upper and lower plots). Logarithmic transformation provides a more detailed visualization of data clusters within the lower value ranges. (Please *see* the companion CD for the color version of this figure.)

between wells, thus hiding important information. Visualization of raw and corrected data together, confirms the effect of employing such methods to cleanse data before analysis (**Fig. 3**).

Log transformation is a transformation method used before visualizing data that exists across a wide range of values. Large outliers can skew the visualization scale creating a condensed cloud of data points at lower ranges. By applying log transformation, the data is spread more evenly across the visualization, providing more information about data within the lower range (**Fig. 4**).

Another data transformation method often used normalizes data based on controls or experimental condition such as time. For example, based on experimental design, samples might be normalized to a zero time-point to identify significant changes over the experimental conditions. This can be done by taking the signed log ratio of each sample to a specific baseline value-typically
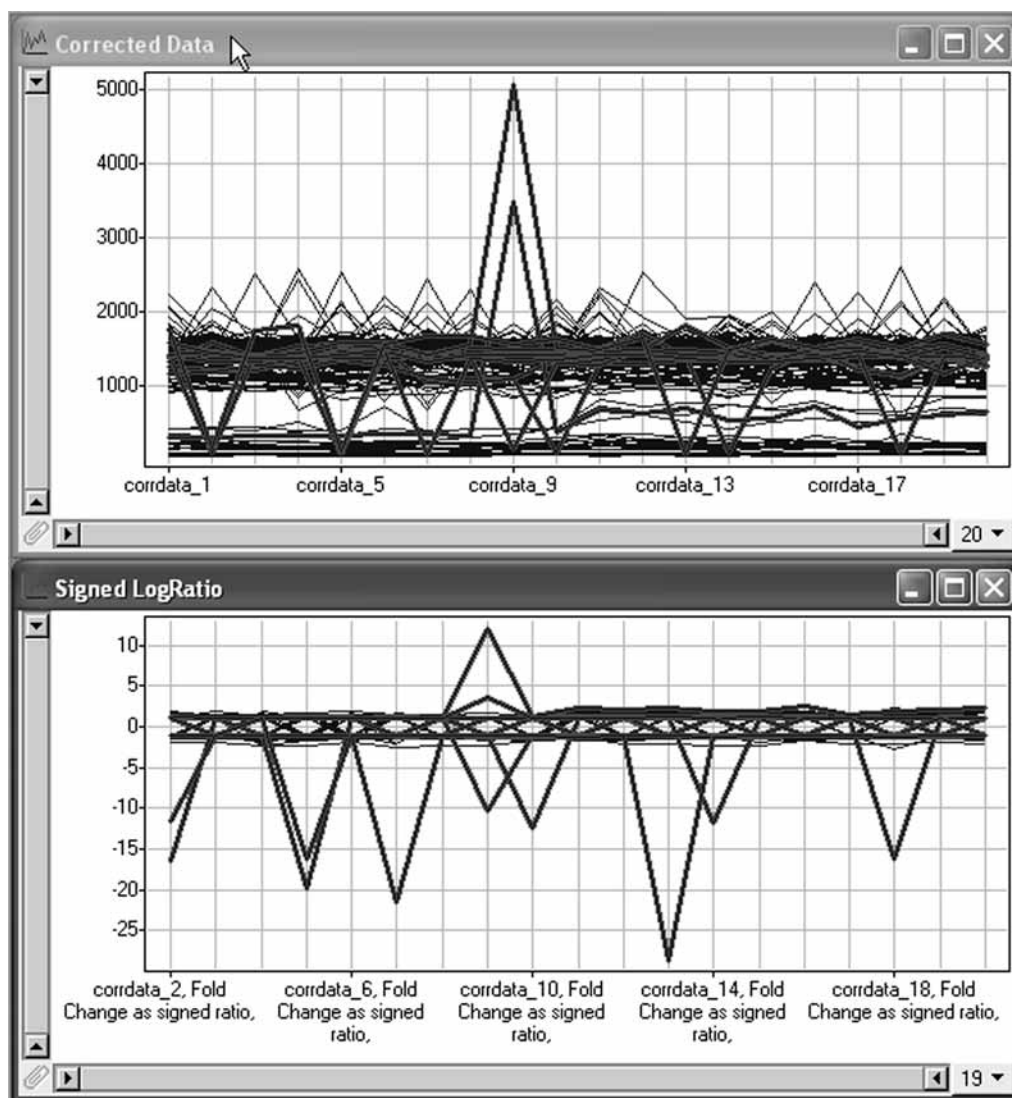
Fig. 5. Fold change: samples might be compared to a baseline value, for example time zero before treatment. By taking the signed log ratio to a baseline, shown in the lower profile plot, the largest changes in a positive or negative direction from baseline across all experiments can be easily identified. Corrected data before fold change calculation is shown in upper plot. (Please *see* the companion CD for the color version of this figure.)

time zero, before addition of drug or start of experimental condition. Visualization of data before and after this normalization allows one to isolate those samples that are most changed from baseline (**Fig. 5**). Normalization might also involve rigorous statistical methods that fit the data most appropriately based on assumptions about the type of data or instrument bias. These and additional quality control methods can provide an assessment of the assay quality, improvement in hit selection and overall data quality *(7–9)*.

Scaling data is another method that affects visualization. For example, *z*-score normalization sets the mean of the data to zero and presents data in standard deviation units. This is another useful way to see relationships between samples that can be masked by large differences in scale (**Fig. 6**).
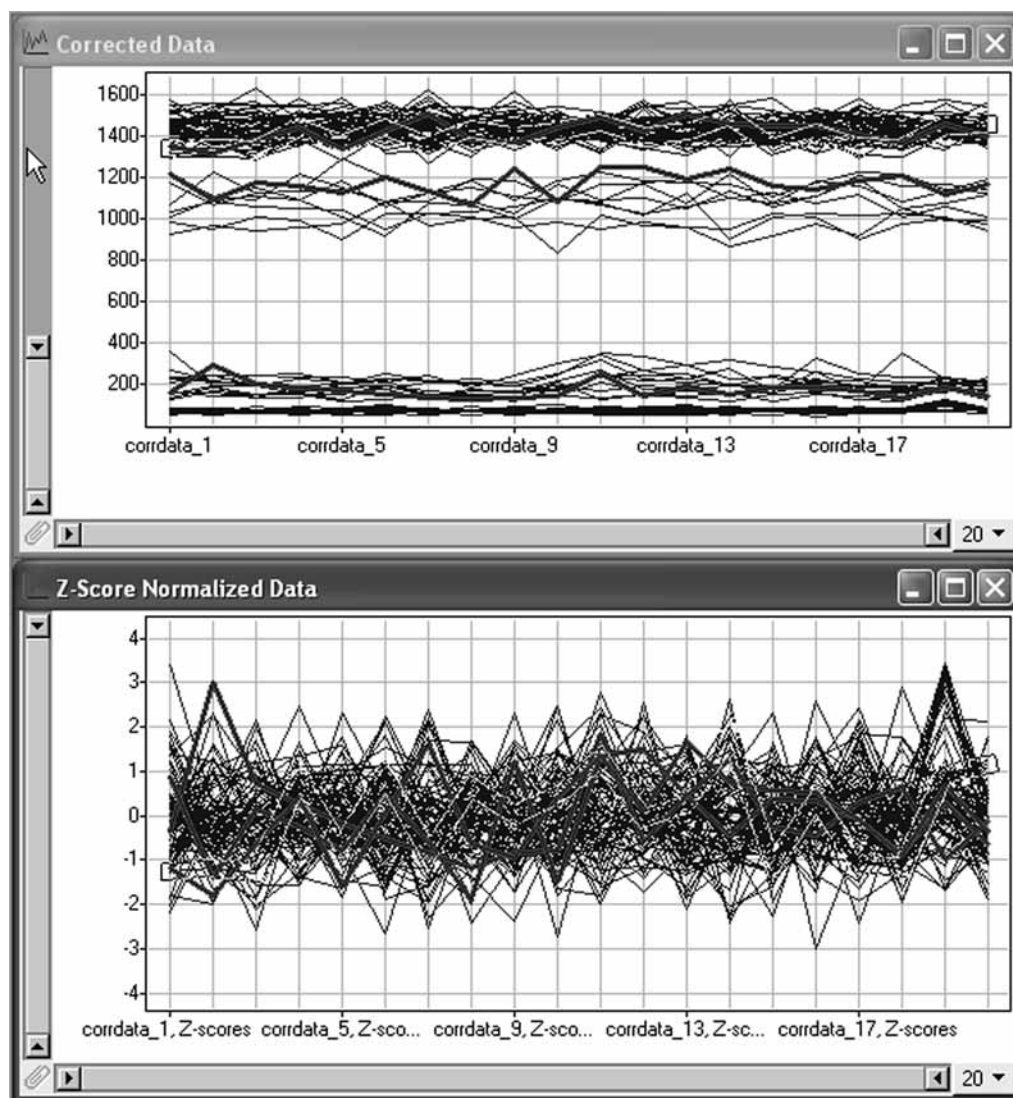
Fig. 6. *z*-Score normalization: data is normalized by setting the mean to zero and displaying the data in standard deviation units. In these profile plots the visual comparison of patterns is facilitated by normalization even though scales differ widely as shown in the upper, nonnormalized plot. (Please *see* the companion CD for the color version of this figure.)

Data transformations can be used for a variety of purposes to display and analyze data more effectively. Some of these methods, which are briefly described here, help to remove bias, and provide a better representation of values over a broad range within the same plot. They might also be used to display the response of cells within an experiment in relation to controls or baseline values, and transform data so that similar patterns can be more easily identified regardless of scale.

### 2.3. Combining Data

Data from different screening runs might be combined to view trends over time. Incorporating new information about each run is useful in detecting problems with a particular run. For example, screens run on different days might have bias introduced by different instruments, a miscalibrated pipet, varying lot numbers of reagent, temperature variations, or even culture conditions
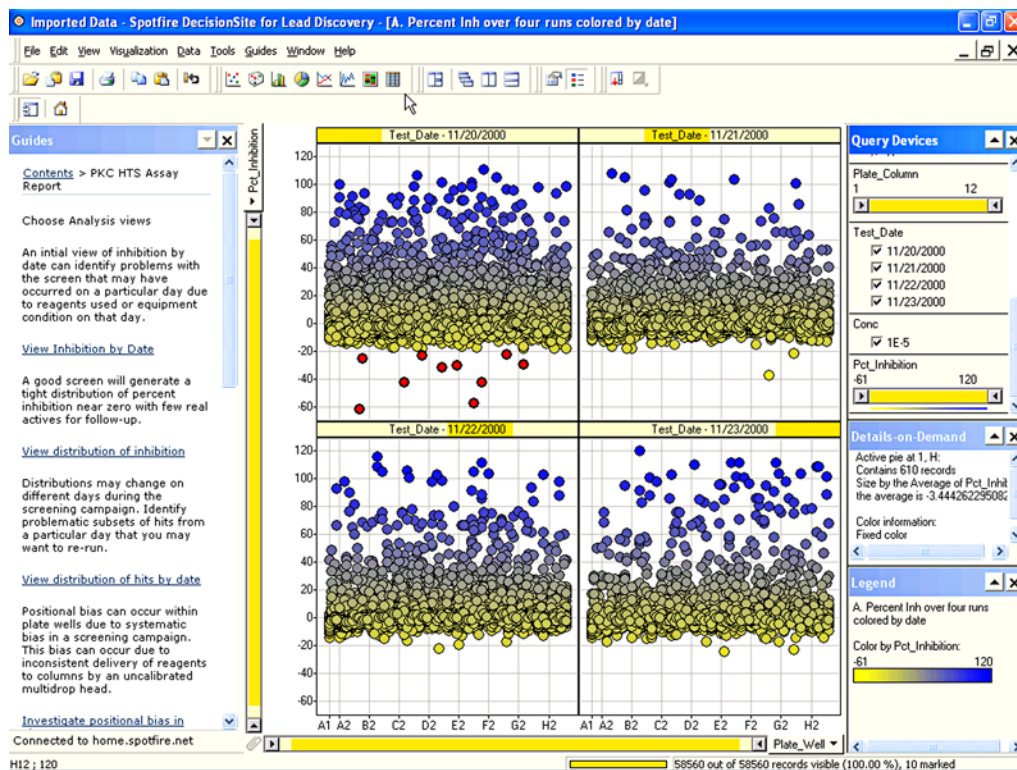
Fig. 7. Screening runs by date: scatterplots of experimental results are displayed in a trellis view by the date the experiment was run. This can identify outliers that might be specific to the instruments or reagents used on that date. The red markers show a number of negative results from experiments run on a particular date that might have resulted from factors or conditions on that day.

(**Figs. 7** and **8**). By including this experimental information researchers can be aware of any extraneous factors that might affect the results of their experiments.

Derived HCS data can be displayed and analyzed within visualizations that combine spatial information, such as well position, with data from specific cells within that well. It might also be valuable to incorporate original image information from which the quantitative data was derived. This allows one to spot specific problems with the image or better understand the relationships within the wells from which derived data is acquired (**Fig. 9**).

Associated data might also include data from HTS, proteomics, genomics, and the chemical structures of the compounds tested within each well. In large-scale experiments, many compounds might be tested and the responses identified within visualizations can be correlated to structural motifs that influence drug design (**Fig. 10**).

Combining derived HCS data with metadata such as dates of experiment or lot number of reagents can help to identify the causes of bias in results because of extraneous variables. Merging original images can also help identify bias and provide additional information such as cell shape that might be lost in quantitation. Incorporating other data such as chemical structures in the analysis of HCS data, assists in identifying structural motifs that might be responsible for cellular response and can drive the direction of synthesis of the most promising drugs.

These examples demonstrate the power of visualization in the analysis of HCS data and related information in order to better understand the complex interactions within biological systems.
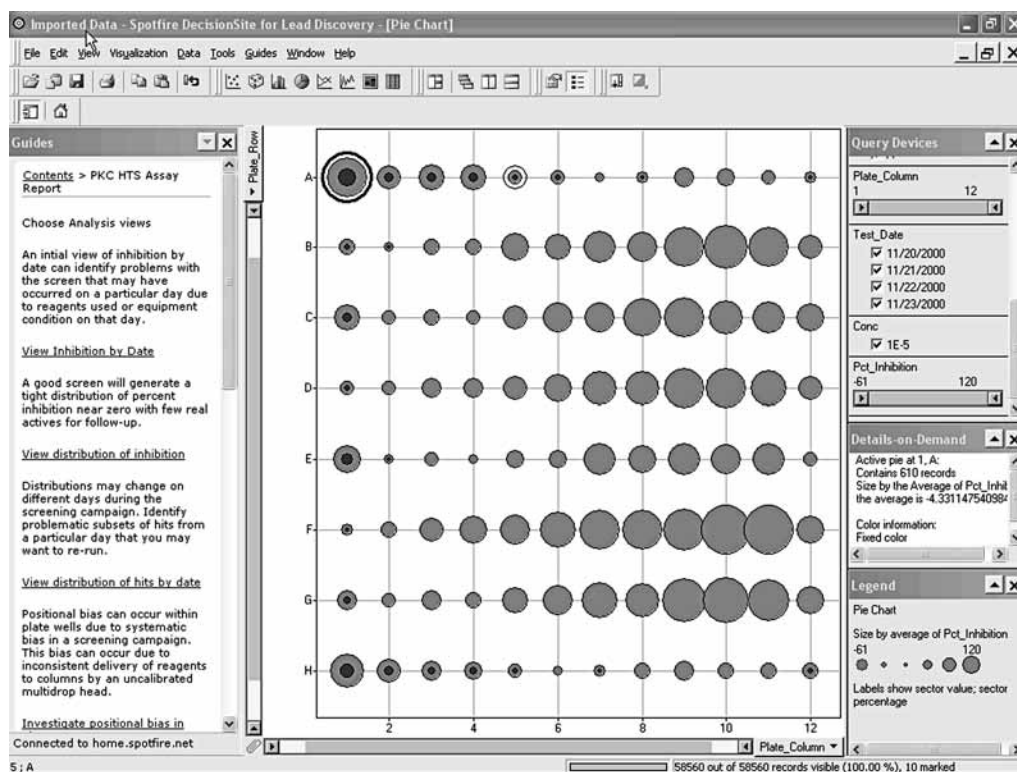
Fig. 8. Investigating positional bias: pie charts with positional information are used to display and summarize all 96-well plates from a screening run. Pie size reflects the average results for each well across all plates and larger pies on the right of the plate indicate potential pipetting bias for instruments moving from left to right in the plate. (Please *see* the companion CD for the color version of this figure.)

## 3. Visualization

### 3.1. Spatial and Temporal Views

Data gathered from HCS often includes information on cell density and morphology, for example, cellular volume or nuclear volume. Spatial relationships within a cell culture might reflect induction of certain pathways by an agonist or repression by antagonists.

Temporal views are often presented as profile charts. This type of view describes relationships between time-points for compounds or cultures. Profiles might represent compounds over time or cell culture over time in the presence of compounds. Together with clustering techniques such as k-means clustering, one can isolate the overall pattern of cellular response as a group of similarly shaped profiles (**Fig. 11**). Temporal views provide dynamic information that can reflect cyclic responses or complex upstream and downstream induction events. These views can also include spatial relationships such as cellular migration over time (**Fig. 12**). This visualization of cell tracking data provides information on both displacement and velocity of motile cells *(2)*.

### 3.2. Well- and Plate-Level Information

HCS captures many values for each cell within individual wells. This well level data is displayed and evaluated along with well summary information yielding both the cell level and summary values for each well (**Fig. 13**). With these types of views one can inspect the background intensity across all wells in the bar chart, evaluate correlation between nuclear elongation and cell area, and
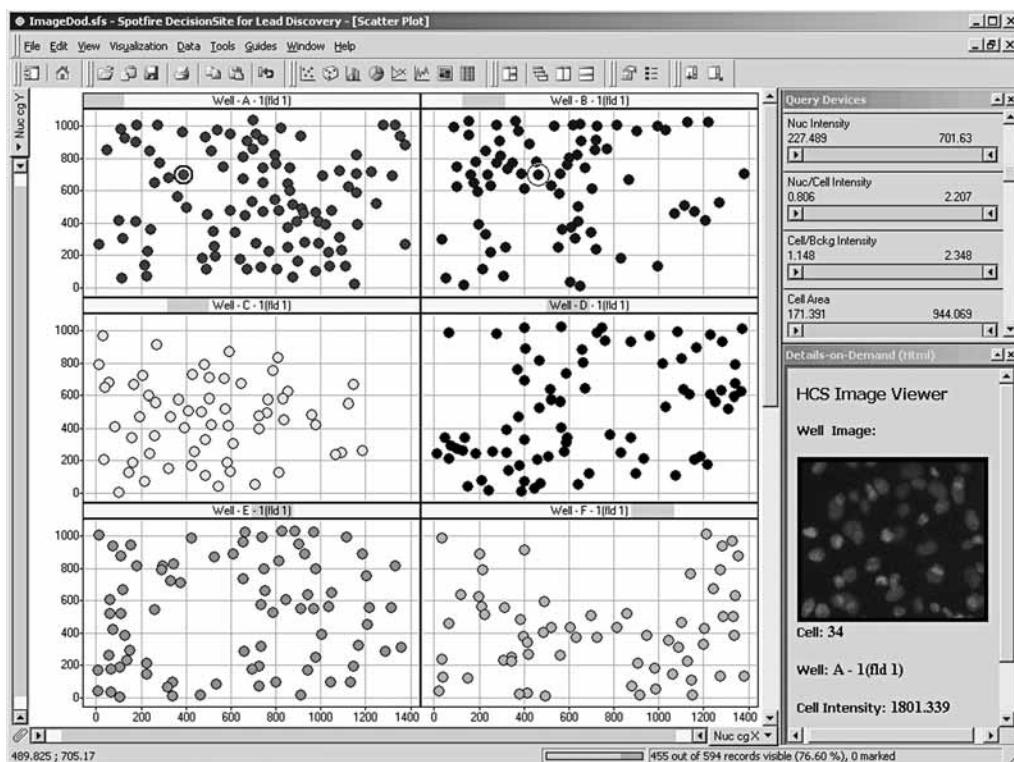
Fig. 9. Combine primary image data: scatterplots can be used to display the *x* and *y* coordinates of data captured from an image. In this visualization, the scatterplots are trellised by well. By clicking on markers representing cellular measurements within a well, the original well image is retrieved and can be reviewed along with the derived data. (Please *see* the companion CD for the color version of this figure.)

isolate outliers that can affect the summary information. Using these techniques, cell populations can be identified and explored within and across wells both at the cell and well level.

HCS data can also be viewed at the plate level to look for trends across a series of test samples or concentrations. This type of view also gives an indication of overall performance, and in which plate biases occur resulting from, for example, edge effects caused by stacking of plates in a cell culture incubator (**Fig. 2**).

### 3.3. Visualizing Related Data

Bringing together the large number of variables measured in HCS can be difficult. Visualization of more than three to four dimensions within a single plot pushes the limits of perception. By plotting key variables in multiple linked views it is possible to select a cluster of data points in one view and identify where these same data points lie in a second visualization of additional variables. Analysis can then proceed stepwise from one view to another. For example, viewing the results of principle component analysis (PCA) along with well level and correlation plots can help to assess a cluster of cells having a similar response that group together in PCA space. From a selected cluster in this view, one might then explore any spatial relationship within the wells in a second view and lastly, how these cells correlate in terms of nuclear and cytomplasmic intensities in a scatterplot (**Fig. 14**). Using color and size to add phase classification and additional information like spot average intensity, adds information to drive the selection process. Thus, complex interactions can be explored using statistical methods such as PCA within multiple visualizations in order to identify features that drive biological responses.
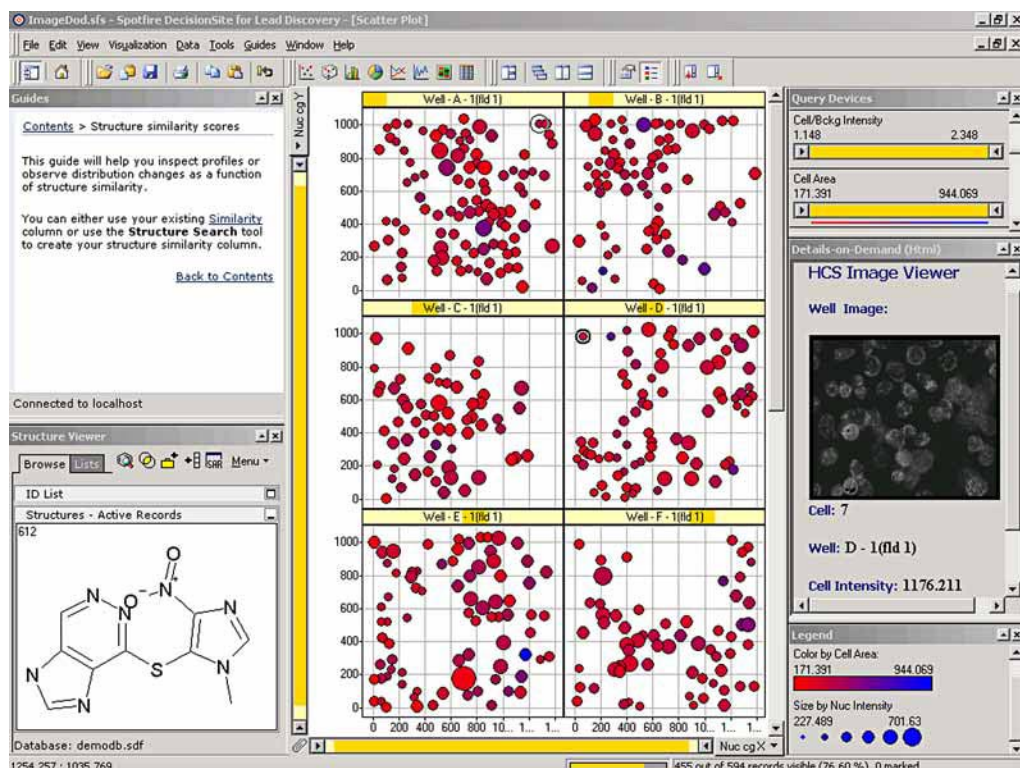
Fig. 10. Adding chemical structure information: well images and chemical structures are retrieved from chemistry and image databases and displayed simultaneously when clicking on data points within scatter-plot visualizations. Scatterplots are split by well number. Markers in the scatterplots are colored by cell area, sized by nuclear intensity, and, as expected, reflect an inverse relationship.

### 3.4. Interactive Visualization

Although exploring the large number of measurements in an HCS data it is often valuable to remove subsets of data that might obscure patterns and trends. This type of interactive analysis requires fast and intuitive methods that facilitate understanding of the data. By making selections within linked views, patterns can emerge that are obscured by the density of data. By selecting a grouping of cellular measurements within the PCA results that cluster together and are from the same phase classification (blue color) of **Fig. 14**, one can isolate three cellular measurements from the entire dataset and observe their spatial and morphological relationships in linked views (**Fig. 15**). The ability to rapidly filter information using visual filtering devices facilitates inquiry and analysis of complex relationships found in HCS data (**Fig. 16**).

### 4. Conclusion

The analysis and interpretation of HCS data poses new and sometimes difficult issues for researchers. Although this technology provides rich biological response data, it is often difficult to access and visualize in an optimal manner to address the questions that it is designed to answer. A key aspect to understanding the information is to design the way in which metadata, primary data, and derived data should be visually presented in order to answer specific questions. Too often, data integration efforts place too much emphasis on data access and transformation, losing sight of the hypothesis. Beginning with the hypothesis one can construct the appropriate questions, define the visualizations and analyses that answer these questions, and lastly, define the data
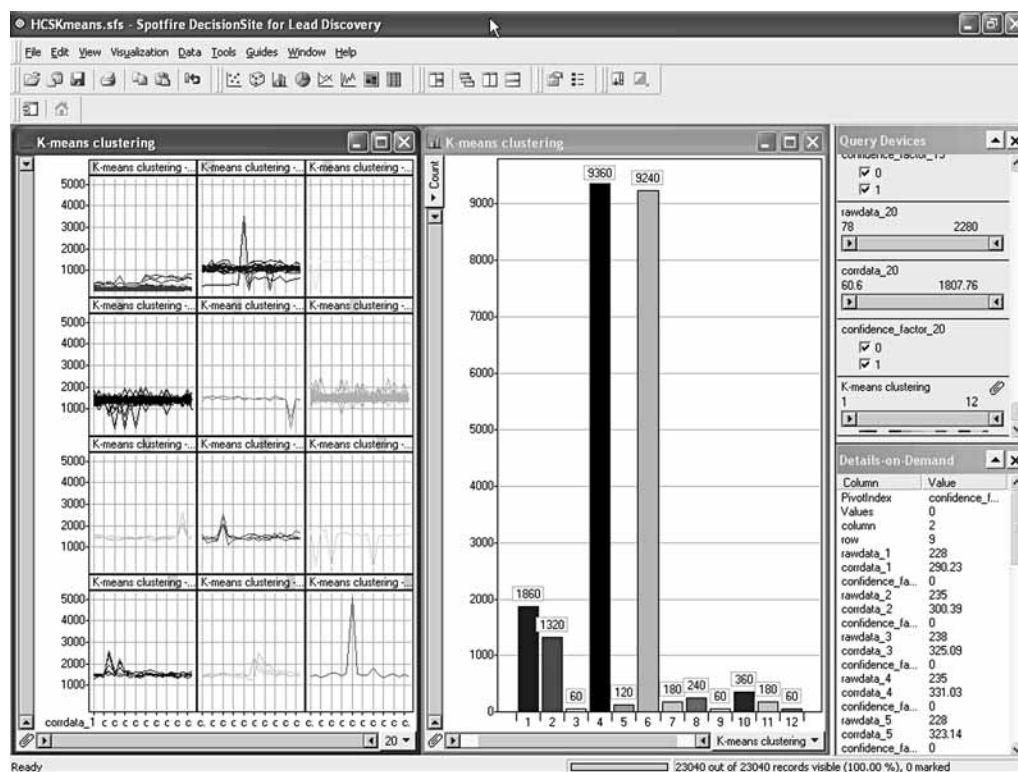
Fig. 11. Cellular response: a temporal view of a cellular response or compound performance over time is captured in profile charts which are trellised by similar pattern. A linked histogram displays the distribution of members in each cluster showing that most fall into two clusters. (Please *see* the companion CD for the color version of this figure.)

access and transformation steps needed to create the appropriate views and analyses. Additional hypotheses can also be explored using an interactive visual environment that immediately responds to repeated inquiries necessary when working with such highly dimensional information.

## Acknowledgments

## Note

1. For more information on interactive visualization and data analytics contact www.spotfire.com.

## References

1. Giuliano, K. A., Haskins, J. R., and Taylor, D. L. (2003) Advances in high content screening for drug discovery. *Assay Drug Dev. Technol.* **1,** 565–577.
2. Abraham, V. C., Taylor, D. L., and Haskins, J. R. (2004) High content screening applied to large scale cell biology. *Trends Biotechnol.* **22,** 15–22.
3. Searls, D. B. (2005) Data integration: challenges for drug discovery. *Nat. Rev. Drug Discov.* **4,** 45–49.
4. Gosh, G. N. and Haskins, J. R. (2004) A flexible large-scale biology software module for automated quantitative analysis of cell morphology. *Business Briefing Future Drug Discov.* 1–4.
5. Kevorkov, D., Makarenkov, V., and Zentilli, P. (2005) New methods for statistical analysis and data correction in HTS. *Cheminformatics, Library Design and Virtual HTS: Poster Session, SBS 11th*
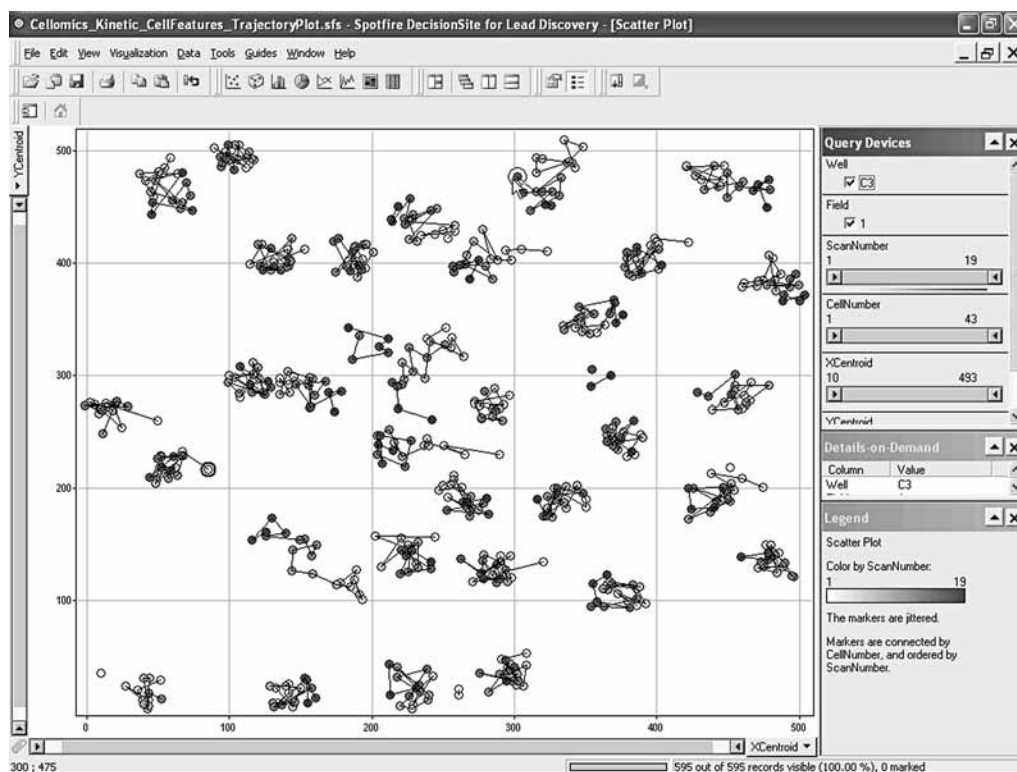
Fig. 12. Cellular migration: spatial coordinates for cells in culture are captured over time and connected by lines with arrows to indicate the migration of the cells in culture. (Please *see* the companion CD for the color version of this figure.)

*Annual Conference and Exhibition Drug Discovery: From Targets to Candidates*, 11–15 September 2005 Geneva, Switzerland.

6. Heuer, C., Haenel, T., and Prause, B. (2002) A novel approach for quality control and correction of HTS data based on artificial intelligence. *The Pharmaceutical Discovery and Development Report, 2003/03, PharmaVentures Ltd.*, [Online] Retrieved from http://www.worldpharmaweb.com/pdd/new/overview5.pdf. Last accessed March 1, 2006.

7. Zhang, J. H., Chung, T. D. Y., and Oldenburg, K. R. (1999) A simple statistic parameter for use in evaluation and validation of high-throughput screening assays, *J. Biomol. Screen.* **4,** 67–73.

8. Brideau, C., Gunter, B., Pikounis, B., and Liaw, A. (2003) Improved statistical methods for hit selection in high-throughput screening. *J. Biomol. Screen.* **8(6),** 634–647.

9. Gunter, B., Brideau, C., Pikounis, B., and Liaw, A. (2003) Statistical and graphical methods for quality control determination of high-throughput screening data. *J. Biomol. Screen.* **8(6),** 624–633.
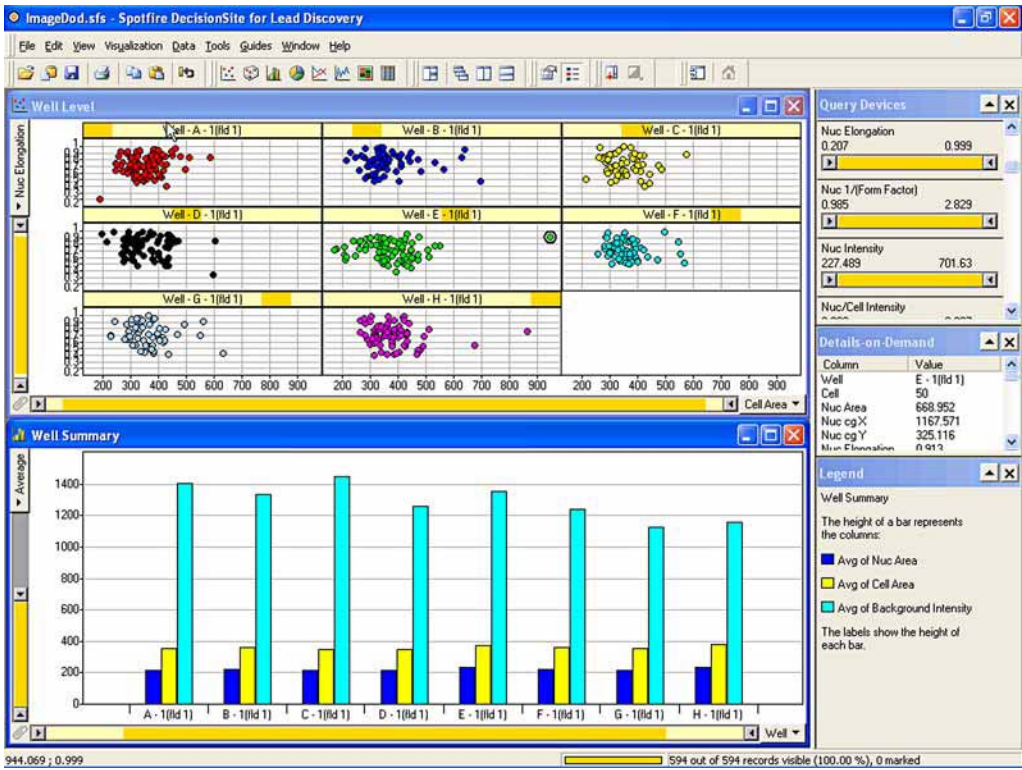
Fig. 13. Well and plate level data: each trellised pane in the upper scatterplot represents a single well with the cells' nuclear elongation on the *y*-axis and total cell area on the *x*-axis. The bar charts below summarize information from each well as average nuclear area, cell area and background intensity.
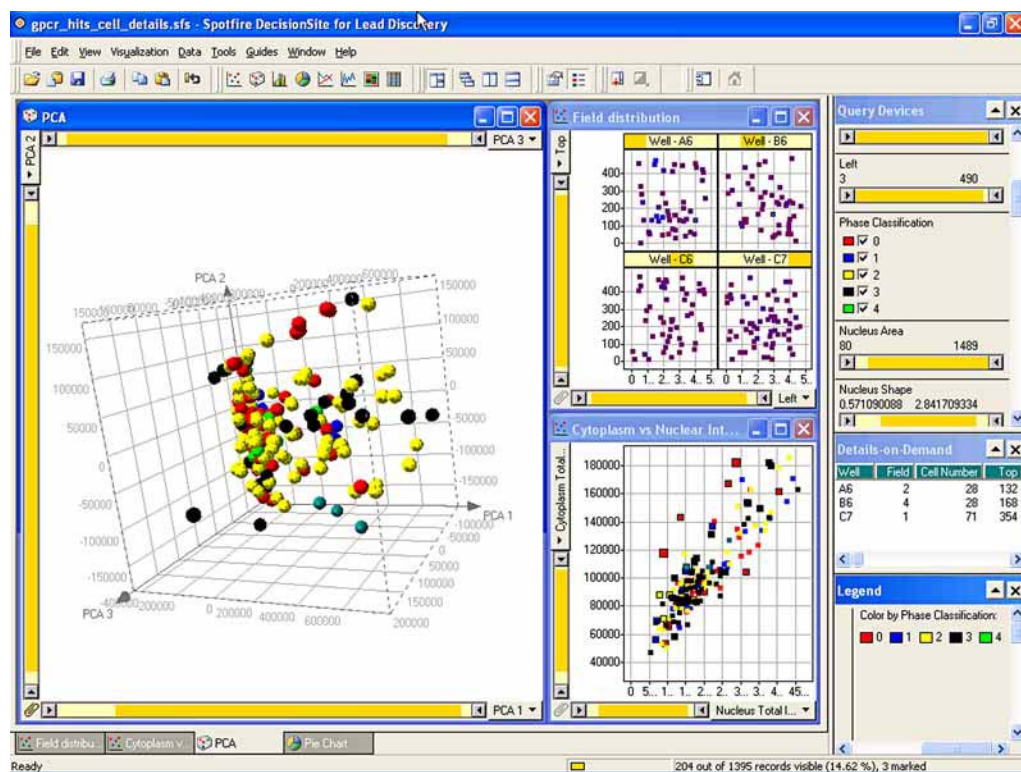
Fig. 14. Linked visualizations: results of principle component analysis groups cells in a new three-dimensional (3D) space designed to capture and summarize the variability of a large number of experimental measurements. Selecting a cluster of these cells in the 3D plot automatically selects the corresponding points in all visualizations. Thus, their location within wells and correlation between cytoplasmic and nuclear intensity can be assessed together in a stepwise fashion from one view to another.
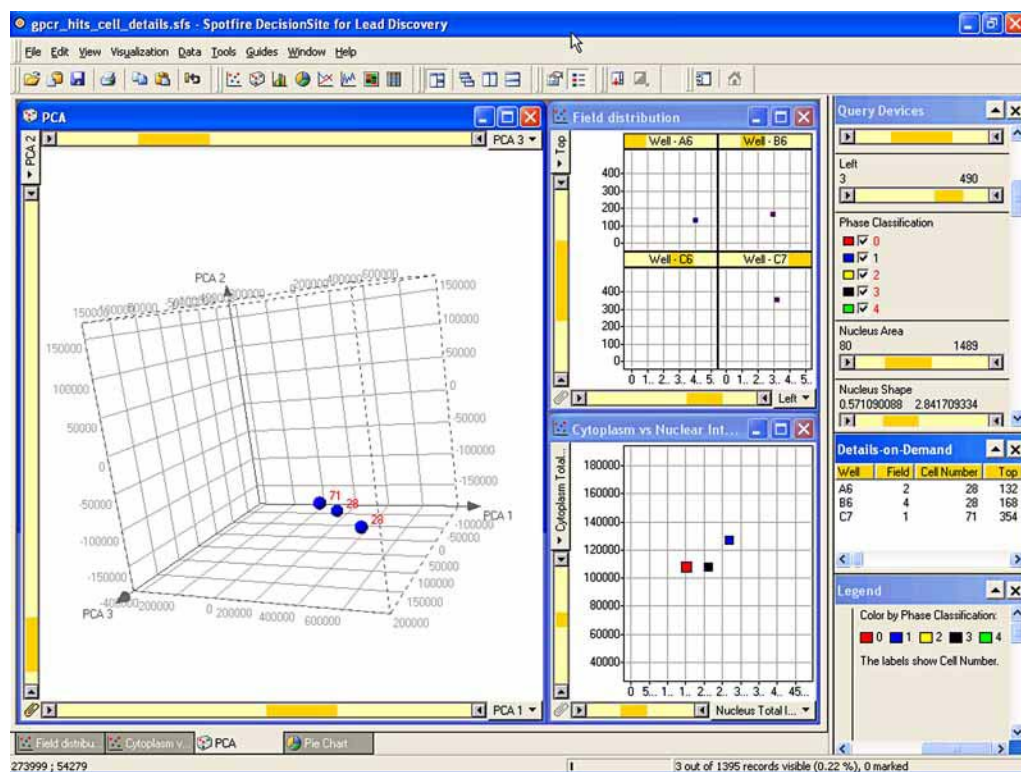
Fig. 15. Interactive visualization: by selecting a subset of cellular measurements clustering in 3D space, the spatial and morphological relationships can be assessed in well level and cytoplasmic vs nuclear intensity correlation plots.
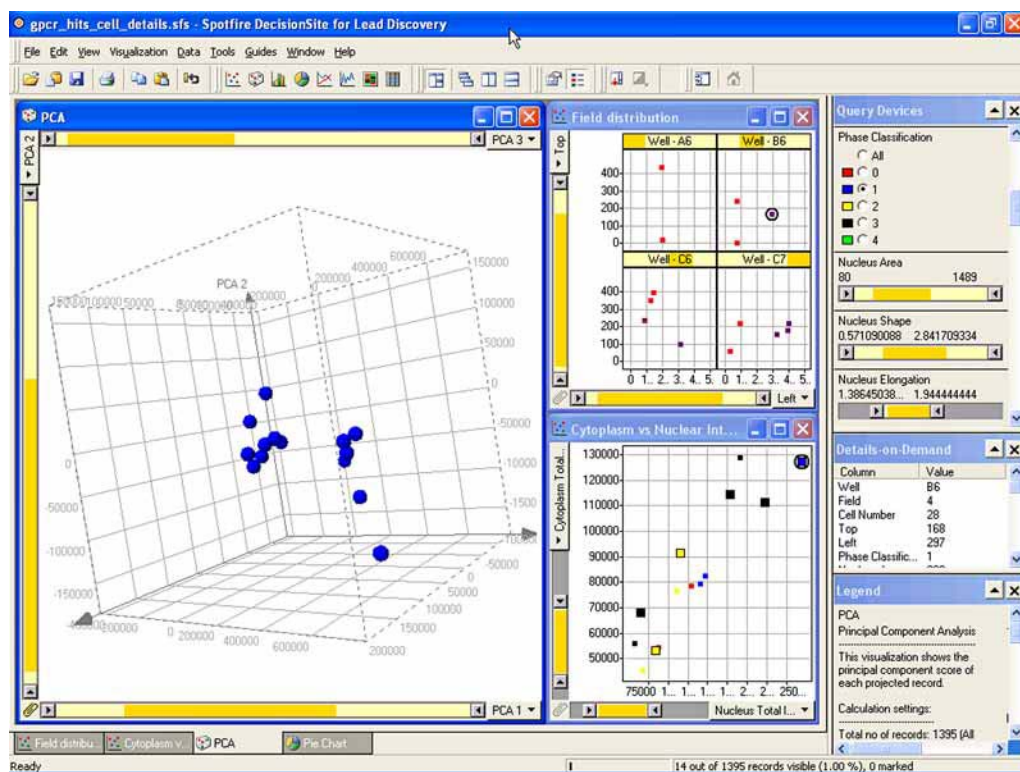
Fig. 16. Interactive visual filtering: by selecting phase classification of 1 and a range of nuclear elonga-tion using visual filtering devices, clusters of cellular measurements are revealed that might indicate response of cell populations in the experiment. A single value highlighted in the 3D plot becomes high-lighted in all views.