# Spotfire Company Update

- 700+ customers
- 25,000+ users
- 150+ employees
- Headquartered in Somerville, MA & Göteborg, Sweden
- Dominant in life science
  - All major Pharmas & 250 Biotechs
  - Discovery
    - Proteomics
    - Genomics
    - Lead Discovery
    - HTS
    - ADME/Tox
  - Development
    - Clinicals, Pilot Manufacturing Supply-chain, Sales/Marketing

# Analytics Proposition

- Drug discovery is also an information business

- Analytics is one of the few remaining areas of competitive differentiation

- Dozens of critical R&D problems where data analysis is key factor for success

- Knowledge management systems should be architected around the people not the data

- Analytical applications require a platform and strategy to go beyond individual productivity gains
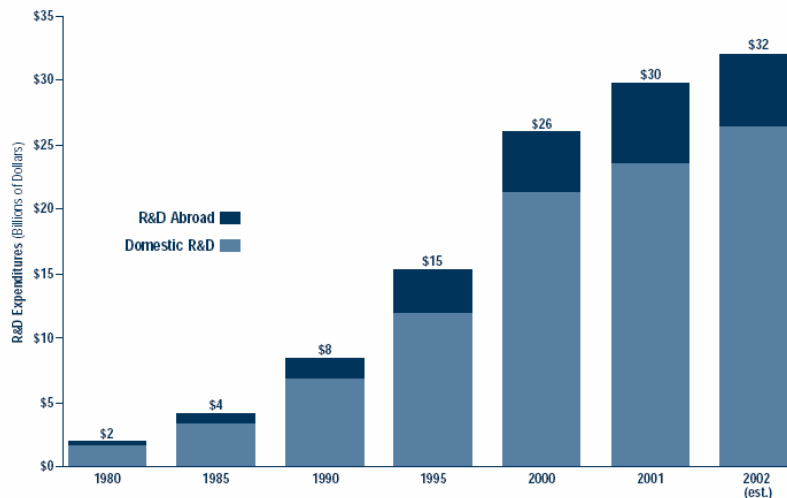
# Creating Competitive Advantage in R&D

## "First Generation Approach"

- High throughput research process investments
  - High throughput screening technology
  - Gene expression technology
  - Combichem technology
  - Genomics databases
- All companies invested, productivity improvement realized?
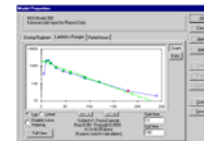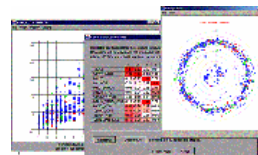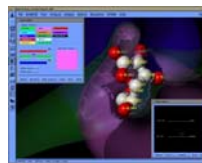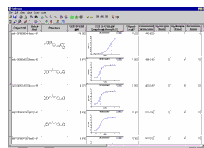
## "Second Generation Approach"

- High quality decision-making investments
  - Focus on taking advantage of data
  - Spread best-practices for choosing compounds and targets
  - Create proprietary decision making capabilities
- Towards effective decision making

**Figure 2.1 | Research and Development Continues to Grow**

R&D Expenditures (Billions of Dollars)

Legend:
- R&D Abroad
- Domestic R&D

| Year | Total |
|------|-------|
| 1980 | $2 |
| 1985 | $4 |
| 1990 | $8 |
| 1995 | $15 |
| 2000 | $26 |
| 2001 | $30 |
| 2002 (est.) | $32 |

Source: Pharmaceutical Research and Manufacturers of America, PhRMA Annual Membership Survey, 2003.

Spotfire®

# First Steps of Data Analysis

| Target Identification | Target Validation | HTS | Lead Optimization | ADME | Toxicology | Pilot MFG | Clinical |



Domain specific analysis software



Instrumentation specific analysis software



Spotfire®

# R&D Wide Analytics



- Faster, high quality R&D requires integrated drug discovery approaches

- Integrated project teams need
  - Integrated data access
  - Integrated data analysis

# Guided Analytics

**Guided Analytics supports processes, analysis workflows, and common end-user tasks required to analyze data from multiple sources**

- Capture decision-making processes
  - Expert knowledge capture – spread the knowledge
  - Best practices – higher accuracy and efficiency

- Tie together data access, analysis, visualization and reporting tools suitable for a particular task
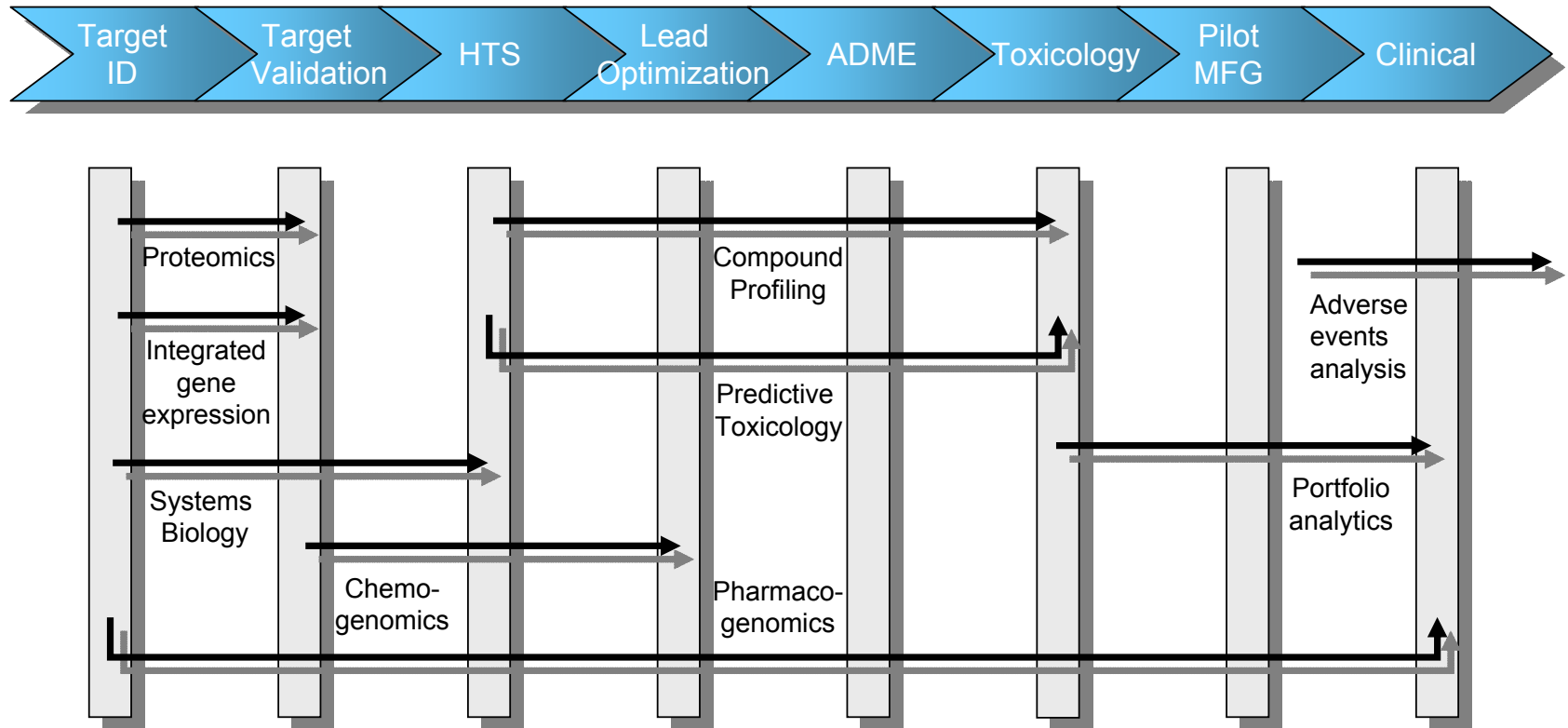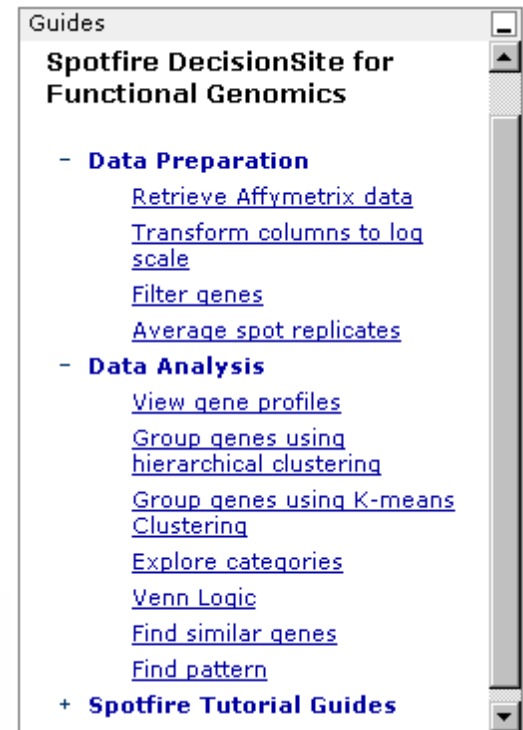  - Get the right information
  - Apply the right calculations
  - Create the right visualizations
  - Create the appropriate reports



Guides

**Spotfire DecisionSite for Functional Genomics**

- **Data Preparation**
  - Retrieve Affymetrix data
  - Transform columns to log scale
  - Filter genes
  - Average spot replicates
- **Data Analysis**
  - View gene profiles
  - Group genes using hierarchical clustering
  - Group genes using K-means Clustering
  - Explore categories
  - Venn Logic
  - Find similar genes
  - Find pattern
- **Spotfire Tutorial Guides**

**Spotfire®**

# Functional Genomics

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|
| •Databases<br>•Text files<br>•Clipboard<br>•Affymetrix AADM<br>•Affymetrix MAS<br>•GenePix files<br>•Rosetta Resolver | •Pivoting<br>•Normalization<br>•Scaling<br>•Combine replicates | •Interactive views<br>•Dynamic filtering<br>•Clustering<br>•PCA<br>•Profile searching<br>•ANOVA/t-test<br>•List comparison | •Annotations<br>•Pathways<br>•Public websites<br>•Internal websites | •Identify active genes<br>•Accept/reject hypothesis<br>•Share thought process<br>•Create a report |

- Steps in analyzing gene expression data
  - Identify the data sources
  - Select methods of data transformation
  - Incorporate appropriate analysis tools to test hypothesis
  - Add known information to extend the analysis
  - Identify subsets of genes for further study
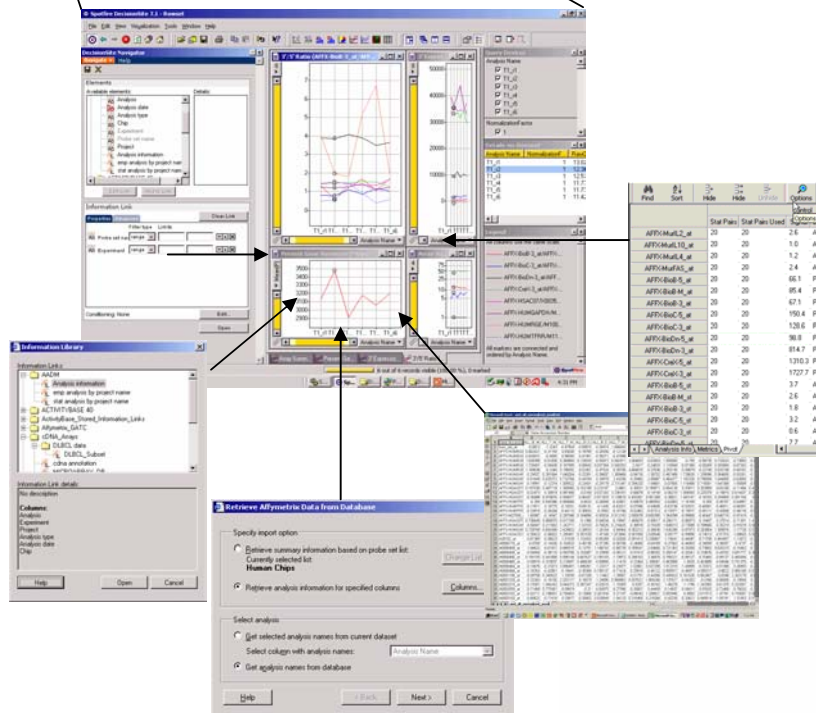  - Accept or reject hypothesis and share conclusions

**Spotfire®**

# Data Import

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|



- **Flexible, easy data import**
  - Spreadsheet, text files, clipboard
  - Relational databases
  - Define columns on the fly
  - Merge additional columns of data
  - Support multiple types of gene expression data: cDNA arrays, GeneChip, A/P, experimental parameters, average and raw data together, p-values
  - data access from standard tools
    - Affymetrix MAS5, databases, QA/QC
    - Rosetta Resolver

**Spotfire®**

# Data analysis demonstration

- Dataset
  - **Publication:** "Influence of the period-dependant circadian clock on diurnal, circadian, and aperiodic gene expression in *Drosophila melanogaster*" PNAS, July 9, 2002.
  - **Data available at:**

    http://circadian.wustl.edu/circadian.html



Influence of the *period*-dependent circadian clock on diurnal, circadian, and aperiodic gene expression in *Drosophila melanogaster*

Yiing Lin*, Mei Han[†], Brian Shimada[‡], Lin Wang[‡], Therese M. Gibler[§], Aloka Amarakone[†], Tarif A. Awad[‡], Gary D. Stormo*, Russell N. Van Gelder[§¶], and Paul H. Taghert[†‖]

Departments of *Genetics, [‡]Anatomy and Neurobiology, [§]Ophthalmology and Visual Sciences, and [¶]Molecular Biology and Pharmacology, Washington University Medical School, St. Louis, MO, 63110; and [‡]Affymetrix, Santa Clara, CA 95051

Communicated by Robert H. Waterston, Washington University School of Medicine, St. Louis, MO, May 6, 2002 (received for review February 2, 2002)

We measured daily gene expression in heads of control and *period* mutant *Drosophila* by using oligonucleotide microarrays. In control flies, 72 genes showed diurnal rhythms in light-dark cycles; 22 of these also oscillated in free-running conditions. The *period* gene significantly influenced the expression levels of over 600 nonoscillating transcripts. Expression levels of several hundred genes also differed significantly between control flies kept in light-dark versus constant darkness but differed minimally between *per*[01] flies kept in the same two conditions. Thus, the *period*-dependent circadian clock regulates only a limited set of rhythmically expressed transcripts. Unexpectedly, *period* regulates basal and light-regulated gene expression to a very broad extent.

**Materials and Methods**

Details of fly stocks, fly collections, and microarray target preparation are provided in *Supporting Materials and Methods*, which is published as supporting information on the PNAS web site, www.pnas.org.

**Data Analysis.** Average difference calls for each gene were calculated by using Affymetrix MICROARRAY ANALYSIS SUITE software. The data on each chip were normalized to the mean expression of that chip. To classify a gene as having circadian expression across two cycles of data, three criteria were applied. First, the gene had to be called present by the analysis software

Spotfire®

# Data analysis demonstration

- Research performed at Washington University on *Drosophila melanogaster* (fruit fly)

- Gene expression data from Affymetrix AADM (Affymetrix Advanced Data Model)

- Experiments were run on 48 chips with some extra repeats and time points, most due to noise

- Samples:          per01  (period gene mutants)

  W33    (wild type)

- Condition:       LD= exposed to light and dark cycles

  DD= exposed to continuous dark cycles

- Time:             0,4,8,12,16,20  hour time points

- Replicates:      2 replicates for every experiment

- Repeated:       Chips having "#2" notation repeated due to pixel noise

# Data Import

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|



- **Flexible, easy data import**
  - Spreadsheet, text files, clipboard
  - Relational databases
  - Define columns on the fly
  - Merge additional columns of data
  - Support multiple types of gene expression data: cDNA arrays, GeneChip, A/P, experimental parameters, average and raw data together, p-values
  - data access from standard tools
    - Affymetrix MAS5, databases, QA/QC
    - Rosetta Resolver

Spotfire®

# Data Transformation

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|



Average replicates, Log transformation using Guides

Pivot by a variable of interest and view gene expression

Normalize to control, scale

Transpose and view genes as profiles

- **Flexible and easy data transformation**
  - Guided data normalization, replicate average, Log transformations
  - Data pivot and transpose during or after import

# Analysis

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|



- **Powerful data analysis**
  - Guided workflow captures knowledge transfer and facilitates standardization
  - Analytical functions that everyone can use routinely and with ease
  - Informative views of the data generated at each step of the analysis
  - Linked interactive visualizations filtered simultaneously
  - Gene lists capture information on the most significant genes that can be flagged in future experiments/analysis

Guided hierarchical clustering prompts user for parameter input, creates appropriate views and extends the analysis through PCA

Spotfire®

# Secondary Data

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|

- Incorporate known information
  - Retrieve annotations from proprietary sources
  - Retrieve information from the web
  - Interact with pathway images to identify relationships
  - Use GO classifications to group genes
  - Dynamically interact with information

Add annotations from files or internal annotation source

Information from the web can be retrieved for selected records based on any value in the dataset including sequence

Pathway viewer updates to mark the pathways corresponding to genes selected in visualizations

Filtering to select a single cluster updates all visualizations and reveals majority of gene are involved in apoptosis

Spotfire®

# Secondary Data – GO Browser

Locate records (Gene IDs) in the gene ontology hierarchy by function (molecular, biological, and cellular) and view associated annotations



- Node information including a customizable live web links to additional sources of information about the GOID, term or a gene

- Search capabilities allowing to search on substring, ID or exact match

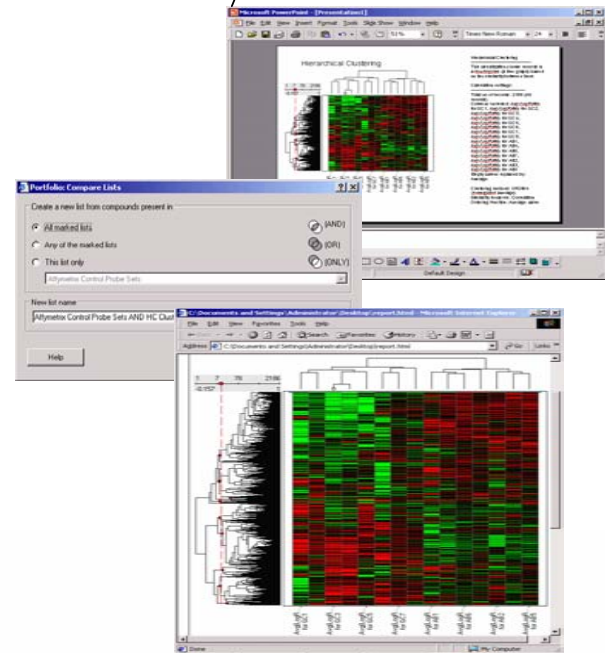- Search for p-values below a certain threshold

# Decision Making

| Data Import | Data Transformation | Analysis | Secondary Data | Decision |
|---|---|---|---|---|

- Report results, save summary information and share thought processes

  - Use visualization and dynamic filtering in project team meetings
  - Identify gene clusters, compare lists and confirm expression patterns in new experiments
  - Save lists of significant, co-regulated genes
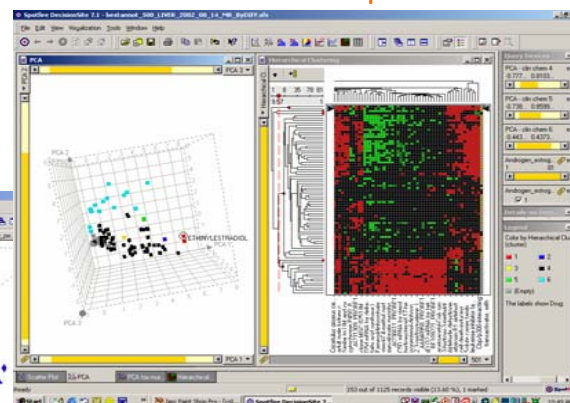  - Share and archive the decision making process



**Spotfire**®

# Extending the analysis

- ## What other information is available?
  - Validate microarray results with rtPCR data
  - Toxicogenomics: integrate Toxicology data
  - Compare gene expression with protein expression
  - Chemogenomics: integrate chemistry support



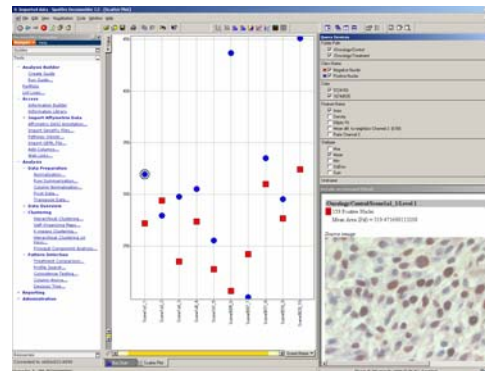Compare microarray and rtPCR data



Toxicogenomics



Protein expression



Chemogenomics

Spotfire®

# Extending the analysis

- ## What other information is available?

  - Investigate statistically significant regions of homozygosity

  - Validate data extracted from histological imaging

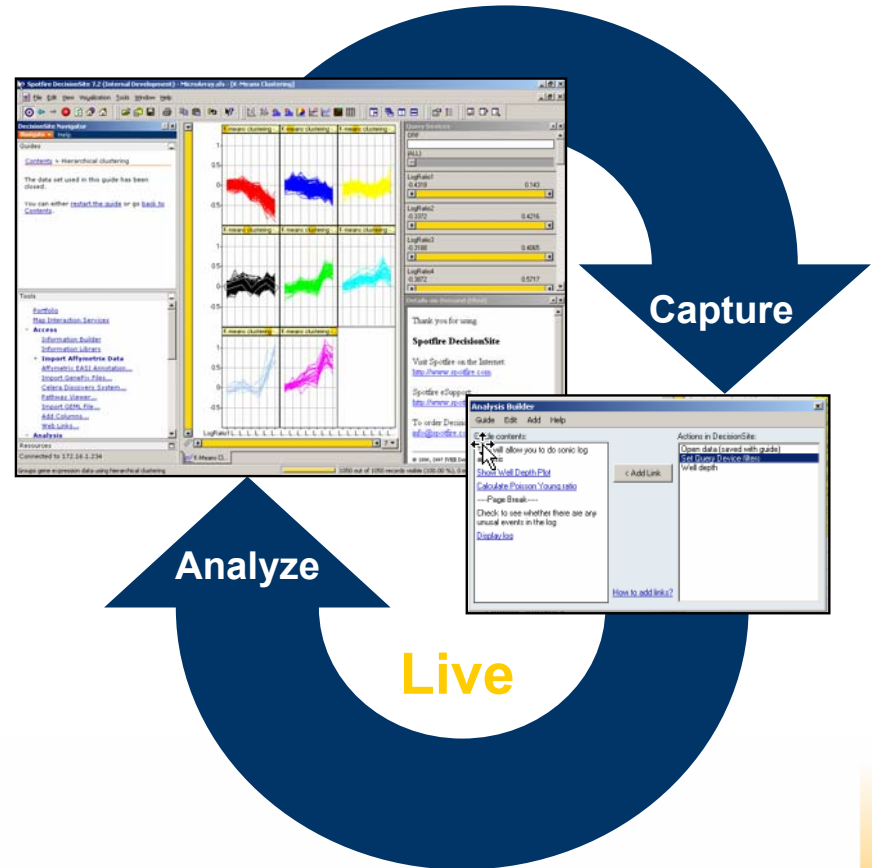  - Adverse event analysis

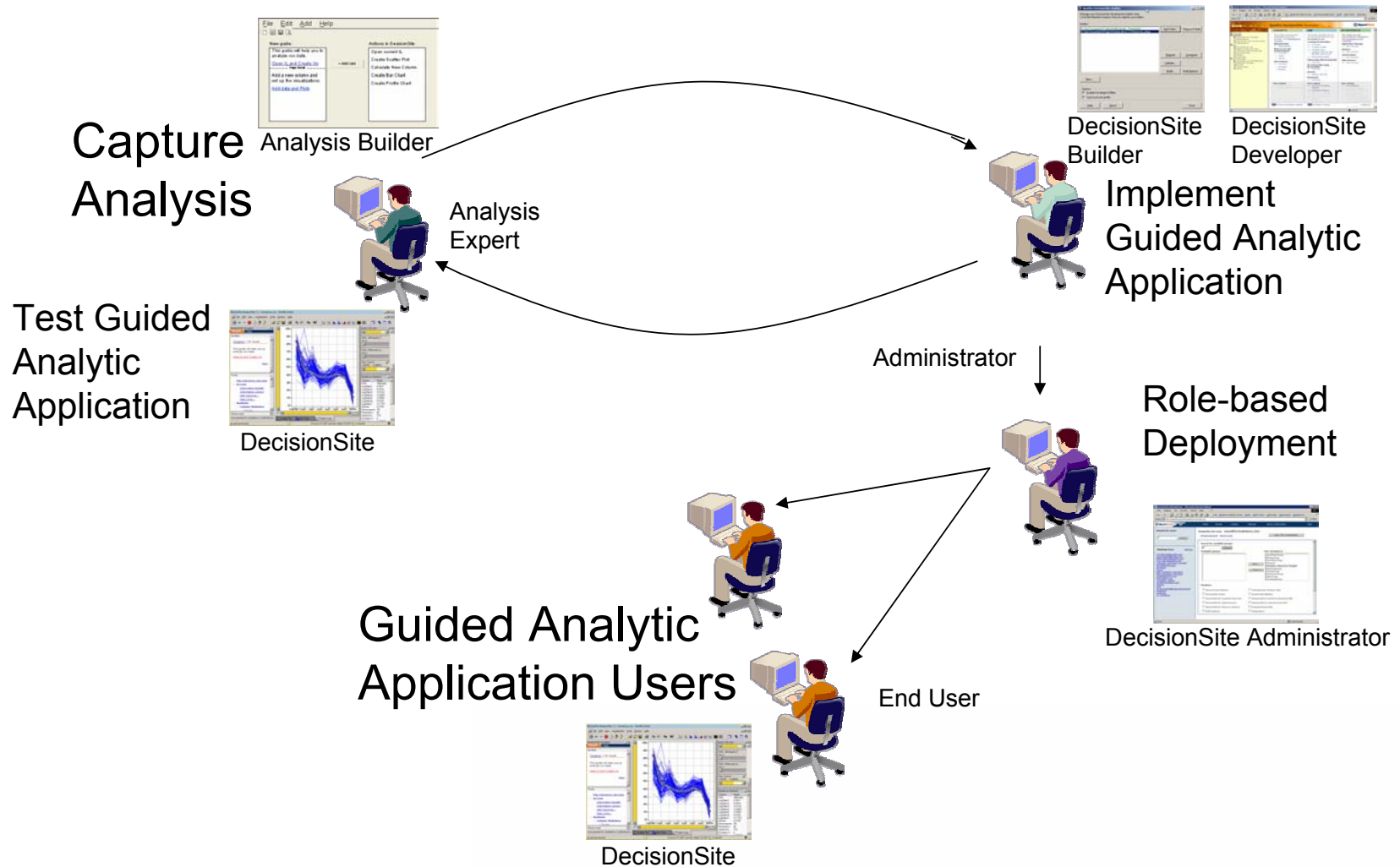SNP Analysis

Histological Imaging

Drug safety/Pharmacovigilance

Spotfire®

# Creating Custom Guides

- Expert analysis processes captured live
- End-user generated guided analytic applications
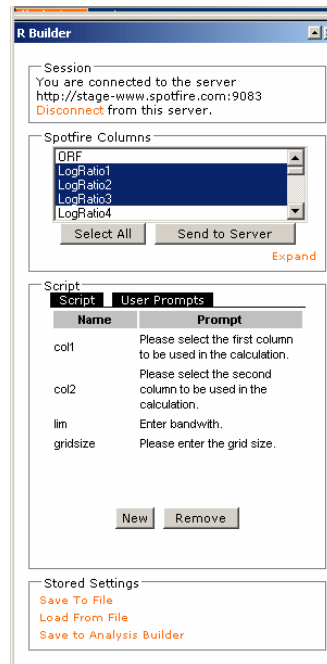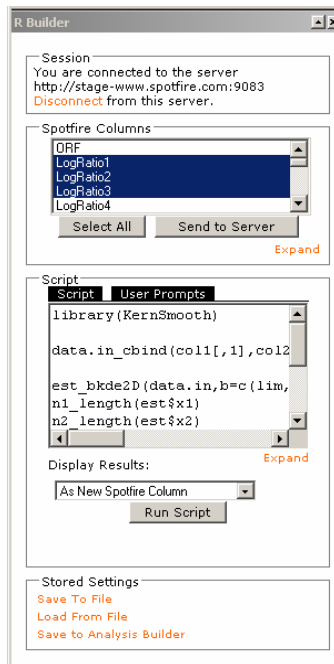  - Ready for deployment with no additional development



**Capture**

**Analyze**

**Live**

# Capturing and Reusing Analysis Processes



**Capture Analysis**

Analysis Builder

Analysis Expert

DecisionSite Builder

DecisionSite Developer

**Implement Guided Analytic Application**

**Test Guided Analytic Application**

DecisionSite

Administrator

**Role-based Deployment**

DecisionSite Administrator

**Guided Analytic Application Users**

End User

DecisionSite

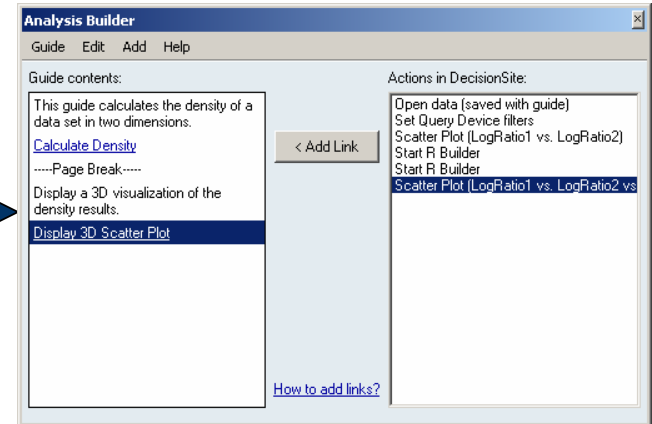**Spotfire®**

# Statistics in R

- Interactive multiplatform statistics environment
- Open-source implementation of the S language
- Easily extendable to include additional algorithms
- Active community developing algorithms and extensions
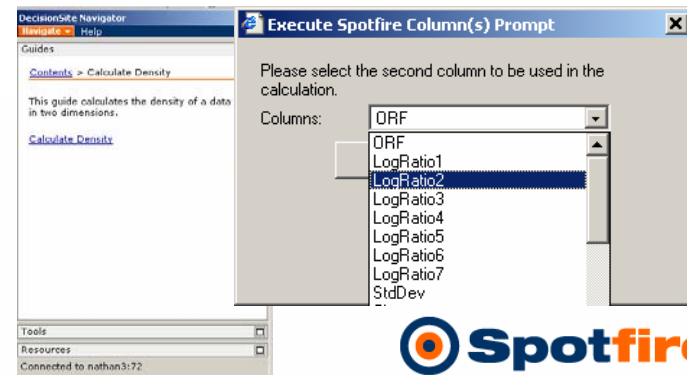- For more info on R statistics and the Bioconductor Project see: **www.bioconductor.org**

**Spotfire**®

# R Integration

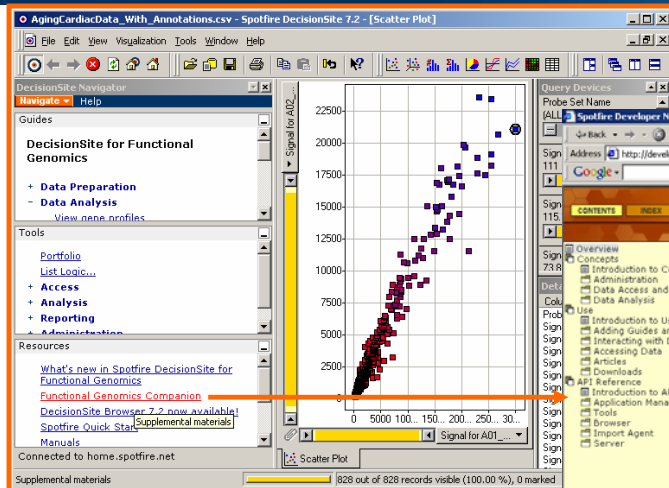## Develop R Code



## Assemble DecisionSite Guide



## Run DecisionSite Application

# Development Resources



**SUCCESSFULLY DEPLOY NEW METHODS**

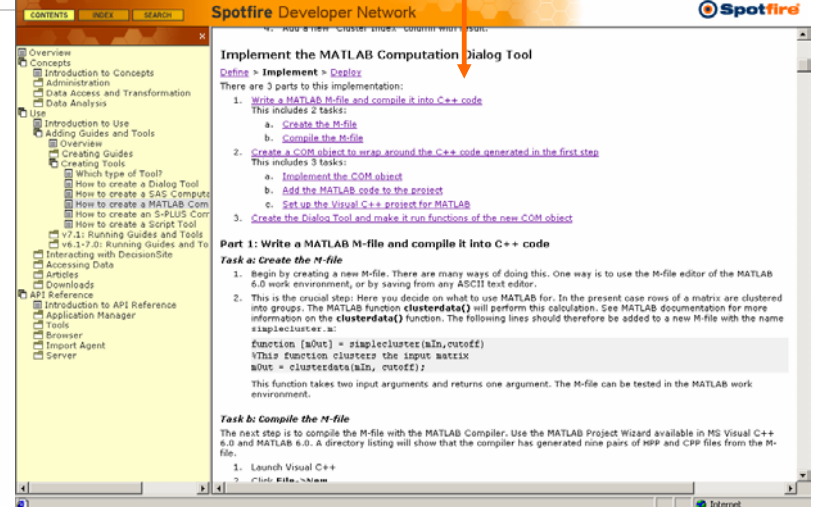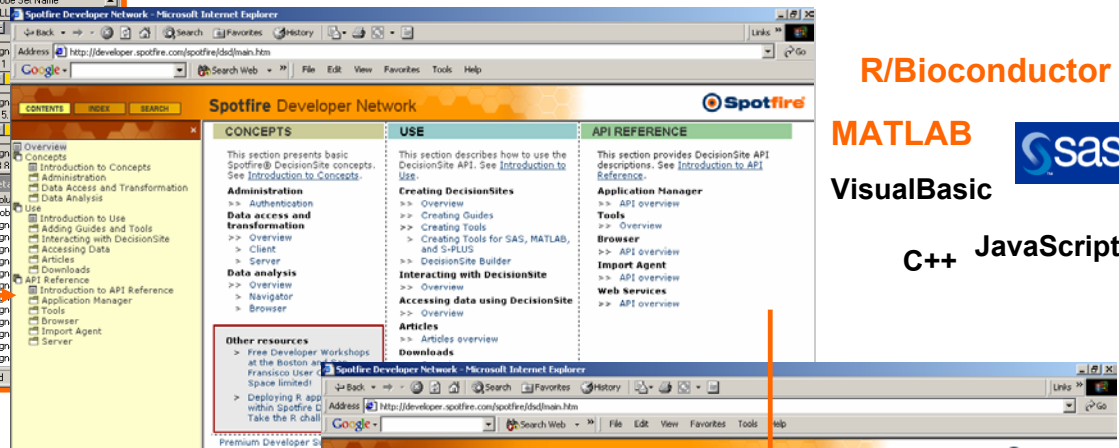R/Bioconductor

MATLAB

VisualBasic

C++   JavaScript

- **Create your Analytic Application**:

  – Extensive API documentation that is searchable, current & online

  – Instruction on how to add Matlab, SAS, S-Plus computations

  – Instruction on how to add your algorithm or specialized tools

# Thank you for your attention!

/freedom to wonder_power to decide

For more information and archived webcast recordings:

**www.spotfire.com**

*Contact: mark@spotfire.com*