# Visual Analysis of Chemical Space with Scaffold Hunter

## Nils Kriege

Dept. of Computer Science, TU Dortmund

OpenTox Euro, Mainz, 02. October 2013

# Chemical Data in Drug Discovery

## Chemical Space

- Theoretical chemical space: $\sim 10^{62}$ molecules
- *De-novo* libraries: several hundreds of millions
- Commercially-available: $21$ million molecules (ZINC)

## Trend

- Increasing amount of available data (public or in-house)
- Need to systematically explore and analyze data to speed up drug discovery process

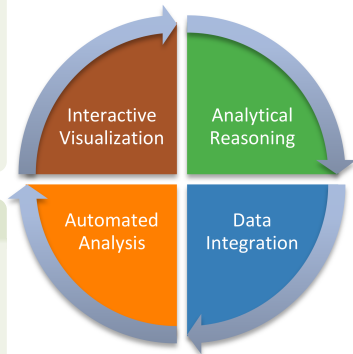# Cyclic Knowledge Discovery by Visual Analysis

**Classical Approach:** Raw data $\rightarrow$ Analysis $\rightarrow$ Visualization

## Visualization

- Analysis results
- Raw data
- Linked views

## Reasoning

- New hypotheses
- Decision making
- Intuition of domain experts



Interactive Visualization · Analytical Reasoning · Data Integration · Automated Analysis

## Analysis
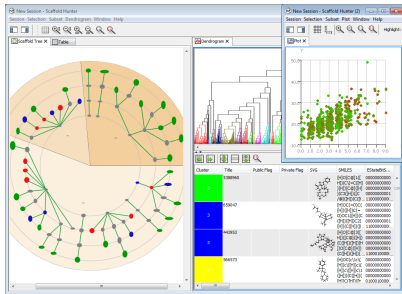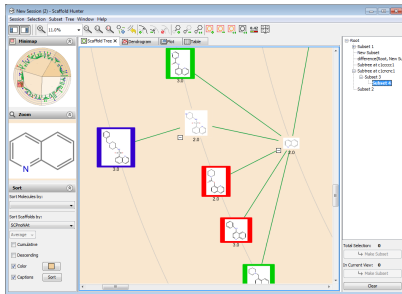
- Clustering
- Classification
- SAR

## Data Integration

- Diverse data sources
- Experimental results

# Scaffold Hunter

- Java-based Open Source tool
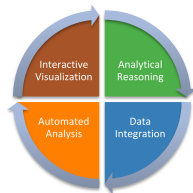- Development started 2007, TU Dortmund

**Scaffold Hunter**



**Goal:**

- Import of data from a variety of sources
- Integrated visualization and analysis
- Interactive exploration in a systematic manner

# Scaffold Hunter for Visual Analysis



- Facilitate cyclic knowledge discovery process
- Refinement of subsets, analysis parameters
- Integration of additional experimentally obtained data

# Scaffold Tree: Concepts & Algorithms

- Hierarchical classification scheme based on core structures
- Rule-based parent scaffold selection
- Scaffolds as representatives for sets of similar molecules
- *Virtual* scaffolds without associated molecules

## Algorithm

- For each molecule:
  1. Prune terminal side chains
     → scaffold
  2. Successively remove rings
     → unique parent scaffolds
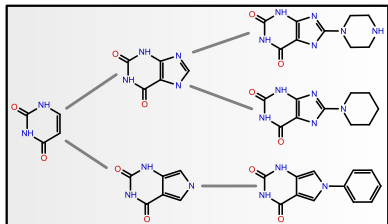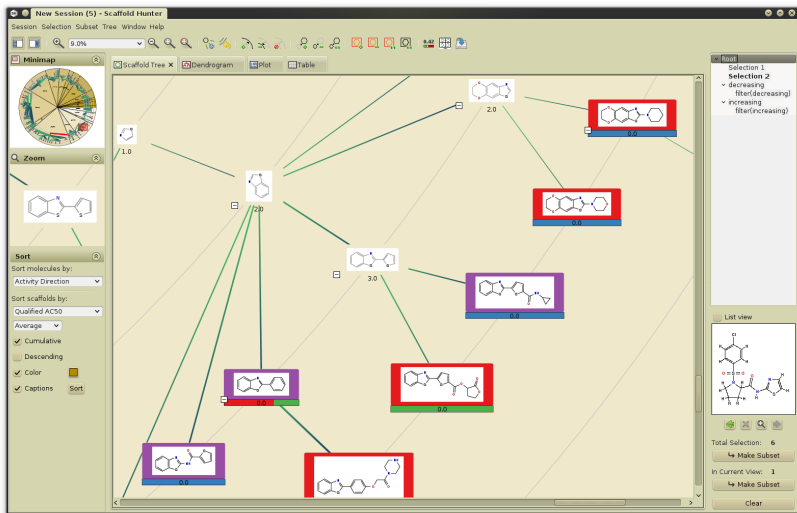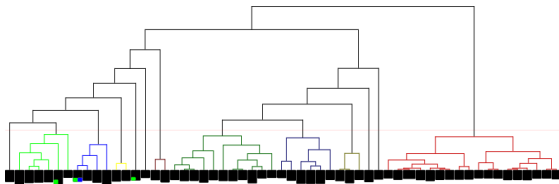- Merge multiple scaffolds
  → scaffold tree



Figure: Branch of a scaffold tree

*The Scaffold Tree - Visualization of the Scaffold Universe by Hierarchical Scaffold Classification*
**Schuffenhauer, Ertl, Roggo, Wetzel, Koch, Waldmann**; J. Chem. Inf. Model., 2007, 47, 47-58

# Scaffold Tree: Visualization



- Details-on-demand: Scaffold depiction adapts to zoom level
- Property Mapping: Representation by visual attributes

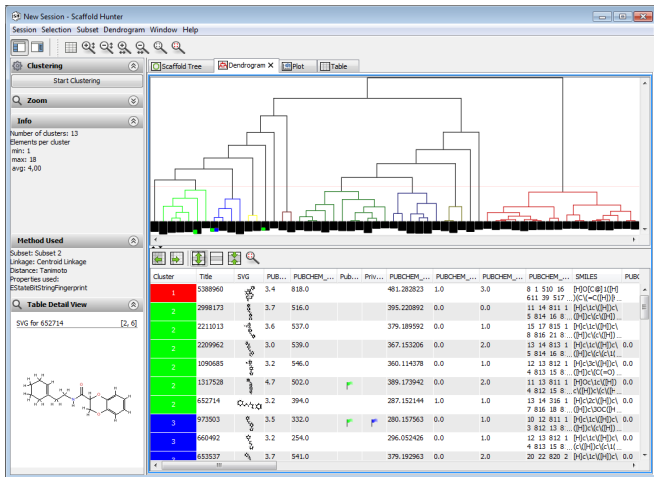# Hierarchical Clustering: Concepts & Algorithms



## SAHN Clustering

- **Distance** between molecules, e.g., Tanimoto & fingerprints
- **Linkage:** Distance between clusters, e.g., Average or Ward
- **Algorithm:**
  1. Start with singleton clusters
  2. Merge pairs of clusters with minimum distance until a single cluster is obtained
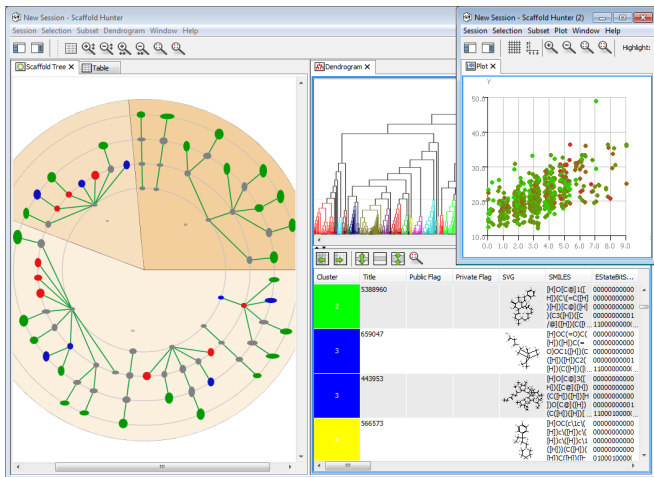
## Heuristic SAHN Clustering

- Subquadratic running time in practice, low memory footprint
- Support for arbitrary metric distance measures
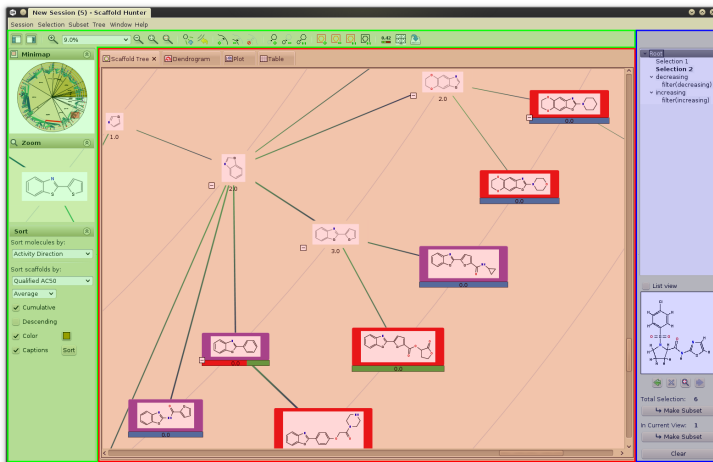
# Hierarchical Clustering: Visualization



- Zoomable user interface with details-on-demand
- **Cluster selection bar:** Interactive refinement of clustering
- **Table View:** Embeddable synchronized spreadsheet

# Plot View



- 2D/3D scatter plot
- Mapping of attributes to axes, color, dot size etc.

# Plot View



- 2D/3D scatter plot
- Mapping of attributes to axes, color, dot size etc.

$\rightarrow$ **Effective use requires intuitive linking of views!**

# Scaffold Hunter Main Window



- **Red**: Currently open views in tabs
- **Green**: View-specific tool- and sidebar
- **Blue**: Global subset and selection management

# Coordination & Linkage of Views

- **Global selection:**
  - Synchronized selection over all views
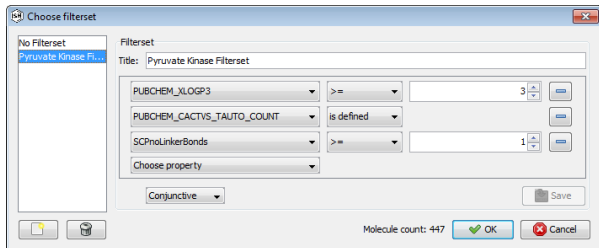  - Selection browser for quick access

# Coordination & Linkage of Views

- **Global selection:**
  - Synchronized selection over all views
  - Selection browser for quick access
- **Subset Management:**
  - Hierarchy of subsets
  - Change underlying subset of view
  - Multiple views on different subsets

# Coordination & Linkage of Views

- **Global selection:**
  - Synchronized selection over all views
  - Selection browser for quick access
- **Subset Management:**
  - Hierarchy of subsets
  - Change underlying subset of view
  - Multiple views on different subsets
- **Filtering:** Selection, Property-based, SSS
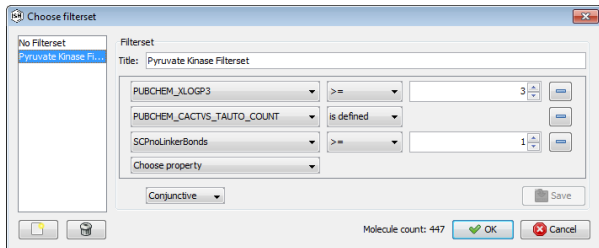
# Coordination & Linkage of Views

- **Global selection:**
  - Synchronized selection over all views
  - Selection browser for quick access
- **Subset Management:**
  - Hierarchy of subsets
  - Change underlying subset of view
  - Multiple views on different subsets
- **Filtering:** Selection, Property-based, SSS
- **Annotations:** Tooltip, comments, …

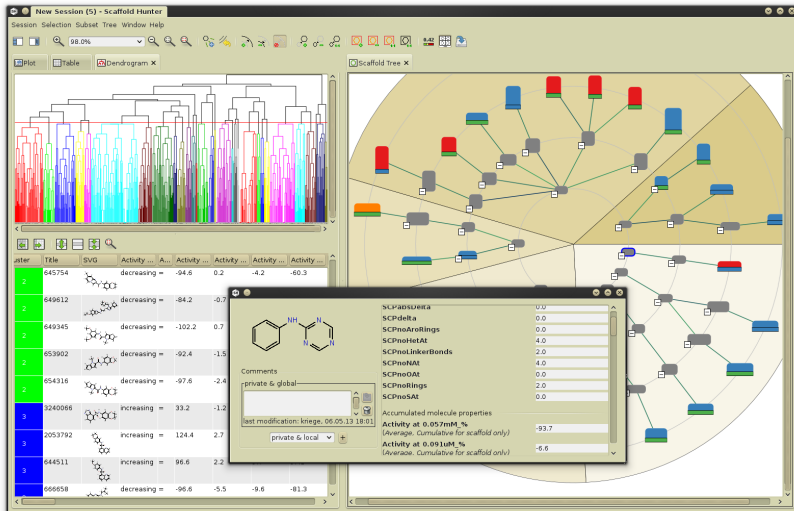# Multiple Views & Tooltip



Figure: Split View: Dendrogram, Scaffold Tree & Tooltip

# Realization & Technical Details

- Freely available under GNU GPL v3
- Implemented in Java for platform independent use
- Modular software architecture:
    - Seamless integration of novel views and analysis features
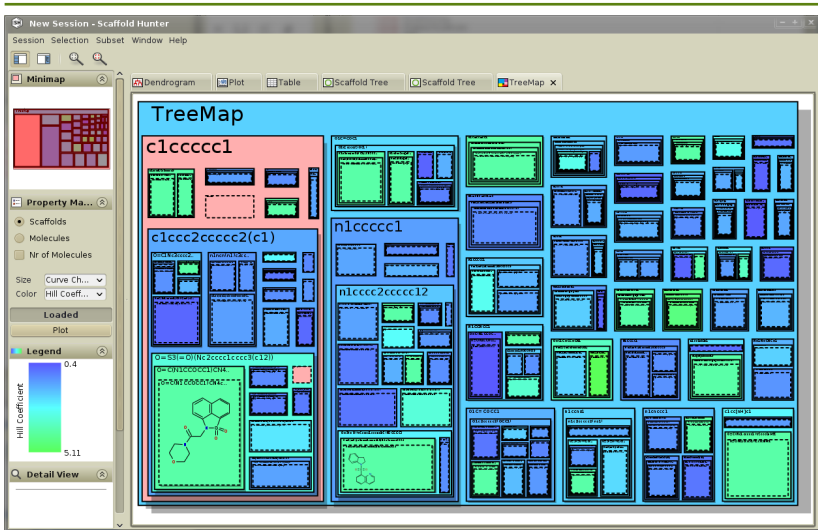    - Plugin system for data sources and property calculation

# Realization & Technical Details

- Freely available under GNU GPL v3
- Implemented in Java for platform independent use
- Modular software architecture:
  - Seamless integration of novel views and analysis features
  - Plugin system for data sources and property calculation

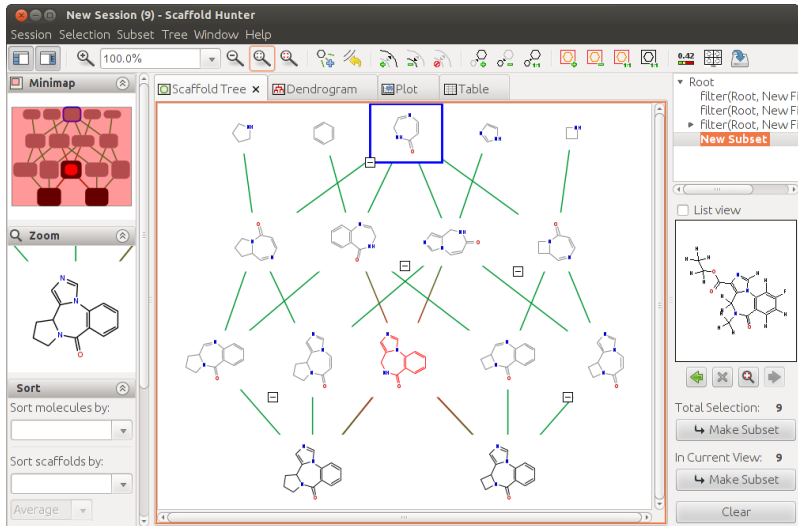## Toolkits & Database Support

- **Chemistry Development Kit (CDK):** Various cheminformatics tasks
- **Piccolo2D:** Zoomable user interfaces
- **Batik:** SVG support
- **Hibernate:** Object-relational mapping
- **MySQL/HSQLDB:** Back-end databases

# Future Work: Scaffold TreeMaps



- Space-filling approach to visualize scaffold trees
- Google Summer of Code project 2013: Jeroen Lappenschaar

# Future Work: Scaffold Networks



- Visualization of multiple parent scaffolds (Sugiyama layout)
- Dynamic filtering of networks

# Conclusion

- Exploratory visual analysis of chemical compound databases
- Clustering and classification of molecular datasets
- Multiple complementary interconnected views

# Conclusion

- Exploratory visual analysis of chemical compound databases
- Clustering and classification of molecular datasets
- Multiple complementary interconnected views

## Development & Acknowledgements

- **TU Dortmund, Prof. Mutzel:** Nils Kriege, Till Schäfer
- **University of Sydney:** Dr. Karsten Klein
- **GSoC 2013:** Jeroen Lappenschaar
- **Cooperation:**
  - MPI for Molecular Physiology Dortmund (Prof. Waldmann)
  - Dr. Koch, Computational Molecular Design, TU Dortmund