

Canvas 2.1

User Manual

Canvas User Manual Copyright © 2014 Schrödinger, LLC. All rights reserved.

While care has been taken in the preparation of this publication, Schrödinger assumes no responsibility for errors or omissions, or for damages resulting from the use of the information contained herein.

Canvas, CombiGlide, ConfGen, Epik, Glide, Impact, Jaguar, Liaison, LigPrep, Maestro, Phase, Prime, PrimeX, QikProp, QikFit, QikSim, QSite, SiteMap, Strike, and WaterMap are trademarks of Schrödinger, LLC. Schrödinger, BioLuminate, and MacroModel are registered trademarks of Schrödinger, LLC. MCPRO is a trademark of William L. Jorgensen. DESMOND is a trademark of D. E. Shaw Research, LLC. Desmond is used with the permission of D. E. Shaw Research. All rights reserved. This publication may contain the trademarks of other companies.

Schrödinger software includes software and libraries provided by third parties. For details of the copyrights, and terms and conditions associated with such included third party software, use your browser to open [third_party_legal.html](#), which is in the docs folder of your Schrödinger software installation.

This publication may refer to other third party software not included in or with Schrödinger software ("such other third party software"), and provide links to third party Web sites ("linked sites"). References to such other third party software or linked sites do not constitute an endorsement by Schrödinger, LLC or its affiliates. Use of such other third party software and linked sites may be subject to third party license agreements and fees. Schrödinger, LLC and its affiliates have no responsibility or liability, directly or indirectly, for such other third party software and linked sites, or for damage resulting from the use thereof. Any warranties that we make regarding Schrödinger products and services do not apply to such other third party software or linked sites, or to the interaction between, or interoperability of, Schrödinger products and services and such other third party software.

August 2014

Contents

Document Conventions	ix
Chapter 1: Introduction	1
1.1 Running Schrödinger Software	1
1.2 Citing Canvas in Publications	2
Chapter 2: The Canvas Graphical Interface	5
2.1 Starting and Quitting Canvas	5
2.2 The Main Window	5
2.2.1 The Menu Bar	6
2.2.2 The Toolbars	7
2.2.3 Keyboard Shortcuts	9
2.3 Canvas Projects	10
2.3.1 Opening and Closing Projects	10
2.3.2 Importing Structures and Data	11
2.3.3 Exporting Structures and Data	15
2.3.4 The Project View Panel	18
2.3.5 The Messages View Panel	19
2.4 The Canvas Spreadsheet	19
2.4.1 Configuring and Navigating the Spreadsheet	19
2.4.2 Shortcut Menus	20
2.4.3 Selecting and Copying Rows, Columns, and Cells	22
2.4.4 Finding Text	23
2.4.5 Coloring Structures	24
2.4.6 Copying and Pasting Structures	24
2.4.7 Renaming Structures	25
2.4.8 Editing Structures	25
2.4.9 Editing Data Cells	28
2.4.10 Adding and Deleting Rows and Columns	29
2.4.11 Displaying Different Views of the Spreadsheet	29

2.4.12	Creating and Using Custom Views	30
2.4.12.1	Applying a Custom View to the Master View.....	32
2.4.12.2	Combining Views with Logical Operations	33
2.5	Organizing the Data	34
2.5.1	Sorting	34
2.5.2	Partitioning Rows Into Classes	35
2.5.2.1	Creating Partitions	35
2.5.2.2	Changing the Classes in a Partition	38
2.5.3	Creating a Heat Map.....	40
2.6	Filtering and Querying the Data.....	41
2.6.1	Filtering by Property Values.....	41
2.6.2	Filtering by Property Classes.....	42
2.6.3	Filtering by Structure.....	43
2.6.4	Filtering by Substructure.....	45
2.6.5	Detecting Duplicate Structures	47
2.7	Making Charts of the Data	47
2.7.1	Scatter Plots.....	48
2.7.2	Histograms.....	50
2.7.3	Pie Charts	51
2.7.4	Chart Settings.....	52
2.7.5	Saving Images of Charts.....	54
2.8	Calculating Statistics.....	55
2.9	Calculating New Properties from the Data	56
2.10	Viewing Structures in PyMOL	58
2.11	Setting Preferences	59
Chapter 3: Running Applications from Canvas.....		61
3.1	General Features	61
3.2	Molecular Properties.....	62
3.2.1	Calculating Molecular Properties	63
3.2.2	Selecting Representative Properties.....	65

3.3 Fingerprints	68
3.3.1 2D Fingerprints	68
3.3.2 3D Fingerprints from Pharmacophores.....	70
3.4 Similarity, Dissimilarity, and Clustering	73
3.4.1 Computing a Similarity or Distance Matrix.....	73
3.4.2 Screening Structures by Similarity or Distance.....	75
3.4.3 Screening Structures by Shape	76
3.4.4 Selecting Diverse Structures.....	78
3.4.5 Hole-Filling and Library Optimization	79
3.4.6 Comparing Libraries.....	81
3.4.7 Clustering Structures by Similarity or Distance	82
3.4.7.1 Hierarchical Clustering	83
3.4.7.2 Leader-Follower and K-Means Clustering.....	86
3.5 Building a Predictive Model	88
3.5.1 Multiple Linear Regression	89
3.5.2 Partial Least-Squares Regression	91
3.5.3 Kernel-Based Partial Least-Squares Regression.....	92
3.5.4 Principal Components Regression.....	95
3.5.5 Bayes Classification	96
3.5.6 Neural Networks	98
3.5.7 Recursive Partitioning	99
3.5.8 Self-Organizing Maps	101
3.6 Other Applications	104
3.6.1 Principal Components Analysis	104
3.6.2 Finding the Maximum Common Substructure.....	105
3.6.3 Scaffold Decomposition	108
3.6.4 R-Group Analysis.....	110
3.6.5 Running a 3D Minimization	110
3.7 Running External Applications	111
3.8 Running Python Scripts	114

Chapter 4: Using Canvas	115
4.1 Selecting Compounds Not Represented in a Library	115
4.2 Removing Duplicate Structures from a Project	116
4.3 Choosing a Fingerprint	117
4.4 Creating Modal Fingerprints	117
4.5 Tips	118
Chapter 5: Running Applications from the Command Line.....	119
5.1 Common Syntax Descriptions	119
5.2 2D Fingerprints	122
5.2.1 canvasFPGen	122
5.2.2 canvasFPCombine.....	127
5.2.3 canvasFPBinary2CSV	128
5.2.4 canvasCSV2FPBinary	129
5.3 3D Pharmacophore Fingerprints	130
5.4 Similarity, Dissimilarity, and Clustering	132
5.4.1 canvasCSV2PW.....	133
5.4.2 canvasCSVMatrix	133
5.4.3 canvasDBCS.....	134
5.4.4 canvasLibOpt	136
5.4.5 canvasFPHist.....	139
5.4.6 canvasFPMatrix	140
5.4.7 canvasHC.....	143
5.4.8 canvasHCBuild.....	146
5.4.9 canvasHCSelect.....	147
5.4.10 canvasTreeDraw	148
5.4.11 canvasKMeans.....	149
5.4.12 canvasLC	150
5.4.13 canvasPW2CSV.....	151

5.5 Model Building and Related Applications	152
5.5.1 canvasMDS	154
5.5.2 canvasMLR	155
5.5.3 canvasMolDescriptors	156
5.5.4 canvasPCA	160
5.5.5 canvasPCAGen	160
5.5.6 canvasPCAProj	161
5.5.7 canvasPCAReg	162
5.5.8 canvasPLS	163
5.5.9 canvasKPLS	164
5.5.10 canvasBayes	165
5.5.11 canvasNnet	167
5.5.12 canvasRP	168
5.5.13 canvasSOM and canvasSOMBits	170
5.6 Utilities	172
5.6.1 canvas_app	172
5.6.2 canvasConvert	172
5.6.3 canvasJob	174
5.6.4 canvasProjectDB	176
5.6.5 canvasSDMerge	179
5.6.6 canvasSearch	180
5.6.7 canvasMCS	182
5.6.8 canvasScaffold	184
5.7 Scripting with Canvas	187
References	189
Getting Help	191
Index	193

Document Conventions

In addition to the use of italics for names of documents, the font conventions that are used in this document are summarized in the table below.

Font	Example	Use
Sans serif	Project Table	Names of GUI features, such as panels, menus, menu items, buttons, and labels
Monospace	<code>\$SCHRODINGER/maestro</code>	File names, directory names, commands, environment variables, command input and output
Italic	<i>filename</i>	Text that the user must replace with a value
Sans serif uppercase	CTRL+H	Keyboard keys

Links to other locations in the current document or to other PDF documents are colored like this: [Document Conventions](#).

In descriptions of command syntax, the following UNIX conventions are used: braces { } enclose a choice of required items, square brackets [] enclose optional items, and the bar symbol | separates items in a list from which one item must be chosen. Lines of command syntax that wrap should be interpreted as a single command.

File name, path, and environment variable syntax is generally given with the UNIX conventions. To obtain the Windows conventions, replace the forward slash / with the backslash \ in path or directory names, and replace the \$ at the beginning of an environment variable with a % at each end. For example, `$SCHRODINGER/maestro` becomes `%SCHRODINGER%\maestro`.

Keyboard references are given in the Windows convention by default, with Mac equivalents in parentheses, for example CTRL+H (⌘H). Where Mac equivalents are not given, COMMAND should be read in place of CTRL. The convention CTRL-H is not used.

In this document, to *type* text means to type the required text in the specified location, and to *enter* text means to type the required text, then press the ENTER key.

References to literature sources are given in square brackets, like this: [10].

Introduction

Canvas is a cheminformatics package that provides a range of applications for structural and data analysis, including fingerprints, similarity searching, substructure searching, selection by diversity, clustering, building regression and classification models. The Canvas graphical interface is project-oriented and provides chemical structure storage and organization, data analysis and visualization, and access to the applications.

You can run Canvas applications from the graphical interface, from the command line, from KNIME with the Schrödinger KNIME Extensions, and from Python scripts with the Canvas Python API.

A good introduction to cheminformatics can be found in the book by Leach and Gillet [1].

1.1 Running Schrödinger Software

Schrödinger applications can be run from a graphical interface or from the command line. The software writes input and output files to a directory (folder) which is termed the *working directory*. If you run applications from the command line, the directory from which you run the application is the working directory for the job.

Linux:

To run any Schrödinger program on a Linux platform, or start a Schrödinger job on a remote host from a Linux platform, you must first set the SCHRODINGER environment variable to the installation directory for your Schrödinger software. To set this variable, enter the following command at a shell prompt:

cshtcsh: `setenv SCHRODINGER installation-directory`

bash/ksh: `export SCHRODINGER=installation-directory`

Once you have set the SCHRODINGER environment variable, you can run programs and utilities with the following commands:

```
$SCHRODINGER/program &  
$SCHRODINGER/utilities/utility &
```

You can start the Canvas interface with the following command:

```
$SCHRODINGER/canvas &
```

It is usually a good idea to change to the desired working directory before starting the Canvas interface. This directory then becomes the working directory.

Windows:

The primary way of running Schrödinger applications on a Windows platform is from a graphical interface. To start the Canvas interface, double-click on the Canvas icon, on a Canvas project, or on a structure file; or choose **Start** → **All Programs** → **Schrodinger-2014-3** → **Canvas**. You do not need to make any settings before starting Canvas or running programs. The default working directory is the Schrodinger folder in your Documents folder.

If you want to run applications from the command line, you can do so in one of the shells that are provided with the installation and have the Schrödinger environment set up:

- Schrödinger Command Prompt—DOS shell.
- Schrödinger Power Shell—Windows Power Shell (if available).

You can open these shells from **Start** → **All Programs** → **Schrodinger-2014-3**. You do not need to include the path to a program or utility when you type the command to run it. If you want access to Unix-style utilities (such as `awk`, `grep`, and `sed`), preface the commands with `sh`, or type `sh` in either of these shells to start a Unix-style shell.

Mac:

The primary way of running Schrödinger software on a Mac is from a graphical interface. To start the Canvas interface, click its icon on the dock. If there is no Canvas icon on the dock, you can put one there by dragging it from the `SchrodingerSuite2014-3` folder in your Applications folder. This folder contains icons for all the available interfaces. The default working directory is the Schrodinger folder in your Documents folder (`$HOME/Documents/Schrodinger`).

Running software from the command line is similar to Linux—open a terminal window and run the program. You can also start Canvas from the command line in the same way as on Linux. The default working directory is then the directory from which you start Canvas. You do not need to set the `SCHRODINGER` environment variable, as this is set in your default environment on installation. To set other variables, on OS X 10.7 use the command

```
defaults write ~/.MacOSX/environment variable "value"
```

and on OS X 10.8 and 10.9 use the command

```
launchctl setenv variable "value"
```

1.2 Citing Canvas in Publications

The use of this product should be acknowledged in publications as:

Canvas, version 2.1, Schrödinger, LLC, New York, NY, 2014.

Please also cite the following reference:

Duan, J.; Dixon, S.L.; Lowrie, J.F; Sherman, W. Analysis and comparison of 2D fingerprints: Insights into database screening performance using eight fingerprint methods. *J. Molec. Graph. Model.* **2010**, *29*, 157-170.

The Canvas Graphical Interface

The Canvas graphical interface is designed for project-oriented work and provides chemical structure storage, data analysis and visualization, and access to a wide range of applications that are also available through the command line. The interface also provides links to Maestro that allow you to easily transfer structures and data between the two applications.

2.1 Starting and Quitting Canvas

Linux: To start the Canvas interface, enter the following command:

```
$SCHRODINGER/canvas [ -proj projName.cnv [ -mae maeFile | -sd sdFile  
[-append|-replace] [-autodelete]] | -zippedproj projName.cnv[.]zip]
```

The `-proj` option opens the specified project. If the project exists, it is opened. If the project does not exist, it is created, and you must provide a Maestro or SD structure file to add to the project with the `-mae` option or the `-sd` option. If you specify a structure file with an existing project, the structures and properties in the file replace those in the existing project unless you specify `-append`. The `-autodelete` option automatically deletes the structure file after it has been imported into the project. You can open a zipped project with the `-zippedproj` option.

Windows: To start the Canvas interface, double-click the Canvas icon on the desktop, or choose `Start → Programs → Schrodinger-2014 → Canvas`. You can also double-click on a Canvas project, to start Canvas and open the project.

Mac: To start the Canvas interface, click the Canvas icon on the dock. If there is no Canvas icon on the dock, drag it from the `SchrodingerSuite2014` folder in your `Applications` folder to the dock. You can also start the interface from the command line in a terminal window, as for Linux.

To quit Canvas, choose `File → Quit` or type `CTRL+Q (⌘Q)` in the main window.

2.2 The Main Window

The Canvas main window has a menu bar, toolbars, a spreadsheet area, a Project View panel, a Messages View panel, and a status area. The two panels are docked into the main window by default but can be undocked. You can show or hide these views from the Project menu or from the context menu (right-click menu, shortcut menu). The status area reports statistics on the number of rows and columns: the total, the number visible, and the number selected.

	Structure	mol MW	#stars	#amine	#amidine	#acid	#amide	#rotor	#H
1	<chem>Cc1ccc2c(c1)nc(=O)c2</chem> mol 0004	Class 3: 200 - 300	0	0	0	0	0	1	0
2	<chem>Cc1ccc2c(c1)nc(=O)c2</chem> mol 0007	Class 2: 100 - 200	1	0	0	0	0	2	0
3	<chem>Cc1ccc2c(c1)nc(=O)c2</chem> mol 0008	Class 3: 200 - 300	0	0	0	0	0	1	0

Figure 2.1. The Canvas main window

The spreadsheet contains the 2D structures in the first column (including the structure name if there is one), followed by columns containing the properties of the structures. It is described in detail in [Section 2.4 on page 19](#).

2.2.1 The Menu Bar

The menu bar has the following menus:

- **File**—Open and close projects, import and export structures and data, set preferences, and quit Canvas. This menu is described in detail in [Section 2.3 on page 10](#).
- **Edit**—Copy cells, clear cell content, add or delete rows or columns, select cells, rows, or columns, find and select text. This menu is described in detail in [Section 2.4 on page 19](#).
- **View**—Create, delete, display, and perform operations on custom views. This menu is described in detail in [Section 2.4.11 on page 29](#).

- **Partition**—Create, store, modify, apply, and delete partitions of the spreadsheet rows. This menu is described in detail in [Section 2.5.2 on page 35](#).
- **Structure**—Configure the structure display and copy and paste structures. This menu is described in detail in [Section 2.4.1 on page 19](#) and [Section 2.4.5 on page 24](#).
- **Data**—Perform operations on the table data, such as sorting, filtering creating a heat map, calculating statistics and new properties, and editing data. This menu is described in detail in [Section 2.5 on page 34](#) and [Section 2.6 on page 41](#).
- **Query**—Query the database and return rows that match the query. This menu is described in detail in [Section 2.6 on page 41](#).
- **Chart**—Create and display scatter plots, histograms, and pie charts. This menu is described in detail in [Section 2.7 on page 47](#).
- **Applications**—Run Canvas applications. This menu is described in detail in [Chapter 3](#).
- **Tasks**—Run Canvas applications. This menu is organized by the kind of task, rather than the list of applications. It opens the same panels as the Applications menu.
- **Project**—Show or hide the Project View and the Message View.
- **Python** —Open a Python interpreter to run Python commands or scripts.
- **Scripts**—Install, manage, and run Python scripts. This menu is described in detail in [Section 3.8 on page 114](#).
- **Help**—Display online help, the user manual, Canvas legal notices, and the Knowledge Base; open the Diagnostics panel. This menu is described in detail on [page 191](#).

If you start Canvas without opening a project, only the File and Help menus are available.

2.2.2 The Toolbars

The Canvas interface has four toolbars: File, Edit, Select Rows, and Structure Size. You can show or hide the toolbars from the shortcut menu, and you can show or hide toolbar buttons in the Preferences panel (File → Preferences → Toolbar Buttons).

The File toolbar has buttons for many of the common actions. A brief description of the File toolbar is given below, with cross-references to fuller descriptions of the actions.



New Project

Create a new Canvas project. The current project is closed. See [Section 2.3.1 on page 10](#).



Open Project

Open an existing Canvas project. [Section 2.3.1 on page 10](#).



Import

Import structures into the project. See [Section 2.3.2 on page 11](#).



Export

Export structures from the project to an external file. See [Section 2.3.3 on page 15](#).



Find

Find text in the spreadsheet. Opens the Find panel. See [Section 2.4.4 on page 23](#).



Collapse to Selected

Collapse the view to show only the selected structures. Same as View → Collapse to Selected.



Hide Selected

Hide the selected structures. Same as View → Hide Selected. See [Section 2.4.11 on page 29](#).



Apply to Master

Apply the current custom view to the master view. Only available in custom views. See [Section 2.4.11 on page 29](#).



Save View

Save the current view. In the master view, the view is saved as a new custom view, and the Save Custom View dialog box opens. See [Section 2.4.11 on page 29](#).



Restore Default View / Revert All Changes

Restore the default view in the master view, or revert all changes to the view in a custom view. In either case, the view is restored to its initial (default) state. See [Section 2.4.4 on page 23](#).



Edit Spreadsheet

Edit cells in the spreadsheet. When you click this button, the Edit toolbar is displayed. See [Section Note: on page 28](#).



Add Row

Add a new empty row to the end of the spreadsheet.



Add Column

Add a new empty column to the right side of the spreadsheet.



Substructure Query

Run a substructure query. See [Section 2.6.4 on page 45](#).



Scatter Plot

Create a scatter plot. See [Section 2.7.1 on page 48](#).



Histogram

Create a histogram. See [Section 2.7.2 on page 50](#).



Pie Chart

Create a pie chart. See [Section 2.7.3 on page 51](#).

The Edit toolbar is only active when you are editing the spreadsheet. It is displayed in place of the File toolbar when you start editing, and has the following two buttons.



Save Spreadsheet Changes

Save the spreadsheet with the changes included. When you click this button, the spreadsheet is saved, and the File toolbar is redisplayed. See [Section Note: on page 28](#).



Exit Without Saving Spreadsheet Changes

Discard all the changes to the spreadsheet. When you click this button, the spreadsheet is reverted to its state before you started editing, and the File toolbar is redisplayed.

The Select Rows toolbar is described in [Section 2.4.3 on page 22](#), and the Structure Size toolbar is described in [Section 2.4.1 on page 19](#).

2.2.3 Keyboard Shortcuts

Many of the menu items have keyboard shortcuts. These shortcuts are listed in [Table 2.1](#).

Table 2.1. Keyboard shortcuts for menu items.

Menu Item	Shortcut
File → New Project	CTRL+N (⌘N)
File → Open Project	CTRL+O (⌘O) [letter O]
File → Close Project	CTRL+W (⌘W)
File → Quit	CTRL+Q (⌘Q)
Edit → Copy	CTRL+C (⌘C)
Edit → Select All	CTRL+A (⌘A)
Edit → Clear Selection	CTRL+U (⌘U)
Edit → Find	CTRL+F (⌘F)
View → Refresh	F5
Partition → New → From Selection	ALT+N (⌘N)
Partition → Open	ALT+O (⌘O) [letter O]
Partition → Add/Move to → Class <i>n</i>	ALT+ <i>n</i> (⌘ <i>n</i>)
Partition → Add/Move to → New Class	ALT+L (⌘L)
Partition → Unclassify	ALT+0 (⌘0) [number zero]
Partition → Close	ALT+W (⌘W)
Data → Undo Last Edit	CTRL+Z (⌘Z)
Data → Redo Last Edit	CTRL+SHIFT+Z (⇧⌘Z)

Table 2.1. Keyboard shortcuts for menu items. (Continued)

Menu Item	Shortcut
Python → Command Window	CTRL+P (⌘P)
Scripts → Manage	CTRL+M (⌘M)
Scripts → Import	CTRL+I (⌘I)

2.3 Canvas Projects

Canvas stores all of its data in a project, which is simply a directory that contains a number of subdirectories and files. A Canvas project directory has the extension `.cnv`. Canvas does not have scratch projects, so you must create a project before you can import structures.

2.3.1 Opening and Closing Projects

The basic operations of opening and closing projects can be performed from the File menu:

- **New Project**—Create a new Canvas project. You can also use CTRL+N (⌘N), or click the New Project toolbar button.



Opens a file selector in which you can navigate to a location and name the project. The default location on Linux is the directory from which Canvas was started, and on Windows is the Schrodinger folder in your documents folder. If a project is already open, the open project is closed first.

- **Open Project**—Open an existing project. You can also use CTRL+O (⌘O), or click the Open Project toolbar button.



Opens a file selector in which you can navigate to a location and select the project or a project archive. Project archives are unzipped in place.

On Windows, you can open a project in a new Canvas session by double-clicking the project or project archive icon.

- **Recent Projects**—Open one of the most recently used projects. The submenu displays a list of up to 10 projects, by default. You can set the number of projects listed in the Preferences panel (File → Preferences → History).

- **Close Project**—Close the current project. You can also use CTRL+W (⌘W). The data in a project is automatically saved whenever you make a change.

Canvas projects are saved whenever a change is made, so there is no Save menu item.

You can create a zipped archive of a project by choosing Archive Project. The archive is given a `.cnvzip` extension. When you open an archived project, changes that are made in the opened project are not automatically re-archived: you must choose Archive Project and overwrite the previous archive.

You can make a project read-only by removing write permission from the project directory (`.cnv`). Users who do not have write permission will be able to open the project in Canvas and import structures from the project, but not perform any operations that would change the project.

2.3.2 Importing Structures and Data

There are several ways of importing structures and data into a Canvas project. From inside Canvas, you can import from a file or from a project. From Maestro, you can export selected entries to a Canvas project—see [Section 3.5](#) of the *Maestro User Manual* for details.

To import structures and data from another Canvas project, choose File → Import Project. A file selector opens, in which you can navigate to a location and select the project. All structures and all data are imported.

To import structures and properties from a file into Canvas, choose File → Import, or click the Import button on the toolbar.



A file selector opens, in which you can navigate to a location, choose the format, and select one or more files. The supported file formats are:

- Maestro, compressed or uncompressed (`.mae`, `.mae.gz`, `.maegz`)
- SD file, compressed or uncompressed (`.sd`, `.sdf`, `.sd.gz`, `.sdf.gz`)
- CSV file with SMILES strings and properties (`.csv`)
- SMILES file with SMILES strings and optional titles (`.smi`).
- Canvas fingerprint file (`.fp`)

When you have selected the files, a dialog box opens ([Figure 2.3](#) or [Figure 2.2](#)) for each file type, so that you can set options for importing the structures and matching the properties. Properties are imported in the order in which they are present in the input files or project. When you close this dialog box, the import progress is shown in the status bar, at the bottom of the main panel. There is also a button you can use to stop the import operation (and discard the results).

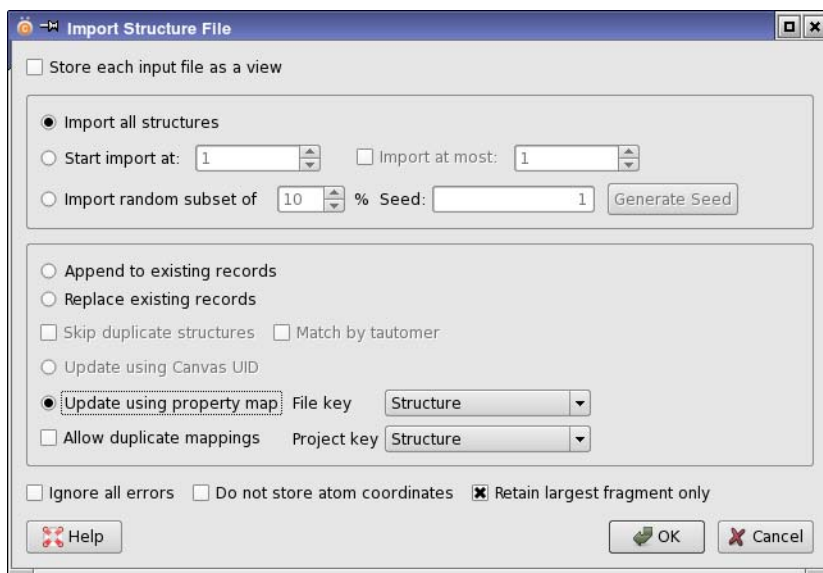


Figure 2.2. The Import Structure File dialog box

You can create a view for each file that is imported by selecting *Store each input file as a view*. The views are named after the input file. By default, no views are created on import.

All structures are imported by default. To restrict the range of structures, select *Start import at*, and enter the index of the first structure that you want to import in the box. The first structure in the file has the index 1. Then to import a given number of structures, select *Import at most* and enter the number of structures to import in the box. If this option is not selected, all structures from the first specified to the end of the file are imported. You can also import a random subset of structures, by selecting *Import random subset of*, specifying the percentage of structures, and generating a seed for the random number generator.

The structures that you import can be appended to the project (*Append to existing records*), used to replace the entire project (*Replace existing records*), or update the existing records where the structures are already in the project and append new structures (*Update using Canvas ID* or *Update using property map*).

You can also set options to skip structures that are already in the Canvas project, by selecting *Skip duplicate structures*. If any of the structures in the set that you selected for import are already in the project, they are not imported. If you want to treat tautomers as duplicates, select *Match by tautomer*.

Likewise, you can automatically skip problematic structures by selecting **Ignore all errors**. By default, you are prompted for an action when an unreadable structure is encountered. Problematic structures are written in the same format as the input file to the directory `failedimports`.

When you import from a structure file (Maestro or SD format), you can choose to import only the connectivity information, and discard the atom coordinates. To do so, select **Do not store atom coordinates**. By default, the coordinates are imported.

If there are structures in the input file that contain multiple fragments, you can import only the largest fragment by selecting **Retain largest fragment** (the default). Otherwise, all fragments are imported.

If you choose to update the project, you must choose an input file property from the **File** key option menu and a project property from the **Project** key option menu, so that the input structures with their data and project records can be matched, as follows:

- If the value of the file key property for a given structure matches a value of the project key property, the project record is updated with the properties from the file, and the structure is not changed.
- If multiple project records are matched for a single file key property value, and **Allow duplicate mappings** is selected, the properties are updated for all matching project records.
- If the file key property has multiple instances of the same value, then the last one in the file is used to update the project.
- If no match is found for the import property, a new record is added.

You might want to use a corporate ID, a unique SMILES string, or a title, as the property to match.

If you want to ensure that you do not add duplicate structures to the project, choose **Structure** for the **Import** property and the **Project** property, and ensure that **Allow duplicate mappings** is selected. The properties from the input file are still used to update those in the project if the input file structure matches a project structure.

When you import Canvas fingerprints, **Update using property map** is the only option available.

If you are importing structures that were previously exported from the same Canvas project, with the Canvas UID included, you can select **Update using Canvas UID** to update the properties in the project from the input file. The UID in the file is matched with the UID in the project, so using this option with another project is likely to corrupt the project data. An example of its use is to export binary fingerprints and reimport the bits as ordinary numeric properties with the value 0 or 1.

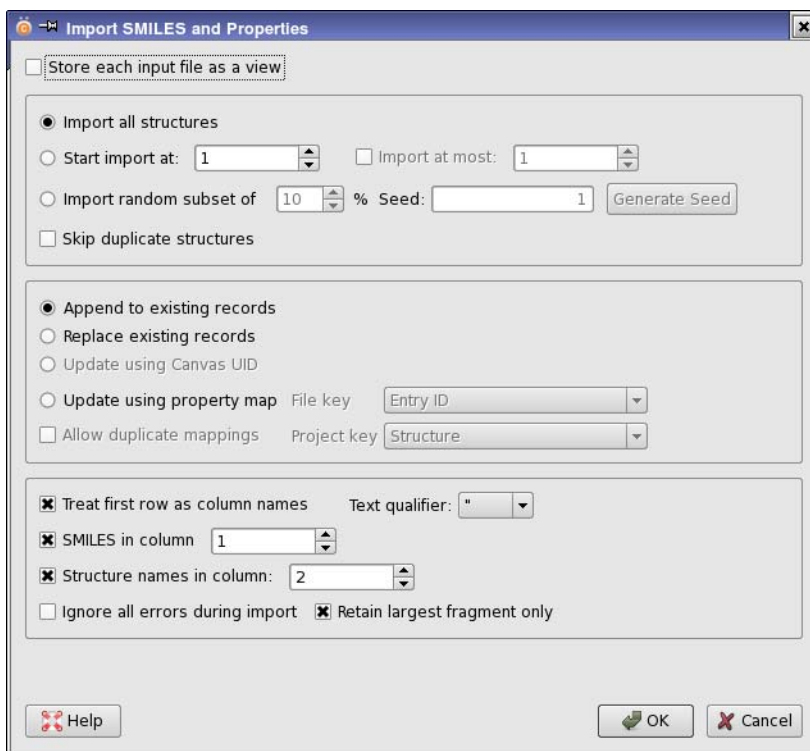


Figure 2.3. The Import SMILES and Properties dialog box

If you are importing structures in SMILES format with properties, from a .csv or .txt file, you have the choice of adding the structures with their properties, or just adding the properties. To add the structures and the properties, select SMILES in column and specify the column number that contains the SMILES strings. The default is 1. If you deselect this option, you can import properties without structures. If you use the Update using property map option, you can add properties from a CSV file to existing structures, provided you have a column in the CSV file to match a project property. Otherwise, the records in the input file are added to the project with an empty structure field.

If the first row of the file contains the names of the properties, select Treat first row as column names. If this option is deselected, the property names are set to Prop1, Prop2, and so on.

For string properties, you can choose the character that defines the string from the Text qualifier option menu. These characters are removed from the string when the string is added as a project property. The choices are a double quote, a single quote, or none. If you choose None, all quotes are preserved in the project. If you choose the single quote, then single quotes are

removed but double quotes are retained. Likewise, if you choose the double quote, then double quotes are removed but single quotes are retained.

To define which column the structure names (or titles) are in, select Structure names in column, and set the column number in the box. The default is column 2. If you deselect this option, the title is left blank. If you are importing from a SMILES file, you can select or deselect this option to indicate that the file contains a title, separated from the SMILES string by a space or a tab. The column number is then set to 2 and cannot be changed.

You can also import 2D structures on the basis of their common name or IUPAC name. To do this, create a new table row (Edit → Add Row or use the toolbar button), then double-click the Structure cell in this row to open the 2D structure editor. Choose Edit → Common or IUPAC Name Search, and enter the name in the text box. When you click OK in the 2D structure editor, the structure is added to the project and appears in the Structure cell.

2.3.3 Exporting Structures and Data

Canvas provides several ways of exporting structures and data. You can export the selected rows (rows in which all columns are selected) to a file or to Maestro. You can also copy and paste from Canvas into a spreadsheet. Properties are exported in the order in which they appear in the view, so you can rearrange the properties in the view into the order you want them to appear in the file before exporting.

To export structures and properties to a file, choose File → Export, or click the Export toolbar button.



A file selector opens, in which you can navigate to a location, choose the format, and name the file. The supported file formats are:

- Maestro, compressed or uncompressed (.mae, .mae.gz, .maegz)
- SD file, compressed or uncompressed (.sd, .sdf, .sd.gz, .sdf.gz)
- CSV file with SMILES strings and properties (.csv)
- SMILES file with SMILES strings and optional titles (.smi).
- Fingerprint file (.fp).

When you click Export, the Export Options dialog box opens. This dialog box allows you to select the properties and fingerprints to export, and set various options. When you dismiss this dialog box, the export starts, and the progress of the export is shown in the status bar, at the bottom of the main panel. There is also a button you can use to stop the export operation (and discard the results).

All properties are initially selected for export. If you want to export a random subset of the rows, select Export random subset, and specify the desired number of rows in the box. You can also specify the random number seed, by entering it in the Seed text box, or clicking Generate Seed.

If you are exporting to a structure file, you can export the coordinates as 2D coordinates by selecting Force 2D coordinates. When you do so, you can choose one of the three options for exporting hydrogens: Suppress hydrogens, Export polar hydrogens only, or Export all hydrogens. The default is to export only polar hydrogens. If the file is an SD file, you can choose to export in the newer version 3.0 format by selecting Use MDL version 3.0 format for SD files.

If you are exporting to a SMILES or CSV file, the default is to export unique (canonical) SMILES. To export a SMILES string that depends on the order of the atoms in the structure (first atom in the structure is first in the SMILES string), deselect Export unique SMILES.

The structures in a Canvas project have a unique identifier (UID), which can be exported by selecting Export Canvas UIDs. Exporting the UIDs allows you to add or change properties for the structures that you export, and then ensure that these properties are associated with the correct structures when you reimport them into the Canvas project.

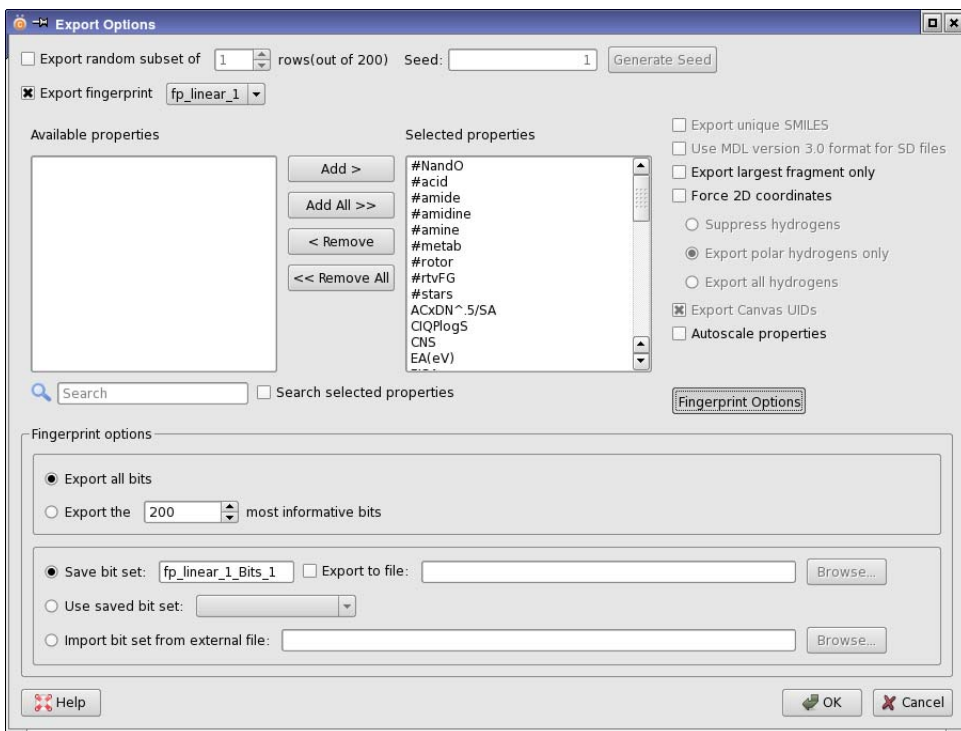


Figure 2.4. The Export Options dialog box showing fingerprint options

If you want to scale the properties so that they have a mean of 0 and a standard deviation of 1, you can do so by selecting Autoscale properties. This can be useful when exporting properties to other applications that require automatically scaled variables. If you want to scale only some of the properties rather than all properties exported, select them in the Selected properties list.

Fingerprints can be exported for two choices of file format: a CSV file (SMILES+Properties) or a binary file (Fingerprint, .fp). When you choose either of these file formats, the Export fingerprint option is available in the Export Options dialog box, and you can select this option to export fingerprints and choose the fingerprint to export. To display options for exporting the fingerprints, click Fingerprint Options.

You can decide whether to export all bits or the most significant n bits. By default, all bits are exported to a binary file, but only the most informative 200 bits are exported to a CSV file. This is because each bit in the fingerprint is stored as a separate property in the CSV file, with values of 0 or 1. The prefix for naming these properties can be entered in the Prefix for bit columns text box, which is absent if exporting to a binary file.

Another alternative is to export a set of bits that has previously been stored (a “bit set”). A bit set stores information on which bits were chosen for a given fingerprint to form the set, and thus indirectly stores information on the features that these bits represent. One reason for choosing a stored bit set is for consistency on export. The most informative 200 bits for one set of molecules might not be the most informative for another set, so if you export the most informative bits, the exported fingerprints will not have information about the same features in the molecules. You can avoid this discrepancy by exporting a saved bit set.

To use a bit set from the project, select Use saved bit set, and choose the bit set from the option menu. This menu lists saved bit sets for the fingerprint you are exporting. To use a bit set from a file, select Import bit set from external file, and click Browse to navigate to the file, or enter the file name in the text box. When you choose either of these options, the options to export all or the most significant bits are not available.

If you choose to export all or the most significant bits, the bit set that you use to export the fingerprint is saved in the project with the fingerprint, and can also be saved in an external file. Save bit set is automatically selected and you can enter a name for the bit set in the text box. The bit set is then saved in the project. To save it in an external file as well, select Export to file, and click Browse to navigate to the file, or enter the file name in the text box.

To export structures and properties from a Canvas project directly to Maestro, choose File → Export to Maestro. The structures and properties are exported from the Canvas project to a temporary Maestro file, then Maestro is started and the file contents are imported.

To copy rows or columns from Canvas to a spreadsheet program, such as Microsoft Excel or Open Office Calc, simply select the rows or columns and copy them to the clipboard with CTRL+C (⌘C) or Edit → Copy. You can then paste them into the chosen application.

2.3.4 The Project View Panel

This panel gives a summary of all the information that is available in the project. It lists the applications and the various classifications of the data (filters, views, partitions) as a tree. When information exists in any of these categories, you can expand the display by clicking the turner (the + in a square) to the left of the category name. To collapse the display, click the turner again (which now shows a – in the square). For applications, the display lists the names of the jobs that have been run. The status of the job is displayed in the status column. For views, partitions, and filters, the named entities for that category are listed. The description of a view is displayed in a tool tip when you pause the pointer over the view name in the tree.

Both the “branches” and the “leaves” in the tree have a shortcut menu that you can use to apply an action. For the branches, the shortcut menu has the following items, some of which are only available on particular branches:

- **Open**—Open an application panel or a filter panel; open the dialog box for selecting a view or a partition.
- **New**—Create a new partition from a property or a random assignment. Only on the Partitions branch.
- **Import Queries**—Import queries from another project. Opens a file selector so you can navigate to and select the project. Only on the Substructure Queries branch.
- **Sort**—Sort views in the default order (order of addition), ascending order or descending order. Only on the Views branch.
- **Delete All**—Delete all items in the branch and the data associated with them.

For the leaves, there is one common item, **Delete**, which deletes the item and its associated data from the project. For applications, the other available actions are:

- **View**—view the results of the job, by opening the relevant panel.
- **Incorporate**—incorporate the results of the application into the project.
- **Clone Job**—run a job with the same settings as the selected job.
- **Kill Job**—kill a running job.

Other actions that are available only for certain branches are:

- **Apply to Master**—Apply the view to the Master View. Only for Views.
- **Rename View**—change the name of the view and its description. Only for Views.
- **Apply Filter**—Apply the filter to the Master View. Only for filters.
- **Run Query**—Run the selected query. Only for Substructure Queries.
- **Delete Query**—Delete the selected query. Only for Substructure Queries.
- **Open**—Open a view or a partition. Only for Views and Partitions.

You can show and hide this panel by choosing **Project** → **Project View**. You can hide the panel by clicking the close button in the title bar. You can undock this panel and place it wherever you want, and you can redock it, using the docking button.



2.3.5 The Messages View Panel

The Messages View panel displays messages generated by actions on the project, including the running of applications. Log messages and job error messages are displayed in this panel.

You can show and hide this panel by choosing **Project** → **Messages View**. You can hide the panel by clicking the close button in the title bar. You can undock this panel and place it wherever you want, and you can redock it, using the docking button.



2.4 The Canvas Spreadsheet

The spreadsheet is the area where structures and properties are displayed. Although it does not have the full capabilities of a spreadsheet program such as Microsoft Excel or OpenOffice Calc, it does function in many respects in the same way.

2.4.1 Configuring and Navigating the Spreadsheet

The spreadsheet can be configured and navigated in much the same way as other spreadsheets, as described below.

You can **resize** individual rows and columns by dragging their borders. To resize multiple rows or columns simultaneously, select the rows or columns and drag one of the borders inside the selection. To resize a column to fit the data, double-click the right border of the column heading. If you have several columns selected, double-clicking the right border of the heading of one of these columns fits each column to the width of its data.

You can **move** a column by dragging it to the new location, or you can move it to either end of the table by right-clicking in the heading row and choosing **Move to Start** or **Move to End**. You can **sort** the columns by name by right-clicking in the heading row and choosing **Sort columns by name**. The sort is case-insensitive.

You can change the number of decimal places displayed for property values in the Preferences panel—see [Section 2.11 on page 59](#).

You can **scroll** the spreadsheet with the mouse wheel or the scroll boxes. If you scroll rapidly in a vertical direction, structures are not drawn until you pause or release the scroll box, because Canvas retrieves structures and properties on demand from an SQLite database. This allows you to scroll through a project containing thousands of rows with no delay. You can also jump to a specific row by entering the row index in the Rows text box on the toolbar and clicking Jump to or pressing TAB, then ENTER.

You can **step** through cells in the table with the arrow keys or the tab key. Stepping through cells changes the selection. The current cell has a border drawn around it in reverse video (black on white, yellow on blue).

The first column in the spreadsheet displays the **structure**. This is a 2D image created from the connection table for the structure. You can hide the structure and just display its title in this column by choosing Structure → Hide. To show it again, choose Structure → Show. You can also use the toolbar button labeled Show (if the structures are hidden) or Hide (if they are displayed). When the structures in the spreadsheet are hidden, the tool tip for the table cell displays the structure as well as the title. You can also show or hide the hydrogen atoms in the structure that are not displayed by default by choosing Structure → Show All Hydrogens.

The structures are scaled automatically so that the bond lengths are a reasonable size and approximately equal for all structures. Structures that are too large for the cell with the scaling are scaled down in size. If you want the structures to always fill the cell, deselect Autoscale on the Structure menu. You can also change the structure size with the Structure Size slider. If you adjust the structure sizes, you can reset them to the default sizes by clicking Reset.

2.4.2 Shortcut Menus

The spreadsheet has shortcut menus for the selected rows, for the column headings, and for the selected structure cell. These are displayed when you right-click and hold.

The selected rows shortcut menu has the following items, given with their menu equivalents:

- Delete—delete the selected rows (Edit → Delete).
- Collapse to Selected—show only the selected rows (View → Collapse to Selected).
- Hide—hide the selected rows (View → Hide Selected).
- Export—export the selected rows to a file (File → Export).
- Export to Maestro—export the selected rows to a new Maestro session (File → Export to Maestro).

The column heading shortcut menu has the following items when you right-click on a single column:

- Sort by *property* (Ascending)—Sort the rows by the values of the property in this column, in ascending order.
- Sort by *property* (Descending)—Sort the rows by the values of the property in this column, in descending order.
- Sort Columns by Name—Rearrange all columns in alphabetical order.
- Restore Natural Column Order—restore the columns to the order in which they were added to the project.
- Rename Column—Change the name of the column. Opens the Edit Column Name dialog box, so you can change the name and description.
- Move to Start—Move the column to the first position (leftmost) in the spreadsheet.
- Move to End—Move the column to the last position (rightmost) in the spreadsheet.
- Hide—Hide the column from the current view. The data remains in the project.
- Delete—Remove the property in this column from the project, and from all views.
- Set Property Value—Set the value of the property in the selected rows to a specified value. A blank value clears the property values.
- Create Modal Fingerprint—Create a modal fingerprint row using the data in the column, which must be a fingerprint column. Opens a dialog box in which you can choose the selected or the visible rows to define the modal fingerprint, and name the fingerprint. A new row is added, with the name displayed in the Structure column. This item is only available if you right-click on a fingerprint column.

When multiple columns are selected, only the Hide and Delete items are present.

The selected structure cell shortcut menu has the following items:

- Edit Structure—Edit the structure in the 2D structure editor.
- Edit Structure Name—Change the name of the structure.
- Copy—Copy the structure to the clipboard in SD format.
- Copy SMILES—Copy the structure to the clipboard as a SMILES string
- Copy Image—Copy an image of the structure to the clipboard.
- Paste—Paste the SD-format structure from the clipboard into the current cell to replace the current structure.

- **Paste SMILES**—Paste a SMILES string from the clipboard into the current structure cell to replace the current structure.

2.4.3 Selecting and Copying Rows, Columns, and Cells

There are several ways of selecting cells in the spreadsheet. The selected cells are highlighted with a blue background.

Selection of cells with the mouse follows the usual rules for tables:

- To select a row, click the row index.
- To select multiple rows, use shift-click and control-click, or drag over the row indices.
- To select a column, click the column heading.
- To select multiple columns, use shift-click and control-click.
- To select a cell, click in the cell.
- To select multiple cells, use shift-click and control-click, or drag over the cells.

You can move the selection to a neighboring cell with the arrow keys. The arrow keys move the current cell, and the new cell is selected. This action clears the selection of any other cells.

You can make selections or change the selection from the Edit menu:

- **Select All**—Select all cells in the spreadsheet.
- **Clear Selection**—Clear the selection of cells.
- **Invert Selection**—Select the cells that are not selected, and deselect the cells that are selected.
- **Expand Selection To Entire Rows**—Expand the selection so that the entire row is selected for any row that has one or more cells selected. Entire rows must be selected to export data.
- **Expand Selection To Entire Columns**—Expand the selection so that the entire column is selected for any column that has one or more cells selected.

To select specific rows, enter the row indices in the **Rows** text box on the toolbar, and click **Select** (or press **TAB** twice, then **ENTER**). The row indices are given as a comma-separated list, and can include row ranges. Ranges can be specified either as $n:m$ or $n-m$. To specify a range that includes the first or the last row, use $:m$ or $n:$. Examples are given in the tool tip for the text box. The current cell is set to the first cell in the row selection.

To select cells whose contents match a text expression, use the **Find** panel, which you open from the **Edit** menu. When the text is found, click **Select All Matches**. The **Find** panel is described in detail in [Section 2.4.4](#).

You can also make some selections from the View menu:

- **Select Rows with Missing Values**—Select the rows for which one or more columns do not have values.
- **Select Columns with Missing Values**—Select the columns for which one or more rows do not have values.
- **Select Rows with Invalid Chemistry**—Select the rows for which a structure is not available. This can happen on import if the structure block is badly formatted or the SMILES string is invalid; in this case a string noting the error is displayed in the structure cell.

Once you have a selection, you can copy it to the clipboard with CTRL+C (⌘C) or Edit → Copy. The selection that is copied is the rectangular region that contains the selected cells. You can then paste the selection into a spreadsheet application such as Microsoft Excel (Windows), OpenOffice Calc (Linux), or iWork Numbers (Mac).

2.4.4 Finding Text

If you want to search for text in the spreadsheet, you can use the Find panel. To open this panel, choose Edit → Find in the main window, or click the Find toolbar button.

To specify the search string, enter it in the Find what text box. You can qualify the search to match the entire contents of the cell and to match the case. If you want to find the exact string, select **Literal string**. To match a string using the wildcard characters * (zero or more characters of any type) and ? (one character of any type), select **Allow wildcards**. If you want to perform more complex pattern matching, select **Interpret as regular expression**, and enter a regular expression in the text box. The regular expression uses metacharacters to represent variable elements in the search, and follows the regular expression syntax for Perl. To match numeric values to the precision given in the table, select **Interpret as floating point value**.

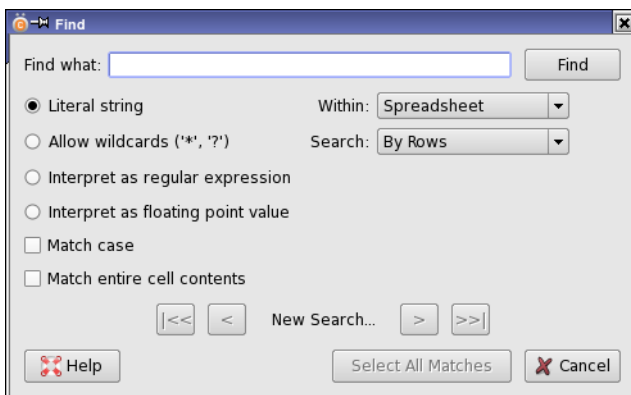


Figure 2.5. The Find panel

To restrict the range of the search in the spreadsheet, choose Spreadsheet or Selected cells from the Within option menu. To find columns rather than data values, choose Column Names.

To set the search order, choose By rows or By columns from the Search option menu.

Once you have set up the text, the range, and the search order, click the Find button to perform the search. All instances of the search string are located with a single click. You can select all the matching cells in the spreadsheet by clicking Select All Matches.

You can navigate through the matches with the navigation buttons. The double-arrow buttons |<< and >>| skip to the first and last matches; the single-arrow buttons step to the next and the previous match. The current match is indicated between the forward and backward navigation buttons. The current cell in the spreadsheet is set to the current match, and the spreadsheet scrolls to the current cell.

2.4.5 Coloring Structures

You can control some aspects of the structure coloring. To color the non-carbon atoms, choose Structure → Color Heteroatoms. If you do a substructure query, you can color the matches by choosing Structure → Highlight Matches. The color used can be set by choosing Structure → Choose Highlight Color.

2.4.6 Copying and Pasting Structures

Structures can be copied and pasted in MDL (SD) format within Canvas. To copy the current structure in MDL format, choose Structure → Copy Structure; to paste a structure into a cell (replacing the existing structure), select the cell and choose Structure → Paste. Pasting cannot be undone.

You can copy and paste structures as SMILES strings, both within Canvas and to other applications. To copy in SMILES format, choose Structure → Copy SMILES. You can then paste the structures into another application, or into a different cell in the Structure column. To paste a structure into a cell, select the cell and choose Structure → Paste SMILES. The structure that you paste can come from any external application that supports SMILES format, such as Maestro, ChemDraw, or ISIS/Draw. SMILES strings can also be pasted into a text box that requires a SMILES string. The structure is pasted as a 2D structure.

Copying and pasting structures can be done from the shortcut menu for the structure cell: right click in the cell and choose the menu item. These items are the same as on the Structure menu.

For creating presentations or other documents, you can copy an image of a structure, by selecting the structure and choosing Structure → Copy Image. The image is copied to the clipboard, and can then be pasted into any application that accepts an image.

2.4.7 Renaming Structures

You can change the structure name by selecting the structure and choosing Structure → Edit Name or right-clicking in the structure cell and choosing Edit Name. In the dialog box that opens, you can type in a new name. You can also change the names used for all the structures by choosing a property for the structure name. To do so, choose Structure → Change Structure Name Property. You can then choose a property from those in the spreadsheet, and store the current structure name in a new property. Properties with real number values are not available. The property you use for the structure name can be kept in the spreadsheet, or it can be removed from the spreadsheet.

2.4.8 Editing Structures

If you want to edit a structure, you can open the 2D structure editor by any of these methods:

- Double-click the cell containing the structure.
- Right-click on the structure and choose Edit Structure.
- Select the structure in the spreadsheet and choose Structure → Edit Structure.

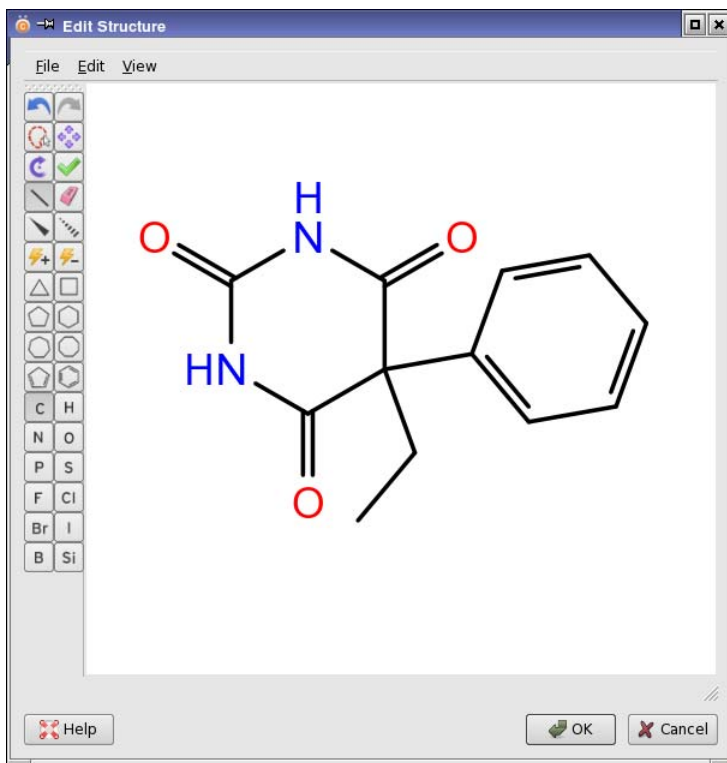


Figure 2.6. The Edit Structure dialog box

The 2D structure editor (Edit Structure dialog box) consists of a menu bar, a toolbar, and a display area.

The toolbar has three sets of buttons. The first set of buttons is for performing actions, which are described below. The second set is a collection of rings that you can add to the structure. The third set is a collection of elements, which you can use to change atoms to a chosen element. The selection of the element can also be done by typing the element symbol (case insensitive).



Undo

Undo the last action. Can be performed multiple times.



Redo

Redo the last action that was undone. Can be performed multiple times.



Lasso

Select atoms by drawing around them (“lassoing”). The selected atoms can be moved, rotated, deleted, copied.



Move

Move (translate) the selected atoms or the entire structure. Drag in the drawing area to move the structure. Drag the selection to move the selected atoms. (This is the default action for dragging a selection, so you do not need to click this button before dragging it.) You can also move the structure by dragging with the right mouse button.



Rotate

Rotate the selected atoms or the entire structure. Drag in the drawing area to rotate the structure. The angle through which the structure has been rotated is displayed near the center of rotation while you are rotating it. If rotation does not seem to be working properly, check whether you have other structures in the drawing area that are not visible. You can also rotate the structure by dragging with the middle mouse button.



Cleanup

Clean up the structure (2D coordinates) and normalize orientation and arrangement of groups. If you have atoms selected, only those atoms are cleaned up.



Draw

Draw a single bond. Click on an atom to draw a bond to a new (carbon) atom, click on a bond to add another bond between two atoms.



Erase

Delete atoms or bonds. Clicking on a carbon atom deletes it; clicking on a non-carbon atom changes it back to carbon, then clicking again deletes it. Clicking on a bond removes one bond.



Wedge Bond

Draw a wedge bond. Clicking on an atom draws a wedge bond to a new (carbon) atom; clicking on a single bond converts it to a wedge bond.



Dashed Bond

Draw a dashed bond. Clicking on an atom draws a dashed bond to a new (carbon) atom; clicking on a single bond converts it to a dashed bond.



Increase Charge

Increase the formal charge on an atom by 1.



Decrease Charge

Decrease the formal charge on an atom by 1.

Some common tasks are listed below:

- To add a singly bonded carbon, click the Draw button, then click on the atom that you want to attach it to. Use the Wedge Bond or Dashed Bond button instead to indicate the stereochemistry.
- To define the stereochemistry of a bond, click the Wedge Bond or Dashed Bond button, then click on the bond you want to change.
- To increase the order of a bond, click the Draw button, then click on the bond.
- To decrease the order of a bond, click the Erase button, then click on the bond.
- To delete an atom, click the Erase button, then click on the atom you want to delete. If it is not a carbon atom, it is first converted to carbon, then clicking again deletes the atom.
- To add a ring fragment, click the fragment button, then click on the atom or bond that you want to replace with the fragment.
- To add a functional group, pause the pointer over the atom you want to add it to, then type – followed by its name or SMARTS pattern, and press ENTER, e.g. -et to add an ethyl group. The named functional groups are listed in [Table 2.2](#), with the corresponding SMARTS patterns; the names are case-insensitive. You can also replace an atom with a functional group, by omitting the initial –.
- To change the element, click the button for the desired element or type the element symbol, then click on the atom you want to change.
- To increase or decrease the charge on an atom, click the Increase Charge or Decrease Charge button, then click on the atom. You can also pause the pointer over the atom and type in the charge.
- To move the structure, click the Move button, then drag in the display area.
- To rotate the structure, click the Rotate button, then drag in the display area.

Table 2.2. Functional group names and SMARTS patterns for use in editing structures.

Name	SMARTS	Name	SMARTS	Name	SMARTS
me	C	ph	C(C=C1)=C-C=C1	coo-	C(=O)[O-]
et	CC	bz	CC(C=C1)=C-C=C1	no2	[N+](=O)[O-]
pr	CCC	chex	C(CC1)CCC1	so2	S(=O)=O
ipr	C(C)C	cp	C(C=C1)C=C1	so2n	S(=O)(=O)N
nbu	CCCC	ome	OC	so3	S(=O)(=O)[O-]
ibu	CC(C)C	oet	OCC	po3	P(=O)([O-])[O-]
sbu	C(C)CC	cooh	C(=O)O	con	C(=O)N
tbu	C(C)(C)C	coome	C(=O)OC	nco	NC=O
nhex	CCCCCC	cooet	C(=O)OCC	cn	C#N

Note: If you edit a 3D structure in the 2D structure editor, all 3D structural information is discarded, and the structure is stored as 2D. To edit a structure in 3D, you can export it to Maestro, edit it, and reimport it.

2.4.9 Editing Data Cells

You can edit the data in the spreadsheet cells directly by choosing Data → Edit Spreadsheet, or clicking the Edit Spreadsheet toolbar button. The spreadsheet is placed in editing mode, and you can click in any cell and change the data in the cell, with the exception of the structure cells (first column) and the cells of a partition. You can undo the last edit with CTRL+Z (⌘Z) or Data → Undo Last Edit, and you can redo an edit with CTRL+SHIFT+Z (⇧⌘Z) or Data → Redo Last Edit. When you have finished making changes, choose Data → Save Spreadsheet Changes to make the changes permanent and exit editing mode, or click the Save Spreadsheet Changes toolbar button. If you want to discard all the changes you have made and return the spreadsheet to the state it was in before editing, choose Data → Exit Without Saving Spreadsheet Changes, or click the Exit Without Saving Spreadsheet Changes toolbar button.

To clear the contents of one or more cells, select the cells and choose Edit → Clear.

If you want to assign a single value to the selected rows for a particular property, right-click on the column heading and choose Set Property Value. A dialog box opens, in which you can set the value in a text box. When you click OK, the property values in the selected rows are set to the value you entered in the text box. If there was no value specified in the text box, the values in the selected rows are cleared.

2.4.10 Adding and Deleting Rows and Columns

To append a row to the spreadsheet, choose **Edit** → **Add Row** or click the toolbar button. A new empty row is added, and the Structure cell is selected, so that you can paste a structure into the cell, or build one in the **Edit Structure** dialog box.



To append a column to the spreadsheet, choose **Edit** → **Add Column**. The **New Column** dialog box opens, in which you can name the column. The name must be unique. When you click **OK**, a new empty column is added.



To delete rows or columns from the spreadsheet, select the rows or columns and choose **Edit** → **Delete**. The rows or columns are removed after confirmation. You can also right-click on a column heading and choose **Delete from Project** to delete the columns.

2.4.11 Displaying Different Views of the Spreadsheet

There are occasions on which you might want to see only a particular range of rows and columns, arranged in a particular way. Such a selection and arrangement is called a “view”. Changes to the view can be made manually, or as a result of some other process. When the view is changed, the title bar indicates that the view is modified. To revert the changes, choose **View** → **Restore Default View**, or click the **Restore Default View** toolbar button.



The **View** menu allows you to hide or show the selected rows and columns. To display a subset of columns and rows, select them and choose **View** → **Collapse to Selected**. This action displays all rows and columns that have any cells selected. If you want to hide particular rows or columns, select them and choose **View** → **Hide Selected**, or click the **Hide Selected** button.



To select the columns that you want to show, choose **View** → **Manage Properties**. The **Manage Properties** dialog box has two lists: **Available properties**, which are the hidden properties, and **Visible properties**, which are the shown properties. By default, all properties are on the **Visible properties** list. You can transfer properties between lists by selecting them and clicking **Add** or **Remove**. You can limit the list to those properties that match a text string by typing the string in

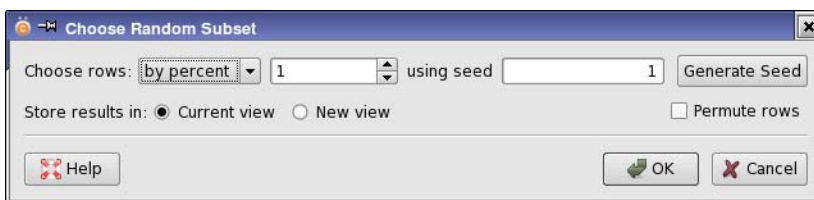


Figure 2.7. The Choose Random Subset dialog box.

the search text box, after setting the Search visible properties option to select the list to limit. When you click OK, the spreadsheet shows only those properties that were in the Visible properties list.

If you want to display a random subset of rows, choose View → Random Subset. You can select the number or percentage of rows in the subset, enter or generate a seed for the random number generator, randomize the row order, and decide whether to save the results in the current view or create a new custom view (see below).

When you run a query or filter the data, only the rows that match are displayed. These operations are described in [Section 2.6 on page 41](#).

2.4.12 Creating and Using Custom Views

Changes made to the view are transient: as soon as another action is taken, the previous arrangement and visibility of the rows and columns is modified. If you want to preserve a particular view of the data, you can create a *custom view* that is saved in the project and can be displayed later. As well as providing a view of the data, custom views can be used to select data to run applications.

To **create** a custom view, choose View → Save As, or click the Save View toolbar button.



The Save Custom View dialog box opens, and you can name the view and provide a description. The view can be based on all rows that are visible in the parent view (Use all rows), or the selected rows (Use selected rows). When you create the view, you can choose to display it immediately by selecting Open view. The named view is stored in the project and is listed under the View node in the Project View panel, and on menus in application panels.

Once you have created a custom view, you can work in this view to run queries, filter the data, and create charts. However, you cannot change the data in a custom view: this must be done in the *master view*, which is the view of the data shown in the main window.

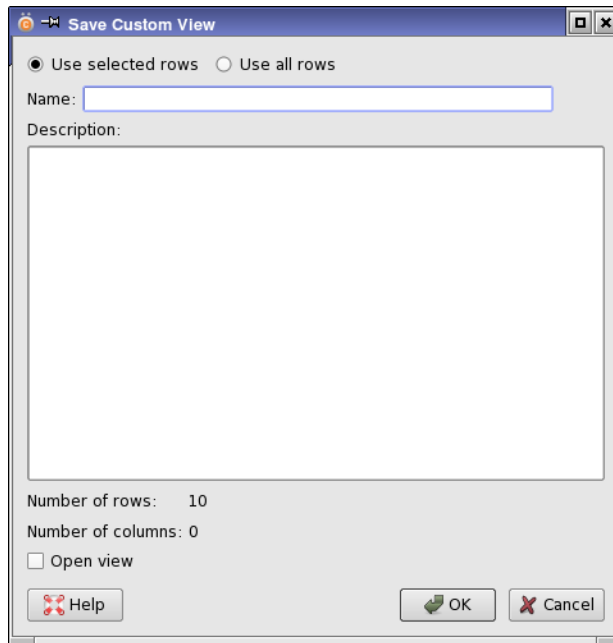


Figure 2.8. The Save Custom View dialog box

If you want to **rename** a custom view, you can do so in the Project View panel. Expand the View node, right-click on the view and choose Rename View. The Save Custom View dialog box opens, and you can enter a new name and change the description. If the view is open, it is closed before renaming.

If you want to **examine** or change the description of the view, choose View → Details. The Details panel for the view opens, and displays the name, description, row and column count, and information on the time it was created, last saved, and last modified. You can edit the description in this panel.

Custom views can be **opened** by choosing View → Open. This action opens a dialog box, in which you can select the view that you want to open. Another way to open a custom view is to expand the View node in the Project View panel, and double-click the desired view, or right-click it and choose Open. The tool tip for the view shows its description.

Custom views are always opened in a **separate window**. This window is a reduced version of the main window: the menu bar does not have the Applications and Project menus, the File menu only has export items, and the Project View and Messages View windows are absent. The absence of these components reflects the fact that the custom view is intended for viewing and rearranging the data, but not making changes in the data.

Changes that you make to a custom view are not automatically saved. To **save** a custom view, choose View → Save in the custom view window. You can also save a custom view as a new custom view, in the same way as from the master view.

You can **undo** changes to the view, by choosing View → Undo change or View → Undo All View Changes. The latter operation undoes changes back to the last time the view was saved.

When you have finished working in a view, you can **close** it by choosing View → Close. If you no longer want to keep the view, choose View → Delete to **delete** it, or right-click on its name in the Project View and choose Delete.

To **sort** the list of views that is displayed in the Project View panel and in application menus, right-click on the View node, choose Sort, then choose a direction, from Ascending, Descending, or Default. The first two items sort the views in alphabetical order; the last restores the default order, which is the order the views were created.

If you want to run **applications** from the command line with the rows in a custom view, you can **export** the row IDs for the view by choosing View → Export Row IDs. The binary file that is written can then be used in conjunction with the project to run an application on the rows in the view. If you run a command line job that generates new properties with `canvasJob` (see [Section 5.6.3 on page 174](#)), you should choose View → Refresh to view the new properties.

Custom views can also be created from the results of an **application**. These views have a File menu, for exporting the view; a View menu, from which you can apply the view to the master, save the view in the project and manage properties; and a Structure menu with items for the scaling and visibility of the structures. These views show only the structures by default, but properties can be added by choosing View → Manage Properties and selecting properties. Once you have selected a property, this property is added to the view whenever you open a view of the results of that application.

2.4.12.1 Applying a Custom View to the Master View

Custom views are independent views of the spreadsheet. However, you can apply a custom view to the master view by choosing View → Apply to Master. This action is also available from the shortcut menu for the view in the Project View panel.

Choosing View → Apply to Master allows you to apply the current row and column selection and ordering to the master view. The Apply to Master dialog box opens, offering choices for the rows, columns, and visibility of the rows in the master view.

To change the visible rows in the master view to match those in the current view, select Retain current for Rows. If there are rows in the current view that are not visible in the master view, and you want to make them visible, select Allow expansion of Master View row set. If you want

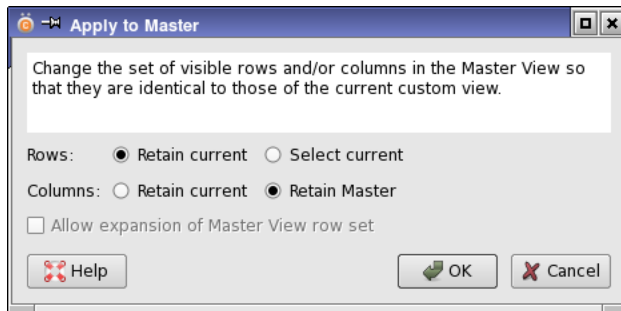


Figure 2.9. *The Apply to Master dialog box.*

to change the visible columns in the master view to match the current view, select Retain current for Columns; to keep the master view visible columns the same, select Retain Master.

If you only want to apply the selection of rows to the master view rather than change the visibility, select Select current for Rows and Retain Master for Columns.

2.4.12.2 Combining Views with Logical Operations

Custom views can be combined with logical operations to update the current view or create new custom views. This capability is useful, for example if you have created custom views with several different filters, and want to combine the results. To combine other custom views with the current view, choose View → Logic, which opens the Apply Logic to View dialog box.

Four logical operations are available:

- And—keep only the rows or columns that are in both views.
- Not—keep the rows or columns that are not in the other view.
- Or—keep the rows or columns that are in either view.
- Exclusive Or—keep the rows or columns that are in either view but not in both.

You can choose whether to retain the rows or the columns, by selecting the appropriate Retain option. These selections restrict the scope of the logical operations: they are not applied to the rows if you select Retain rows in current view, and likewise for columns.

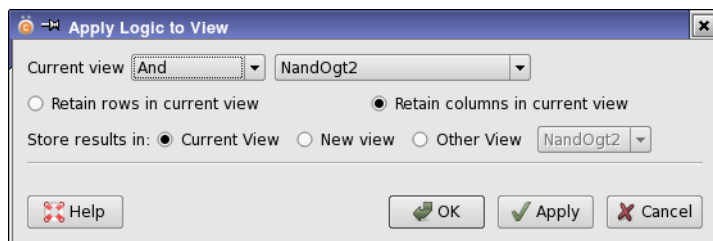


Figure 2.10. *The Apply Logic to View dialog box.*

If you want to create a new view, select **New view**. When you click **Apply** or **OK**, the **Save Custom View** dialog box opens, in which you can name, describe, and save the view.

If you want to store the results in an existing view (overwriting that view), select **Other view** and select the view from the option menu.

To apply multiple logical operations to a view, first select **Current view**. After selecting the operation, the other view, and the row and column retention options, click **Apply**. You can then repeat the selection and click **Apply** as many times as you like.

If you want to apply a sequence of logical operations on a view and create a new view with the result, but keep the original view, choose **View → Save** before you start the process. After you have created the new view, choose **View → Undo All View Changes** in the original view. This operation undoes the operations back to the last save.

2.5 Organizing the Data

Canvas provides various ways of organizing the rows in the spreadsheet. Some of them, such as clustering, involve running an application. These means of organizing data are treated in [Chapter 3](#). The simplest ways of organizing the data are sorting and partitioning, which are described in the first two subsections below. Another way of highlighting the data values is by coloring the table cells using a heat map.

2.5.1 Sorting

One of the simplest ways of organizing the data is sorting. To sort the rows in any view of the spreadsheet, choose **Data → Sort**, which opens the **Sort** dialog box. This dialog box allows you to sort the data (**New sort**), to reverse the order of the previous sort (**Reverse current order**), or restore the order of the rows as they are stored in the database (**Restore natural order**).

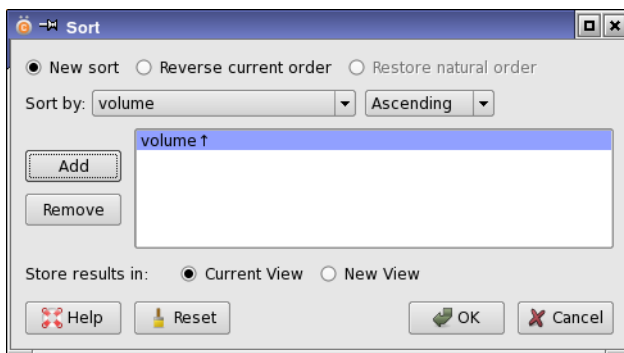


Figure 2.11. The Sort dialog box.

To start a new sort of the data, choose **New sort**, then start the selection of the sort keys. You can sort by multiple property values, in the order in which you add them to the sort list. To add a sort key to the sort list, choose the property from the **Sort by** option menu, choose the sort direction, from **Ascending** or **Descending**, and click **Add**. The property and the direction is added to the end of the sort list. You can then add another sort key.

Sorting is done in the order in which the keys are specified in the sort list. If you want to change the order, you can drag a sort key to a new position in the list.

You can apply the results of the sort to the current view or create a new view with the results, by choosing one of the **Store results** in options. If you create a new view, the **Save Custom View** dialog box opens, and you can name, describe, and save the sort results as a new view.

For a simple sort on the values of a single property, right-click in the column heading and choose **Sort by *property* (Ascending)** or **Sort by *property* (Descending)**.

2.5.2 Partitioning Rows Into Classes

Another way of organizing the structures (or rows) is to partition them into classes. The classes can be based on the value of some property, or they can be assigned manually or randomly.

Partitioning the rows creates a new property, which is shown as a column in the spreadsheet whose cells are colored according to the class that the row belongs to. The column heading has **Class::** prepended to the name, so it is distinguished from any actual property of the same name. The value of this property for any row is the class value. You can then sort the rows by the class value to display classes contiguously.

2.5.2.1 Creating Partitions

There are three ways of creating a partition from scratch: from the selected rows, from values of a property, and by random assignment. You can also save an existing partition with a new name (**Partition** → **Save As**) and then edit the partition (**Partition** → **Edit**).

To create a new partition from the selection, choose **Partition** → **New** → **From Selection** or press **ALT+N**. The **New Partition from Selection** dialog box opens. In this dialog box, you can name the partition and provide a description. When you click **OK**, a new partition is created, and the **New Class** dialog box opens with the selected rows as the members of the class. If you selected **Create classes for selected and unselected rows**, the unselected rows are assigned as members of a second class. You can edit the class values in the table, and assign a color by right-clicking in the cell in the **Color** column and choosing **Change Color**. A color selector opens, in which you can choose the desired color.

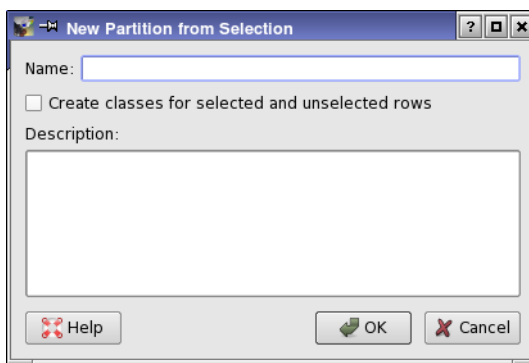


Figure 2.12. The New Partition from Selection dialog box.

To create a partition based on property values, choose Partition → New → From Property. The New Partition from Property dialog box opens. You can choose the property from the Property option menu. The property name is the default name for the partition, but you can edit it in the Name text box. You can enter text describing the partition in the Description text area.

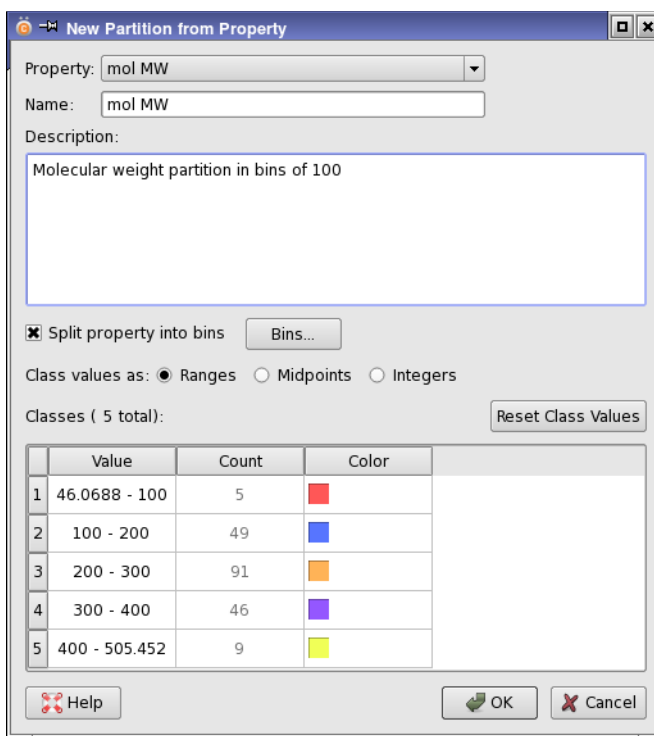


Figure 2.13. The New Partition from Property dialog box.

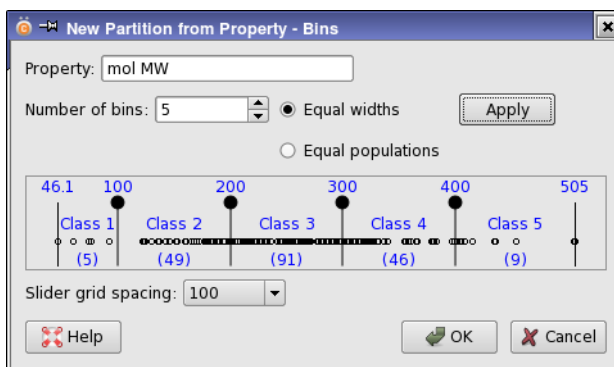


Figure 2.14. The Bins dialog box.

If you want to create a class for each property value, ensure that Split property into bins is not selected. If you want to create classes that cover a range of property values, select this option and click Bins, to specify the widths of the classes. Property values that are less than the minimum value are placed in the first class, and property values that are greater than the maximum value are placed in the last class. The lower boundary of each of these classes is included in the class; the upper boundary is in the next class.

You can choose how to set the class values by selecting one of the Class values as options:

- Ranges—Use a string m – n that indicates the range of property values in the class.
- Midpoints—Use the midpoint of the range.
- Integers—Use the bin index for the class value.

For classes based on individual property values, the class value is the property value, and these options are not available.

The Classes table shows the class values, the count, and the color used in the partition column in the spreadsheet. You can edit the class values if you wish by editing the table cells. For classes based on property value ranges, you can reset the class values to the defaults as described above by clicking Reset. You can change the color for a class by right-clicking in the cell in the Color column and choosing Change Color. A color selector opens, in which you can choose the desired color.

To create a random partition, choose Partition → New → Random Split. The New Partition from Random Split dialog box opens. A standard name is supplied in the Name text box, which you can edit. You can enter text describing the partition in the Description text area. You can specify the number of classes, and set the seed used to generate random numbers for assignment of the rows to classes. To assign the rows to classes after making a change, click Update.

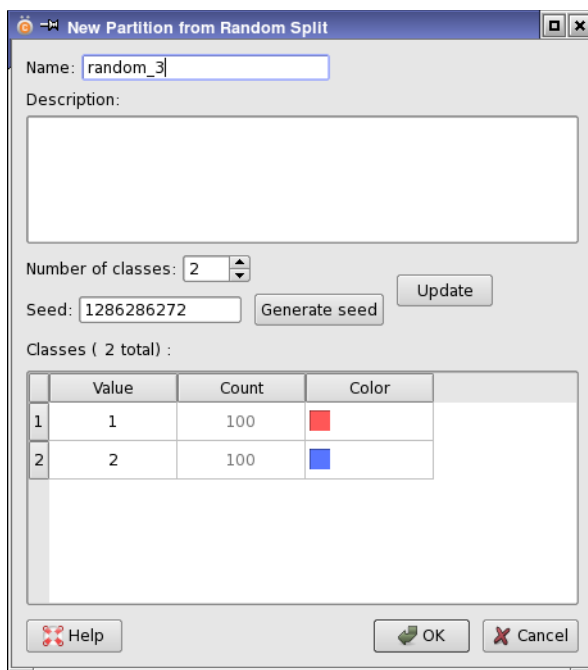


Figure 2.15. The *New Partition from Random Split* dialog box.

As for the *New Partition from Property* dialog box, the *Classes* table shows the class values, the count, and the color used in the partition column in the spreadsheet. The class values here are a sequence of integers, which you can change by editing the table cells. You can change the values to text strings if you want labels rather than numeric values. You can change the color for a class by right-clicking in the cell in the *Color* column and choosing *Change Color*, which opens a color selector.

2.5.2.2 Changing the Classes in a Partition

When a partition has been created, you can add classes, move rows to new classes or declassify them, edit the partition, or delete the partition.

To create a new class, first select the rows in the spreadsheet. The rows can be unclassified rows, or rows that are already in one or more classes. Next, choose *Partition* → *Add/Move to* → *New Class*, or press ALT+L (⌘L). The *New Class* dialog box opens, showing the class table. You can then edit the class value and change the color, as described above.

To move rows to an existing class, select the rows and choose *Partition* → *Add/Move to* → *Class n*, or type ALT+n (⌘n), where *n* is the class index. Classes are indexed sequentially from 1, and the class value is displayed on the *Add/Move to* submenu.

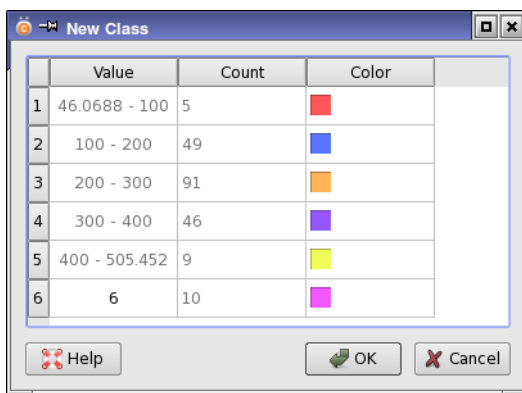


Figure 2.17. The New Class dialog box.

To remove rows from their classes, so that they are not in any class, select the rows and choose Partition → Unclassify, or press ALT+0 (⌘0). You do not need to select entire rows to add or remove them from a class: any row with a selected cell is included in operations on a class.

To make changes to the partitions, choose Partition → Edit. In the Edit Partition dialog box, you can change the name and description of the partition, change the class values and colors, and delete classes. To delete a class, right-click on it in the Class table and choose Delete.

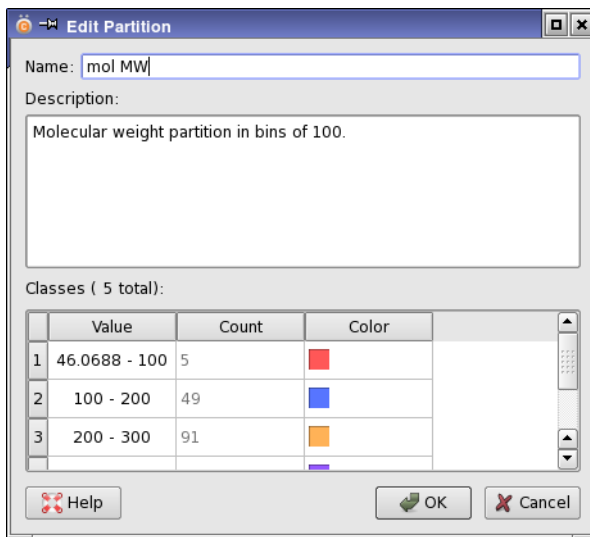


Figure 2.16. The Edit Partition dialog box.

2.5.3 Creating a Heat Map

To highlight the values of a property, you can apply a heat map for the property, which colors the cells in the spreadsheet according to the property value. The colors provide a visual organization of the property values. To set up the heat maps, choose **Data** → **Heat Map**, which opens the Heat Map dialog box.

You can choose the properties to which you want to apply a heat map, by selecting them in the Properties list and clicking **Add**. Properties to which you no longer want to apply a heat map can be removed by selecting them in the heat map table and clicking **Remove**.

Four color spaces are available:

- **RGB**—Color is interpolated linearly between the RGB colors for the maximum and minimum values.
- **Rainbow (HSV)**—Color follows the spectrum (rainbow) between the maximum and minimum values.
- **Transparency**—Transparency increases from the maximum value and the minimum value towards the mean value (which is white).
- **Color Wheel (discrete)**—A selection of colors from the spectrum is applied to individual integer values. For real-valued properties, the value is rounded down to the nearest integer and the color assigned based on the integer value.

Only one of these spaces can be used at a time, and is applied to all properties. You can set the colors for the minimum and maximum values of the property range, by right-clicking in the **Min** or **Max** column of the heat map table, and choosing **Edit Color**.

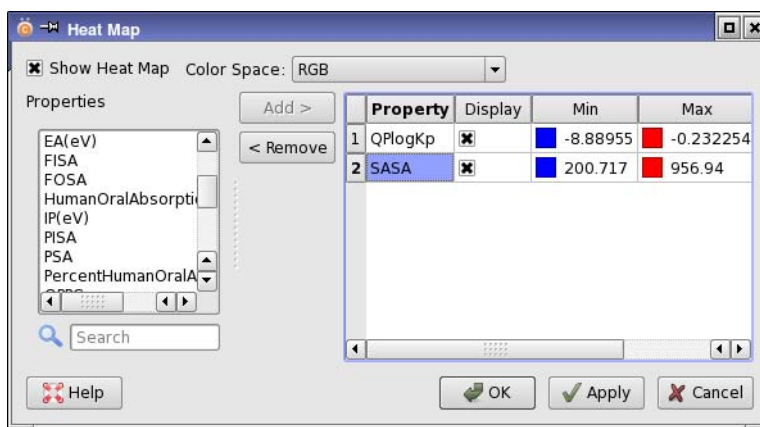


Figure 2.18. The Heat Map dialog box.

Overall display of the heat maps is controlled by selecting or deselecting Show Heat Map. Display of heat maps for individual properties is controlled by selecting or deselecting the check box in the Display column for the property.

The property ranges are divided evenly to create the heat map. If you want to customize the property ranges used for the colors, you can instead create a partition for the property and color the classes with the colors that you choose.

2.6 Filtering and Querying the Data

Canvas allows you to filter the data on the values of properties or on counts of functional groups, which are represented as SMARTS patterns. You can also run a query to find all molecules that have a particular substructure. Filters and queries can be stored and be reused. All types of filters and queries can be run from the Project View panel with default settings, using the shortcut menu.

2.6.1 Filtering by Property Values

You can filter the structures in the spreadsheet by a combination of property values. This task is performed in the Property Filter dialog box, which you can open by choosing Data → Property Filter or Query → Properties.

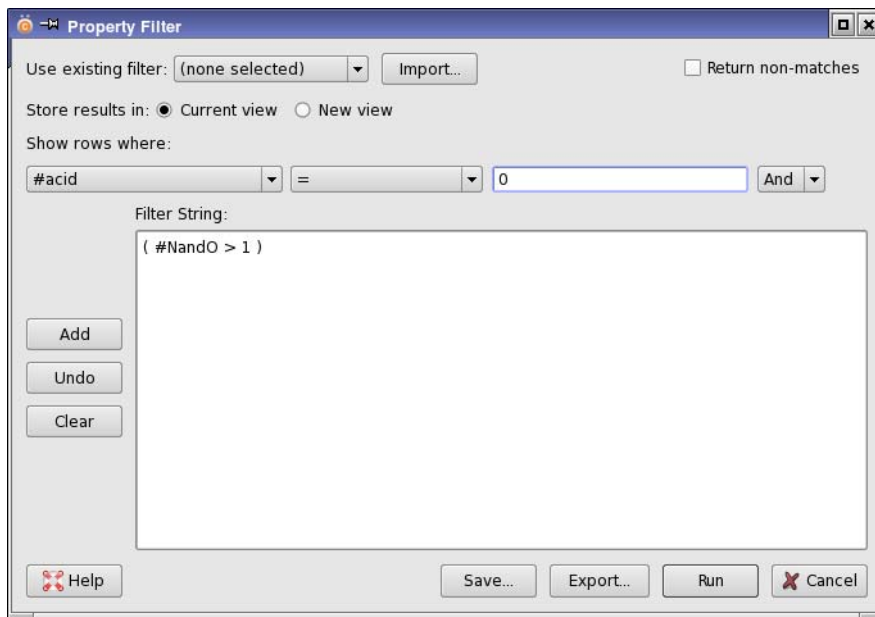


Figure 2.19. The Property Filter dialog box.

To construct a new filter, you first set up a condition on a property by selecting a property and a relational operator from the Show rows where option menus, and entering a value in the text box. The operator selection includes arithmetic operators, operators for matching strings, and operators for highest and lowest counts and percentages. Once the condition is defined, click Add to add it to the Filter String text area, where the entire filter string is displayed. To add more conditions, construct each condition in the same way, select And or Or to define its relationship to the previous conditions, then click Add. To remove the last condition added, click Undo. If you need to start again, click Clear, and the filter string is removed. The string can also be edited directly.

You can load, edit, and apply existing filters. To load a filter from the project, choose the filter from the Use existing filter option menu. The Filter String is replaced with the string for the selected filter. To load a filter from outside the project, click Import. A file selector opens, in which you can navigate to and choose the filter file (.flt).

When you have created a filter, you can save it in the project by clicking Save, and entering a name in the Save Property Filter dialog box. The filter is then listed on the Use existing filter option menu, and under Property Filters in the Project View window. You can also export the filter to a .flt file so that it can be used in other projects, by clicking Export.

To delete a filter, right click on it in the Project View window and choose Delete.

The filter can be applied to the current view, or the results can be stored as a new view. If you select New view, the Save Custom View dialog box opens when you apply the filter, so that you can name and describe the view.

To apply the filter, click Run. If you have not saved the filter, you are prompted to save it. If you click Yes, the Save Custom View dialog box opens; if you click No, the filter is applied without saving it.

To apply the inverse of the filter, select Return non-matches before applying the filter.

2.6.2 Filtering by Property Classes

In a similar way to filtering by property values, you can filter the data by the property classes stored in a partition. This feature allows you to perform more complex filtering on a single property. To do so, first open the partition, then choose Partition → Filter. In the Partition Filter panel, you can select the classes whose rows you want to show (or hide). To hide the selected classes instead of showing only the selected classes, select Hide rows for selected classes. You can also choose whether or not to hide rows that are not in any class with the Hide unclassified rows option. As for property filters, the filter is applied to the visible rows, so the number of rows never increases when you apply the filter. The filter can only be applied to the current view.



Figure 2.20. The Partition Filter dialog box.

If you click Apply, the panel remains open, and you can interact with the spreadsheet normally. This allows you to make different selections of classes and apply the filter. For example, if you want to display different classes in turn, you can select the first class and click Apply, then perform whatever actions you want to do in the spreadsheet. You can then select the next class and click Apply, and perform actions.

When you click OK, the most recent selection is applied and the panel closes.

2.6.3 Filtering by Structure

You can filter the structures based on their chemical content, which includes the usual functional groups. This task is performed in the Structure Filter dialog box, which you open by choosing Structure → Structure Filter.

Each condition on the occurrence of a chemical feature is specified by a row of the table. Complex filters can be set up by using multiple rows. The structures that pass all the filters are returned when you run the job. Thus, there is an implicit AND applied between the rows. You can invert the filter to return the structures that fail any of the filters, by selecting Invert logic (retain failures).

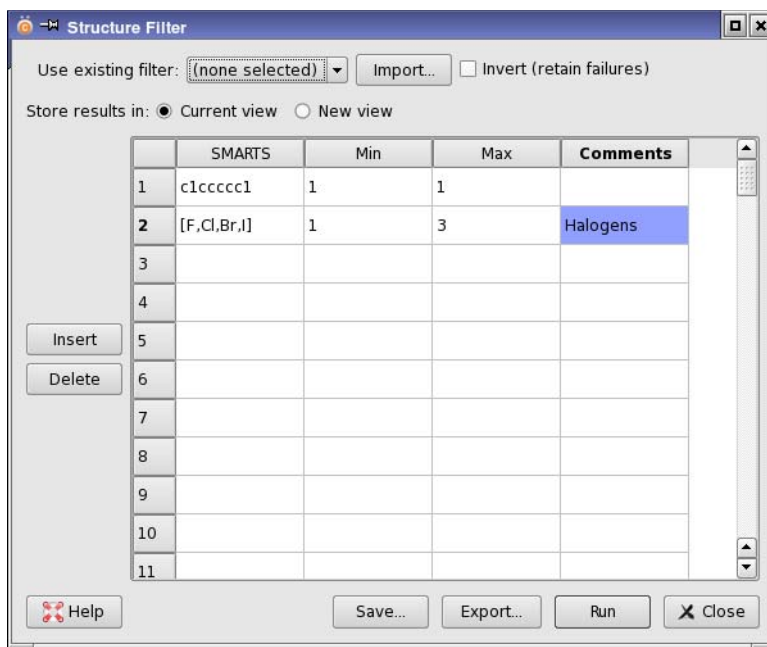


Figure 2.21. The Structure Filter dialog box.

The chemical features are defined in terms of SMARTS patterns. In the SMARTS column you can enter the SMARTS patterns to filter on, and in the Min and Max columns, you can enter the minimum and maximum number of occurrences of the SMARTS pattern in a structure that are required for the structure to pass through the filter. You must enter at least a minimum or a maximum, and you can enter both. In the Comments column you can annotate the SMARTS pattern—for example, to give names to functional groups. You can insert a row above the selected row by clicking Insert, and you can delete the selected rows by clicking Delete.

Structure files can be used to generate SMARTS patterns. To do this, import a file in Maestro, SD, or SMILES only format with the Import button. The structures in the file are converted into SMARTS patterns, and added to the rows of the table.

When you have created a filter, you can save it in the project for later use by clicking Save. A dialog box opens, in which you can name the filter. The filter is then listed on the Use option menu and under Structure Filters in the Project View window. If you want to store the filter outside the project, for example, to import into another project, click Export. A file selector opens, in which you can save the filter in a file with a .flt extension.

To delete a filter, right click on it in the Project View window and choose Delete.

You can also load, edit, and apply existing filters. To load a filter from the project, choose the filter from the Use option menu. The table is replaced with the data for the selected filter. To apply multiple filters, select Apply, then choose the filter from the option menu. The selected filter is appended to the patterns in the table. Four standard filters are available: REOS (“rapid elimination of swill”), PAINS1, PAINS2, and PAINS3 (pan assay interference filters [15] for different classes of assays, corresponding to filter C, filter B, and filter A, respectively). To load a filter from outside the project, click Import. A file selector opens, in which you can navigate to and choose the filter file (.flt).

The filter can be applied to the current view, or the results can be stored as a new view. If you select New view, the Save Custom View dialog box opens when you apply the filter, so that you can name and describe the view.

To apply the filter, click Run. If you have not saved the filter, you are prompted to save it. If you click Yes, the Save Custom View dialog box opens; if you click No, the filter is applied without saving it. If Highlight Matches is selected on the Structure menu, the substructures are highlighted on the structures that match.

2.6.4 Filtering by Substructure

You can filter the structures in the database by running a substructure query. This task is performed in the Substructure Query dialog box, which you open by choosing Query → Substructure Query or clicking the Substructure Query toolbar button.



The query can be run on the selected rows in the spreadsheet, on all rows, or on an external project. To make the choice, select the appropriate Screen option.

By default, the query returns all structures that contain the specified substructure. If you want to return structures that match the query exactly, select Require exact match. This choice will find all instances of a single molecule, and distinguish it from structures that contain other molecules, such as water. You can also return structures that do not match the query by selecting Invert (return non-matches). If you are querying an external project, the matching structures are appended to the current project.

The query can be specified either by providing a SMARTS string, or by building or importing a 2D structure.

- To use a SMARTS string, choose SMARTS from the Query Source option menu. You can enter or paste a SMARTS string into the Enter SMARTS query text box, or choose a previously used SMARTS string from the option menu.

- To use a 2D structure, choose Builder from the Query Source option menu. You can build a structure, by clicking New and drawing the structure in the Edit Structure panel; import a query from a MOL file or an SD file, by clicking Import and navigating to the file; or select a previously used query from the option menu. Only the first structure is used if an SD file is chosen.

The display area shows the query as a 2D structure. You can set the number of queries on the option menus in the Preferences panel, in the History category.

By default, the structures are prescreened with the 2D index for the structure (a special 512-bit fingerprint), to skip structures that cannot match the query, before checking for a match. This feature speeds up the search. It can be used for prescreening if the query is an exact SMARTS pattern; it cannot be used if the query contains SMARTS primitives that disqualify it as an exact substructure, such as [C,N], or [*]. To turn off the prescreening, deselect Use 2D index.

You can create or edit a query with the 2D builder or with ChemDraw (if ChemDraw is installed) by clicking the Build button in the Mol File tab. If you have already specified a file, it is loaded into the 2D editing application. When you close the application, the query is saved and imported into Canvas. You have the option of saving the query in a named file for future use.

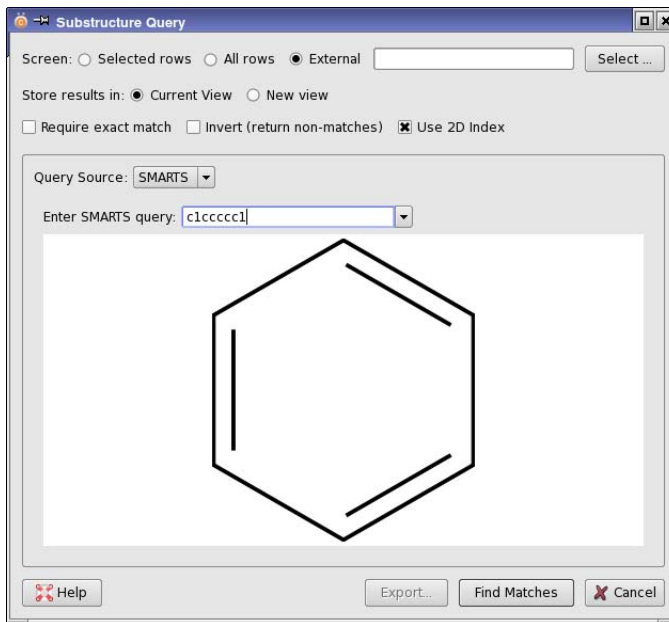


Figure 2.22. The Substructure Query dialog box.

The query results can be applied to the current view, or the results can be stored as a new view. If you select **New view**, the **Save Custom View** dialog box opens when you run the query, so that you can name and describe the view. The matches are highlighted on the structure if **Structure** → **Highlight Matches** is selected.

To run the query, click **Find Matches**.

When you have run the query, you can align the 2D structure of the matches, as shown in the spreadsheet, to a reference structure of your choice. The alignment is done with a least-squares method on the matching atoms. The reference structure is the first match by default, but you can choose any single matching structure as the reference. To perform the alignment, choose **Structure** → **Align Matches**. If you deselect this option, the matching structures are returned to their default orientation. The alignment is removed if you delete the reference structure, modify it so that it no longer matches the query, or hide it in a custom view and then save the view. or if you undo the query filter. If you want to align to a different structure, deselect the **Align Matches** option, then choose the new reference structure and select the option again.

2.6.5 Detecting Duplicate Structures

There is no restriction in Canvas on the structures that you can import, other than any filters you apply in the process of importing. This means that the project can contain duplicates. If you want to locate the duplicates, you can do so by choosing **Structure** → **Detect Duplicates**. Two columns are added to the spreadsheet, named **Has duplicates** and **Duplicate of**. These properties contain the UID for the first of the set of identical structures. **Has duplicates** is set only for the first occurrence of the structure; **Duplicate of** is set only for the second and subsequent occurrences of a structure.

You can use this column to select the unique structures or the duplicates with a property filter. In the **Property Filter** dialog box, choose the **Duplicate of** property under **Show rows where**, and choose **has a value** for the operator to select the duplicates; and choose **has no value** to select the unique structures (including ones that have no duplicates). To select just the unique structures from the sets of duplicates, choose **has a value** for the **Has duplicates** property.

2.7 Making Charts of the Data

Canvas provides several kinds of charts for visualizing the data: scatter plots, histograms, and frequency pie charts. To create a chart, choose the chart type from the **Chart** menu, or click the toolbar button for the chart type.



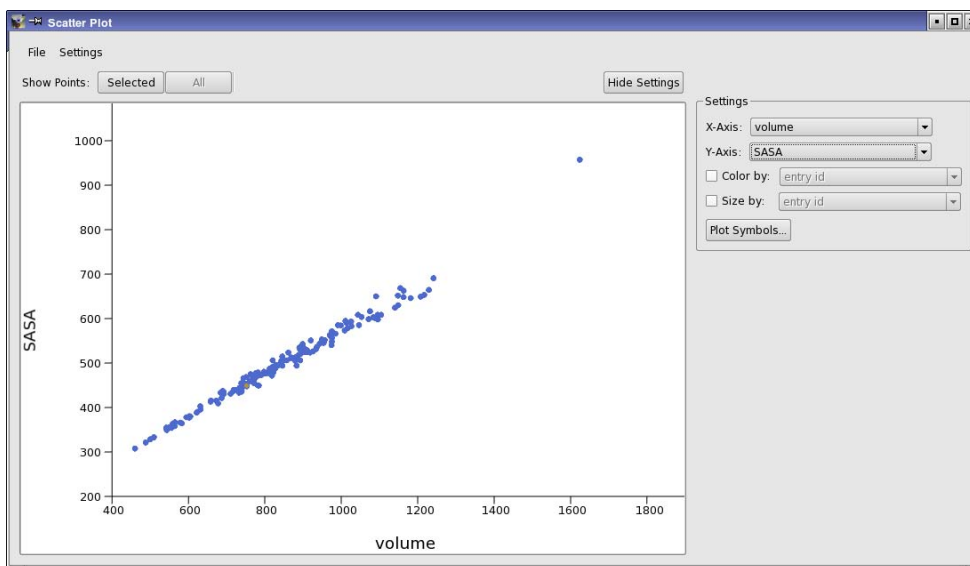


Figure 2.23. The Scatter Plot panel.

Each chart panel has a File menu, a Settings menu, and one or more option menus to choose the properties that are used for the chart. The properties listed on the option menus are the properties that are selected in the current view.

2.7.1 Scatter Plots

Scatter plots can be used to represent up to four different variables. The default plot shows two variables, on the x and y axes. The other two variables can be represented by the color and the size of the plot symbols. The property values for each point are displayed in a tool tip.

The properties that are displayed can be chosen from the X-Axis, Y-Axis, Color by, and Size by option menus in the Settings section of the panel. If this section is not displayed, click Show Settings. The color and size properties can also be chosen in the Plot Symbols dialog box, which you open by choosing Settings → Plot Symbols. To make your choice, select By property and choose the property from the option menu. For numerical properties, you can enter the limits of the range in the Limits text boxes. For coloring by a numerical property, you can set the colors for the ends of the property range with the Min and Max buttons. The colors between the end points are a linear interpolation in RGB color space. For coloring by a categorical property, the colors are chosen from a color cycle, and the limits and ranges are not available. For size by property, you can set the minimum and maximum point size of the plot symbols. You can also set the color and point size of the plot symbols if you do not want to map them to properties. To do this, choose Fixed for Color and Size, and set the color or the point size.

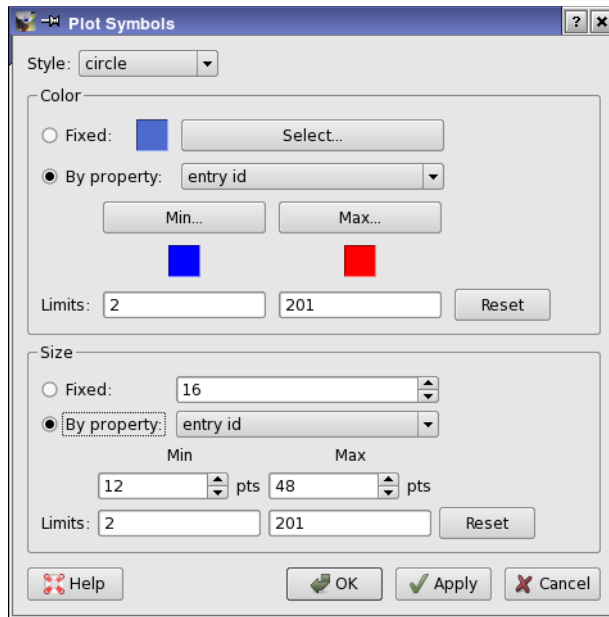


Figure 2.24. The Plot Symbols dialog box.

To configure the axes, choose Settings → Axes. The configuration options are described in [Section 2.7.4 on page 52](#).

Scatter plots can be used to interact with the view. If you select points in the scatter plot, the corresponding rows in the view are selected. Likewise, if you select rows in the view, the corresponding points in the scatter plot are selected. You can select points by clicking or dragging. Shift-clicking or shift-dragging adds to the selection without removing any points; control-clicking or control-dragging changes (inverts) the selection of points without affecting the selection of other points.

Once you have selected some points, you can show only those points in the chart by clicking Selected. If you then change the selection, click Selected again to show only the new selection. To redisplay all points, click All.

You can also label the selected points, by right-clicking and choosing Show Labels from the shortcut menu. The labels include the structure and title, and the names and values of the properties that are represented by that point (x , y , color, size). The contents change when any of the properties plotted changes. The labels can be dragged to any point in the plot. They are connected with the plot point by a line. To hide individual labels, click the X in the label. To hide all the labels, choose Hide Labels from the shortcut menu.

2.7.2 Histograms

Histograms show the frequency of occurrence of values of a single variable, which you can choose from the Property option menu. The variable can be numeric, and the bars defined in terms of ranges of the variable; or it can be categorical, and the bars record the frequency of occurrence of the values of the variable. String variables and partitions are always treated as categorical, and real variables as numeric. Integer variables can be treated as either, by selecting Property is numeric or Property is categorical. If they are treated as categorical, the bars represent the frequency of occurrence of each integer value.

The bars can be shown as frequency counts or as percentages. To set up the histogram bars, choose Settings → Bars. To configure the vertical axis, choose Settings → Y Axis. Many of the settings are common to other plots, and are described in [Section 2.7.4 on page 52](#).

Clicking on a bar selects the table rows that have the property values covered by the bar. The bar is highlighted with a pattern. The proportion of the bar that is patterned represents the fraction of rows in that bar that are selected in the spreadsheet. You can add a bar to the selection with shift-click, and you can change the selection of a bar with control-click. To clear the selection, click in an empty part of the chart.

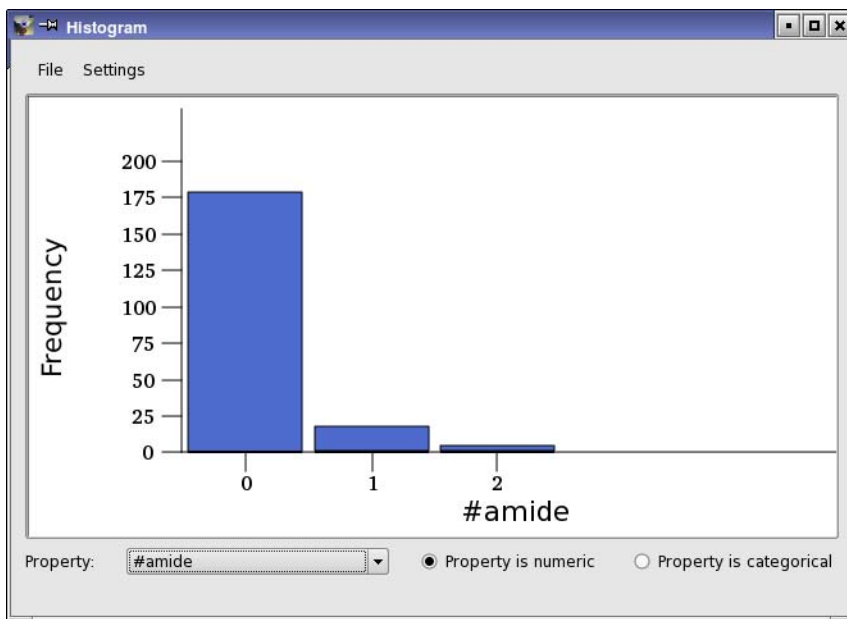


Figure 2.25. The Histogram panel.

2.7.3 Pie Charts

Pie charts also show the values of a single variable. You can choose the property that is represented from the Property menu. The variable can be numeric, and the sectors defined in terms of ranges of the variable; or it can be categorical, and the sectors record the frequency of occurrence of the values of the variable. String variables and partitions are always treated as categorical, and real variables as numeric. Integer variables can be treated as either, by selecting Property is numeric or Property is categorical. If they are treated as categorical, the sectors represent the frequency of occurrence of each integer value.

Clicking on a sector selects the table rows that have the property values covered by the sector. The sector is highlighted with a pattern. The proportion of the sector that is patterned represents the fraction of rows in that sector that are selected in the spreadsheet. You can add a sector to the selection with shift-click, and you can change the selection of a sector with control-click. To clear the selection, click in an empty part of the chart.

To set the ranges represented by the chart, the legend, and other options, choose Settings → Sectors. The settings are described in Section 2.7.4 on page 52.

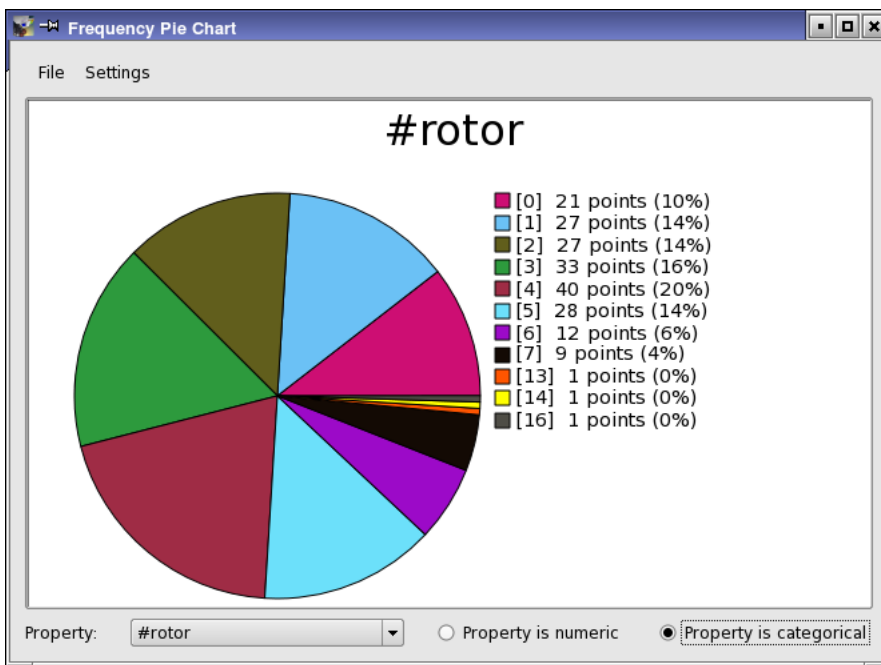


Figure 2.26. The Frequency Pie Chart panel.

2.7.4 Chart Settings

The common settings for each chart are made in the **General** dialog box, which you open by choosing **Settings** → **General** in the chart panel. You can set the title of the chart and change its font size and color, and set the background and foreground colors for the chart.

For scatter plots and histograms, you can change various attributes of the axes. To make these settings, choose **Settings** → **Axes** in a scatter plot panel or **Settings** → **Y Axis** in a histogram panel. Most of these settings are common to the two chart types:

- Minimum and maximum values
- Axis line color and weight
- Tick mark spacing, start value, color and weight
- Tick mark label spacing, start value, font size, color, and number of decimals
- Axis title, font size, and color.

For scatter plots, you can decide where the *y* axis line crosses the *x* axis, and display or hide grid lines. For histograms, you can choose whether to display raw counts or percentages.

For histograms, the horizontal (*x*) axis can be configured by choosing **Settings** → **Bars**. The dialog box depends on whether the property is numerical or categorical. You can set the following attributes for numeric properties:

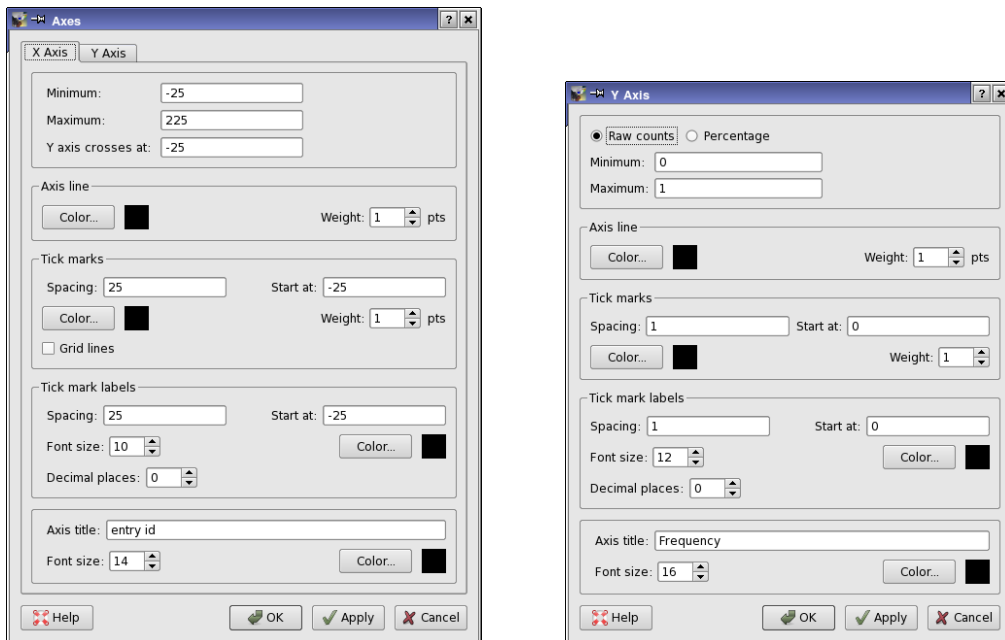


Figure 2.27. Axis settings dialog boxes.

- Axis line color and weight
- Tick mark color, weight, and centering
- Tick mark label font size, color and number of decimal places
- Axis title, font size, and color
- Bin width, lower and upper limits of the data
- Bar width relative to the data range, bar color

For categorical properties, you can set the following attributes:

- Axis line color and weight
- Tick mark color and weight
- Tick mark label font size and color
- Axis title, font size, and color
- Bar width and color

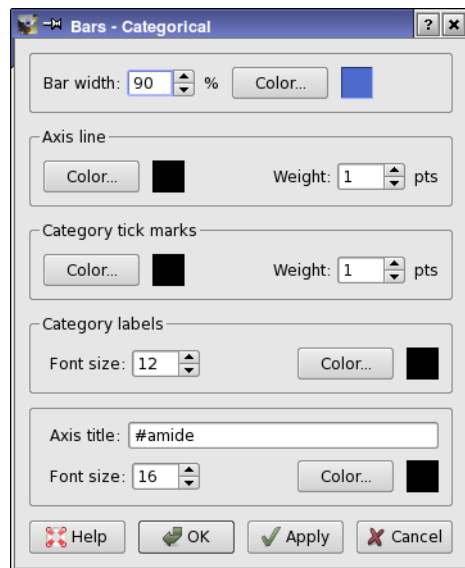
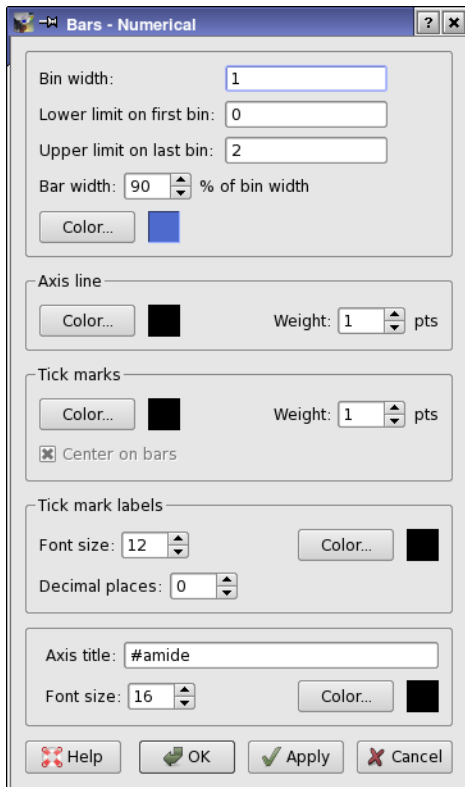


Figure 2.28. The Bars dialog box for numerical and categorical properties.

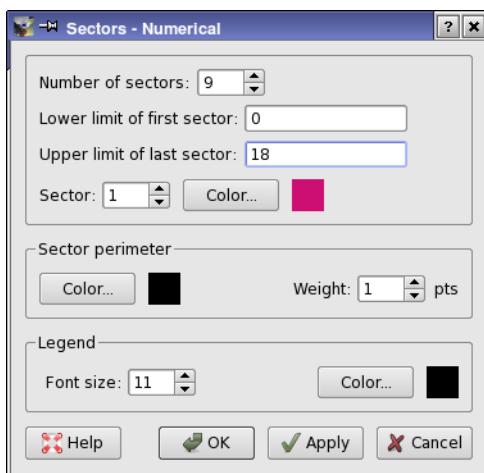


Figure 2.29. The Sectors dialog box.

For pie charts, you can make settings for the sectors and the legend, by choosing Settings → Sectors, which cover the following attributes:

- Number of sectors and lower and upper limits for the sectors (numeric only)
- Color of each sector and color and weight of the perimeter
- Legend font size and color.

2.7.5 Saving Images of Charts

You can save an image of any chart by choosing File → Save Image in the chart panel. A file selector opens, in which you can navigate to a location, name the image, and choose the image format. The available formats are PNG, TIFF, and JPEG.

2.8 Calculating Statistics

You can compute various univariate and bivariate statistics for properties in the Statistics dialog box, which you open by choosing Data → Statistics.

The statistics are calculated for the row and column selection of the current view if you select Selected Data, or for all rows and columns if you select All Data. When the statistics are calculated, you can select Ignore properties with missing values if you want to exclude such properties from the calculation (and the results). If you do not exclude them, a warning is posted and the results table contains question marks (?) for this property.

Univariate and bivariate statistics are computed and displayed separately. You can compute multiple univariate statistics by selecting them from the list in the Univariate tab, and clicking Compute. The results are displayed in a table. If you make a new selection, the statistics that were not previously computed are computed and added to the table.

In addition to the statistics, a Count column is included in the table that lists the number of values that are set for the property. This enables you to identify properties that do not have values for all rows in the spreadsheet.

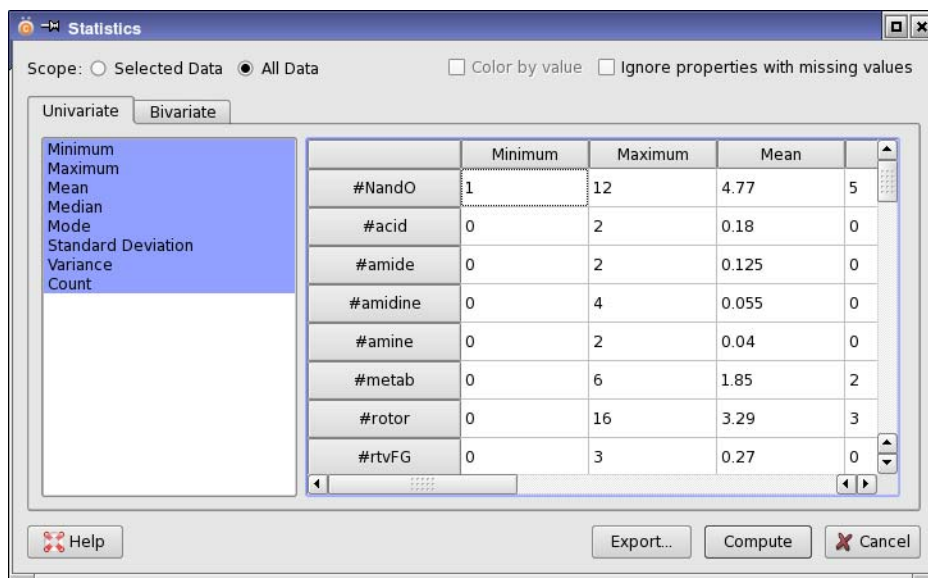


Figure 2.30. The Univariate tab of the Statistics dialog box.

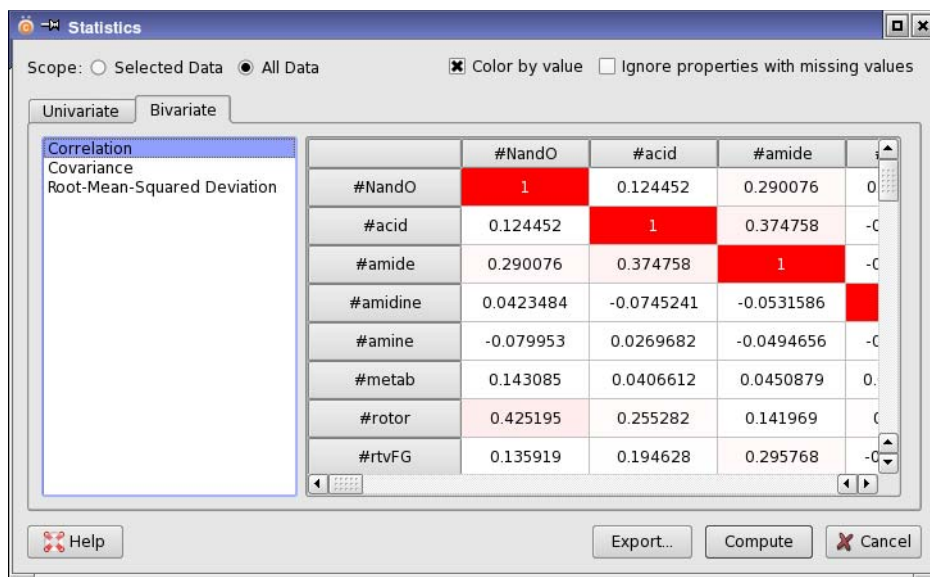


Figure 2.31. The Bivariate tab of the Statistics dialog box.

You can compute one bivariate statistic at a time, by selecting it in the list in the Bivariate tab and clicking Compute. The results are displayed in a table. If you select a statistic that has been previously computed in this session, the results are displayed without having to click Compute. For bivariate statistics, you can select Color by value to color the cells by value, using a blue-white-red color map. You can also display a scatter plot of the two variables (properties) used to compute the bivariate statistic by clicking the table cell.

To export the statistics in the current tab to a CSV file, click Export. A file selector opens, in which you can navigate to the location and name the file.

2.9 Calculating New Properties from the Data

Canvas provides a calculator that can be used to combine properties by constructing a formula. The calculator provides mathematical functions as well as basic arithmetical operations. It is designed to generate entire new properties from existing properties by operating on property columns, rather than to operate on individual cells, as in a spreadsheet program like Microsoft Excel or Open Office Calc. To open the calculator, choose Data → Calculator.

The calculation can be performed on all rows, or on the selected rows. If you want to use the selected rows, you must select them before opening the dialog box. The result is stored in a new property, for which you must provide a name in the Create column text box.

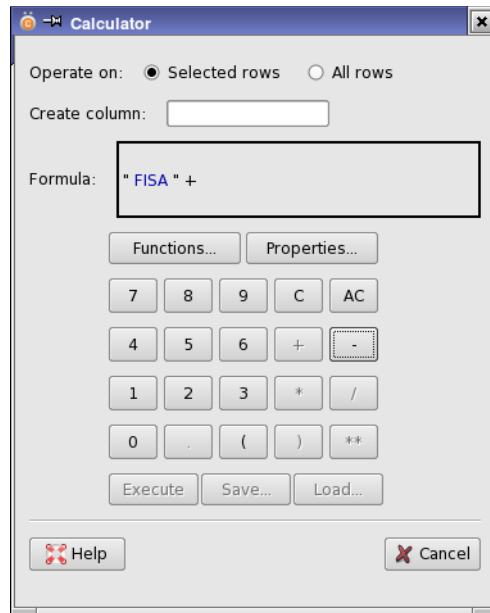


Figure 2.32. The Calculator dialog box.

To construct a formula, use the keys on the keypad, which are set out like a traditional calculator, and the Functions and Properties buttons. The Functions button opens a dialog box, in which you can choose a mathematical function. The Properties button opens a dialog box, in which you can choose an existing property from the spreadsheet.

When you have finished constructing the formula, click **Execute** to perform the calculation and add the new property to the spreadsheet. If you chose to operate on the selected rows, values for the unselected rows are not defined.

You can save a formula for later use by clicking **Save**. A dialog box opens, in which you can name the formula. To use a saved formula, click **Load**, and select the formula in the Formulas dialog box. This dialog box also allows you to delete saved formulas.

2.10 Viewing Structures in PyMOL

The Canvas spreadsheet only displays structures in 2D. To view them in 3D, you would normally have to export them and then import them into a 3D viewer (like Maestro). Canvas provides an interface to PyMOL that allows you to select structures in the Canvas spreadsheet and display them directly in PyMOL. To do so, the location of the PyMOL executable must be included in your PATH environment variable. For information on setting this environment variable on Windows, see [Appendix A](#) of the *Installation Guide*.

When you choose Structure → Show Structures in PyMOL, the PyMOL interface is opened, and check boxes are added in the lower left corner of the structure cells in the spreadsheet. Selecting a check box sends the structure to PyMOL for display. You can select multiple check boxes with the usual shift-click and control-click actions, and the structures are displayed in PyMOL.

Custom views have an independent set of checkboxes that can be selected and deselected. When you create a new view, the view inherits the structure selection for PyMOL from the parent view. When you filter a view (master or custom), any structures that are filtered out are not considered to be selected for viewing.

The structures that you see in PyMOL are a combination of the structures that are marked for display in all open views. That is, the displayed structures correspond to rows that are visible in any of the open views and are checked for viewing in PyMOL. So if a particular structure is selected and visible in two open views, deselecting it in one view does not remove it from PyMOL, because it is still selected and visible in the other view.

To finish viewing structures in PyMOL, choose Structure → Close PyMOL.

If you want PyMOL to open automatically when you open a project, you can set a preference in the Preferences dialog box, which you open with File → Preferences. Click General in the list on the left, then select Start PyMOL.

2.11 Setting Preferences

Some customization of the interface features and other settings can be made in the Preferences dialog box, which you open with File → Preferences. The panel has a list of preference categories on the left. When you select an item from this list, the preferences in that category are displayed on the right. If you want to restore the default preferences, click Reset. The preferences in each category are described below.

General preferences

This category has two settings.

- **Start PyMOL**—Start a PyMOL session when a project is opened, so that you can view structures in 3D. If this option is not selected, you can start PyMOL at any time from the Structure menu.
- **Canvas max memory**—Set the maximum memory that can be used by Canvas, in bytes. You can use commas in the value, e.g. 1,000,000. You can use suffixes for kilobytes, megabytes, and gigabytes, in any form, case-insensitive, e.g. 2 Mb, 2 Mbytes, 2 megabytes all mean 2097152 bytes. This limit is used by any applications that are started from Canvas. Sets the `SCHRODINGER_CANVAS_MAX_MEM` environment variable for the Canvas session.
- **Incorporate jobs automatically**—When a job finishes, incorporate the results immediately into the Canvas project. The project is unavailable while the incorporation takes place. This preference sets the default incorporation option in the Start Application dialog box.
- **Number of decimal digits to display**—Set the number of digits past the decimal point to display in the spreadsheet. This value applies to all properties.

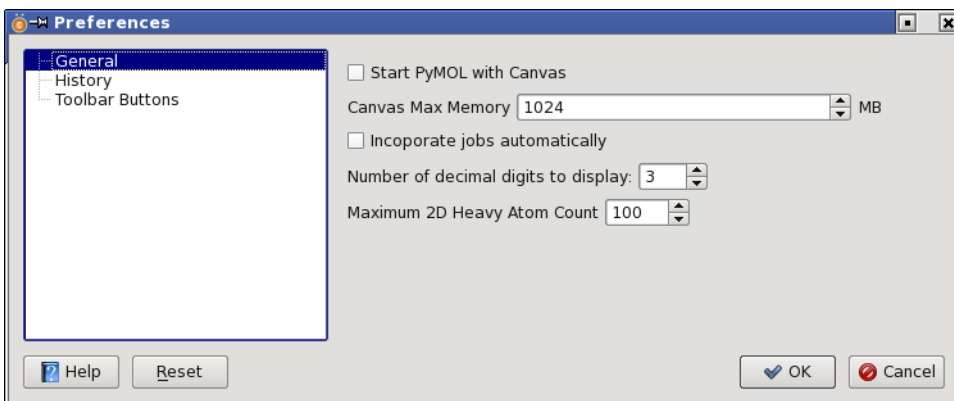


Figure 2.33. The Preferences dialog box.

- Maximum 2D heavy atom count—Specify the maximum allowed number of heavy atoms in any structure when importing and when drawing structures. On import, only the first 200 structures are checked for the heavy atom limit, and a dialog box is displayed if a structure exceeds the limit. Thereafter, all structures are imported. Any structure that is imported and exceeds the limit is drawn with a projection of the 3D structure onto the xy plane, rather than as a 2D structure.

History preferences

This category allows you to set preferences for the number of items stored in various history lists.

- Recent project list—Specify the number of projects to show on the Recent Projects submenu of the File menu. The default is 10. Click Clear to clear the project history list.
- Substructure query list—Specify the number of SMARTS queries or query files to show in the Substructure Query dialog box, in the lists that are displayed for the Query File text box or the SMARTS Query text box.

Toolbar buttons preferences

With these preferences, you can check or clear the boxes for the toolbar buttons that you want to have on the File toolbar, or click Select All to check all the boxes. The buttons are added or removed when you click OK in the dialog box. The buttons are described in [Section 2.2.2 on page 7](#).

Running Applications from Canvas

The Canvas interface provides access to a range of applications for molecular descriptors, fingerprints, similarity and clustering, calculating statistics, and building regression and other predictive models based on properties of the structures. The applications can only be run from the master view.

3.1 General Features

Nearly all of the applications run on a data set that includes both structures and properties. The *Data set* option menu is used to choose the data set. You can choose the selected rows or the visible rows in the master view, which makes all properties in the spreadsheet available for selection. You can choose to use a saved view and select the view from an option menu; the rows or columns or both can be restricted to those of the saved view by choosing from the *Domain* option menu, while the other rows or columns are taken from the visible part of the master view. If you are examining the results of a job or cloning a job, you can choose the original data set from that job to start a new job.

Some applications require the selection of multiple properties. These panels contain property selection tools, which consist of two lists, *Available properties*, and *Selected properties*, a search box to restrict one or other list to properties that contain the text entered in the box, and buttons for transferring the properties that are selected in one list or all properties to the other list. The *Selected properties* list may have a different name in some panels.

For applications that are run under Job Control, when you click the action button (usually *Compute*) in the application panel, the *Start Application* dialog box opens. This dialog box allows you to choose the host on which to run the application, and the number of CPUs, for applications that can be distributed across CPUs. It also allows you to choose whether to incorporate the results immediately into the project or to wait until you explicitly incorporate the results. The default incorporation behavior can be set in the *Preferences* panel (*File* → *Preferences*). For information on setting up the host list and access to hosts, see [Chapter 7](#) of the *Installation Guide*. The job is run under the Job Control facility, which is described in detail in the *Job Control Guide*. If you want to bypass Job Control, you can do so from the command line—see [Chapter 5](#).

The default amount of memory used by the job is about 500 MB. You can change this value in the *Preferences* dialog box (*File* → *Preferences*), in the *General* category.

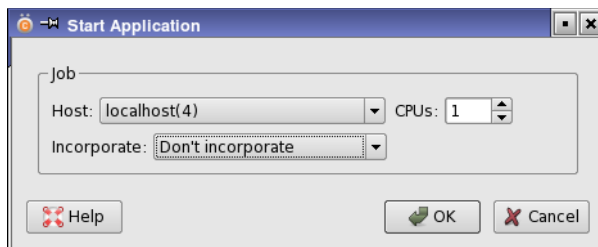


Figure 3.1. The Start Application dialog box.

When you run an application, messages about its progress are displayed in the Messages View panel, and a record is added under the application in the Project View panel. The record is color-coded by job status. When a job finishes, the application node is expanded. The results are not automatically incorporated into the spreadsheet unless you set the preference to do so or selected the option to do so in the Start Application dialog box. To add them, display the jobs for the application in the Project View panel, right-click on the record and choose Incorporate. If you want to be notified when a job finishes, select Notify when jobs complete in the Project View panel.

To run another job with the same options or using the job's options as a basis, right-click on the job in the Project View panel and choose Clone Job. The Clone Job dialog box opens, and offers an option to use the original data set. When you click OK, the application dialog box opens with the options set as they were for the job you cloned. You can then modify the options and start the job.

Some applications have panels for viewing and using the results of the job. For these applications, right-click on the job in the Project View panel and choose View.

3.2 Molecular Properties

Molecular properties are used as the basis of modeling applications such as QSAR or QSPR. Canvas provides a means of calculating a wide range of molecular properties and descriptors based on the 2D structure. 3D properties can be generated by other applications and imported into Canvas. With a large number of properties it is often necessary to reduce the number of properties to a relatively independent set. This section describes applications for molecular property generation and molecular property selection.

3.2.1 Calculating Molecular Properties

Canvas can calculate a range of molecular properties from the 2D structures that can be used as descriptors. Some of these are simple counts of features, others are numerical properties evaluated from a set of rules. The calculation of these properties can be done in the Molecular Properties dialog box, which you open by choosing Applications → Molecular Properties.

You can choose whether to calculate the properties for the selected rows or the visible rows, a saved view, or the data set from a previous run, and choose a name to use to save the results in the project.

There are four classes of properties, Physicochemical Descriptors, Topological Descriptors, LigFilter Descriptors, and QikProp Descriptors, which you can select by choosing the class from the tree on the left. The property selection tools are displayed on the right. For each class of properties, you can select a host and a number of CPUs for the job. The jobs for each class run independently and are incorporated independently. The job status is listed in the Project View under the class name.

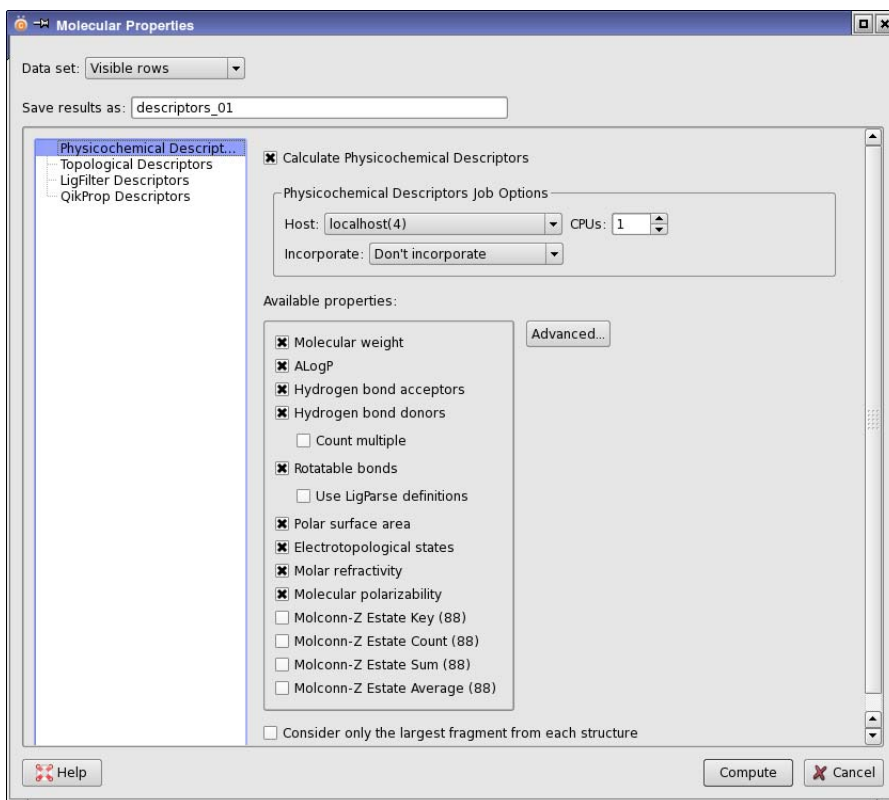


Figure 3.2. The Molecular Properties dialog box.

The physicochemical descriptors that can be calculated are:

Molecular weight	Molar refractivity [10]
AlogP [9]	Molecular polarizability [11]
Hydrogen bond acceptors	Molconn-Z Estate Key [12]
Hydrogen bond donors	Molconn-Z Estate Count [12]
Rotatable bonds	Molconn-Z Estate Sum [12]
Polar surface area [13]	Molconn-Z Estate Average [12]
Electrotopological states (Estates) [8]	

References to methods are given in square brackets. Each of these corresponds to an option in the Available properties section. By default all but the Molconn-Z Estate properties are selected. The latter correspond to 88 properties for each of the options.

For the hydrogen bond acceptor and donor counts and the rotatable bond counts, you can specify alternative definitions. The definitions used by LigParse (in LigPrep) for rotatable bonds are available as an option, Use LigParse definitions. The default (Canvas) definitions differ from the LigParse definitions by excluding amide bonds and bonds that connect to a terminal heavy atom (one that otherwise only has H atoms attached, such as methyl, amino, or hydroxyl). You can also count groups that can act as donors of more than one H-bond as multiple donors, e.g. an NH₂ group would be counted as two donors. For all three counts, you can specify custom definitions in the Molecular Properties - Advanced Options dialog box, which you open by clicking Advanced. To use a custom definition, select the appropriate option and click its Browse button to load the definition file. The file must contain a list of SMARTS patterns for the feature. More information on creating custom definitions is given in Section 5.5.3 on page 156.

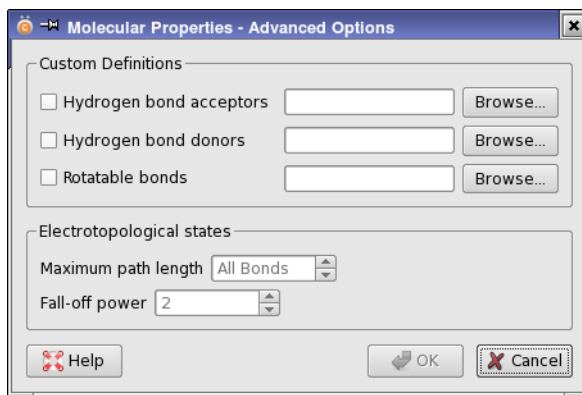


Figure 3.3. The Molecular Properties - Advanced Options dialog box.

When computing the sum and average for the Estates, you can set the maximum path length and the fall-off power in the Molecular Properties - Advanced Options dialog box.

If any of the structures contains more than one molecule (fragment), you can limit the property calculation to the largest fragment by selecting Consider only the largest fragment from each structure.

The topological descriptor list is long, so a property selection tool is provided in which you can search for properties in the list, and then transfer them from the Available properties list to the Properties to calculate list.

The LigFilter descriptor list is also long, but it is broken into sections with a tree tool. For more information on the LigFilter descriptors, see [Section 2.4](#) of the *General Utilities* manual.

The QikProp descriptors are described in [Chapter 1](#) of the *QikProp User Manual*. By default, “fast mode” is used, in which the dipole moment, ionization potential and electron affinity are not calculated, as they require a PM3 semiempirical quantum mechanical calculation. To include these properties, select Compute PM3 properties.

3.2.2 Selecting Representative Properties

Often, the set of molecular properties or descriptors that you obtain from various sources can have a high degree of correlation. For model-building, it is useful to remove the linear dependence in the properties, to speed up the calculation and to reduce numerical error. One way of doing this is to select a smaller set of properties that adequately represents the large set, in some measurable way.

Canvas provides the ability to select a representative set of properties. The set is chosen by clustering the properties with a hierarchical clustering method, based on the correlation matrix. The distance between two properties is calculated as one minus the absolute value of the correlation matrix element. A default number of clusters is chosen using Kelley criteria [4], and the property that is nearest the centroid of each cluster is chosen as the representative property. The job to perform the clustering is set up in the Feature Selection dialog box, which you open from the Applications menu.

To set up the job, choose a structure source from the Data set option menu, enter a name for the results, and choose the property set that you want to select from. You should check the list for properties that are not relevant, such as an entry ID or possibly results from other jobs. The list of properties that you choose can be reduced by removing those that have mostly identical values. To do this, you can select Exclude properties whose values are identical in more than N % of molecules, or Exclude properties with Shannon entropy lower than X . Click Compute to run the job.

When the job finishes, you can view the results in the Feature Selection Viewer panel, which you open by selecting View from the shortcut menu for the feature selection job, under Feature Selection in the Project View panel. This panel lists the original set of properties and the reduced (representative) set of properties, and provides tools to change the selection of representative properties.

To determine how well the selected properties (reduced list) represents the original properties, multiple linear regression is performed for each property to calculate an R^2 value for the fit of the property by the reduced list. First, a principal components analysis is done on the reduced list and principal components are selected to account for 95% of the variance. The X variables in the PLS model are taken from the projection of the principal components onto the property data from the reduced list, and the Y variable is the original property. These R^2 values are reported in the second column of the Original properties table.

If when setting up the job, you chose to remove properties for which a specified percentage of the values were identical, these properties are dimmed in the Original properties list, and the cells in the R^2 column are colored yellow and contain text describing the elimination criterion rather than a numeric value. Likewise, cells for properties that did not meet the Shannon entropy cutoff are colored magenta, with descriptive text. If you did not choose to remove such properties, the cells in the R^2 column in which all values of the property are identical are colored red and contain the text “Zero variance”.

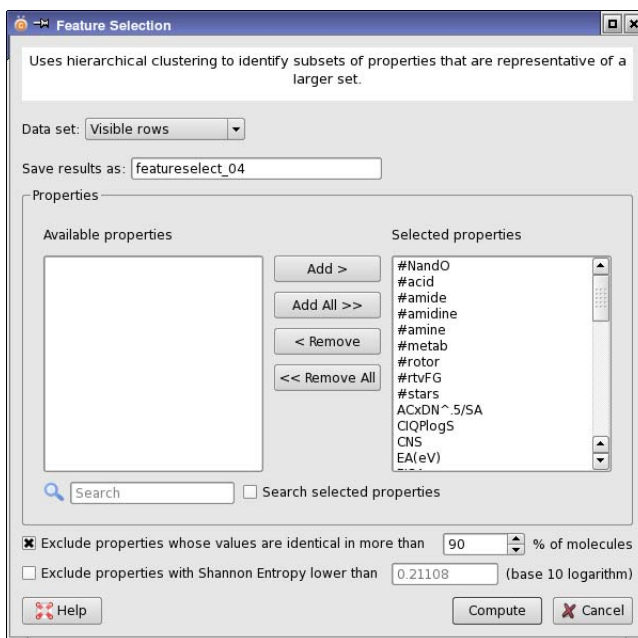


Figure 3.4. The Feature Selection dialog box.

Properties on the original list that are also on the reduced list are dimmed, and have a line of - characters instead of a numeric value in the R² column. Properties that are manually transferred between lists (as described below) are colored cyan.

You can change the number of properties in the reduced list in several ways:

- Select Number of clusters and change the number of clusters, then click Update Properties to update the lists and the regression results.
- Select Maximum correlation and change the value of the maximum allowed correlation between properties on the reduced list, then click Update Properties.
- Select a property in either of the lists and click Add or Remove, then click Update R².

You might want to change the property manually if you have a preference for a particular property. The property chosen from each cluster is the one that is closest to the centroid. You might want to replace it with another property from the same cluster. To do this, it is useful to sort the list of original properties by the clustering order (right-click and choose Sort by Cluster), so that properties in the same cluster are adjacent. It might also be useful to view the correlation matrix (click Correlation Matrix) for the original properties, ordered by cluster, so that you can identify the properties that are in the same cluster as the one you want to replace. With cluster ordering, the heat map of the correlation matrix should have red diagonal blocks that mark the clusters.

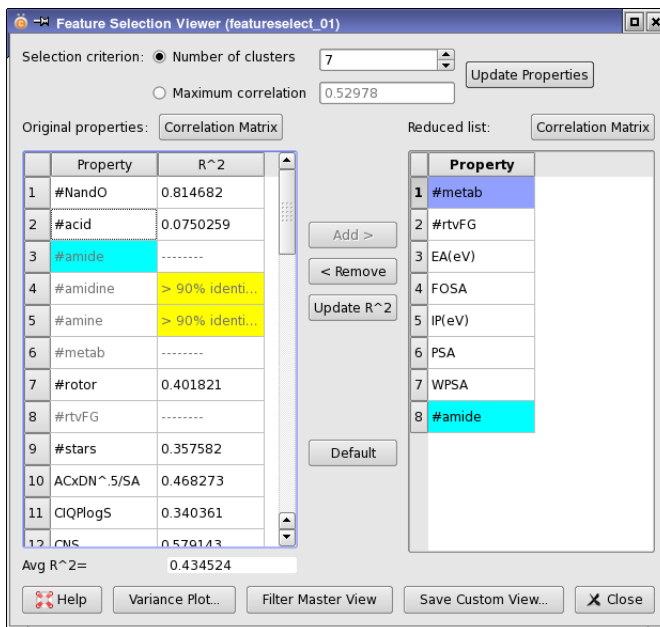


Figure 3.5. The Feature Selection Viewer panel.

3.3 Fingerprints

Canvas can calculate a variety of fingerprints from the 2D structure, both hashed fingerprints and structural keys, which are used by other Canvas applications. Canvas can also calculate fingerprints from 3D all-atom structures, based on 3-point or 4-point pharmacophore models. Once you have a fingerprint, you can create modal fingerprints [6] (averages over several structures) using the shortcut menu for the fingerprint property—see [page 21](#) or [page 117](#) for details.

3.3.1 2D Fingerprints

Binary fingerprints from the 2D structure are calculated using the Binary Fingerprints from Structures dialog box, which you open by choosing Applications → Binary Fingerprints from Structures.

You can choose whether to calculate the properties for the selected rows or the visible rows, a saved view, or the data set from a previous run, and choose a name for the new fingerprint column. If any of the structures contains more than one molecule (fragment), you can limit the property calculation to the largest fragment by selecting Consider only the largest fragment from each structure.

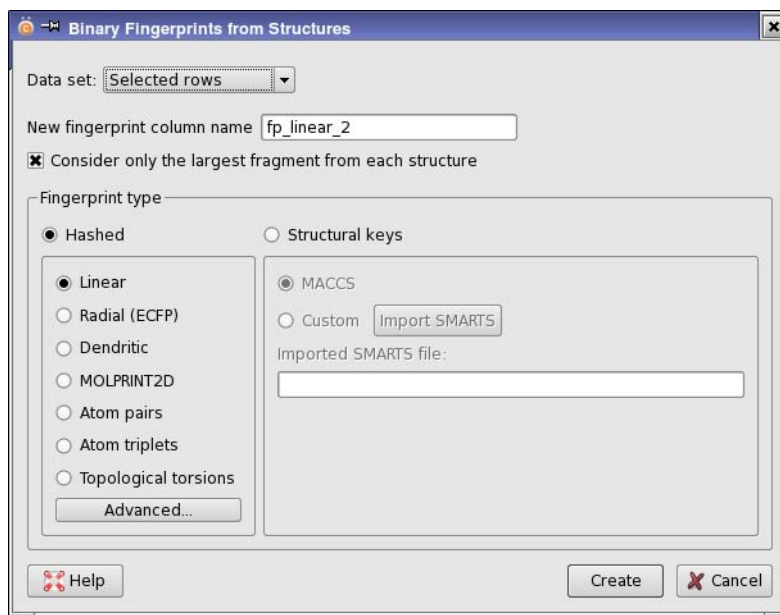


Figure 3.6. The Binary Fingerprints from Structures dialog box

To calculate a hashed fingerprint, select Hashed, then choose the particular fingerprint type, from Linear, Radial, Dendritic, MOLPRINT2D, Atom pairs, Atom triplets, or Topological torsions.

If you want to change any of the parameters used in the fingerprint calculations, click Advanced, and make your selections in the Hashed Fingerprints - Advanced Options dialog box. This dialog box has four main sections.

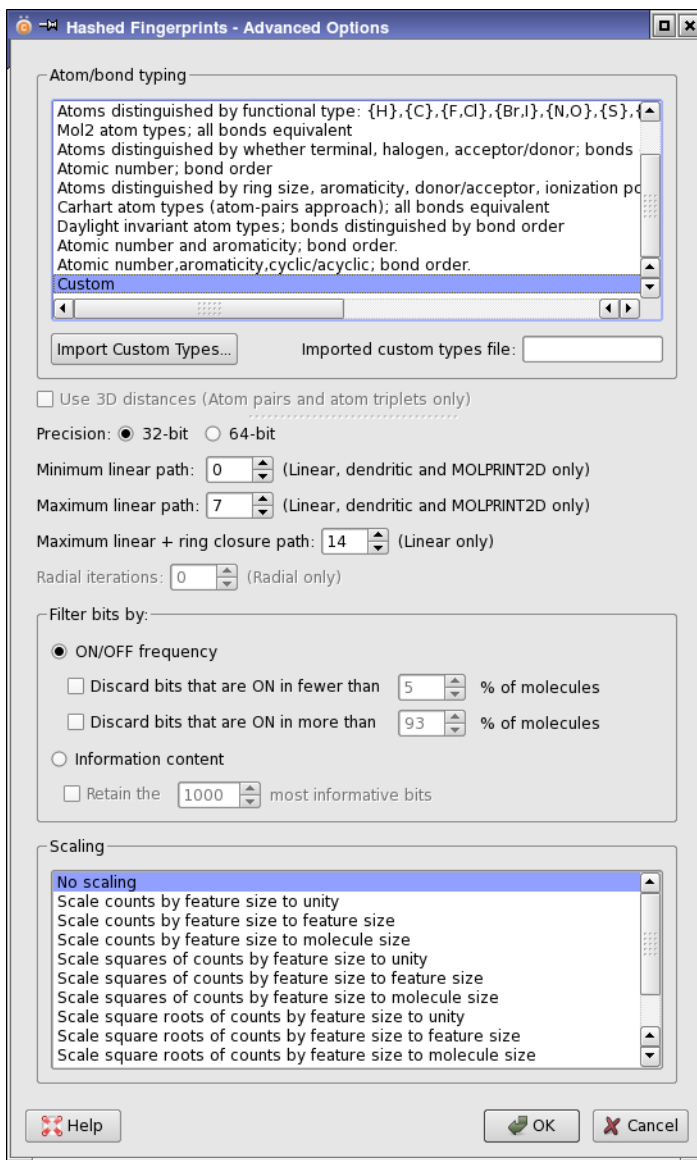


Figure 3.7. The Hashed Fingerprints - Advanced Options dialog box

- Atom/bond typing—Choose the atom and bond typing scheme to be used to classify atoms and bonds in the structures. If you choose Custom, click Import Custom Types to locate the file that contains the custom atom and bond types. The file must contain SMARTS patterns for the fingerprints, one per line.
- Various options—These include options for individual fingerprint types and precision options:
 - You can set the precision to 32 or 64 bits. Using 64 bits reduces collisions of ON bits, but doubles the space required to store each key.
 - For atom pairs and atom triplets, you can choose to use 3D distances if they are available.
 - For linear, dendritic, and MOLPRINT2D fingerprints, you can set the minimum and maximum linear path, and for linear fingerprints, you can also set the maximum linear+ring closure path.
 - For radial fingerprints, you can set the number of radial iterations.
- Filter bits by—Select an option for filtering out bits. If you select ON/OFF frequency, you can then choose to discard bits that are ON in fewer than or more than a specified percentage of molecules, and set the percentages. If you select Information content, you can choose to retain the N most informative bits, and set the value of N .
- Scaling—Select a method from this list to scale the binary fingerprint data to real numbers. The default is not to scale the fingerprint data.

As well as hashed fingerprints, you can calculate MACCS or custom structural keys. If you choose Custom, you must import a SMARTS definition file for the structural keys, by clicking Import SMARTS. Each line in the file must contain one SMARTS pattern, then the name of the key, separated by one or more spaces.

3.3.2 3D Fingerprints from Pharmacophores

Fingerprints based on pharmacophore models of the 3D structure are calculated using the 3D Pharmacophore Fingerprints dialog box, which you open by choosing Applications → 3D Pharmacophore Fingerprints. The fingerprint bits are assigned based on the feature type and the distances between the features.

You can choose whether to calculate the properties for the selected rows or the visible rows, a saved view, or the data set from a previous run, and choose a name for the new fingerprint column. A default name is supplied. If any of the structures contains more than one molecule (fragment), the largest fragment is used.

The fingerprints can be generated from 3-point pharmacophores (the default), 4-point pharmacophores, or both. 4-point pharmacophores usually set an order of magnitude more bits than 3-point pharmacophores, and do not provide much additional benefit.

Pharmacophores can be generated for a single conformer of each structure, or for a set of conformers that are generated temporarily for this purpose. To use the existing structure and not generate conformers, select **Use existing 3D structures**. To generate conformers, select **Generate conformers in memory**. The conformer set does not explicitly include the original structure, but is likely to contain a similar structure. The conformers are discarded after use.

When you start this job, you can choose to distribute it over multiple processors by setting the CPUs value in the **Start Application** dialog box.

If you want more control over the process, click **Advanced**, and make settings in the **3D Pharmacophore Fingerprints - Advanced Options** dialog box.

- **Custom feature definitions**—Select **Use custom feature definitions** to supply your own definitions of the pharmacophore features in terms of SMARTS patterns. The definitions are in the format of a Phase feature definition file—see [Appendix B.3](#) of the *Phase User Manual* for details. Click **Import** to import the definition file.
- **Precision**—Select **32-bit** or **64-bit** for the fingerprint precision. Using 32 bits is usually adequate and results in few collisions. Using 64 bits reduces collisions, but requires more space to store the fingerprints.
- **Inter-feature distance bins**—Set the smallest distance in the **Lower limit** text box, and the largest distance in the **Upper limit** text box. These distances set the range of the bins: distances smaller than the lower limit are placed in the first bin, and distances larger than the upper limit are placed in the last bin. Set the width of the bin in the **Bin width** text box, and

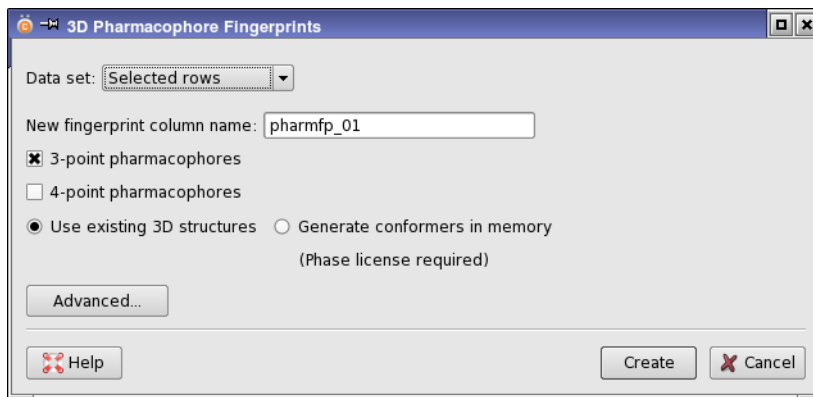


Figure 3.8. The **3D Pharmacophore Fingerprints** dialog box

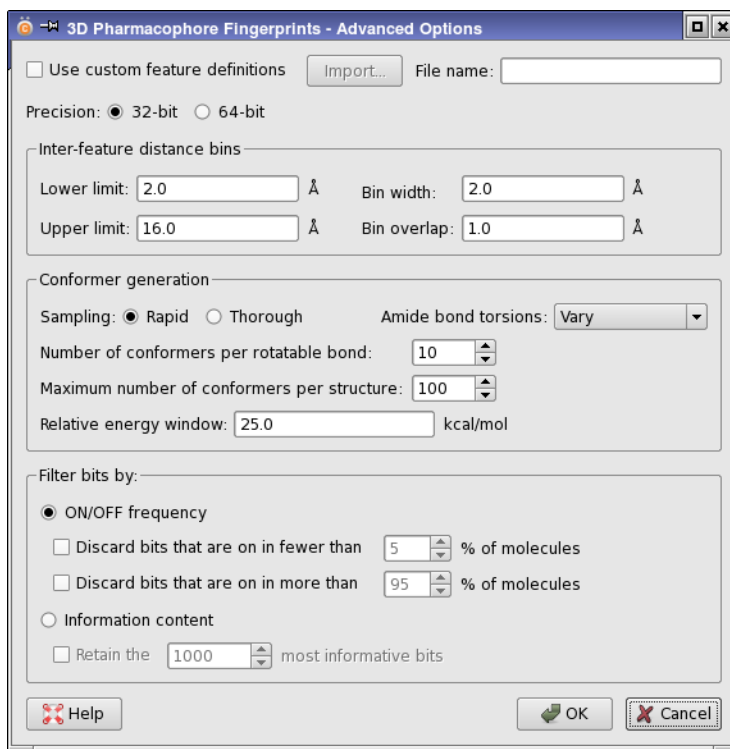


Figure 3.9. The 3D Pharmacophore Fingerprints - Advanced Options dialog box

the overlap in the Bin overlap text box. Distances are assigned to a neighboring bin as well if the distance is within the overlap distance of the bin boundary. This means that the effective bin width is the specified width plus twice the overlap width.

- Conformer generation—Set options and thresholds for generating conformers. This section is only available if you selected Generate conformers in memory in the main panel.
 - Select Rapid or Thorough for the sampling of conformational space. Rapid sampling varies rotatable groups independently; thorough sampling varies them together.
 - Choose how to treat amide bonds from the Amide bond torsions option menu. Vary allows full rotational freedom; the other two options fix the value at the input value or at the trans conformation.
 - Limit the number of conformers generated by setting values in the Maximum number of conformers per structure and Number of conformers per rotatable bond text boxes. Both limits are applied. You can also limit the number of conformers by filtering out conformers whose energy is higher than that of the lowest conformer by the amount specified in the Relative energy window text box.

- Filter bits by—Choose a filtering method to reduce the number of bits. Select ON/OFF frequency to discard bits that are on in most molecules or only a few molecules, and select the relevant options and the percentage of molecules. Select Information content and set the number of informative bits to keep to filter by the amount of information provided by the bits.

3.4 Similarity, Dissimilarity, and Clustering

Molecular similarity or dissimilarity is measured by a distance between molecules in the space of a property set. The similarity or distance is defined by a metric. Similarity or distance can be used for clustering and screening of structures, or selection of diverse structures.

3.4.1 Computing a Similarity or Distance Matrix

Similarity and distance matrices can be computed in Canvas based on fingerprints or on molecular properties, and over a range of metrics. The computed matrices can be examined in table form, and they can be exported to a file, where they can be used by other applications. The structures for any table cell are shown in a tool tip if you choose Structure tooltips.

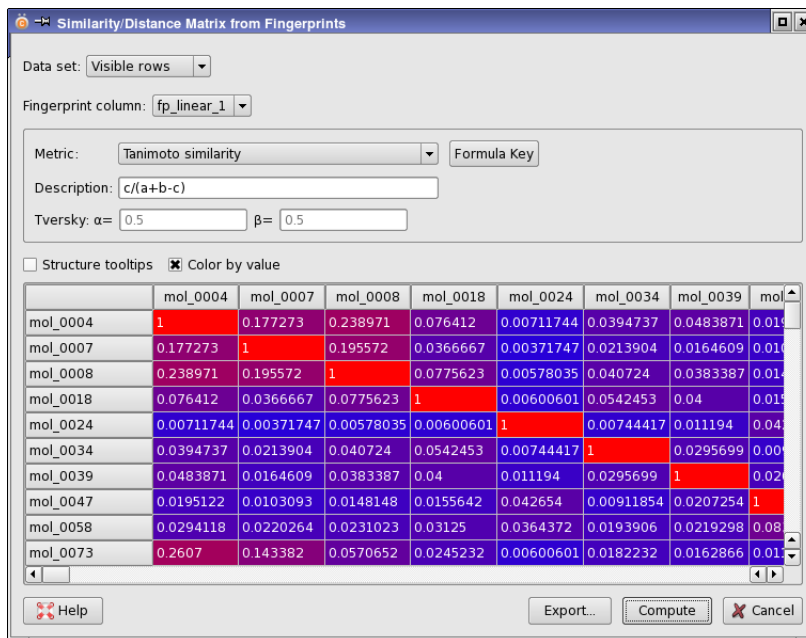


Figure 3.10. The Similarity/Distance Matrix from Fingerprints dialog box

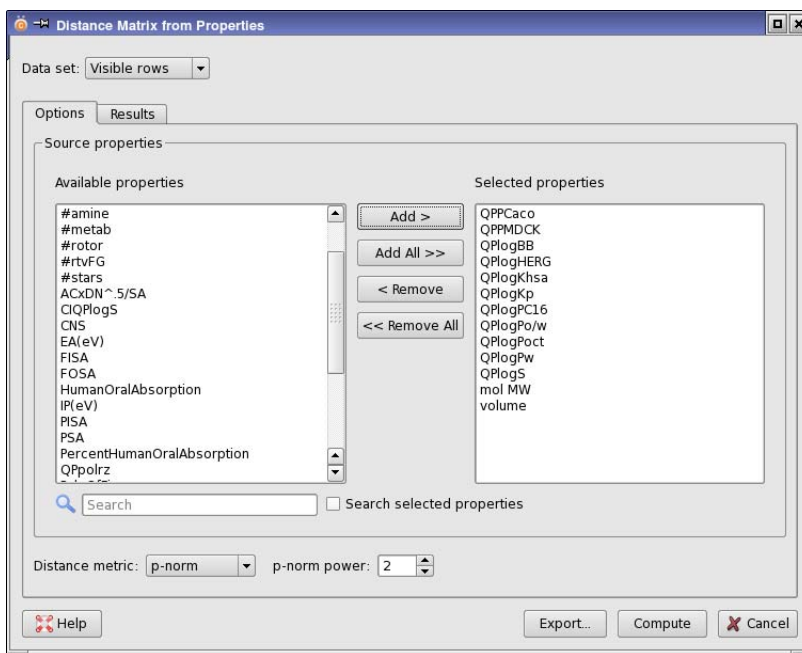


Figure 3.11. The *Distance Matrix from Properties* dialog box, *Options* tab.

To calculate similarity or distance matrices based on fingerprints, choose Applications → Similarity/Distance Matrix → Create from Fingerprints. In the Similarity/Distance Matrix from Fingerprints panel, you can choose whether to calculate the matrix for the selected rows or the visible rows, a saved view, or the data set from a previous run, choose the fingerprint column and the metric, and choose whether to color the results table cells by value. Twenty-four metrics are available—see Table 5.22 on page 142 for a list. When you click Compute, you are prompted to either to display the results in a table or export them to a file.

To calculate distance matrices based on properties, choose Applications → Similarity/Distance Matrix → Create from Properties. In the Distance Matrix from Properties panel, you can choose whether to calculate the matrix for the selected rows or the visible rows, a saved view, or the data set from a previous run, select the properties and the metric, and choose whether to color the results table cells by value. Two metric types are available: p-norm and Mahalanobis. The p-norm is defined by $(\sum x^p)^{1/p}$. The Mahalanobis distance matrix is calculated using the covariance matrix to correct non-orthogonality of the data. When you click Compute, you are prompted to choose whether to display the results in a table or export them to a file.

	mol...04	mol...07	mol...08	mol...18	mol...24	mol...34	mol...39	mol...47	mol...58	mol...
mol_0004	0	547.817	1450.59	2489.19	1040.09	1712.05	156.397	348.706	1423.4	82...
mol_0007	547.817	0	910.829	1964.58	639.999	1192.24	615.493	610.213	1001.97	479...
mol_0008	1450.59	910.829	0	1064.57	817.232	333.245	1525.56	1491.65	785.311	137...
mol_0018	2489.19	1964.58	1064.57	0	1767.97	777.65	2577.26	2553.76	1511.63	241...
mol_0024	1040.09	639.999	817.232	1767.97	0	1088.25	1062.9	977.359	431.964	989...
mol_0034	1712.05	1192.24	333.245	777.65	1088.25	0	1802.76	1790.29	951.478	163...
mol_0039	156.397	615.493	1525.56	2577.26	1062.9	1802.76	0	220.137	1465.09	225...
mol_0047	348.706	610.213	1491.65	2553.76	977.359	1790.29	220.137	0	1399.49	389...
mol_0058	1423.4	1001.97	785.311	1511.63	431.964	951.478	1465.09	1399.49	0	136...
mol_0073	82.7971	479.007	1375.3	2410.23	989.116	1632.87	225.387	389.097	1365.13	0
mol_0075	170.823	379.709	1287.44	2333.15	897.769	1557.02	251.345	345.418	1285.06	118...
mol_0100	368.788	673.652	1465.41	2441.31	1096.55	1675.16	503.439	687.406	1405.39	341...
mol_0104	214.749	755.97	1662.04	2702.12	1209.87	1925.27	182.796	381.707	1598.42	295...
mol_0106	287.011	801.947	1675.77	2687.5	1270.8	1912.59	359.092	573.155	1628.35	331...
mol_0115	261.412	514.651	1412.6	2473.3	952.222	1704.59	174.486	125.961	1368.42	287...

Figure 3.12. The Distance Matrix from Properties dialog box, Results tab.

3.4.2 Screening Structures by Similarity or Distance

Similarity or distance matrices computed from fingerprints can be used to screen structures by similarity to one or more reference structures. The results are returned as properties that contain the similarity to the reference structure.

To set up a screen, choose Applications → Similarity/Distance Screen. The screen can be run on the selected rows or the visible rows, a saved view, or the data set from a previous run, a saved view, or the data set from a previous run, and you can choose the fingerprint from those that are already available in the project.

The names of the similarity properties are generated from a base name, which you can enter, with a suffix added that distinguishes the reference structures. The suffix can be taken from either a generated sequence of names, or the reference structure name. A default is supplied.

The reference structures can be taken from the selected rows, a saved view, or they can be selected by hand. If you choose Select by hand, the structures in the screen are shown in the reference structures table, and you can select the desired table rows. The column name is displayed for the selected structures when you select them.

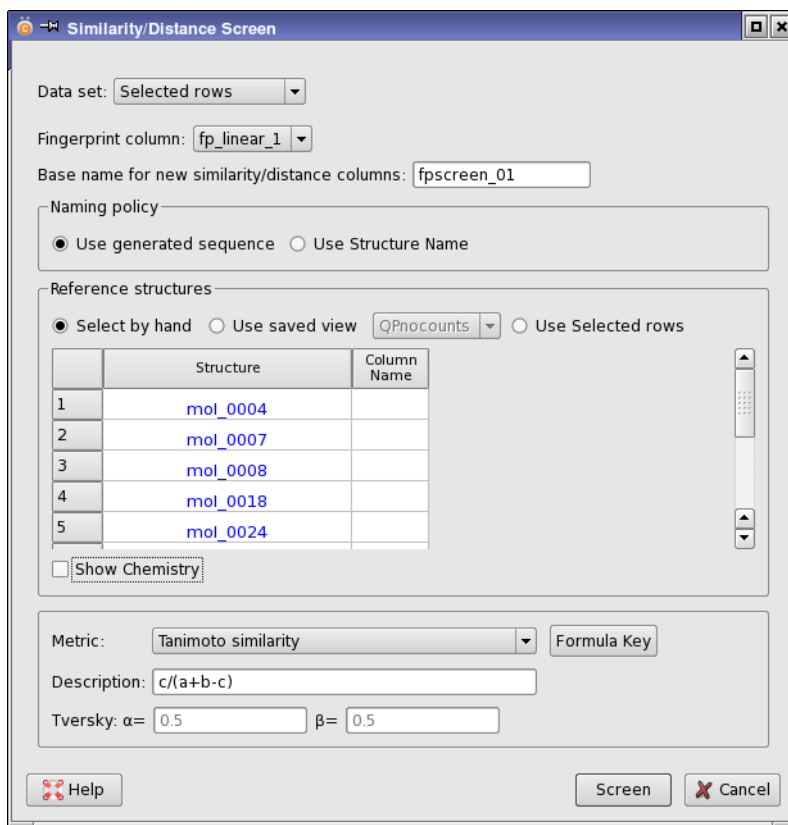


Figure 3.13. The Similarity/Distance Screen dialog box

The metric used to calculate similarity or distance can be chosen from the 24 available metrics. The description of the metric is given in the Description text box, and the definition of the symbols can be viewed by clicking Formula Key.

When you have selected all the options, click Screen to generate the results. After the results are incorporated into the spreadsheet, you can run a property query to select structures based on similarity to one or more of the references.

3.4.3 Screening Structures by Shape

If you have 3D structures, you can select structures that have a similar shape to a query molecule in the Shape Screen dialog box, which you open from the Applications menu. The screening is done with the phase_shape program, and requires a Phase license. For details on shape screening, see Chapter 14 of the *Phase User Manual*.

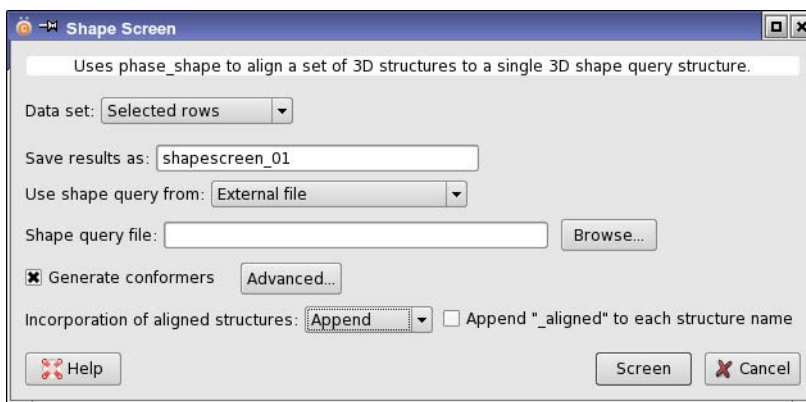


Figure 3.14. The Shape Screen dialog box

The screen can be run on the selected rows or the visible rows, a saved view, or the data set from a previous run, and the results are saved with the name you provide. The shape query must be a single molecule and can come from an external file (Maestro or SD), or the structure that is displayed in the PyMOL viewer.

If you want more control over the screening process, you can set options in the Shape Screen - Advanced Options dialog box, which you open by clicking Advanced.

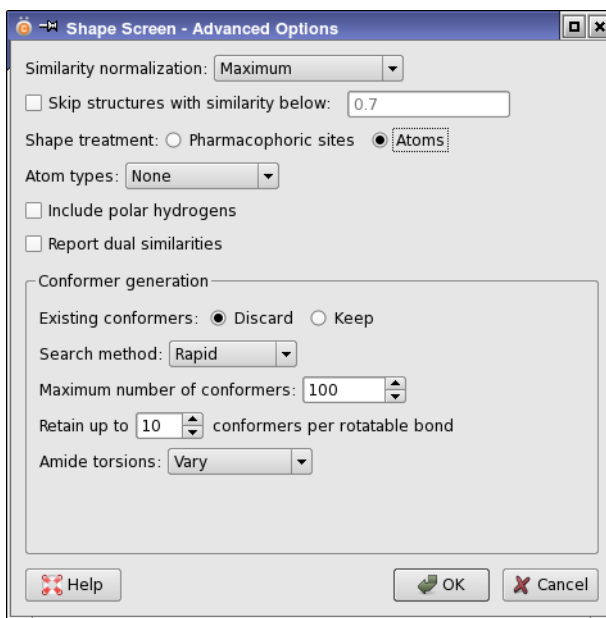


Figure 3.15. The Shape Screen - Advanced Options dialog box

- You can generate conformers for the molecules that are screened. The conformational search varies only the torsional angles of rotatable groups in the molecules.
- You can screen molecules based on pharmacophore types (the default) or atom types, and select different atom type schemes for determining which atoms are treated as similar. For a simple shape screen, choose Atoms with Atom type: None. When using atom types, you can choose whether to include polar hydrogens in the determination of similarity, and report shape similarities with and without atom typing.
- You can select a similarity normalization (default is Maximum), and skip structures with low similarity. The similarity to the query is returned as a property.

3.4.4 Selecting Diverse Structures

It is often useful to be able to select structures that are as different as possible from each other. You can do this in Canvas in the Diversity-Based Selection dialog box, which you open from the Applications menu.

The analysis can be run on the selected rows of the spreadsheet or on all rows, and the results are saved with a name which you can specify. The similarity or distance matrix is computed from fingerprints, which you can select from those already available in the project. You can also choose the metric used for the similarity or distance from the 24 available metrics.

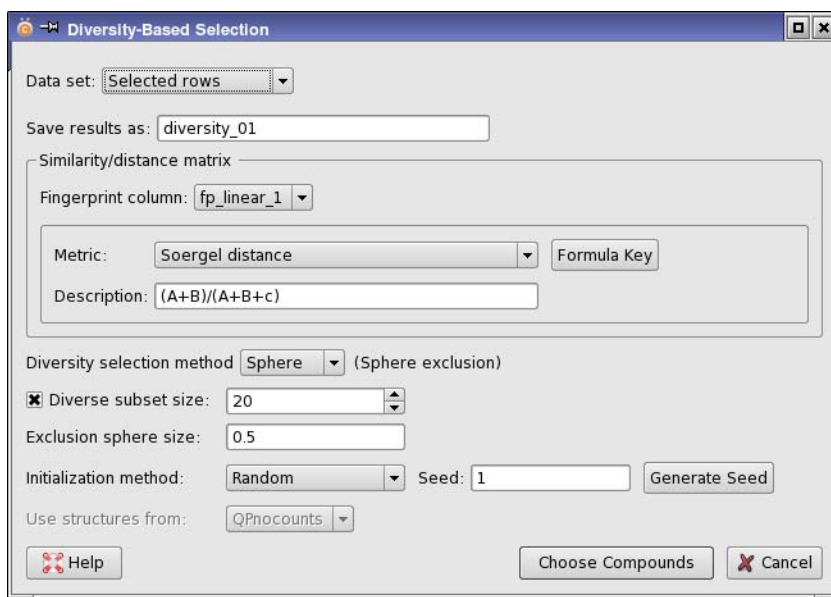


Figure 3.16. The Diversity-Based Selection dialog box

Four methods are available from the Diversity selection method option menu: Sphere (sphere exclusion), DISE (directed sphere exclusion [2]), MAXSUM (maximum sum of pairwise distances), and MAXMIN (Maximum nearest-neighbor distance). The size of the subset can be specified for all four methods, but is only applied for the Sphere and DISE methods if the exclusion sphere size produces a larger diverse subset.

The selection must be initialized with one or more structures. The initialization method depends on the selection method. The DISE method has its own initialization, so no method is available. For the other three methods, you can choose random initialization and specify the or generate the seed. For the Sphere method, you can choose existing structures, and specify the structures as all structures that are in the view of your choice. These structures serve as an excluded set: the diverse set does not include any structure within the exclusion sphere of the input structures. For the MAXSUM and MAXMIN initialization method, you can choose Representative (structure with the greatest similarity to the other structures) or Dissimilar (structure with the least similarity to the other structures).

When you have made your choice, click Choose Compounds. After the job has finished, you can open a restricted custom view of the selected structures, by right-clicking on the record in the Project View panel under Diversity and choosing View. You can save this view as a normal custom view, or apply it to the master view. When the job is incorporated, a new property is added to the spreadsheet, with 1 for the structures that were selected and 0 for those that were not selected. The property is named *jobname* set, with underscores replaced by blanks.

3.4.5 Hole-Filling and Library Optimization

If you want to choose a set of compounds based on their dissimilarity and optionally by optimizing the values of a set of properties, you can do so with in the Hole-Filling and Library Optimization dialog box, which you open from the Applications menu. With this application you can choose a set of dissimilar structures from a pool, and you can use them to fill holes in an existing library. Like Diversity-Based Selection, this application uses nearest-neighbor distances to choose the most diverse structures, but with a different approach. Here, the desired number of structures is selected at random, and then the similarity of the structures to each other and to the library is minimized. In addition, the process can be guided by selection of structures whose properties lie in specified ranges. This is done by minimizing the fraction of properties for each selected structure whose values lie outside the ranges.

The structure pool can be chosen from selected or visible rows in the master view, a saved view, or the data set from a previous run, and the results are saved with a name that you can specify. If you want to use the selected rows, you must select them before opening the dialog box.

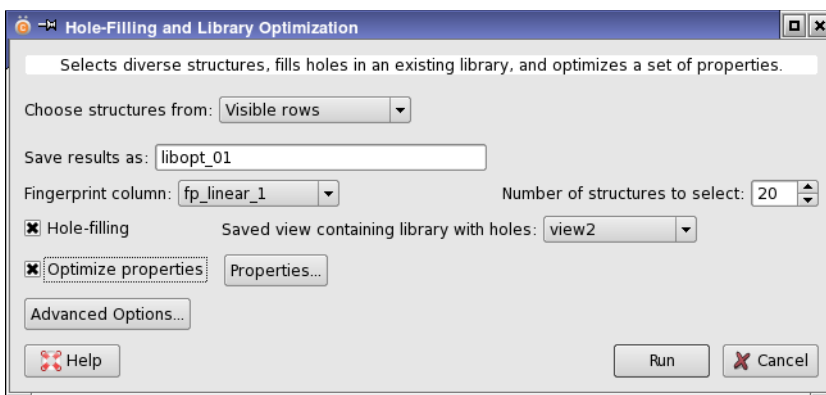


Figure 3.17. The Hole-Filling and Library Optimization dialog box

If you select Hole-filling, the library whose holes are to be filled is taken from a saved view, which you can choose. It is best to ensure that the structure pool and the library do not have structures in common. The similarity matrix is computed from fingerprints, which you can select from those already available in the project.

If you want to optimize the values of properties, select Optimize properties, and click Properties to choose the properties and set the desired ranges. To save the property set and their ranges as a property filter, click Export Filter. Filters can be imported (Import) or read from another job (Use filter from).

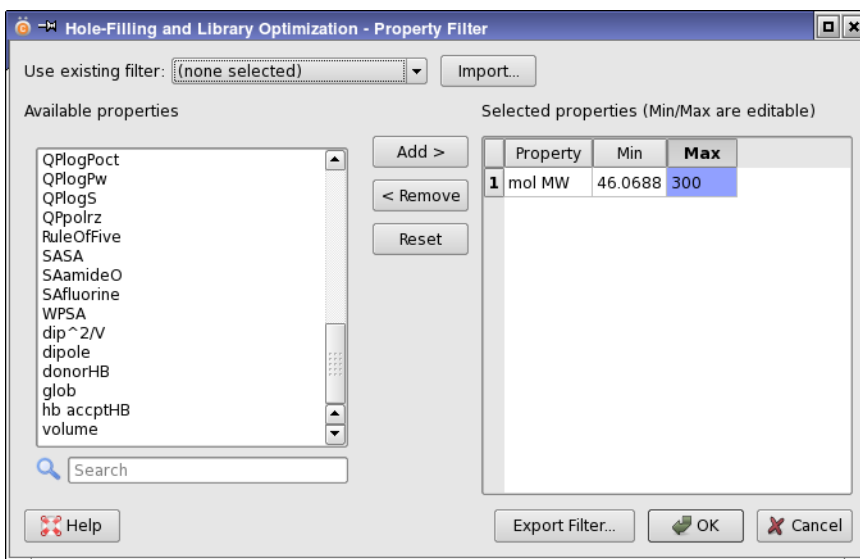


Figure 3.18. The Hole-Filling and Library Optimization - Property Filter dialog box

You can optimize properties without filling holes in a library, but the structures are still selected by minimizing their similarity to one another.

When you have specified the structures and set options, including the number of structures to select, click Run. The time taken scales as the square of the number of structures selected, so selecting a large number of structures could take some time. If you want to set parameters to control the optimization process, click Advanced Options. In the dialog box that opens, you can set the maximum number of optimization cycles, thresholds for stopping the optimization, a limit on the run time, the random seed for selecting structures and other aspects of the optimization, and a threshold for a simulated annealing process used when the score used in the optimization increases. See [Section 5.4.4 on page 136](#) for more information on these parameters and the algorithm used for the optimization.

3.4.6 Comparing Libraries

A histogram of nearest neighbor similarities is a convenient way to assess the diversity of a single collection of compounds, or the similarity between two collections of compounds. You can create the data for this histogram in the Library Comparison dialog box, which you open from the Applications menu.

The library comparison works by finding the nearest neighbor in the reference library for each compound in the query library, using fingerprint similarity based on a chosen metric. The nearest-neighbor similarities are returned in a column named MaxSim. The results of the library comparison can be used to create a histogram from the nearest neighbor similarities.

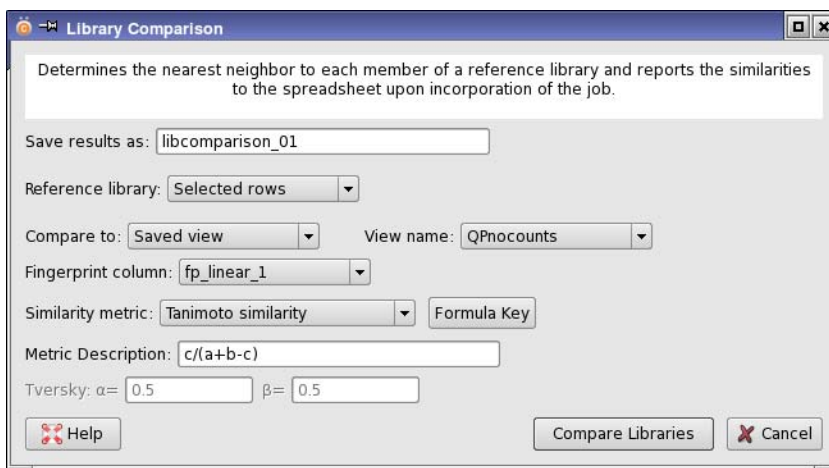


Figure 3.19. The Library Comparison dialog box

The reference library can be chosen from the selected rows or the visible rows of the spreadsheet, a saved view, or the data set from a previous run, and the results are saved with a name which you can specify. If you want to use the selected rows, you must select them before opening the dialog box. The library for comparison can be taken from a saved view, which you can choose, or it can be the reference library itself. The similarity matrix is computed from fingerprints, which you can select from those already available in the project. You can also choose the metric used for the similarity from the available similarity metrics.

When you have made all the desired settings, click **Compare Libraries**. After the job finishes, you can create a histogram based on the **MaxSim** property, by clicking the **Histogram** button.



The histogram for a library with a high degree of similarity to the reference library will have most of the library at the greatest similarity, whereas a library with a low similarity to the reference library will have a large part of the library at smaller similarities.

3.4.7 Clustering Structures by Similarity or Distance

Canvas provides three methods for clustering structures: hierarchical clustering, leader-follower clustering, and k-means clustering. Leader-follower and k-means clustering can be used on much larger data sets than hierarchical clustering. Clustering by these methods can be done in the Hierarchical Clustering dialog box, the Leader-Follower Clustering dialog box, and the K-Means Clustering dialog box, which you open from the Applications menu.

For all three methods, you can choose the dataset for the clustering as described in [Section 3.1 on page 61](#), and the results are saved with a name that you can specify.

All three methods can use fingerprints for clustering, which you can select from those already available in the project. For hierarchical and leader-follower clustering, you can choose the metric used for the similarity or distance from the 24 available metrics. For leader-follower and k-means clustering, you can use properties for clustering, which you select from the properties available in the project. The properties can be automatically scaled to a common range, by selecting **Autoscale**.

The remaining options in the three dialog boxes are unique to the method. When you have finished making your selections, click **Create Clusters** to run the clustering job. A job record is added in the Project View panel under the clustering method, and the progress of the job is shown in the Messages View panel.

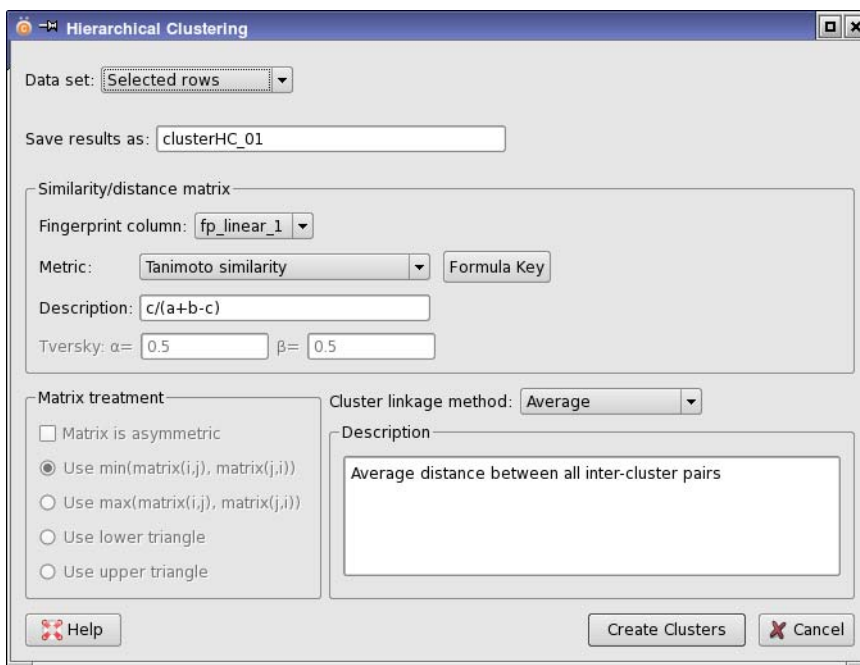


Figure 3.20. The Hierarchical Clustering dialog box

3.4.7.1 Hierarchical Clustering

The specific options for hierarchical clustering include the following:

- If the matrix is asymmetric, you can select **Matrix is asymmetric** and then choose how to treat the off-diagonal elements.
- You must choose a clustering linkage method, from the nine available methods. The linkage method describes how to calculate the distance between clusters. When you choose a linkage method, the description is displayed in the **Description** text area. The linkage methods are listed in [Table 3.1](#).

The results of hierarchical clustering are not incorporated directly into the project. Instead, you can view a dendrogram of the clusters, and choose a set of clusters to export either to a file or to the spreadsheet.

To view the dendrogram, right-click on the job record in the **Project View** panel and choose **View**. In the **Hierarchical Clustering Dendrogram** panel, you can choose the cluster set by the number of clusters, by the Kelley criterion [4], or by the merging distance. The merging distance is displayed as a tool tip when you pause the pointer over a split point. For the last of

Table 3.1. Description of linkage methods for hierarchical clustering.

Method	Description
Single	Closest inter-cluster pair (nearest neighbor)
Complete	Farthest inter-cluster pair (farthest neighbor)
Average	Average distance between all inter-cluster pairs
Centroid	Euclidean distance between cluster centroids
McQuitty	Average distance to the two clusters merged in forming a given cluster
Ward	Sum of squared distances to merged cluster centroid (minimum variance)
Weighted centroid	Weighted center of mass distance, also known as median
Flexible beta	Weighted average intra-cluster and inter-cluster distances (Lance-Williams) with beta=0.25
Schrödinger	Closest distance between terminal (right-to-left) points in 1D cluster orderings.

these, you must click **Update** after changing the merging distance to view the new set of clusters. To change the scaling of the dendrogram, use the mouse wheel. You can export an image of the dendrogram in PNG, JPEG, or TIFF format by clicking **Save Image**.

By default, the dendrogram is colored by cluster membership. You can color the leaf nodes by the values of a property with a heat map, by selecting **Color by property**, and clicking **Edit Heat Map** to choose the properties and setting the color range for the property (see [Section 2.5.3 on page 40](#)). The properties that you select are added to the option menu in the **Hierarchical Clustering Dendrogram** panel, so you can switch between properties. To switch back to coloring by cluster membership, select **Color by cluster**.

When you have selected a set of clusters by any of the three methods, you can click **Export** to export the clusters or their centroids to a file or to the spreadsheet. The **Export Selection** dialog box opens, in which you can specify the destination and make other related settings. In addition to the structures, three properties are exported: **Cluster**, which is the cluster index, **Centroid**, which is 1 for the centroid of each cluster, and 0 otherwise, and **Total**, which is the number of structures in the cluster. If there is a tie for the centroid, the first structure is chosen.

If you choose **To file**, click **Browse** to navigate to the location for the file, select the file format from **Maestro**, **SD**, or **CSV**, and name the file. You can export the structures in each cluster to a separate file, labeled with the cluster index, by selecting **Place results into separate files**. CSV files include the structure name, a SMILES string for the structure, and the three properties. If you chose to export only the centroids, they are written to a single file.

If you choose **To spreadsheet**, the column names for the properties are derived from the property name by adding a prefix, which you can specify in the **Name prefix** text box. The resulting



Figure 3.21. The Hierarchical Clustering Dendrogram panel

column names are *prefix:Cluster*, *prefix:Centroid*, and *prefix:Total*. The default prefix is *HClust::jobname*. The clusters can be exported to separate views, named *prefix::cluster N*, by selecting Place clusters in separate views. If you chose to export only the centroids, they are written to a single view named *prefix::centroids*.

You can view the structures that belong to a cluster by clicking the “leaf” (the square at the bottom of the dendrogram) for one of the cluster members. A limited custom view is displayed that contains only the structures by default, with the background color of the cluster from the dendrogram. You can add properties to the view by choosing View → Manage Properties and selecting the properties. When you do, these properties are added by default to any cluster view opened by clicking its leaf. To view the structures for more than one cluster, click in the dendrogram panel, then control-click the leaf.

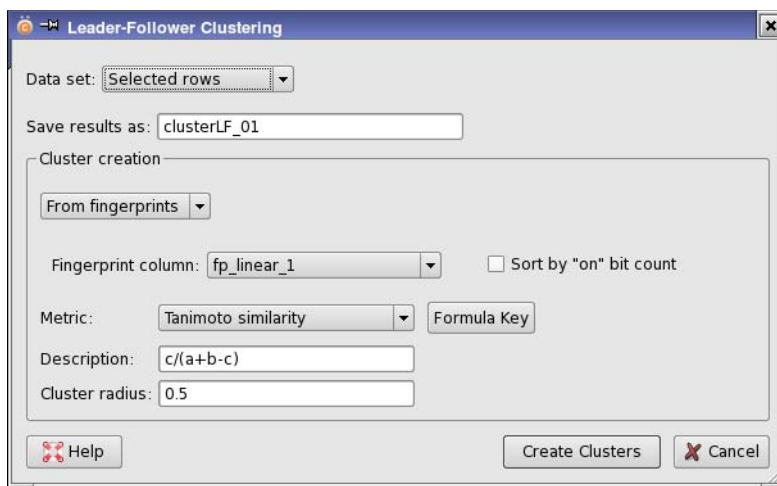


Figure 3.22. The Leader-Follower Clustering dialog box

3.4.7.2 Leader-Follower and K-Means Clustering

For leader-follow clustering, you can sort the fingerprints by the count of “on” bits, in descending order. This option is only available with some of the metrics. Sorting can speed up the clustering calculation considerably.

For k-means clustering, the parameters of the method can be set. You must specify the number of clusters, but the remaining parameters can be left at their defaults.

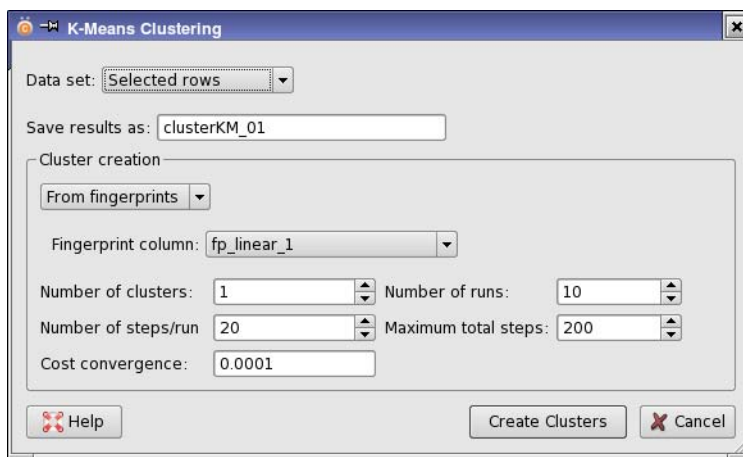


Figure 3.23. The K-Means Clustering dialog box

Information on the algorithms for k-means and leader-follower clustering is given in [Section 5.4.11 on page 149](#) and [Section 5.4.12 on page 150](#).

Incorporating the results of leader-follower and k-means clustering adds a new column to the spreadsheet, named *prefix:jobname:Cluster*, where prefix is LFClust or KMClust. This column contains the index of the cluster to which each structure belongs. For leader-follower clustering, another column is added that marks the leader in each cluster, named LFClust:*jobname:Leader*. The column contains 1 for the leader in each cluster and 0 otherwise. Similarly, for k-means clustering a column is added that marks the centroid of each cluster, named KMClust:*jobname:Centroid*.

You can view the clusters by right-clicking the job record in the Project View panel, and choosing View. A restricted view of the spreadsheet opens, named Leader-Follower Clustering Results, or K-Means Clustering Results. This view has only the structure column and the File, View and Chemistry menus, but you can add properties from the View menu. You can export the clusters, select the cluster to view, save the view as a custom view, or apply the view to the master view. You can open more than one of these panels to display multiple clusters. For leader-follower clustering, the cluster leader is marked by a gray background.

3.5 Building a Predictive Model

Canvas provides a range of methods for building a model that can be used to predict properties. These models include the common regression models—multiple linear regression, partial least-squares regression, and principal components regression—and also Bayes classification, recursive partitioning, and neural network models. The models can be built, applied, exported and plotted from their dialog boxes, which you open from the Applications menu.

The dialog boxes have a common structure, as they are performing similar tasks. Each has an Options tab for setting up the input and options, and a tab for selecting the training set and displaying the results. This tab is labeled Training Set initially, and is relabeled Results after the model is built. The common features are described below. See also the general descriptions in [Section 3.1 on page 61](#).

- The models can be built or applied for the selected rows or the visible rows of the spreadsheet, for a saved view, or for the original data set for a run, by selecting from the Data set option menu.
- The results are saved with the name that you specify in the Save results as text box. This name is also the job name and appears in the Project View panel.
- Existing models can be applied to the data set. Models from within the project can be loaded by selecting Apply saved model and choosing the model from the option menu. External models can be loaded by selecting Apply imported model, clicking Import, and navigating to the file that contains the model. To apply the selected model, click Apply Model. The results are displayed as columns in the structures table, and include the predicted value, and the statistics are displayed in the tabs next to it (Training and Test).
- A new model can be built when you select Build new model. You can select the independent variables with the X variable property selection tools or fingerprint option menu (if available), and the dependent variable (or variables, for neural networks) from the Y variable option menu or selection tools. The values of a Y variable are displayed in the data set structures table when you select it.
- The training set and the test set can be assigned in three ways:
 - Randomly: select Assign randomly, specify a percentage in the text box, specify a seed by entering it the Seed text box or clicking Generate Seed, then click Apply. The training and the test set are not completely random: the random selection process ensures that they have comparable distributions of activities.
 - From a partition: select Assign from partition, choose the partition name from the option menu, select partition values from the Training set value and Test set value option menus, then click Apply. The percentage of the values in the partition is displayed in text boxes to the right of the option menus.

- Manually: set the values in the Model Set column of the structures table. To change the value for a single row, double-click in the table cell and choose a value from the option menu that appears. To change the value for multiple rows, select the rows, right-click on any of the selected cells and choose the value from the Change Selected submenu of the shortcut menu.
- The model is built when you click Build Model. For most of the models, statistics on the training set and the test set are displayed in the tabs to the right of the structures table.
- You can export the model to a file, by clicking Export Model, and navigating to the desired location and naming the file in the file selector that opens.
- You can display a scatter plot of the results by clicking Scatter Plot (if the panel has this button). A dialog box opens, in which you can specify the data set to be plotted, choose colors for the training and test set points, and choose whether to draw the 45-degree line of perfect fit. Clicking Plot then opens a standard scatter plot panel, as described in [Section 2.7.1 on page 48](#).
- You can export the data from selected rows in the structures table, by right-clicking and choosing Export to File or Export to Spreadsheet. For export to a file, data in all columns for the selected rows are written to a CSV file with column headings. Exporting to the spreadsheet adds the predicted values to the spreadsheet.

The features that differ between the dialog boxes are for setting up the parameters of the method, and are described in the following sections.

The final section describes Kohonen self-organizing maps. Although not quite the same as the other predictive models, it can also be built and applied.

3.5.1 Multiple Linear Regression

For multiple linear regression, you can build a model from the selected X variables, or you can choose the best subsets. If you select Choose best subsets, click Options to set parameters for finding the best subsets in the Multiple Linear Regression - Best Subsets Options dialog box.

The best subsets are determined by a simulated annealing process, in which you specify a number of Monte Carlo steps, an initial temperature and a final temperature. The temperatures are expressed in terms of the standard deviation in the dependent variable. The number of X variables in the subsets must be specified. To average the best models, select Use average of and specify the number of models in the box. The average can be weighted by the R^2 value.

You can force the regression line to pass through the origin by selecting Suppress y intercept.

When the model is built, the coefficients for each independent variable and the standard deviation for the variable is given in a table below the statistics tabs.

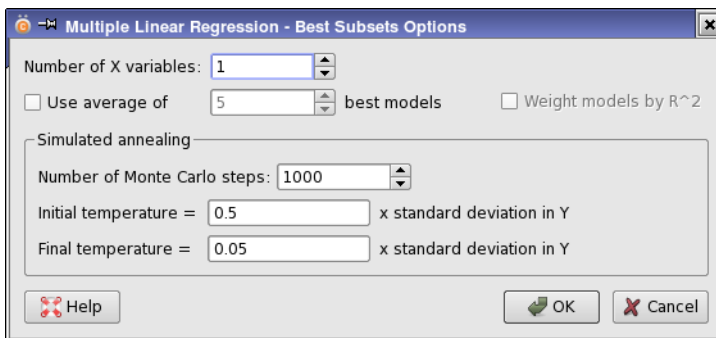


Figure 3.25. The Multiple Linear Regression - Best Subsets Options dialog box

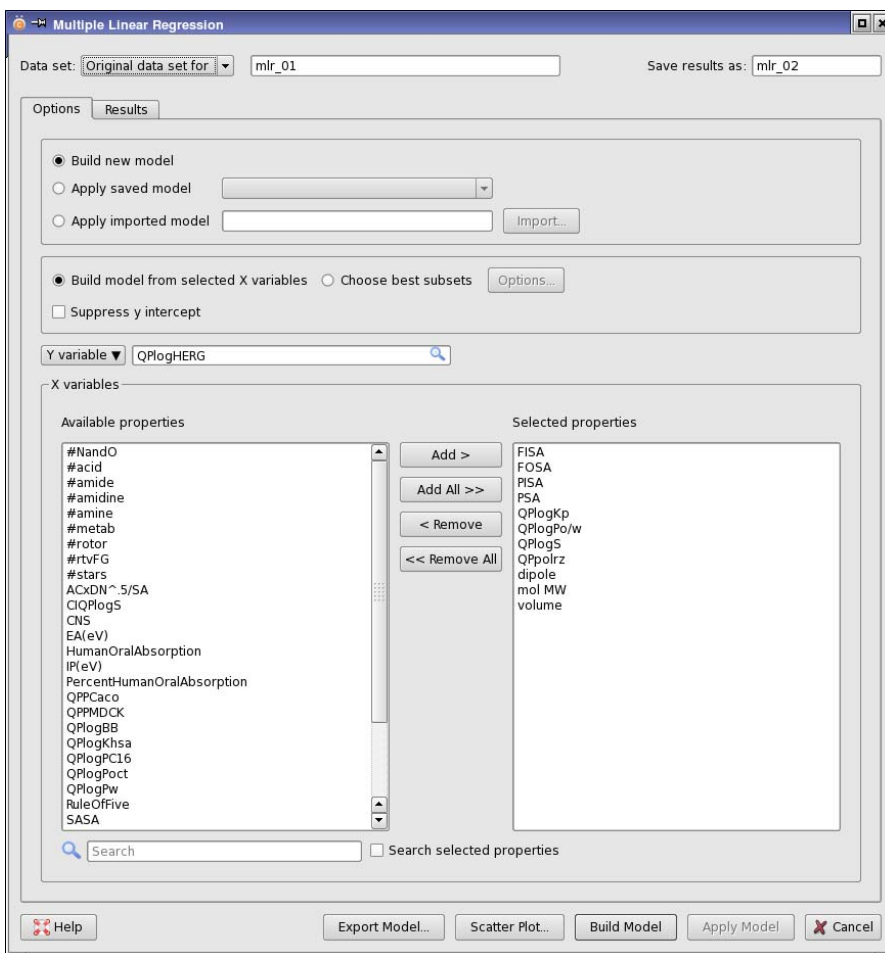


Figure 3.24. The Multiple Linear Regression dialog box

3.5.2 Partial Least-Squares Regression

For partial least-squares (PLS) regression, you can set the maximum number of PLS factors, and you can stop adding PLS factors when the standard deviation drops below a limit that you set. The dependent variables can be scaled automatically so that they lie on the same range. You can also set a minimum t-value for selecting significant independent variables.

When you build the model, it includes models for each number of PLS factors from 1 to the maximum, or the point at which the standard deviation falls below the threshold.

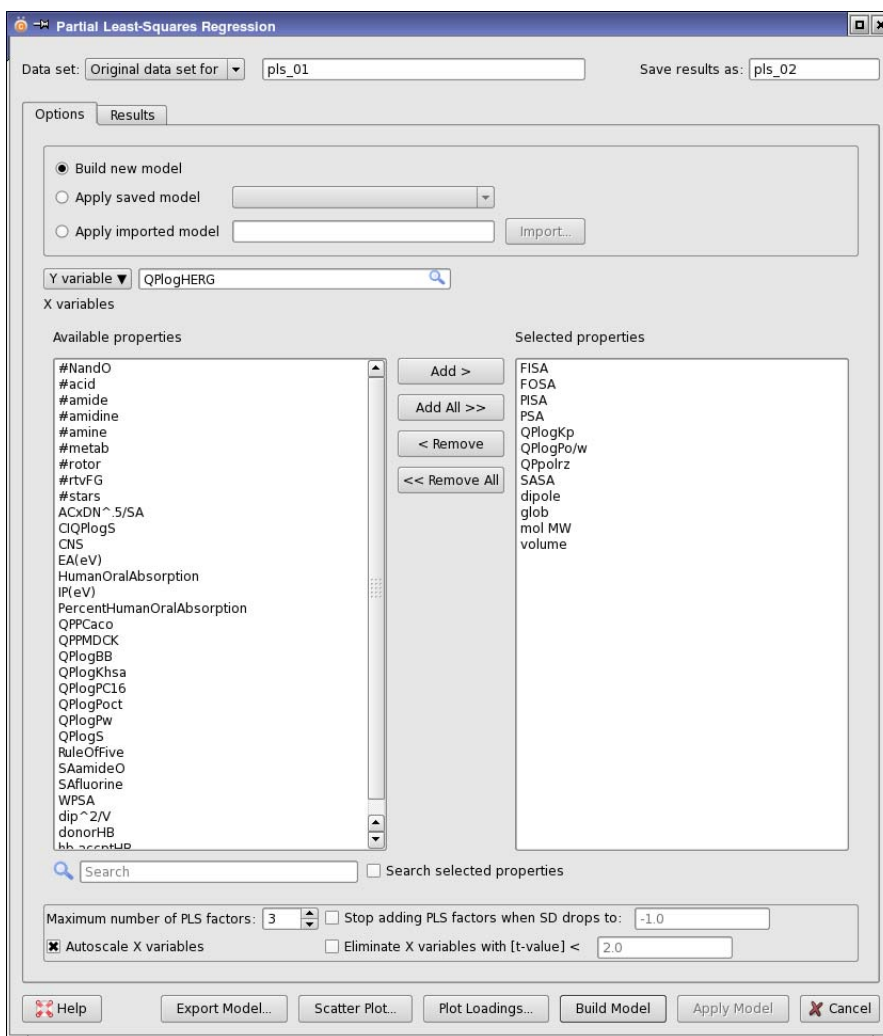


Figure 3.26. The Partial Least-Squares Regression dialog box

After the model is built, the loadings for each PLS factor for each X variable are displayed in the Factor loadings tab. You can plot the loadings by clicking Plot Loadings. A standard scatter plot panel opens (see [Section 2.7.1 on page 48](#)), and you can choose the factors that you want to plot on the x axis, the y axis, by color, and by symbol size.

3.5.3 Kernel-Based Partial Least-Squares Regression

Kernel-based partial least-squares regression is an extension of partial least-squares regression that introduces some nonlinearity into the scalar products of X variables used in the regression via a “kernel”, which is some nonlinear function of these scalar products [7]. In Canvas, the kernel is a Gaussian function,

$$K(i, j) = \exp(-d_{ij}^2/\sigma^2_{ij})$$

where d_{ij} is the Euclidean distance between X variables i and j , and σ is the non-linearity parameter. This kernel replaces the simple scalar products of the X variables in the regression. In Canvas, no automatic tuning of σ is done.

You can set the maximum number of KPLS factors, and you can stop adding KPLS factors when the standard deviation drops below a limit that you set. You can manually adjust the kernel nonlinearity using the slider or the box. The value set by this slider is $1/\sigma$, so values close to zero mean almost linear, and large values mean very nonlinear. Higher nonlinearity typically leads to tighter fitting, but it also tends to give poorer predictions on new compounds.

The uncertainty in the predictions for the test set can be estimated using bootstrapping. Select Calculate uncertainty on test set predictions, and set the number of cycles used in the bootstrapping method. Bootstrapping is done by sampling the training set randomly with replacement to generate a new test set of the same size (that may include duplicates), building a model and making predictions of the test set, then repeating the procedure a specified number of times. The standard deviation from the original test set is then calculated as the uncertainty.

When you build the model, it includes models for each number of KPLS factors from 1 to the maximum, or the point at which the standard deviation falls below the threshold. The error and the predicted uncertainty (if calculated) is displayed in the structures table for the number of KPLS factors chosen in the Test or Training results tab.

Models can be built on properties or on fingerprints. If you build a model based on fingerprints, you can visualize the contributions of each atom to the model. When the job finishes, right-click on the job name in the Project View and choose View to open the Kernel-Based Partial Least-Squares Regression dialog box, then click Visualize Model.

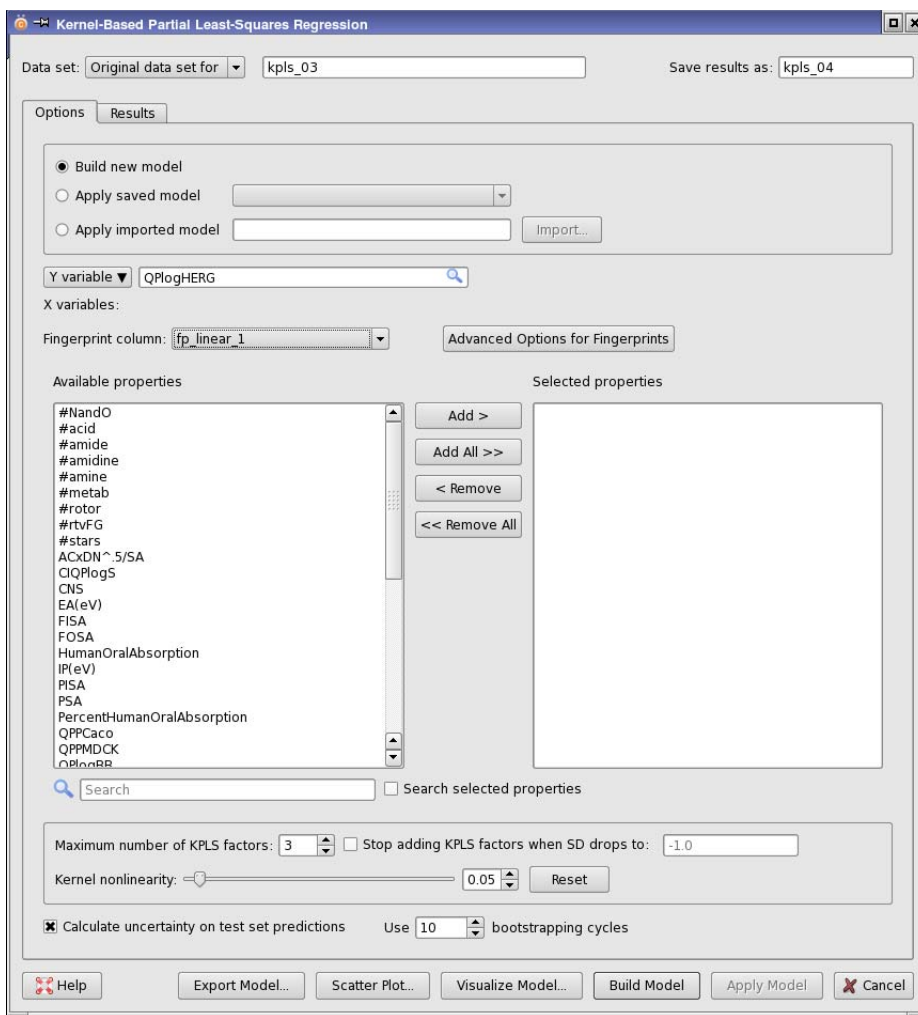


Figure 3.27. The Kernel-Based Partial Least-Squares Regression dialog box

The Model Visualization panel displays the structures in a spreadsheet with the observed and predicted values of the property, the uncertainty in the predicted value for test set molecules (if calculated) and the set to which they were assigned in the building of the model. You can sort the spreadsheet rows in ascending order of property values.

Each atom that contributed to a fingerprint used in building the model is marked with a colored disk that represents the value of the contribution to the property due to that atom. Two colors are used: one for negative values and one for positive values. The color saturation indicates the

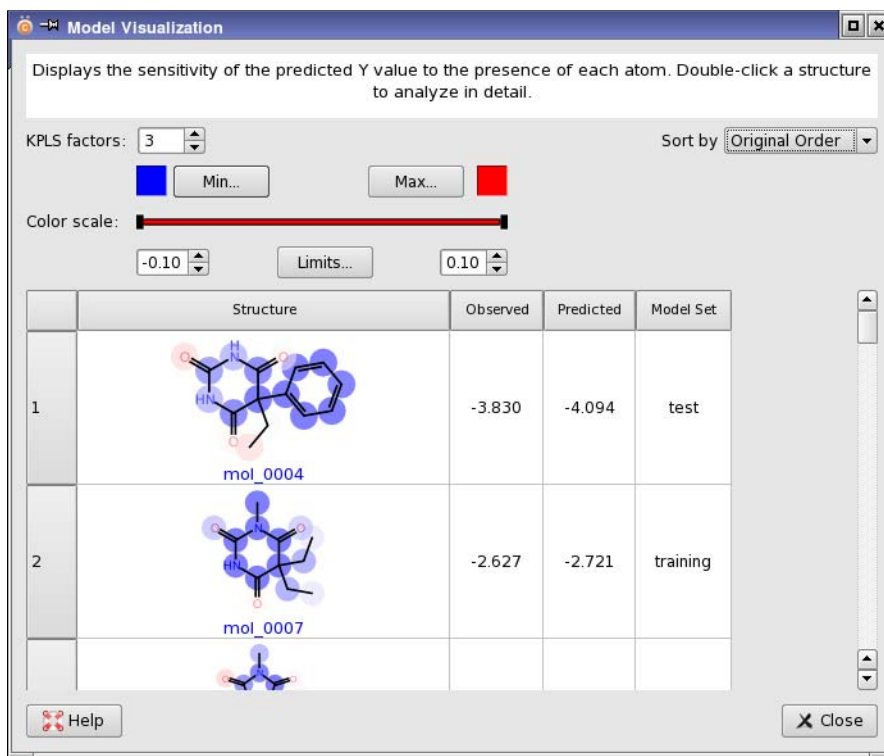


Figure 3.28. The KPLS Model Visualization panel.

magnitude of the contribution. Atoms that did not appear in any fingerprint are not marked with a disk.

You can change the number of KPLS factors for the display, to examine the effect of adding KPLS factors. You can change the color scale, either with the slider or the text boxes, to represent the range of values of interest, and you can adjust the limits on the range of values in the Color Scale Limits dialog box, which you open by clicking Limits. The Min and Max buttons open color selectors, so you can change the colors used in the display.

For more detailed information on the results for a particular structure, double-click it, to open the Analyze Atomic Contributions dialog box. In this panel, tool tips display the numerical values of the contributions for each atom. Double-clicking on an atom opens another dialog box that shows copies of the 2D structure molecule with each fragment in the fingerprint that contains the atom marked on the structure, and the contributions from the fragment and from the atom in the fragment.

As well as examining contributions, you can edit the structure to create a new structure, by clicking Edit Structure. When you click OK, the new structure is displayed, its fingerprint is calculated, and the contributions are reanalyzed and displayed. The observed value is still that of the original, but the predicted value is the value for the new molecule. You can then save the new molecule to the project by clicking Save to Project. It is added as a new row, with the predicted property value. This allows you to interactively optimize the property value.

3.5.4 Principal Components Regression

For principal components (PCA) regression, you can set the maximum number of principal components, and the dependent variables can be scaled automatically so that they lie on the same range.

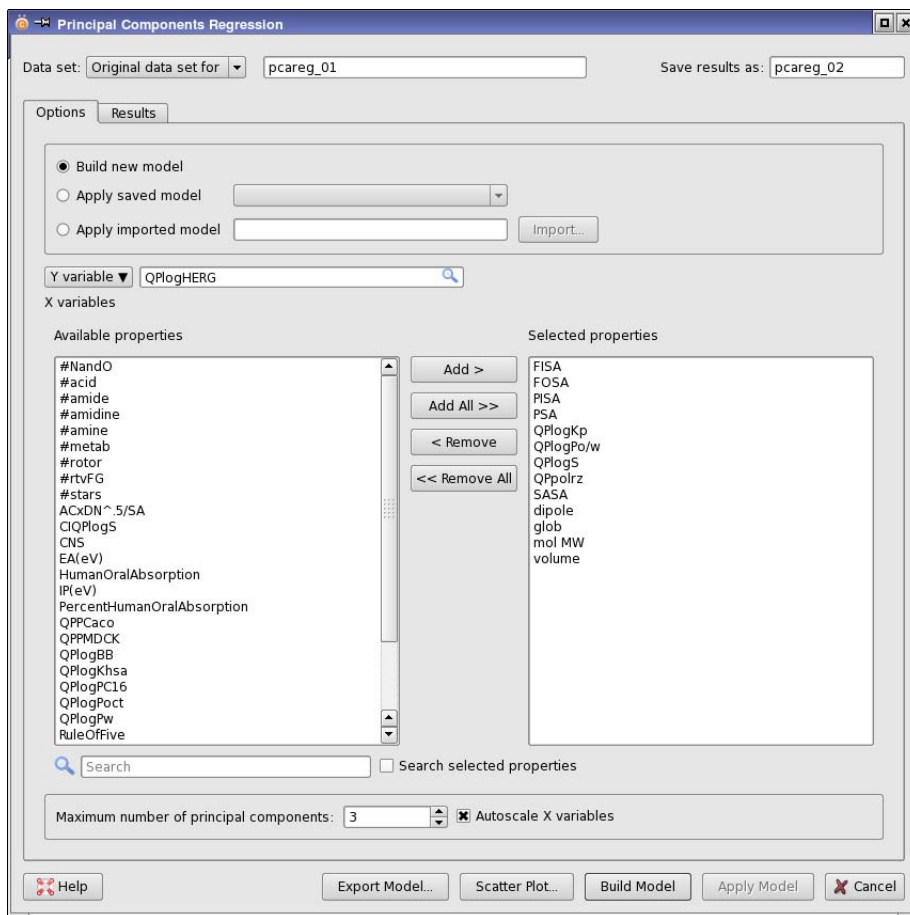


Figure 3.29. The Principal Components Regression dialog box

3.5.5 Bayes Classification

For Bayes classification, you can choose a Canvas fingerprint for the independent variables instead of selecting properties, and set options for partitioning the dependent variable. When a model is built, the tabs to the right of the structures table show the classes and the number of members of each class that were correctly classified and incorrectly classified.

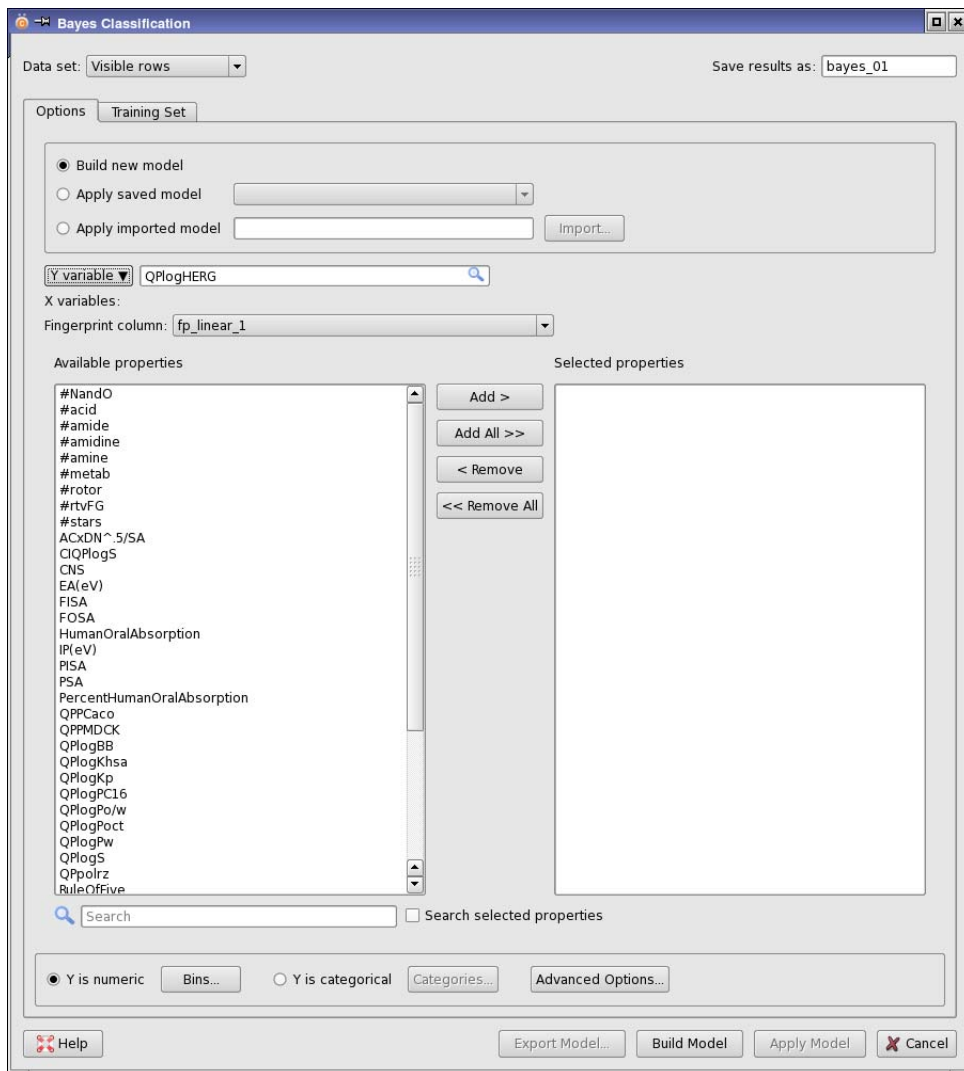


Figure 3.30. The Bayes Classification dialog box

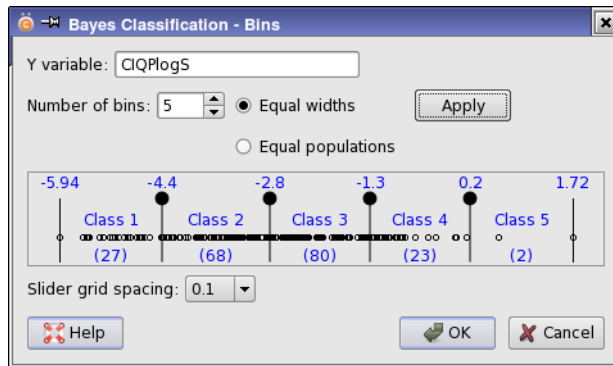


Figure 3.31. The Bayes Classification - Bins dialog box

If the dependent variable is numerical, you can choose whether to divide it into ranges that may cover more than one value by selecting Y is numeric, or you can treat each variable as a category by selecting Y is categorical. To assign the values to bins, click Bins, and make settings in the Bayes Classification - Bins dialog box. You can set the number of bins, and partition the variable so that the bins are of equal width or of approximately equal population. After making changes, click Apply. The display area shows information on the bins. The data points are plotted horizontally. The bin boundaries are represented by vertical lines, with a dot on the top. The value of the Y variable at the bin boundary is displayed above the boundary line. Labels for each bin, or class, are displayed above the data points and populations of each bin are displayed below the data points. You can add boundaries by clicking between the dots, and you can move boundaries by dragging the dots. The amount by which the boundary changes when you drag the dots is determined by the choice made from the Slider grid spacing option menu.

Parameters of the method can be set in the Bayes Classification - Advanced Options dialog box, which you open by clicking Advanced Options. You can set the smoothing coefficient, and for fingerprint data you can choose whether to apply distance cutoffs, or to keep bits in the top N percent. The methods are described in ref. 14.

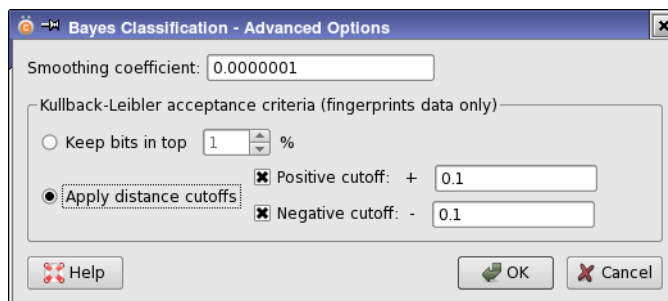


Figure 3.32. The Bayes Classification - Advanced Options dialog box

Note that running a Bayes classification with more than about 10 categories is unlikely to produce meaningful predictions.

3.5.6 Neural Networks

For neural networks, you can specify the number of networks to train, the number to keep, the number of training cycles, and the fraction of the training set to use for cross-validation, selected randomly.

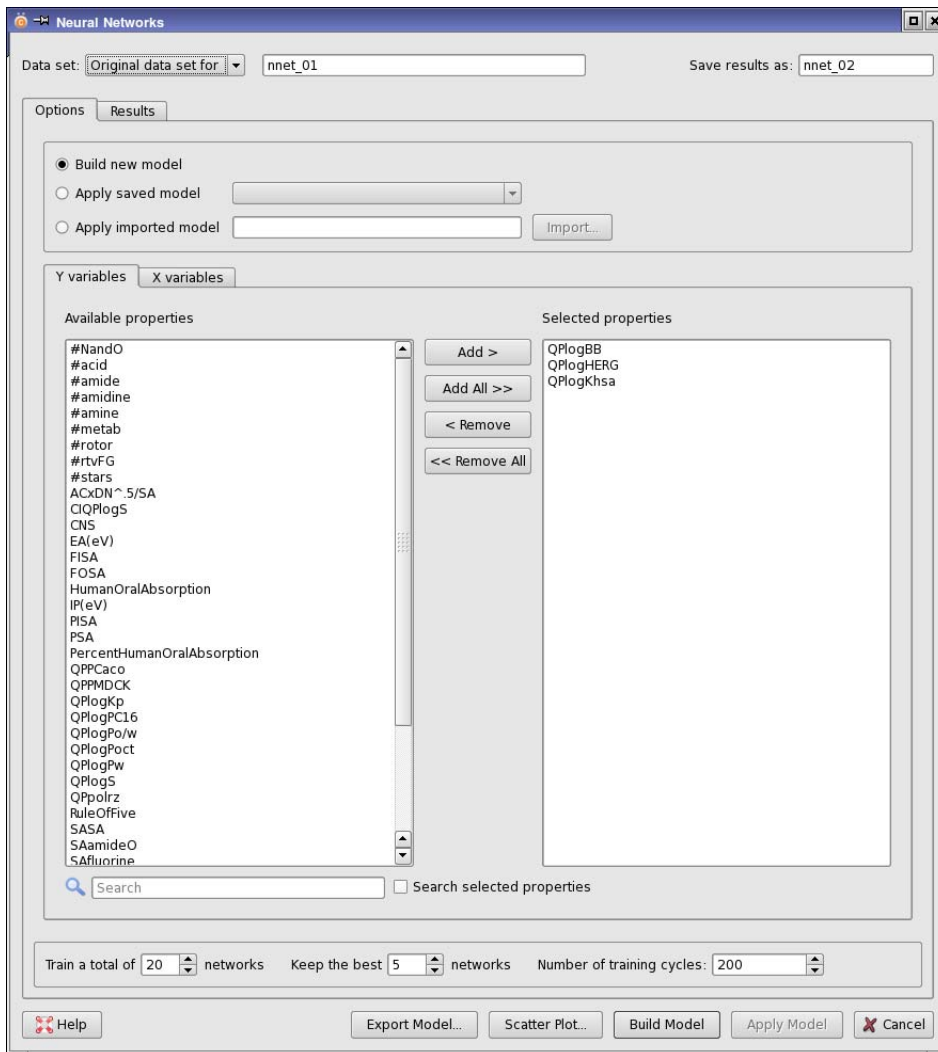


Figure 3.33. The Neural Networks dialog box

3.5.7 Recursive Partitioning

Recursive partitioning is a method like Bayes classification that can be used to build a model for the prediction of property ranges or categories, using a decision tree.

You can set options for partitioning the dependent variable and for building a single tree or an ensemble model, which generates multiple trees. Building an ensemble model filters out noise and corrects the biases in a single tree. When a model is built, the tabs to the right of the structures table show the classes and the number of members of each class that were correctly classified and incorrectly classified.

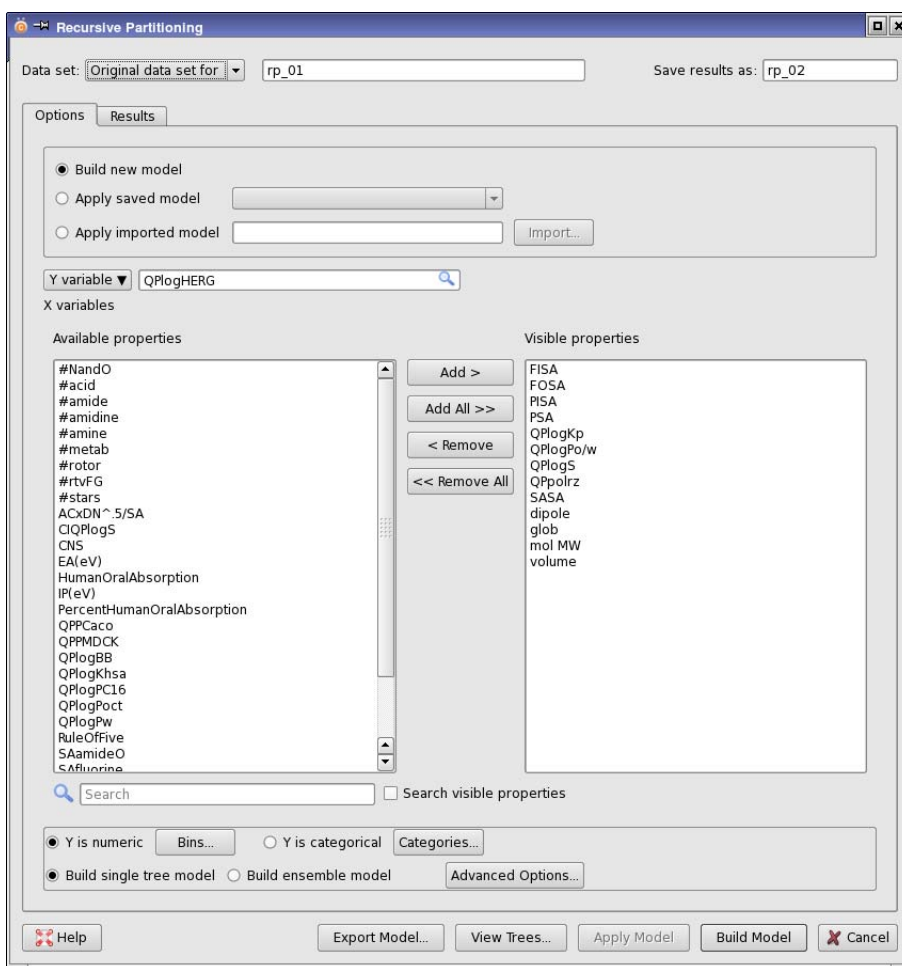


Figure 3.34. The Recursive Partitioning dialog box

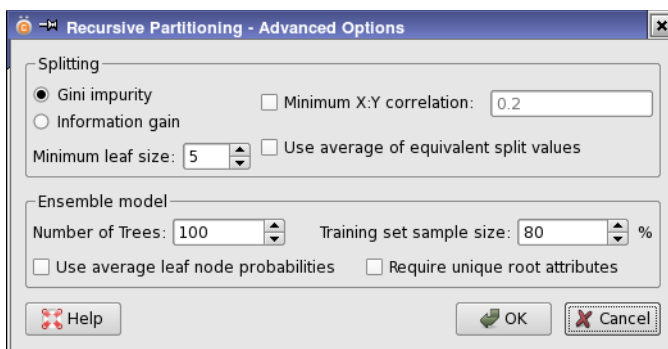


Figure 3.35. The Recursive Partitioning - Advanced Options dialog box

If the dependent variable is numerical, you can choose whether to divide it into ranges that may cover more than one value by selecting *Y is numeric*, or you can treat each variable as a category by selecting *Y is categorical*. To assign the values to bins, click *Bins*, and make settings in the *Recursive Partitioning - Bins* dialog box (This is the same as for Bayes classification—see [Figure 3.31](#).) You can set the number of bins, and partition the variable so that the bins are of equal width or of approximately equal population. After making changes, click *Apply*. The display area shows information on the bins. The data points are plotted horizontally. The bin boundaries are represented by vertical lines, with a dot on the top. The value of the *Y* variable at the bin boundary is displayed above the boundary line. Labels for each bin, or class, are displayed above the data points and populations of each bin are displayed below the data points. You can add boundaries by clicking between the dots, and you can move boundaries by dragging the dots. The amount by which the boundary changes when you drag the dots is determined by the choice made from the *Slider grid spacing* option menu.

Parameters of the method can be set in the *Recursive Partitioning - Advanced Options* dialog box, which you open by clicking *Advanced Options*. You can choose the splitting method, filter out variables that have little correlation with the dependent variable, and set parameters for generating an ensemble model. The ensemble in Canvas can be compared to the random forest method, with the difference being that Canvas builds each tree from a different random sample of the training set, whereas random forest uses the full training set, but chooses a different random pool of *X* variables at each decision node.

3.5.8 Self-Organizing Maps

Kohonen self-organizing maps can be built and applied in the Self-Organizing Map panel, which you open from the Applications menu. The panel has some features in common with those of the regression and other models.

- The models can be built or applied for the selected rows or the visible rows, a saved view, or the data set from a previous run.
- The results are saved with the name that you specify in the Save results as text box. This name is also the job name and appears in the Project View panel.
- Existing maps can be applied to the data set. Maps from within the project can be loaded by selecting Apply saved map and choosing the map from the option menu. Only maps that match the X variables choice are listed on the option menu. External maps can be

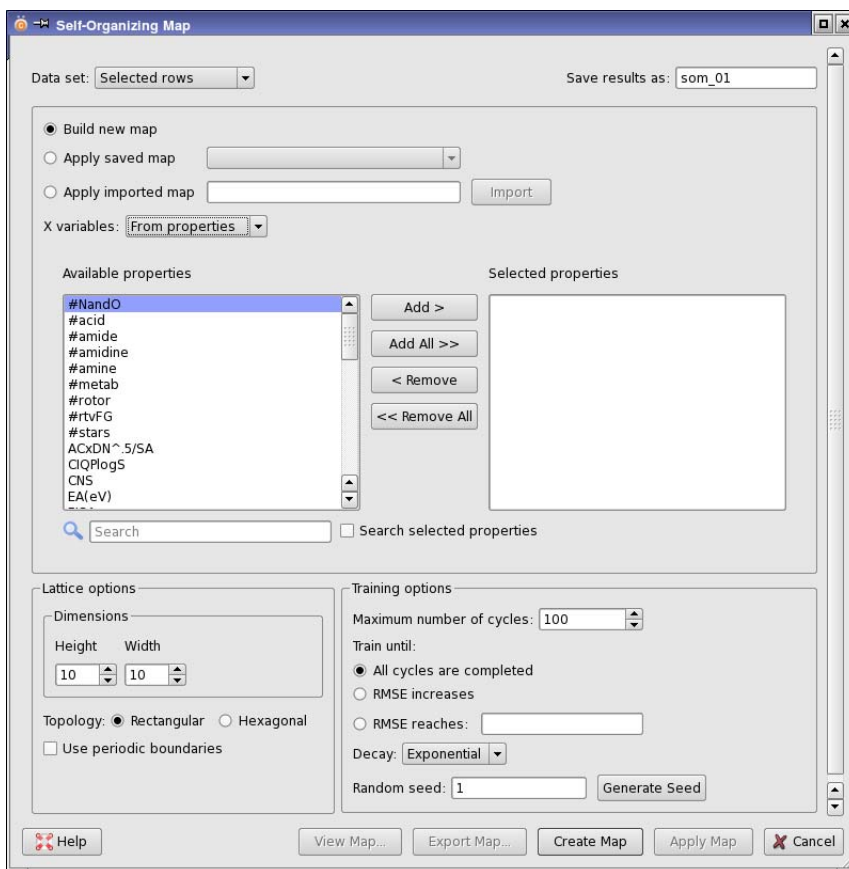


Figure 3.37. The Self-Organizing Map dialog box

loaded by selecting **Apply imported map**, clicking **Import**, and navigating to the file that contains the map. To apply the selected map, click **Apply Map**. The **Self-Organizing Map Viewer** panel opens and displays the results of applying the map to the current data set. When a map is loaded, you can export the map to a file by clicking **Export Map**.

To build a map, first select **Build new map**, then choose the type of X variable you want to use, from **Numeric** or **Fingerprints**. The panel is updated with the appropriate properties or fingerprints. You can choose a single fingerprint, or multiple properties.

In the **Lattice options** section, you can set the dimensions of the lattice and select the topology. You can also choose to use periodic boundary conditions.

In the **Training options** section, you can set the number of cycles, and specify a condition on the RMSE for stopping the training: either stop if the RMSE increases, or stop when the RMSE reaches a specified value. You can select the type of decay used to decrease the gain term as the training proceeds, and specify or generate a random seed for the initialization of the map.

Click **Create Map** to create the map. The dialog box closes, and the progress is logged in the **Messages View** panel.

To display the map, right click on the record in the **Project View** panel and choose **View**. The **Self-Organizing Map Viewer** panel opens with the map displayed. Alternatively, open the **Self-Organizing Map** panel again, select **Apply saved map** and choose the map from the option menu, and click **View Map**. You can have multiple viewer panels open at the same time.

The **Self-Organizing Map Viewer** panel provides several choices for coloring the map cells:

- **Distance to selected cell**—Color the map by the distance to the selected cell, which is outlined in red. Click on a cell in the map to select it.
- **Cell population**—Color the map by the number of structures in each cell.
- **Cell property value**—Color the cell according to the property chosen from the option menu.
- **Average property value**—Color the cell according to the average value of the property chosen from the option menu. This menu is populated by choosing **Numeric** from the **Properties** menu, and selecting the properties you want to use from the dialog box that opens.
- **Category**—Color the cell according to the number of structures in the cell with the selected value of the selected categorical property. The menu of categorical properties is populated by choosing **Categorical** from the **Properties** menu, and selecting the properties you want to use from the dialog box that opens.

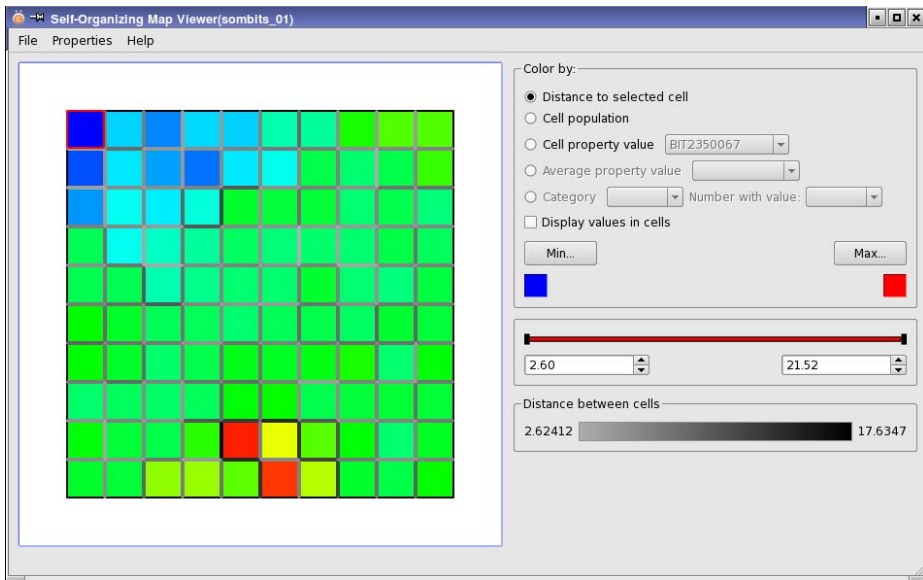


Figure 3.38. The Self-Organizing Map Viewer panel

The value that was used to determine the color for the cell can be displayed in the cell by selecting Display values in cells.

You can choose the colors for the minimum and maximum values of any quantity used to color the map by clicking the Min and Max buttons. Each button opens a color chooser. The color is displayed below the button. Intermediate values are interpolated on a spectrum.

The boxes below the bar display the minimum and maximum values of the property used to color the map. You can adjust these values so that the cells with values outside the range specified in the boxes are colored by the colors for the minimum and maximum. The bar above these boxes indicates the range of possible values for which a color ramp is applied. You can drag the sliders on the bar to adjust the minimum and maximum property values used.

The shade of gray in the cell borders indicates the distance between adjacent cells: the darker the border, the larger the distance. A legend of the gray shades is displayed in the Distance between cells section.

When viewing properties other than the distance to a cell, you can right-click a cell to display a custom view of the structures that are members of that cell. If you want to add members of other cells to the view, control-right-click the cell you want to add. Cells for the current view are outlined in yellow. You can add properties to the view by choosing View → Manage Properties, and selecting the desired properties.

Cell memberships can be exported to the spreadsheet or a file with File → Export. The cells are numbered from left to right, top to bottom, starting with 1. The Cell property contains the cell number for each structure, and the Total property contains the number of structures in the cell.

An image of the map can be saved in PNG, JPEG, or TIFF format by choosing File → Save Image. A dialog box opens that allows you to name the file and set the dimensions of the image, in pixels.

3.6 Other Applications

The remaining applications on the Applications menu are described in this section.

3.6.1 Principal Components Analysis

Analyzing the principal components of a set of independent variables can help in the selection of a subset of variables for use in other applications. The analysis can be done in the Principal Components Analysis dialog box, which you open from the Applications menu.

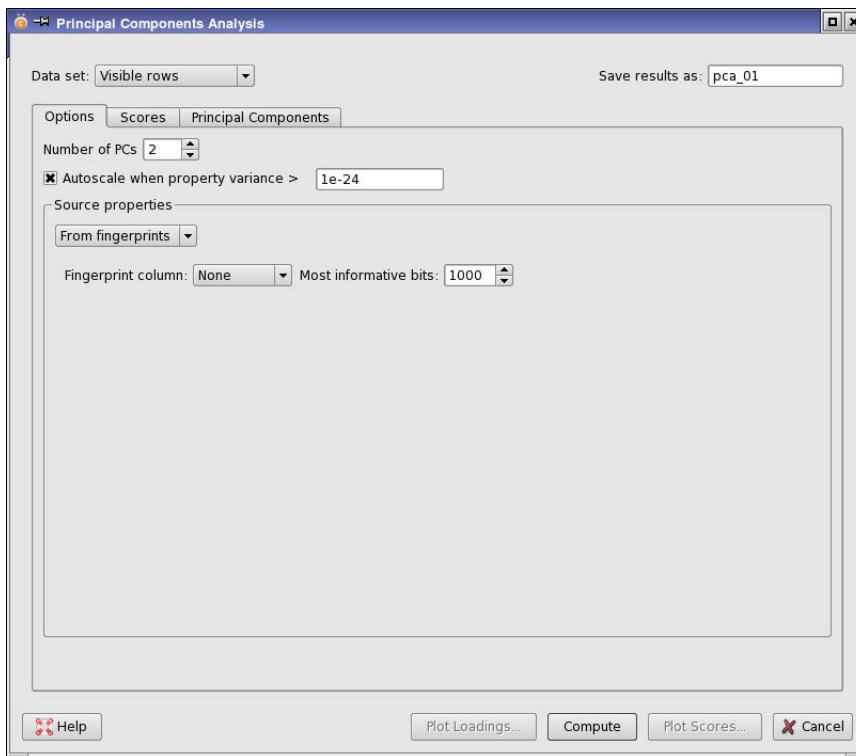


Figure 3.39. The Principal Components Analysis dialog box

To set up the analysis, choose the **Data set** option for the rows you want to use in the analysis, provide a name in the **Save results as text box**, specify the number of principal components (PCs), and choose a set of numeric properties or a fingerprint in the **Source properties** section.

You can select the **Autoscale** option to scale the property values by dividing by the standard deviation if the variance is greater than the specified value.

For fingerprint analysis, you can specify the number of most informative bits to keep. This reduces the time taken for the analysis.

The job to generate the principal components and compute the scores and loadings is run when you click **Compute**.

When the job finishes, the components are displayed in the **Principal Components** tab, showing the eigenvalue, cumulative variance, and coefficients for each property. The principal component scores for each of the structures are shown in the **Scores** tab. If you pause the pointer over the structure name, the 2D structure is displayed in a tool tip. To export the scores to a file or the spreadsheet, select the columns that you want to export and right-click in the table. The **Export Selection** dialog box opens, in which you can choose the destination, and for export to the spreadsheet, the property prefix.

You can plot the scores or the loadings by clicking **Plot Scores** or **Plot Loadings**. In both cases, a **Scatter Plot** panel opens, displaying the scores or loadings for two principal components. You can represent values for two more principal components by the color and size of the plot symbols. See [Section 2.7.1 on page 48](#) for more information on **Scatter Plot** panels. The name and structure of the molecule is displayed in a tool tip for each point in the scores plot, and the name of the property is displayed in a tool tip for each point in the loadings plot.

3.6.2 Finding the Maximum Common Substructure

The task of finding the maximum common substructure (MCS) of a number of structures can be run from the **Maximum Common Substructure** dialog box, which you open from the **Applications** menu. As for other applications, the structures that are processed can be the selected or the visible rows in the spreadsheet, a saved view, or structures from a previous run. If you choose **Selected rows** or **Visible rows**, you must ensure that the rows are selected or visible before you open this dialog box. You can name the job in the **Save results as text box**.

The search does not need to encompass all the structures in the set. You can specify a minimum number of structures that must match the substructure and a maximum number that must match. If the minimum is greater than the number of structures in the set, all structures must match. If the minimum is less than the maximum, a series of substructures (MCS groups) is found for each number of structures from the minimum to the maximum. You can restrict the groups so that any structure is only in one group, which is chosen to be that of the largest MCS.

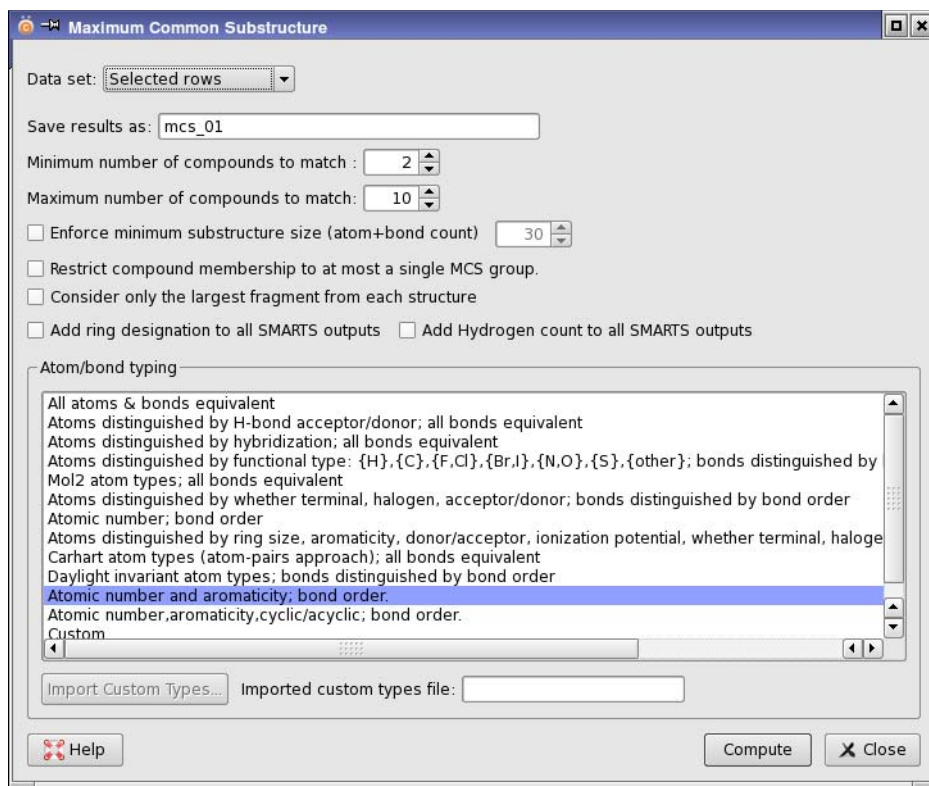


Figure 3.40. The Maximum Common Substructure panel

You can set a lower limit on the size of the MCS, in terms of the combined number of atoms and bonds. If the structure contains more than one molecule (fragment), you can restrict the search to the largest fragment.

The matching is done using a set of definitions for the atom types and the bond types. Twelve definitions are provided in the Atom/bond typing list, and you can add your own scheme by choosing Custom, then clicking Import Custom Types to load the definition file, which must consist of SMARTS strings that define the atom or bond type. The schemes are described in [Table 5.7 on page 125](#); all except E can be used for MCS searches.

The results include the SMARTS strings for the substructures, the number of substructures in each group, and the membership of the groups. You can add ring and hydrogen atom designations to the SMARTS strings by selecting the appropriate options. By default they are not included.

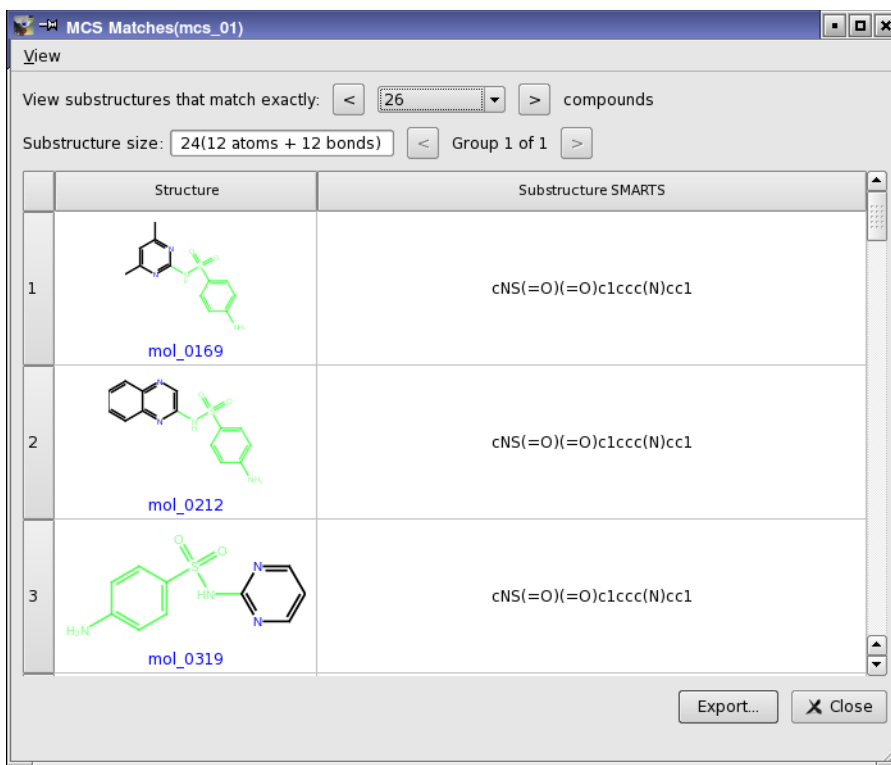


Figure 3.41. The MCS Matches view

To view the results, right-click on the job record under MCS in the Project View panel, and select View. The MCS Matches panel opens, displaying the structures and the substructure SMARTS patterns. The substructures are marked in green on the structures. You can view the substructures that match different numbers of structures, and if there is more than one group for a given number of molecules, you can display each group.

If you want to add properties to the view, choose View → Manage Properties, and select the desired properties in the dialog box that opens. The property is added to all subsequent views of MCS matches.

The set of structures displayed can be applied to the master view by choosing View → Apply to Master. The master view displays only the structures in the MCS Matches panel, and these structures are selected.

You can export the structures to a Maestro, SD, or CSV (SMILES + Properties) file by clicking Export. A file selector opens, in which you can choose the file format, navigate to a location, and name the file.

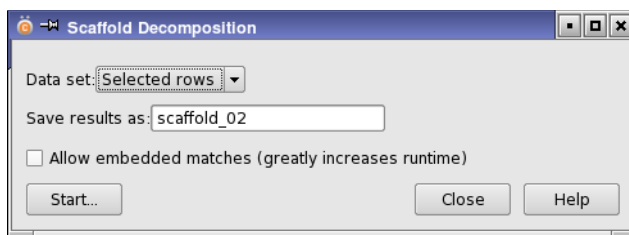


Figure 3.42. The Scaffold Decomposition dialog box.

3.6.3 Scaffold Decomposition

Another way of finding common structural features in a set of structures is to decompose each molecule into a set of scaffolds, then make a list for each scaffold of the structures that contain that scaffold. You can perform this scaffold decomposition in the Scaffold Decomposition panel, which you open from the Applications menu. In this panel you can choose the structure source from the Data set option menu, and give a name to the results.

The decomposition is done by first removing side chains from all ring systems, to define the largest scaffold in the structure. This structure is then split into all possible subscaffolds by breaking bonds and removing linkers between rings, other than those in fused rings. The list of scaffolds is collated, and then for each scaffold, the structures that contain that scaffold are located and added to a list.

When finding the structures that contain a scaffold, you can choose to match scaffolds to a fused ring system in a structure that contains the scaffold. This is called an “embedded match”, and you can allow them by selecting Allow embedded matches. For example, if the list of scaffolds contains imidazole and benzimidazole, a structure that contains benzimidazole is included in the list of structures that contain the imidazole scaffold, if embedded matches are allowed. Identifying embedded matches takes substantially more time than not doing so.

For more information on the process, see [Section 5.6.8 on page 184](#).

After the job finishes, you can view the results in the Scaffold Decomposition of *jobname* panel, by right-clicking on the job and choosing View.

The Scaffold Decomposition of *jobname* panel shows the 2D structures of the scaffolds in a table. The scaffolds can be classified by the number of rings in each scaffold, or by the number of ring systems. The scaffolds for a given classification are shown on the same table row. You can select scaffolds in the table, and limit the display to show only those scaffolds. Each time you limit the display (or remove the limits), that view of the table is stored, and you can use the arrow buttons to step through these views. If you want to see all the subscaffolds or superscaffolds of a particular scaffold, right-click in its table cell and choose the appropriate command.

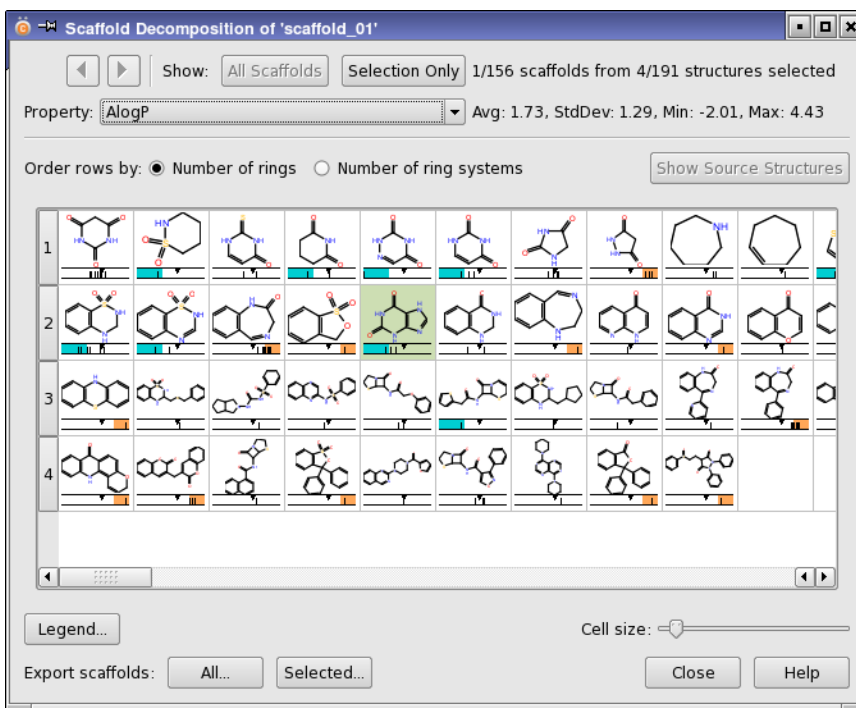


Figure 3.43. The Scaffold Decomposition of jobname panel.

The table also shows the property distribution of the structures that contain each scaffold, for a property that you can choose from the Property option menu. The average, standard deviation, minimum, and maximum values of the property are displayed to the right of the option menu. The property distribution is displayed for each scaffold as a “stick spectrum” in a bar that represents the full range of the property. For each structure there is a vertical line at the property value for that structure. The average property value is marked, and the region outside one standard deviation from the average is colored if it contains vertical lines for property values.

You can export the scaffolds to a structure file. The scaffolds are stored internally as canonical SMILES strings, so if you choose Maestro or SD format for the file, they are exported as 2D structures without hydrogens, except where the hydrogens are part of the SMILES string.

If you want to see the source structures for a set of scaffolds, select them in the table and click Show Source Structures. The Source Structures panel opens, and displays the structures for those scaffolds, arranged in rows with the same classification as for the scaffolds, and including the property distribution display. You can export selected structures or all structures to a file. You can run an R-group analysis on the selected structures or on all structures in the

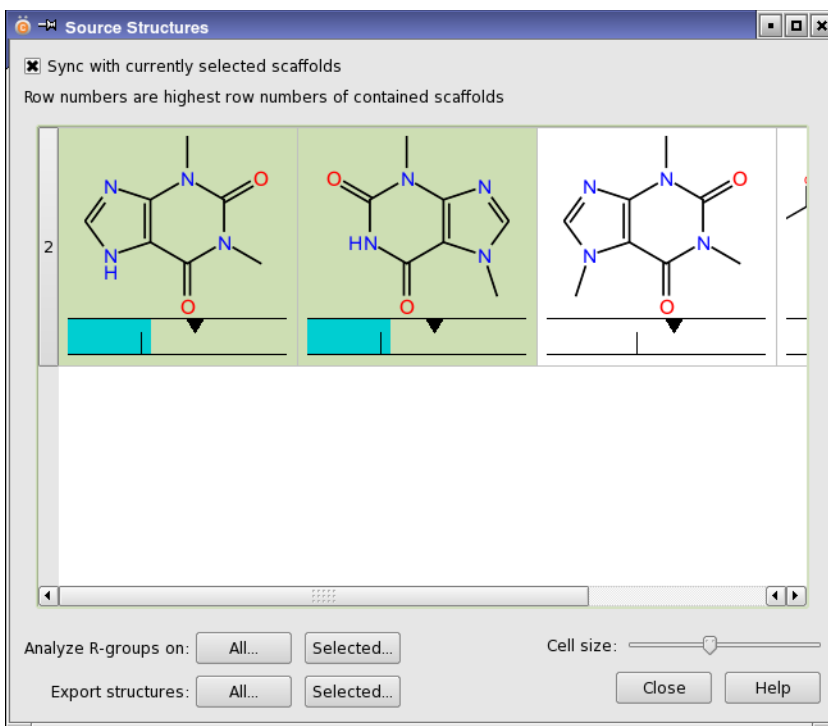


Figure 3.44. The Source Structures panel.

panel. If the structures have more than one scaffold in common, you are prompted to choose the scaffold you want to use for R-group analysis.

3.6.4 R-Group Analysis

R-group analysis identifies the R groups that are attached to a scaffold for a set of structures, and presents analyses of properties as a function of the attachment points and the R groups at each point. This application is described in a separate document, *R-Group Analysis*.

3.6.5 Running a 3D Minimization

The structures that are imported into Canvas or created in Canvas do not have to be 3D structures. You can convert them into 3D structures and minimize them with the 3D Minimization dialog box, which you open from the Applications menu.

You can minimize the structures in the selected rows or the visible rows, a saved view, or the data set from a previous run. When the minimized structures are returned, you can replace the existing structures or append the minimized structures to the project, with an optional suffix to

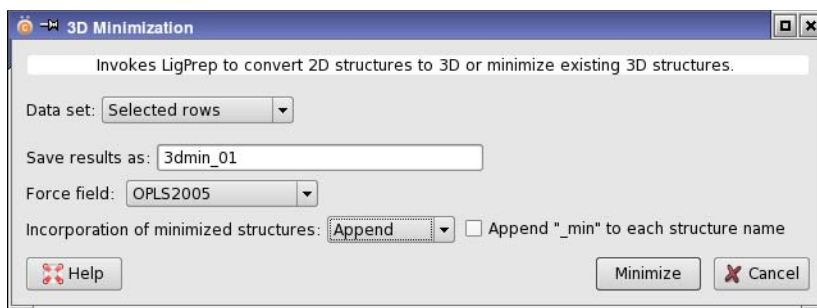


Figure 3.45. The 3D Minimization dialog box.

the name. You can export them to an external file in Maestro or SD format. The extension you give to the file name determines whether it is compressed or not. If no extension is given, the files are compressed, with extension `.maegz` for Maestro format and `.sdf.gz` for SD format.

The minimization process uses LigPrep, for which a license is required. LigPrep can generate stereoisomers, tautomers, and protonation states, and search for ring conformations, as part of the 2D-3D conversion process. In this application, it retains the tautomeric and protonation state, and searches for the lowest-energy ring conformation. If stereochemical information is provided with the 2D structure, it uses that information, but if the information is not available, it generates a stereoisomer derived from the stereochemistry observed in a database of natural products.

Two force fields are available for the minimization of the final structure: MMFFs and OPLS_2005. OPLS_2005 supports a wider range of elements than MMFFs.

3.7 Running External Applications

If you want to run an external application and incorporate the results into Canvas, you can do so with the External Application dialog box, which you open from the Applications menu.

In common with other applications panels, you can use the selected rows or the visible rows, a saved view, or the data set from a previous run for the structure input to the application. If the application requires certain properties as input, you can select the properties to use with the property selection tools in the Application input section. The structures and properties are written to a file, which you can name in the Application input section. The file is written to a special directory in the Canvas project, and the full path to the file is used when the command for the application is executed. You can choose a file format from the Format option menu. The application must be able to read one of the formats on this option menu.

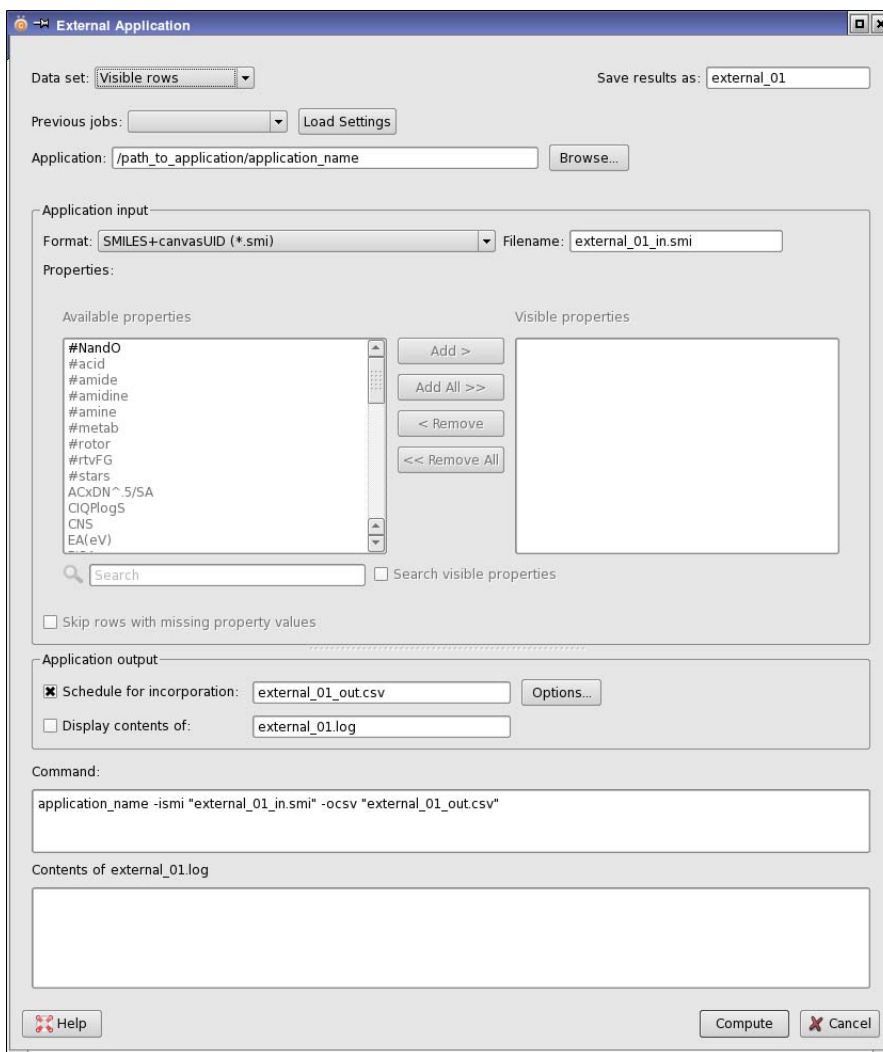


Figure 3.46. The External Application dialog box

The application can be specified in the Application text box as a full path, or you can click Browse and navigate to the application. If the application makes use of Schrödinger software, libraries, or APIs, you can select Invoke with \$SCHRODINGER/run so that the Schrödinger environment is set up automatically.

If you want the results to be incorporated back into the Canvas project, select Incorporate in the Application output section and specify the name of a CSV file in the text box. In order to incor-

porate results into Canvas, the application must generate a CSV file that has a header row, the SMILES string for the structure in the first column, and the structure name in the second column. The following columns can contain the properties generated by the application. If the application cannot process one of the structures, it must still write a row for that structure in the CSV file, with the SMILES string, the structure name, and the correct number of commas. You can also display a text file in the Contents of text area when the job finishes, by selecting Display contents of and specifying the file name in the corresponding text box. As for the input file, the output file and the displayed file are written to a directory in the Canvas project, and the full path is used when the command is run.

The default file names are constructed from the label given in the Save results as text box. This label is used to identify the calculation in the Project View, under External Application.

The command for running the application is specified in the Command text area. You do not need to specify the path to the application, because it is specified in the Application text box, and the path is prepended to the application name when the command is executed. Likewise, you do not need to specify the path to the input file or the output file, because the path is prepended to the file names. A sample command is given in the Command text area, which you can edit to provide the correct syntax for running the command, including the method of specifying the input and output file. The names of these two files are automatically supplied and updated from the relevant text boxes in the Application input and Application output sections.

All files must be specified as command arguments, because the command is not run in a shell. The application must therefore process all the arguments you give for the command, and not use shell redirection for input or output.

If you choose to display the contents of a file in the panel, that file must appear in the command. The command can include other arguments used by the application, and it can specify other files needed by the application. Any files other than those specified in the panel must include the full path.

When you are ready to run the application, click Compute. The controls in the panel are made unavailable, except for the Stop button. The panel must remain open until the run has finished. If you chose incorporation, the results are automatically incorporated into the Canvas project.

If you want to run a similar job, or correct the settings, you can select Use settings from and choose a run from the option menu. The panel is set up with the settings from the selected run, which you can change before clicking Compute. This option does the same as cloning a job.

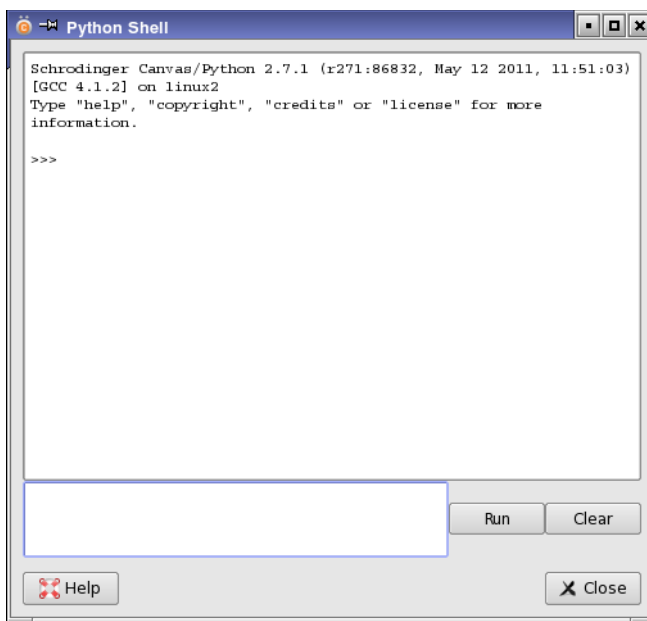


Figure 3.47. The Python Shell panel.

3.8 Running Python Scripts

Canvas provides two ways of running Python scripts from within the GUI.

You can install your own Python scripts on the Scripts menu with Scripts → Install and manage them (delete, reload) with Scripts → Manage. Once installed, you can run them by choosing the item on the Scripts menu. Scripts can be downloaded from the Script Center at <http://www.schrodinger.com/scriptcenter>, or you can create your own. In general, the scripts should run without the need to supply arguments. The scripts are run in the Schrödinger environment, so they can use any of the *Canvas Python API* modules. The mechanism for installing and managing the scripts is the same as used in Maestro—see Section 15.1 of the *Maestro User Manual* for more information. However, scripts that run in Maestro will not in general run in Canvas, because they do not have access to the Maestro libraries.

You can run Python scripts and Python commands in a Python interpreter, which you open by choosing Python → Open Shell or typing CTRL+P. This interpreter uses the version of Python that is in the Schrödinger software distribution, and various Schrödinger and Canvas modules are imported by default.

An overview of the general Schrödinger Python modules and a detailed description of those modules is available from the Help menu.

Using Canvas

Canvas has many possible uses. In this chapter, instructions are given for some of the scenarios in which Canvas can be used to solve a particular problem.

4.1 Selecting Compounds Not Represented in a Library

For purchasing new compounds, it can be useful to determine which compounds in a given set (available for purchase, for example) are not represented in your existing library. You can use diversity-based selection to find these compounds and export the structures to a file as follows. It is assumed that the two libraries are in separate files.

1. Import the existing library and the pool of compounds to choose from into a new Canvas project.

When you import the libraries, ensure that you create a view from each library (select Store each input file as a view).

2. Create a fingerprint for all compounds, with Applications → Binary Fingerprints. The fingerprint is used for assessing the diversity.
3. Select all the compounds in the pool that you want to choose from.

You can do this by opening the view for the pool of compounds, choosing View → Apply to Master, and selecting the Rows option Select current.

4. Open the Diversity-Based Selection dialog box from the Applications menu.
5. For Data set, choose Selected rows.
6. Choose the fingerprint you generated from the Fingerprint column option menu.
7. For the Diversity selection method, choose Sphere.
8. Enter the number of compounds you want in the Diverse subset size text box.
9. For the Initialization method, choose Existing structures.
10. For Use structures from, choose the custom view that contains the existing library.

11. Click Choose Compounds.

A diverse collection is selected that avoids the structures in the existing library, i.e., all diverse compounds will lie outside the spheres of the existing library structures. The collection is displayed in a new custom view.

12. Choose File → Export, or click the Export button in the new custom view, to export the diverse collection to a file.

If the existing library is already in a project, you can open the project, create a view from all the compounds in the project, then import the pool of compounds you want to choose from, ensuring that you create a view when you import it. Then you can proceed from step 2 above. At the end, you can delete the compound pool if you want to keep your library project. Alternatively, you could import the existing library project into a new Canvas project, create a view, and then import the compound pool.

4.2 Removing Duplicate Structures from a Project

1. Choose Structure → Detect Duplicates.

A new column, labeled Duplicate of is added to the spreadsheet.

2. Choose Data → Property Filter.

The Property Filter dialog box opens.

3. Ensure that Current view is selected for Store results in.

4. Choose Duplicate of for the property under Show rows where.

5. Choose has a value for the operator.

6. Click Add.

The filter string Duplicate of has a value is added to the list.

7. Click Run.

The dialog box closes and the rows for duplicate structures are now selected in the Master View.

8. Choose Edit → Delete.

The duplicates are deleted from the project.

4.3 Choosing a Fingerprint

Choosing which fingerprint to use and the options to set can be difficult, as no single setting is best for all targets. Some large scale screening studies [5,6] reported the following findings:

- MOLPRINT2D, Dendritic, Radial, Triplet, and Pairwise fingerprints performed better than others.
- More specific atom typing schemes, such as Daylight, Mol2, Carhart, were better.
- Collisions in the fingerprints significantly degraded performance, so using more bits is better.
- Combining complementary fingerprints can minimize the risk of choosing the wrong fragment.
- Combining actives into a modal fingerprint almost always increases the retrieval rate.
- Not much sensitivity was found to bit scaling or the similarity or distance metric.

4.4 Creating Modal Fingerprints

Modal fingerprints [6] are fingerprints that are averaged over several structures, thus combining the information from several query molecules into a single fingerprint.

To create a modal fingerprint:

1. Create a binary fingerprint with the type of your choice.
2. Select the rows in the spreadsheet for the structures whose fingerprints you want to average to create a modal fingerprint.
3. Right-click on the fingerprint column heading, and choose Create Modal Fingerprint.
4. Choose Selected rows from the Domain option menu.
5. Enter a name for the row in which the modal fingerprint will be stored
6. Click OK.

A new table row is created, with the name you chose displayed in the Structure column, and the modal fingerprint in the fingerprint column you used.

4.5 Tips

This section contains some tips for common operations.

To define training and test sets before building a model:

- Create a partition from a property with classes for the training set and the test set, and select Assign from partition to assign the sets in the model-building panel.

To copy model predictions to the spreadsheet:

- Right-click in the structures table and select To spreadsheet.

To run an application using a custom view:

- Choose View → Apply to Master in the custom view and select Select current, then run the application on the selected rows.

To add properties to an application view:

- Choose View → Manage Properties, and select the properties you want to add.

Running Applications from the Command Line

All the Canvas command-line programs are available in `$(SCHRODINGER)/utilities`. The syntax statements below are therefore given with respect to this directory.

Canvas programs generally run as standalone processes in the foreground. Some programs have Job Control options, which allow you to specify a job name and run the program under Job Control. This gives you access to all the Job Control features, including running the job on a remote host.

Canvas supplies two utilities, `canvasJob` and `canvas_app`, that allow you to run jobs that read from and write to a Canvas project. The list of programs you can run in this way is given in [Section 5.6.3 on page 174](#) with the description of `canvasJob`.

The default amount of memory used by Canvas applications is approximately 500 MB. To change this value, set the environment variable `SCHRODINGER_CANVAS_MAX_MEM` to the desired number of bytes (MB and GB units are not permitted). For information on setting environment variables on Windows, see [Appendix A](#) of the *Installation Guide*.

5.1 Common Syntax Descriptions

Many of the programs share common syntax, especially those that perform similar tasks. This section describes some elements of common usage.

A range specification is a comma-separated list of indices or ranges, with no spaces. A range is defined by the endpoints of the index range, separated by a colon. Examples of valid range specifications are:

1,4	structures 1 and 4
1:10,14	structures 1 through 10 and 14
2:	structures 2 through the end of file
:5,13:18	structures 1 through 5 and 13 through 18

In the option syntax, ranges are often used for selecting structures from a file, and are denoted *structureRange*.

Structure files can often be specified in either compressed or uncompressed format. These files are specified on the command line with `-ifmt` for input or `-ofmt` for output. The file types that are generally accepted by Canvas, listed with the value of *fmt* and the allowed extensions, are:

mae	Maestro file (.mae, .maegz, .mae.gz)
sd	SD file (.sdf, .sd, .sdf.gz, .sd.gz)
csv	Comma-separated value file including the structure in SMILES format, molecule name and properties (.csv, .csv.gz); can use comma, tab or space for the delimiter
smi	Space or tab-separated SMILES and molecule name (.smi, .smi.gz)
proj	Canvas project (.cnv).

When CSV files are used by a program, it is usually necessary to specify the delimiter that separates the fields on each row. You should use ' ' for space, and '\t' for tab. The single quotes are required for these two specifications. If the delimiter is a space, consecutive spaces are ignored.

Programs that take a range of structure files as input have a common input file argument syntax, which is described in Table 5.1. The output file syntax varies more, according to the purpose of the program. However, there are some common options associated with particular file types, which are listed in Table 5.2.

Table 5.1. Common input structure file arguments.

Argument	Description
-ifmt <i>inFile</i>	Structure file containing the input molecules. Required. The supported formats <i>fmt</i> are: smi SMILES string and name, separated by tab or space sd SD file mae Maestro file csv CSV file containing SMILES, name, and properties Valid file extensions are listed on page 119.
-fieldAsName <i>field</i>	Property in the SD or Maestro input file to use as the molecule name. For Maestro files, <i>field</i> must start with <i>s_</i> or <i>i_</i> . Cannot be used with <i>-ism_i</i> , or <i>-icsv</i> . Default: <i>s_m_title</i> property for Maestro files, first line of CT block for SD files.
-noHeader	CSV input file does not have a header line.
-d <i>delimiter</i>	Delimiter used in CSV input file. Valid delimiters are comma (','), space (' '), and tab('\t'). Default: comma.
-smi <i>SMILESCol</i>	Field in CSV input file that contains SMILES strings, either by name if there is a header row, or by column index starting at 1. Default: column 1.
-name <i>nameCol</i>	Field in CSV file to use as the molecule name, either by column name or by column index. Default: column 2 if SMILES is present, column 1 otherwise.

Table 5.2. Common output structure file arguments

Argument	Description
-ofmt <i>outFile</i>	Structure file containing the output molecules. The supported formats <i>fmt</i> are: smi SMILES string and name, separated by tab or space sd SD file mae Maestro file csv CSV file containing SMILES, name, and properties Valid file extensions are listed on page 119 .
-od <i>delimiter</i>	Delimiter used in CSV output file. Valid delimiters are comma (','), space (' '), and tab('\t'). Default: comma.
-v3	Use MDL version 3 format for output SD files.

Many of the programs can run under Job Control. For these programs, the syntax descriptions include *job-options*. Where the options are supported, the common set of supported options are listed in [Table 5.3](#). With the exception of -JOB and -MINREC, these options are described in full in [Table 2.1](#) and [Table 2.2](#) of the *Job Control Guide*.

Most programs run on a single processor. Those that can be distributed over multiple processors are indicated in the description of the program.

Table 5.3. Common Job Control options.

Option	Description
-JOB <i>jobName</i>	Job name. If present, run the job under Job Control. If omitted, no other job control options are permitted.
-HOST <i>host</i>	Run job on <i>host</i> . The format <i>host:n</i> requests use of <i>n</i> CPUs on <i>host</i> .
-LOCAL	Store temporary job files in current directory.
-MINREC <i>nrec</i>	Minimum number of records per CPU. Prevents submission of a large number of subjobs that each contains only a small number of records. Only available for programs that support use of multiple CPUs Default: 100.
-TMPDIR <i>dir</i>	Store temporary job files in <i>dir</i> .
-WAIT	Do not return control to the shell until the job finishes.
-NICE	Run job at reduced priority.

5.2 2D Fingerprints

The programs available for 2D fingerprints are listed in [Table 5.4](#), and are described in detail in the following sections.

Table 5.4. Tools for 2D fingerprints

Tool	Description
canvasFPGen	Generate fingerprints for a set of molecules.
canvasFPCombine	Combine fingerprints for distinct or overlapping sets of molecules.
canvasFPBinary2CSV	Extract fingerprints to a CSV file.
canvasCSV2FPBinary	Convert a CSV file to canvas binary fingerprint file.

5.2.1 canvasFPGen

This program generates fingerprints for molecules in a structure file. The syntax of the command is as follows.

```
canvasFPGen inputFileArgs outputFileArgs [job-options] [options]
```

The input file arguments are the common arguments listed in [Table 5.1](#). The alternatives for specifying the output file are listed in [Table 5.5](#). The job options are standard Job Control options, listed in [Table 5.3](#). This program can be distributed over multiple processors. The options are listed in [Table 5.6](#).

Table 5.5. Output file arguments for the canvasFPGen command.

Argument	Description
-o <i>fpFile</i>	Binary output file for generated fingerprints that stores only molecule name and fingerprint bits.
-odata <i>fpFile</i>	Binary output file for generated fingerprints that also includes additional data fields. If the input file is in SMILES format, there are no such additional fields.
-ocsv <i>fpFile</i>	Write MACCS or custom fingerprints directly to a CSV file. Column headers consist of either the SMARTS patterns in the fingerprint definition file or the optional label that follows each SMARTS pattern in the fingerprint definition file.
-fieldOnly <i>fields</i>	Space-separated list of property names (fields) from a Maestro or SD input file to write to the output file. Only valid with -odata.
-uniform	Each structure block in a Maestro or SD input file contains the same set of properties.

Table 5.6. Options for the *canvasFPGen* command.

Option	Description
-fptype <i>fpType</i> [<i>customfile</i>]	Fingerprint type. Allowed values are: linear, maccs, radial, molprint2D, torsion, pairwise, triplet, dendritic, and custom. custom must be followed by the name of the file that contains the SMARTS patterns for the fingerprints (1 per line). Default: linear.
-xp	Represent fingerprints using 64 bit precision. This reduces collisions of “on” bits, but doubles the space required to store each key.
-fill	Fill in the line for a molecule that fails to generate fingerprints in the output files. No data values are included on this line, apart from the molecule name. This option is useful to create placeholders that preserve positional alignment to external data. Default: skip failed molecule in the output.
-compress	Use frequency-based compression to reduce required storage by approximately tenfold.
-3D	Use actual 3D coordinate distances instead of topological distances. Applicable only to pairwise and triplet fingerprint types. Requires a Maestro or SD input file with 3D coordinates.
-atomtype <i>scheme</i> [<i>customfile</i>]	Atom typing scheme. Must be an integer value between 1 and 12 or C or E. See Table 5.7 for details. For C, <i>customfile</i> must be given. Defaults for each fingerprint are linear:10, radial:4, torsion:10, pairwise:9, triplet:10, molprint2D:5, dendritic:10.
-path <i>pathLength</i>	Maximum path length for linear, dendritic, and molprint2D fingerprint types. Defaults: 7 (linear), 5 (dendritic), 2 (molprint2D).
-ring <i>ringSize</i>	Maximum linear ring path. Valid only with linear fingerprint type. 0 means no ring closure. Default: 14.
-minpath <i>pathLength</i>	Minimum path length for linear, dendritic, and molprint2D fingerprint types. Defaults: 0 (linear, dendritic), 2 (molprint2D).
-halfstep	When growing linear fragments, additionally include growth by bonds alone. Valid only for linear fingerprint type.
-binwidth <i>value</i>	Bin width for distances (pairwise, triplet fingerprint and E atom typing scheme). Default: 1.
-binoverlap <i>value</i>	Distance threshold for assignment of item to multiple bins. Useful for “fuzzy” binning. Default: 0.
-minbin <i>value</i>	Minimum distance for binning. Default: 0.
-maxbin <i>value</i>	Maximum distance for binning. Default: no limit.
-iter <i>iterations</i>	Number of iterations for radial fingerprints. Valid only with radial fingerprint type. Usually between 2 and 6. Default: 4.

Table 5.6. Options for the *canvasFPGen* command. (Continued)

Option	Description
-estatePath <i>length</i>	Path length for Estate atom typing. Default 2.
-estateWidth <i>width</i>	Binning step for Estate atom typing. Default: 0.25.
-miniter <i>iterations</i>	Set minimum for radial iterations below which features are discarded. Valid only with radial fingerprint type. Default: 0.
-startH	Calculate molprint2D codes for hydrogens. Valid only with molprint2D fingerprint type. Default: do not calculate.
-endH	Include terminal hydrogens in molprint2D codes. Valid only with molprint2D fingerprint type. Default: do not include.
-min <i>fracMin</i>	Omit bits that are set by less than the specified fraction of molecules. Default: use all bits.
-max <i>fracMax</i>	Omit bits that are set by more than the specified fraction of molecules. Default: use all bits.
-noone	Omit bits only set in a single molecule. Overrides -min.
-noall	Omit bits set in all molecules. Overrides -max.
-reduce <i>bits</i>	Reduce the precision of fingerprints by the specified number of bits. May be applied after min/max filters. This option increases the chance of feature collisions. For example, a value of 22 will reduce each single precision key (32 bits) into a range of 1024 (10 bits).
-mostSig <i>nbits</i>	Keep only the <i>nbits</i> most informative bits across the chosen input set. Default: use all bits.
-n <i>structureRange</i>	The set of input structures to process. <i>structureRange</i> is a range specification, as defined on page 119 . Default: process all structures.
-obad <i>badMolFile</i>	Save the molecules that failed to generate a fingerprint to a file. This option is only available for SMILES, SD, or CSV input. Default: write information on failures to standard output.
-scaling <i>option</i>	Options to rescale binary fingerprint data to reals. Must be an integer—see Table 5.8 for details. Default: 0.
-strip	Process only the largest fragment by atom count. Default: process all fragments.
-stamp <i>name</i>	Override molecular title with the supplied string for all molecules. Useful in conjunction with <i>canvasFPCombine</i> for creating an ensemble or modal fingerprint.

Table 5.7. Atom typing schemes for the `-atomtype` option of the `canvasFPGen` command.

Scheme	Description
1	All atoms equivalent; all bonds equivalent.
2	Atoms distinguished as hydrogen bond acceptors or donors; all bonds equivalent.
3	Atoms distinguished by hybridization state; all bonds equivalent.
4	Atoms distinguished by functional type: {H}, {C}, {F,Cl}, {Br,I}, {N,O}, {S}, {other}; bonds by hybridization.
5	Mol2 atom types; all bonds equivalent.
6	Atoms distinguished by whether they are terminal, halogens, hydrogen bond acceptors or donors; bonds distinguished by bond order.
7	Atomic number and bond order.
8	Atoms distinguished by ring size, aromaticity, hydrogen bond acceptor or donor, ionization potential, whether terminal or halogen; bonds distinguished by bond order.
9	Carhart atom types (atom-pairs approach); all bonds equivalent.
10	Daylight invariant atom types; bonds distinguished by bond order.
11	Atoms distinguished by atomic number and bond order, aromaticity.
12	Atoms distinguished by atomic number and bond order, aromaticity, if aliphatic whether cyclic or acyclic.
C	Custom. Must be followed by location of a type definitions file. See the definition files (.typ) in <code>\$(SCHRODINGER)/mmshare-vversion/data/canvas</code> for examples. Cannot be used with fingerprint types <code>maccs</code> and <code>custom</code> .
E	Estate atom types. Cannot be used with fingerprint types <code>maccs</code> and <code>custom</code> .

Table 5.8. Fingerprint scaling options.

Option	Description
0	No scaling (default)
1	Scale counts by feature size to unity
2	Scale counts by feature size to feature size
3	Scale counts by feature size to molecule size
4	Scale squares of counts by feature size to unity
5	Scale squares of counts by feature size to feature size
6	Scale squares of counts by feature size to molecule size
7	Scale square root of counts by feature size to unity
8	Scale square root of counts by feature size to feature size
9	Scale square root of counts by feature size to molecule size
10	Use raw feature counts
11	Use square of raw feature counts
12	Use square root of raw feature counts
13	Use constant value of one

5.2.2 canvasFPCombine

This program combines fingerprints for distinct or overlapping sets of molecules. When a molecule appears more than once, the bit sets are combined using a logical OR. If fingerprints were created using the `-scaling` option of `canvasFPGen`, real-valued data for duplicate molecules are averaged. The syntax of the command is as follows:

```
canvasFPCombine -i file-list -o fpFile [options]
```

The command arguments are described in [Table 5.9](#).

Table 5.9. Arguments for the `canvasFPCombine` command.

Argument	Description
<code>-i <i>file-list</i></code>	Space-delimited list of input fingerprint files. If a single file is specified, fingerprints are combined for consecutive structures with the same molecule ID.
<code>-o <i>fpFile</i></code>	Output fingerprint file.
<code>-strict</code>	Input files must have the same molecule IDs in the same order.
<code>-cat</code>	Makes no attempt to merge records, Concatenates files in the order specified in <i>file-list</i> . Within each input file, the natural order is preserved.
<code>-compress</code>	Compress the output file.
<code>-d</code>	Delete input files after successful completion.
<code>-min <i>fracMin</i></code>	Omit bits that are set by less than the specified fraction of molecules. Default: use all bits.
<code>-max <i>fracMax</i></code>	Omit bits that are set by more than the specified fraction of molecules. Default: use all bits.
<code>-noone</code>	Omit bits only set in a single molecule. Overrides <code>-min</code> .
<code>-noall</code>	Omit bits set in all molecules. Overrides <code>-max</code> .
<code>-reduce <i>bits</i></code>	Reduce the precision of fingerprints by the specified number of bits. May be applied after min/max filters. This option increases the chance of feature collisions. For example, a value of 22 will reduce each single precision key (32 bits) into a range of 1024 (10 bits).
<code>-mostSig <i>nbits</i></code>	Keep only the <i>nbits</i> most informative bits across the chosen input set. Default: use all bits.

5.2.3 canvasFPBinary2CSV

This program extracts fingerprints to a CSV file. The syntax of the command is as follows:

```
canvasFPBinary2CSV -i fpFile {-o|-odata} csvFile [options]
```

The input file must be a binary fingerprint file generated by canvasFPGen. The output file is in CSV format, and contains the bit data if specified with `-o`; if `-odata` is specified, additional descriptor data is written out. By default, each row represents a molecule, and each column represents a bit or a descriptor. The options are described in [Table 5.10](#).

Table 5.10. Options for the `canvasFPBinary2CSV` command.

Option	Description
<code>-dense</code>	Include bits that are off in all molecules. Ignored unless input fingerprint file has been prepared with <code>-reduce</code> . Warning: using this option may require large amounts of disk space.
<code>-dtab</code>	Use tab as delimiter.
<code>-force</code>	If binary fingerprint has been rescaled to continuous values, force binary (0,1) representation. Default: print rescaled values.
<code>-max fracMax</code>	Omit bits that are set by more than the specified fraction of molecules. Default: use all bits.
<code>-min fracMin</code>	Omit bits that are set by less than the specified fraction of molecules. Default: use all bits.
<code>-noall</code>	Omit bits set in all molecules. Overrides <code>-max</code> .
<code>-noone</code>	Omit bits only set in a single molecule. Overrides <code>-min</code> .
<code>-notot</code>	Suppress printing of total set bit counts to <i>csvFile</i> . Useful for generating raw data for input into other programs. Default: print total set bit counts.
<code>-off value</code>	Use <i>value</i> instead of a blank when a bit is not set by a molecule.
<code>-on value</code>	Use <i>value</i> instead of a 1 when a bit is set by a molecule.
<code>-p</code>	Pivot output: bits as rows, molecules as columns. Rescaled values are not used in this mode (implies <code>-force</code>).
<code>-pf prefix</code>	Prefix to the bit column names in the header. Default is BIT. Cannot be used with <code>-sparse</code> .
<code>-sparse</code>	Emit comma separated list of keys for each molecule. Preceded by molecule name, without headers or footers. Incompatible with <code>-odata</code> , <code>-p</code> , <code>-pf</code> , <code>-off</code> , <code>-on</code> and implies <code>-force</code> and <code>-notot</code> .
<code>-totonly</code>	Suppress printing of individual rows. Useful for viewing only summary data. Default: print individual rows.

5.2.4 canvasCSV2FPBinary

This program converts a CSV file into a Canvas fingerprint binary file. If binary columns share a common prefix such as BIT, then non-binary data can be embedded alongside. If binary-specified columns contain real-valued data, then a Canvas scaled fingerprint will be generated. The command syntax is as follows:

```
canvasCSV2FPBinary {-icsv csvFile|-isparsed file} -o binaryFile [options]
```

The input CSV file must include a header, and may contain both binary and nonbinary data. If `-isparsed` is used for the input file, the file must contain the molecule ID in the first column, followed by the list of “on” bit indices. The options are described in [Table 5.11](#).

Table 5.11. Options for the `canvasCSV2FPBinary` command.

Option	Description
<code>-prefix string</code>	Prefix for binary columns in the header. The default prefix of a binary column is BIT. A blank string (<code>-prefix " "</code>) means to treat all columns as binary.
<code>-fieldAsName field</code>	Field in the CSV file to use as the molecule name. Default: first column.
<code>-noheader</code>	No header is supplied. The first column is presumed to be the molecule name and the rest binary data.
<code>-d delimiter</code>	Delimiter used in the CSV or sparse input file. Default: comma.
<code>-n rows</code>	Convert the specified rows in the CSV file. The first line in <code>csvFile</code> is the header; the second line in <code>csvFile</code> is the first data row, and is row 1. <code>rows</code> is a range specification, as described on page 119 . Default: convert all rows.
<code>-type string</code>	Supply custom type information for Canvas. If you are using this utility for more than one type of CSV file, where the same bits have different meanings, it is recommended you supply a distinct string for each type using this option. This utility is useful when multiple sets of fingerprints are generated independently from Canvas and you want to distinguish between them in the Canvas GUI. Fingerprints generated by Canvas and KNIME are tracked automatically, and do not need to be differentiated. Default: custom.
<code>-offset N</code>	Add this value to each binary key found. Useful for concatenation of multiple fingerprints.
<code>-bit on, off</code>	Identifiers to use for ON and OFF states delimited with a comma. Example: ON, OFF

5.3 3D Pharmacophore Fingerprints

The program `canvasPharmFP` generates fingerprints from 3D pharmacophores for the input structures. The fingerprint bits are set on the basis of the feature types and inter-feature distances.

The command syntax is as follows:

```
canvasPharmFP [job-options] program-options
```

The job options are described in [Table 5.3 on page 121](#). This program can be distributed over multiple processors. The program options are described in [Table 5.12](#).

Table 5.12. Options for the canvasPharmFP command.

Option	Description
-iproj -imae -isd <i>source</i>	Source of input structures. -iproj: Canvas project (.cnv) -imae: Maestro file (.mae, .maegz, .mae.gz) -isd: SD file (.sdf, .sd, .sdf.gz, .sd.gz) If <i>source</i> is a project, it must include the absolute path, and it must be accessible to all processors on which the job is run.
-fieldAsName <i>prop</i>	Use an alternate property as the source of structure names. For a list of project properties, run <code>canvasPharmFP -helpColumns</code> .
-n <i>list</i>	Process the specified subset of rows or structures. The format of the list is given in Section 5.1 on page 119 .
-file	Interpret <i>list</i> as the name of a binary file containing row numbers exported from the Canvas GUI. Valid only when <i>source</i> is a Canvas project.
-o -odata <i>fpFile</i>	Binary output file for fingerprints. If -o is used, the file contains fingerprints and structure names; if -odata is used the file contains fingerprints and other properties from the input source, as specified by -fieldOnly.
-fieldOnly <i>fields</i>	Write out only the specified properties to <i>fpFile</i> . The list is a space-separated list of property names. Use -helpColumns to get a list of project properties. Valid only with -odata.
-uniform	All structures in the input Maestro or SD file contain the same properties. Valid only with -odata. Not valid with -iproj.
-fill	Insert a blank fingerprint as a placeholder when a given structure cannot be processed.
-3pt	Generate 3-point pharmacophore fingerprints. This is the default.
-4pt	Generate 4-point pharmacophore fingerprints. To include 3-point and 4-point pharmacophores in each fingerprint, use -3pt and -4pt.

Table 5.12. Options for the `canvasPharmFP` command. (Continued)

Option	Description
-dmin <i>dmin</i>	Minimum distance between pharmacophore features in angstroms. Distances smaller than this value are placed in the first bin. Default: 2.0 Å.
-dmax <i>dmax</i>	Maximum distance between pharmacophore features in angstroms. Distances larger than this value are placed in the last bin. Default: 16.0 Å.
-width <i>width</i>	Distance bin width in angstroms. Default: 2.0 Å.
-overlap <i>overlap</i>	Set the bit for a neighboring bin if a distance is within <i>overlap</i> of that bin. Default: <i>width/2</i> .
-fd <i>fdFile</i>	Use pharmacophore feature definitions in <i>fdFile</i> . This file must be a feature definition file in Phase format (see Appendix B.3 of the <i>Phase User Manual</i>). If omitted, default Phase feature definitions are used.
-flex	Generate conformers and combine bits from different conformers using logical OR. Conformers are discarded after use. Requires a Phase license and uses one PHASE_DBCREATE token per subjob.
-sample {rapid thorough}	Sampling method for conformational space. Rapid sampling varies the rotatable groups independently; thorough sampling varies them together. Only valid with <code>-flex</code> . Default: <code>rapid</code> .
-cmax <i>maxConfs</i>	Maximum number of conformers. Only valid with <code>-flex</code> . Default: 100.
-bf <i>numPerBond</i>	Number of conformers to retain per rotatable bond in the structure. The total number of conformers is bounded by the product of <i>numPerBond</i> and the number of rotatable bonds, so if <i>maxConfs</i> is increased, it may be necessary to increase <i>numPerBond</i> in order to retain additional conformers for more rigid structures. Only valid with <code>-flex</code> . Default: 10.
-ewin <i>deltaE</i>	Energy window for keeping conformers in kJ/mol. Conformers whose energy is higher than this amount above the minimum energy conformer are discarded. Only valid with <code>-flex</code> . Default: 104.6 kJ/mol (25 kcal/mol).
-amide {vary orig trans}	Specify the method for amide torsion sampling. The allowed values are: <i>vary</i> —vary the torsion freely <i>orig</i> —use the input conformation <i>trans</i> —set the conformation to trans Only valid with <code>-flex</code> . Default: <code>vary</code> .
-xp	Store fingerprints using 64-bit precision. Default: 32-bit.
-compress	Use frequency-based compression to reduce required storage by approximately tenfold.
-fmin <i>fmin</i>	Omit bits that are set by less than the specified fraction of input structures. Not valid with <code>-mostSig</code> . Default: 0.0.

Table 5.12. Options for the `canvasPharmFP` command. (Continued)

Option	Description
<code>-fmax fmax</code>	Omits bits that are set by more than the specified fraction of input structures. Not valid with <code>-mostSig</code> . Default: 1.0.
<code>-mostSig n</code>	Keep only the <i>n</i> most informative bits over the set of input structures. Not valid with <code>-fmin</code> or <code>-fmax</code> .
<code>-reduce bits</code>	Reduce fingerprint precision by the specified number of bits. May be applied after the bit filters.
<code>-h[elp]</code>	Print this message and exit.
<code>-helpColumns</code>	Print list of Canvas project column names and exit. Valid only with <code>-iproj</code> . No other options are required.

5.4 Similarity, Dissimilarity, and Clustering

The programs for generating similarity information, dissimilarity-based selections of molecules, and for clustering are listed in Table 5.13, and are described in the following sections.

Table 5.13. Tools for similarity, dissimilarity, and clustering.

Tool	Description
<code>canvasCSV2PW</code>	Generate a binary pairwise similarity or distance matrix using a CSV input
<code>canvasCSVMatrix</code>	Generate a distance matrix based on CSV input. The result is a symmetric row-by-row matrix.
<code>canvasDBCS</code>	Dissimilarity-based compound selection using pairwise matrix or Canvas fingerprint files as input
<code>canvasFPHist</code>	Generate a histogram of nearest neighbor similarities from one or two fingerprint files
<code>canvasFPMatrix</code>	Generate a pairwise similarity or distance matrix using fingerprints from one or two sets of molecules.
<code>canvasHC</code>	Hierarchical clustering and report details for specific levels.
<code>canvasHCBuild</code>	Hierarchical clustering.
<code>canvasHCSelect</code>	Report details for specific levels of clustering.
<code>canvasKMeans</code>	K-means clustering.
<code>canvasLC</code>	Leader-follower clustering.
<code>canvasPW2CSV</code>	Convert binary pairwise matrix to CSV.
<code>canvasTreeDraw</code>	Draw dendrogram from a tree file.

5.4.1 canvasCSV2PW

This program generates a binary pairwise similarity or distance matrix from a CSV input file.

```
canvasCSV2PW -i csvFile -o binaryFile [-rowNames ] [-colNames ]
```

The arguments are described in [Table 5.14](#).

Table 5.14. Arguments for the *canvasCSV2PW* command.

Argument	Description
-i <i>csvFile</i>	File containing CSV data to be converted.
-o <i>binaryFile</i>	Output file for matrix in binary format.
-rowNames	Input CSV matrix has molecule name to begin each line. May be omitted if CSV file was created by Canvas.
-colNames	Input CSV matrix lists molecule names as first line. May be omitted if CSV file was created by Canvas.

5.4.2 canvasCSVMatrix

This program generates a distance matrix based on CSV input data, resulting in a symmetric row-by-row matrix.

```
canvasCSVMatrix -icsv csvFile { -o binaryFile | -ocsv csvFile } [-order n]
  [-cov [-scale [-t threshold]]]
```

The input CSV file, specified by `-icsv`, must contain the desired real-valued data. The output file can be a binary file, specified with `-o`, or a CSV file, specified by `-ocsv`. Both formats can be used in subsequent steps. The options are described in [Table 5.15](#).

Table 5.15. Options for the *canvasCSVMatrix* command.

Option	Description
-order <i>n</i>	Order of distance matrix as defined in $(\sum x^n)^{1/n}$. <i>n</i> must be a positive integer. The default value of <i>n</i> is 2, which corresponds to Euclidean distance.
-cov	Generate Mahalanobis distance matrix using the covariance matrix to correct non-orthogonality of the data.
-scale	Scale to unit variance, if it is within threshold.
-t <i>threshold</i>	Threshold for scaling. Default value is 1E-3.

5.4.3 canvasDBCS

This program performs dissimilarity-based compound selection based on a pairwise matrix or a Canvas fingerprint file. The command syntax is as follows:

```
canvasDBCS {-im matrixFile|-ifp fpFile [-ifp2 fpFile2 [-ifp3 fpFile3]]}
           [options]
```

The arguments for this command are described in [Table 5.16](#). The metrics are listed in [Table 5.17](#); a description of the metrics is given in [Table 5.22](#) for `canvasFPMatrix`.

Table 5.16. Arguments for the canvasDBCS command.

Argument	Description
-im <i>matrixFile</i>	Pairwise matrix file generated by <code>canvasFPMatrix</code> or another program. Must be symmetric and contain either similarities or distances. Can be used with all methods except <code>dise</code> .
-ifp <i>fpFile</i>	Canvas fingerprint file generated by <code>canvasFPGen</code> . Requires <code>sphere</code> or <code>dise</code> method.
-ifp2 <i>fpFile2</i>	Canvas fingerprint file representing the chemical space to avoid when choosing diverse compounds. A compound is not selected if it lies within the exclusion sphere of any compound in this file. Requires <code>sphere</code> method.
-ifp3 <i>fpFile3</i>	Canvas fingerprint file representing the chemical space to cover when choosing diverse compounds. A compound is selected only if it lies within the exclusion sphere of a compound in this file. Requires <code>sphere</code> method.
-o <i>outputFile</i>	Output file for results. Default is <code>stdout</code> .
-n <i>size</i>	Diverse subset size. Required for methods <code>maxsum</code> and <code>maxmin</code> . Allowed for methods <code>sphere</code> and <code>dise</code> where it is applied as a secondary filtering step if the exclusion distance produces a diverse subset larger than <i>size</i> .
-d <i>distance</i>	Exclusion distance for methods <code>sphere</code> and <code>dise</code> . The meaning of this parameter is the same whether a similarity or a distance metric is used. For example, if using Tanimoto similarity, <code>-d 0.7</code> selects compounds that exhibit pairwise distances of least 0.7, which means the pairwise similarities will be no higher than 0.3. Default: 0.5.
-method <i>type</i>	Choice of selection method. The allowed values are: <code>sphere</code> Sphere exclusion <code>dise</code> Directed sphere exclusion [2]. <code>maxsum</code> Maximum sum of distances. Each round adds the compound with the largest total distance from the current subset. <code>maxmin</code> Largest minimum distance. Each round adds the compound with the largest minimum distance from the current subset.

Table 5.16. Arguments for the `canvasDBCS` command. (Continued)

Argument	Description
<code>-metric name</code>	Metric type. <i>name</i> can be either the name or the index. Allowed values are listed in Table 5.17 by index and name. Applies only if a fingerprint file is specified as input. Default: <code>soergel (20)</code> .
<code>-noselfcheck</code>	When using <code>-ifp</code> and <code>-ifp2</code> to choose structures from the first set (<code>-ifp</code>) that are dissimilar to those in a second set (<code>-ifp2</code>), do not check whether the structures chosen from the first set are dissimilar to each other. Requires <code>sphere</code> method, <code>-ifp</code> and <code>-ifp2</code> . The default is to exclude structures from the first set that are similar to each other.
<code>-helpMetric</code>	Prints the definition of each metric then exits.
<code>-forceBinary</code>	Ignore any scaled fingerprint values in the input file. Binary values are always used. Only relevant to an input fingerprint file.
<code>-rowNames</code>	Input CSV matrix has molecule name to begin each line. Only applies to input CSV file.
<code>-colNames</code>	Input CSV matrix lists molecule names as first line. Only applies to input CSV file.
<code>-init type</code>	Initialization method. Allowed method types are: <code>random</code> , <code>first</code> , <code>representative</code> , <code>dissimilar</code> . <code>representative</code> and <code>dissimilar</code> can only be used with the <code>maxsum</code> and <code>maxmin</code> methods. Default: <code>random</code> .
<code>-seed int</code>	Seed for random number generation if the random initialization method was specified. Defaults: use the computer clock.
<code>-xpdise</code>	Run DISE method in a mode that produces even better coverage at the expense of speed.
<code>-guiprop</code>	Add a property column to the output file that contains 1 for selected structures and 0 for unselected structures.

Table 5.17. Available metrics for the `canvasDBCS` command.

Index	Name	Index	Name
1	buser	13	patternDifference
2	cosine	14	pearson
3	dice	15	petke
4	dixon	16	rogersTanimoto
5	euclidean	17	shape
6	hamann	18	simpson
7	hamming	19	size
8	kulczynski	20	soergel
9	matching	21	tanimoto
10	mcConnaughey	22	variance
11	minmax	23	yule
12	modifiedTanimoto		

5.4.4 canvasLibOpt

This program selects a subset of compounds from a pool by minimizing their similarity, and optionally by optimizing the ranges of specified properties. It can be used to fill holes in an existing library of compounds. The command syntax is as follows:

```
canvasLibOpt fpIn nstr csvOut [options]
```

where *fpIn* specifies the Canvas fingerprint file containing the pool of structures from which selections will be made, *nstr* specifies the number of structures to select, and *csvOut* specifies the output CSV file containing the selected structures. The options are described in [Table 5.18](#). The columns in the CSV file are described in [Table 5.19](#). The algorithm used is described below.

Given a pool of structures in a Canvas fingerprint file, `canvasLibOpt` uses a greedy, stochastic algorithm to select a diverse subset, with the goal of minimizing the sum of the average nearest neighbor Tanimoto similarity within the subset and the average property filter score, depending on which task is chosen.

If hole-filling was chosen, the nearest neighbor to a given structure is determined by considering both the structures in the subset and the structures in the library (fingerprint file *fpHoles*). If property filtering was chosen, the optimization attempts to minimize the average number of property filters that failed within the subset. A property filter is simply a range of

Table 5.18. Options for the `canvasLibOpt` command.

Option	Description
<code>-range rows</code>	Comma-delimited series of colon-separated ranges that define the subset of rows to select structures from in <i>fpIn</i> . By default, all rows are considered.
<code>-fill fpHoles</code>	Fill holes in the library represented by <i>fpHoles</i> . The fingerprints in this file must have been created using the same settings as <i>fpIn</i> . The range of structures considered to be in the library can be specified with <code>-range2</code> . <i>fpHoles</i> can be the same file as <i>fpIn</i> if you use <code>-range2</code> to define the library. Computational effort increases as the product of <i>nstr</i> and the size of this library, and memory increases linearly with the size of this library.
<code>-range2 rows2</code>	Consider only the specified subset of rows in <i>fpHoles</i> to be in the library. <i>fpHoles</i> and <i>fpIn</i> may be the same file if <code>-range</code> and <code>-range2</code> are used. Only valid with <code>-fill</code> .
<code>-filter file</code>	File containing a set of property filters, with one filter per line. Structures that fail fewer filters are preferred. Each filter consists of a property name, followed by minimum and maximum allowed values for that property. For example, <pre>AlogP -1.5 4.5 HBA 0 5 MW 200 500</pre> If the property name contains spaces, it must be surrounded by quotes. Use <code>-help_props</code> to see a list of the properties stored in <i>fpIn</i> .
<code>-cycles ncyc</code>	Number of optimization cycles. A given cycle consists of making <i>nstr</i> passes through the selected structures, with an attempt to replace the poorest scoring structure at each pass. A candidate structure is chosen randomly from the pool and the poorest scoring structure is replaced by that candidate if doing so improves the average Sim_NN and Filter values among all structures in the subset. If improvement does not occur, the candidate structure is accepted with a probability determined by a Monte Carlo test. The default number of cycles is 10.
<code>-stop dx, n</code>	Stop optimization when the total score changes (up or down) by less than <i>dx</i> for <i>n</i> consecutive cycles. Default: 0.001, 3.
<code>-tmax sec</code>	Run optimization for no longer than the specified number of CPU seconds. There is no time limit by default.
<code>-mc tol</code>	Monte Carlo acceptance criterion (<code>etest</code>). A total score increase of <i>tol</i> is accepted with a 50% probability in the first cycle and a 1% probability in the last cycle. Default: 0.001.
<code>-rand seed</code>	Integer random seed. Controls the choice of initial subset and all random operations throughout the course of the optimization. Default: 1.
<code>-col prefix</code>	Add the supplied column name prefix to the CSV output file columns Row, Sim_NN, Row_NN, Hole_NN, and Filter. For example, <code>-col "job1::"</code> would yield the column names <code>job1::Row</code> , <code>job1::Sim_NN</code> , etc. See Table 5.19 for a description of the columns.

Table 5.18. Options for the `canvasLibOpt` command. (Continued)

Option	Description
<code>-help</code>	Print usage message and exit.
<code>-help_props</code>	Print the names of the properties stored in <code>fpIn</code> .

Table 5.19. Description of columns in the CSV output file for `canvasLibOpt`.

Column	Description
SMILES	SMILES string for the structure. If <code>fpIn</code> does not contain a SMILES property, a dummy SMILES "C" is written to each row.
Name	Name or ID of the structure.
Row	Row number of the structure in <code>fpIn</code> .
Sim_NN	Tanimoto similarity to nearest neighbor in the optimized library.
Row_NN	Nearest neighbor row number.
Hole_NN	Indicator of whether the nearest neighbor is in the original library (with holes). The value is 1 if it is in the original library, 0 otherwise.
Filter	Number of property filters failed by this structure divided by the total number of property filters. Set to 0 if no filters are used.

allowed values for the property. Failure means that the property value lies outside the specified range. The filter score for a structure is computed as the number of failures divided by the total number of filters, i.e. the fraction of failures. The average filter score is the average of the filter scores for all structures in the subset.

The total score is the average nearest neighbor similarity plus the average filter score. Both of these quantities lie on the interval $[0, 1]$, so the total score lies on the interval $[0, 2]$.

Optimization is done by choosing the requested number of structures (*nstr*) at random, then performing a series of optimization cycles that attempt to improve both the average nearest neighbor similarity and the average filter score.

Each optimization cycle consists of making *nstr* passes through the subset, with an attempt to replace the poorest scoring structure at each pass. The poorest scoring structure is the one that has the highest nearest neighbor similarity. If there are ties in the nearest neighbor similarity, the filter score is used to break the tie. A candidate structure is chosen at random as a possible replacement, and the replacement is done if doing so improves the total score for the subset. The total score is considered to improve if one of the following happens:

- Average nearest neighbor similarity and average filter score decrease.
- Average nearest neighbor similarity decreases and average filter score remains the same.
- Average nearest neighbor similarity remains the same and average filter score decreases.

If the total score does not improve, a Monte Carlo test is used:

```
etest = deltaSim + deltaFilter
```

where `deltaSim` is the increase in average nearest neighbor similarity, and `deltaFilter` is the increase in average filter score, both of which are greater than or equal to 0 if the score does not improve. The replacement is done if the probability, defined by

```
prob = exp(-etest/temperature)
```

is greater than a random number between 0 and 1.

The temperature decreases linearly over the course of the optimization, and the cooling schedule is chosen such that an `etest` value of 0.001 will yield `prob = 0.5` in the first cycle, and `prob = 0.01` in the last cycle. The default `etest` tolerance of 0.001 can be overridden using the option `-mc tol`.

When the optimization has completed the requested number of cycles, or when other criteria are met, a CSV file is written with the optimized subset of structures. A SMILES column is included, so the CSV file can be imported into the Canvas GUI. The SMILES corresponds to the actual structure if the fingerprint file contains a SMILES property, which is always the case if the fingerprint file is exported from the GUI. And in that case, the Canvas UID is also included in the CSV file, so it is possible to read it back into the same project and update on Canvas UID. If there is no SMILES in the fingerprint file a dummy SMILES "C" is written for each structure.

5.4.5 canvasFPHist

This utility generates a histogram of nearest-neighbor similarities, by locating the largest off-diagonal value in each row of a similarity matrix for a single fingerprint file, or the largest value in each row of the similarity between two fingerprint files. The command syntax is as follows:

```
canvasFPHist [options] fpFile histFile
```

The input file `fpFile` must be a Canvas fingerprint file. If the matrix between two sets of fingerprints is required, specify the second input file with `-ifp2`. The output file `histFile` is a CSV file containing the similarity bin boundaries and the frequencies, normalized to sum to 1.

The options are described in [Table 5.20](#).

Table 5.20. Options for the `canvasFPHist` command.

Option	Description
<code>-range range</code>	Specify the range of molecules to use from the input fingerprint file <code>fpFile</code> . <code>range</code> is a range specification, as described on page 119 . Default: use all.
<code>-metric name</code>	Metric type. The allowed values of <code>name</code> are <code>buser</code> , <code>cosine</code> , <code>dice</code> , <code>hamann</code> , <code>kulczynski</code> , <code>matching</code> , <code>mcConnaughey</code> , <code>minmax</code> , <code>modifiedTanimoto</code> , <code>pearson</code> , <code>petke</code> , <code>rogersTanimoto</code> , <code>simpson</code> , <code>tanimoto</code> , <code>tversky</code> , <code>yule</code> . Default: <code>tanimoto</code> .
<code>-osim simFile</code>	Write histogram similarities to the CSV file <code>simFile</code> . The format is SMILES,Name,Sim, where SMILES is a dummy column with the string “C” in each row, Name is the molecule name, and Sim is the maximum similarity. The similarity column name can be set with the <code>-col</code> option.
<code>-col simCol</code>	Similarity column name in <code>simFile</code> . Default: <code>Sim</code> .
<code>-bin dbin</code>	Histogram bin spacing. The default spacing is determined by the number of rows and ranges from 0.5 to 0.05, with smaller spacing for larger numbers of rows.
<code>-ifp2 fpFile2</code>	Specify a second fingerprint input file, for the columns of the similarity matrix. If specified, only similarities between <code>fpFile</code> and <code>fpFile2</code> are calculated.
<code>-range2 range</code>	Specify the range of molecules to use from the second input file, <code>fpFile2</code> . <code>range</code> is a range specification, as described on page 119 . Default: use all.
<code>-block value</code>	Maximum number of row fingerprints to load in memory at a time. Default: 1000.
<code>-alpha value</code>	Tversky alpha parameter. Default: 0.5.
<code>-beta value</code>	Tversky beta parameter. Default: 0.5.

5.4.6 canvasFPMatrix

This utility generates a pairwise similarity or distance matrix using binary or scaled fingerprints from one or two sets of molecules. The command syntax is as follows:

```
canvasFPMatrix -ifp fpFile [-ifp2 fpFile2] {-o binaryFile|-ocsv csvFile}
  [options]
```

The input file specified with `-ifp` must be a Canvas fingerprint file. If the matrix between two sets of fingerprints is required, specify the second input file with `-ifp2`. The output file can be in binary format, if specified with `-o`, or in CSV format, if specified with `-ocsv`. Both formats can be used in subsequent operations. The options are described in [Table 5.21](#). The metrics are listed in [Table 5.22](#), along with their type and formula. The quantities in the formula are defined in [Table 5.23](#), except for the Tversky α and β parameters.

Table 5.21. Options for the `canvasFPMatrix` command.

Option	Description
<code>-filter minSim</code> [-all]	Report only rows with at least one similarity greater than or equal to <i>minSim</i> . Valid only with <code>-ocsv</code> and when a similarity metric is used. If <code>-all</code> is used, report only rows where all similarities are greater than or equal to <i>minSim</i> .
<code>-sort avg best</code>	Sort the rows of the output matrix in order of the largest similarity or smallest distance in each row (<i>best</i>) or the largest average similarity or smallest average distance of the row (<i>avg</i>). When using a square matrix (<code>-ifp2</code> not used), the diagonal elements are ignored.
<code>-caprow N</code>	Write out only the top <i>N</i> rows, after sorting. Requires <code>-sort</code> .
<code>-capcol N</code>	For each row, write out only the top <i>N</i> columns. The output consists of the row name, followed by a list of <i>N</i> column names and the corresponding <i>N</i> metric values.
<code>-metric name</code>	Metric type. <i>name</i> can be either the name or the index. Allowed values are listed in Table 5.22 by index and name. Applies only if a fingerprint file is specified as input. Default: <code>soergel (20)</code> .
<code>-helpMetric</code>	Print the definition of each metric then exit.
<code>-forceBinary</code>	Use binary values and Ignore any scaled fingerprint values in the input files.
<code>-range first:last</code>	Specify range of molecules in the input file. <i>first</i> and <i>last</i> are the indices of the first and last molecules to be included.
<code>-ifp2 fpFile</code>	Specify a second input file. If specified, only pairs between the first and the second files are calculated.
<code>-range2 range</code>	Specify range of molecules in the second input file. <i>first</i> and <i>last</i> are the indices of the first and last molecules to be included.
<code>-blocksize value</code>	Maximum number of fingerprints to load in memory at a time. Default: 5000.
<code>-alpha value</code>	Tversky alpha parameter. Default: 0.5.
<code>-beta value</code>	Tversky beta parameter. Default: 0.5.
<code>-flatten value</code>	Gaussian parameter to make output matrix sparse. Default is no flattening.
<code>-limitOffBits</code>	Limit the set of possible off bits. By default, the number of off bits is limited only by the fingerprint address space size (2^{32} or 2^{64}), which may yield undesirable behavior for metrics that incorporate off bits. If this option is used, the set of possible off bits is limited to those bits that are set by at least one compound in the fingerprints provided. Applies only to <code>buser</code> , <code>hamann</code> , <code>matching</code> , <code>modifiedTanimoto</code> , <code>patternDifference</code> , <code>pearson</code> , <code>rogersTanimoto</code> , <code>size</code> , <code>shape</code> , <code>variance</code> , and <code>yule</code> metrics.

Table 5.22. Available metrics for the `canvasFPMatrix` command.

Index	Name	Type	Formula
1	buser	similarity	$(\sqrt{cd}+c)/(\sqrt{cd}+a+b-c)$
2	cosine	similarity	c/\sqrt{ab}
3	dice	similarity	$2c/(a+b)$
4	dixon	distance	$(A+B)^2/(a+b-c)$
5	euclidean	distance	$\sqrt{A+B}$
6	hamann	similarity	$(c+d-A-B)/N$
7	hamming	distance	$A+B$
8	kulczynski	similarity	$0.5(c/a + c/b)$
9	matching	similarity	$(c+d)/L$
10	mcConnaughey	similarity	$(c^2-(a-c)(b-c))/(ab)$
11	minmax	similarity	$\text{sum}\{\min(a,b)/\max(a,b)\}$
12	modifiedTanimoto	similarity	$\alpha *T1 + (1.0-\alpha)*T0$
13	patternDifference	distance	AB/N^2
14	pearson	similarity	$(cd-AB)/\sqrt{ab(A+d)(B+d)}$
15	petke	similarity	$c/\max(a,b)$
16	rogersTanimoto	similarity	$(c+d)/(2(a+b)-3c+d)$
17	shape	distance	$(A+B)/N - ((A-B)/N)^2$
18	simpson	similarity	$c/\min(a,b)$
19	size	distance	$((A-B)/N)^2$
20	soergel	distance	$((A-B)/N)^2$
21	tanimoto	similarity	$c/(a+b-c)$
22	tversky	similarity	$c/(\alpha(a-c)+\beta(b-c)+c)$
23	variance	distance	$(A+B)/(4N)$
24	yule	similarity	$(c*d - A*B)/(c*d + A*B)$

Table 5.23. Variables used in metric formulae

Variable	Description
a	Number of bits that are on in structure 1
b	Number of bits that are on in structure 2
c	Number of bits that are on in both structure 1 and structure 2
d	Number of bits that are off in both structure 1 and structure 2
A	Number of bits that are on in structure 1 but not in structure 2. $A = a - c$
B	Number of bits that are on in structure 2 but not in structure 1. $B = b - c$
L	Total number of bits. $L = a + b - c + d$
N	Restricted total number of bits. $N = a + b - c + \min(d, 10000)$

5.4.7 canvasHC

This program performs full hierarchical clustering and reports details for a specific number of clusters. The command syntax is as follows:

```
canvasHC -im matrixFile [-n clusters|-ns clusters|-d distance|-kelley]
        [-rowNames] [-colNames] [-og groupFile [-noties]] -or runFile]
        [job-options] [common-options]
canvasHC -ifp FPFile [-saveMatrix matrixFile] [-metric metric]
        [-ob baseFile] [-os statFile [-noties]] [-alpha alpha] [-beta beta]
        [job-options] [common-options]
```

where the common options are:

```
[-usemin|usemax|-useupper|-uselower] [-linkage {name|number}]
        [-helpLinkage] [-ot treeFile]
```

The input file specified by `-im` must be a similarity or distance matrix file in either binary or CSV format. The input matrix is assumed to be symmetric. If it is not symmetric, you should use one of the four `-use` options to specify how to handle the asymmetry. Alternatively, you can specify a fingerprint file with `-ifp` and calculate the similarity or distance matrix. The job options are listed in [Table 5.3](#). The program options are described in [Table 5.24](#).

Table 5.24. Options for the *canvasHC* command.

Option	Description
<i>Common options:</i>	
-helpLinkage	Print details about linkage methods and exit.
-linkage <i>name</i>	Method of calculating distance between clusters. The method can be referred to by name or by index. The allowed values are listed by index and name in Table 5.25 . Default: average (3).
-ot <i>treeFile</i>	Intermediate file required by <i>canvasTreeDraw</i> to generate a dendrogram.
-uselower	Asymmetric distance/similarity matrix. Use the lower half for clustering.
-usemax	Asymmetric distance/similarity matrix. Use the minimum of D_{ij} and D_{ji} for clustering.
-usemin	Asymmetric distance/similarity matrix. Use the minimum of D_{ij} and D_{ji} for clustering.
-useupper	Asymmetric distance/similarity matrix. Use the upper half for clustering.
<i>Options for matrix file input:</i>	
-colNames	Input CSV matrix lists all molecule names on the first line.
-rowNames	Input CSV matrix has the molecule name at the beginning of each line.
-d <i>distance</i>	Print details for the set of clusters formed at or below the specified merging distance. The number of clusters decreases as <i>distance</i> increases.
-kelley	Print details for the set of clusters formed at the number of clusters specified by a minimum in the Kelley cost function [4].
-n <i>clusters</i>	Print details for the formation of the specified number of clusters. Default: 1.
-ns <i>clusters</i>	Print details for the formation of the specified number of non-singleton clusters. Singletons are not counted. Default: 1.
-og <i>groupFile</i>	Print cluster memberships in CSV format to <i>groupFile</i> . Default: print to standard output.
-or <i>runFile</i>	Print final embedding stress value to <i>runFile</i> . Default: print to standard output.
<i>Options for fingerprint file input:</i>	
-alpha <i>alpha</i>	Tversky alpha parameter. Only valid with <i>-metric</i> . Default 0.5.
-beta <i>beta</i>	Tversky beta parameter. Only valid with <i>-metric</i> . Default 0.5.
-metric <i>type</i>	Metric used to calculate similarity or distance matrix from a fingerprint file. Valid values of <i>type</i> are listed in Table 5.22 . The value can be the name or the index. Default: <code>tanimoto</code> (21).

Table 5.24. Options for the `canvasHC` command. (Continued)

Option	Description
<code>-noties</code>	Disregard ties when assigning items nearest and farthest from the cluster centroid.
<code>-os statsFile</code>	Report file designed to help select an appropriate number of clusters.
<code>-saveMatrix file</code>	Save the similarity or distance matrix generated from fingerprints to the specified binary file. Default: do not save.

Table 5.25. Description of linkage methods for hierarchical clustering.

Method	Description
1: single	Closest inter-cluster pair
2: complete	Farthest inter-cluster pair
3: average	Average distance between all inter-cluster pairs
4: centroid	Euclidean distance between cluster centroids
5: mcquitty	Average distance to the two clusters merged in forming a given cluster
6: ward	Sum of squared distances to merged cluster centroid (minimum variance)
7: weightedcentroid	Weighted center of mass distance, also known as median
8: flexiblebeta	Weighted average intra-cluster and inter-cluster distances (Lance-Williams) with $\beta=0.25$
9: schrodinger	Closest distance between terminal (right-to-left) points in 1D cluster orderings.

5.4.8 canvasHCBUILD

This program performs full hierarchical clustering, but provides no details. You can use `canvasHCSelect` to obtain details for a specific level of clustering. The command syntax is:

```
canvasHC -im matrixFile [options]
```

The input file, specified by `-im`, must be a similarity or distance matrix file, and can be in either binary or CSV format. The input matrix is assumed to be symmetric. If it is not symmetric, you should use one of the four `-use` options to specify how to handle the asymmetry. The options are described in [Table 5.26](#).

Table 5.26. Options for the `canvasHCBUILD` command.

Option	Description
<code>-rowNames</code>	Input CSV matrix has the molecule name at the beginning of each line.
<code>-colNames</code>	Input CSV matrix lists all molecule names on the first line.
<code>-linkage <i>name</i></code>	Method of calculating distance between clusters. The method can be referred to by name or by index. The allowed values are listed by index and name in Table 5.25 . Default: average (3).
<code>-helpLinkage</code>	Print details about linkage methods and exit.
<code>-helpStats</code>	Print definitions of statistics reported in <i>statsFile</i> and exit.
<code>-n <i>clusters</i></code>	Print details for the formation of the specified number of clusters. Default: 1.
<code>-noties</code>	Disregard ties when assigning items nearest and farthest from the cluster centroid.
<code>-ns <i>clusters</i></code>	Print details for the formation of the specified number of non-singleton clusters. Singletons are not counted. Default: 1.
<code>-d <i>distance</i></code>	Print details for the set of clusters formed at or below the specified merging distance. The number of clusters decreases as <i>distance</i> increases.
<code>-ob <i>baseFile</i></code>	Intermediate file required by <code>canvasHCSelect</code> to generate cluster memberships.
<code>-os <i>statsFile</i></code>	Report designed to help select an appropriate number of clusters.
<code>-or <i>runFile</i></code>	Print final embedding stress value to <i>runFile</i> . Default: print to standard output.
<code>-ot <i>treeFile</i></code>	Intermediate file required by <code>canvasTreeDraw</code> to generate a dendrogram.
<code>-usemin</code>	Asymmetric input matrix. Use the minimum of D_{ij} and D_{ji} for clustering.
<code>-usemax</code>	Asymmetric input matrix. Use the maximum of D_{ij} and D_{ji} for clustering.
<code>-useupper</code>	Asymmetric input matrix. Use the upper half for clustering.
<code>-uselower</code>	Asymmetric input matrix. Use the lower half for clustering.

5.4.9 canvasHCSelect

This program reports details for a specific level of clustering for a prior run of canvasHCBuild. The command syntax is as follows:

```
canvasHCSelect -ib baseFile -im matrixFile [options]
```

The input base file is the file specified with the `-ob` option to canvasHCBuild, and the input matrix file is likewise the matrix file specified with the `-im` option to canvasHCBuild. The options are described in [Table 5.27](#).

Table 5.27. Options for the canvasHCSelect command.

Option	Description
-rowNames	Input CSV matrix has the molecule name at the beginning of each line.
-colNames	Input CSV matrix lists all molecule names on the first line.
-n <i>clusters</i>	Print details for the formation of the specified number of clusters. Default: 1.
-ns <i>clusters</i>	Print details for the formation of the specified number of non-singleton clusters. Singletons are not counted. Default: 1.
-d <i>distance</i>	Print details for the set of clusters formed at or below the specified merging distance. The number of clusters decreases as <i>distance</i> increases.
-og <i>groupFile</i>	Print cluster memberships in CSV format to <i>groupFile</i> . Default: print to standard output.

5.4.10 canvasTreeDraw

This program draws dendrograms from a clustering output file. The command syntax is:

```
canvasTreeDraw -i treeFile -o psFile [options]
```

The input tree file, *treeFile*, must be generated by running `canvasHC` or `canvasHCBuild` with the `-ot` option. The output file of the dendrogram, specified by `-o`, is a postscript file.

Table 5.28. Options for the `canvasTreeDraw` program.

Option	Description
-c	Color links (with number of groups from 2 to 19).
-C	Color labels (with number of groups from 2 to 19).
-h	Use rainbow colors, light and dark, allowing unlimited number of groups.
-H	Use rainbow colors, light only, allowing unlimited number of groups.
-s <i>lines</i>	Number of lines to skip within groups (can be non-integer).
-S <i>lines</i>	Number of lines to skip between groups (can be non-integer).
-2	Use PostScript level 2. Default: level 1.
-L	Suppress printing of labels.
-n <i>groups</i>	Number of groups to color by. Invalid with <code>-d</code> . Default: 1.
-d <i>distance</i>	Distance used to separate groups. Invalid with <code>-n</code> .
-f <i>fontSize</i>	Font size, in points. Default: 8.

5.4.11 canvasKMeans

This program performs k-means clustering using Lloyd’s algorithm [3]. The algorithm consists of a set of runs, each of which starts with the random selection of k data points (structures) as the initial centroids of the clusters. For each run, a series of iterations is performed in which each of the points in the data set is assigned to a cluster based on the centroid closest to that point, and each centroid is updated by setting it equal to the average position of the points assigned to that cluster. The cost function that is minimized is the sum of intra-cluster variances. A run is terminated when the maximum number of steps is reached, or a convergence condition on the cost is reached. The clusters that are returned are taken from the run with the smallest cost.

The command syntax is as follows:

```
canvasKMeans [job-options] {-ifp fpFile|-icsv csvFile [csv-options]} -o outFile
-k clusters [options]
```

The job options are listed in Table 5.3. The *csv-options* are the input file options for CSV files listed in Table 5.1 (except `-smi`). The remaining arguments are described in Table 5.29.

Table 5.29. Arguments for the `canvasKMeans` command.

Argument	Description
<code>-ifp fpFile</code>	Input fingerprint file.
<code>-icsv csvFile</code>	Input CSV file.
<code>-autoScale</code>	Apply autoscaling to the data in the CSV file. Only valid with <code>-icsv</code> .
<code>-o outFile</code>	Output file. Required.
<code>-k clusters</code>	Number of clusters.
<code>-runs maxRuns</code>	Maximum number of runs. Default: 10.
<code>-runSteps maxRunSteps</code>	Maximum number of steps per run. Default: 20.
<code>-totSteps maxTotSteps</code>	Maximum number of steps over all runs. Default: 200.
<code>-conv dCost</code>	Stop iterating when the relative change in the cost function drops below <i>dCost</i> . Default: 0.001.

5.4.12 canvasLC

This program performs leader-follower clustering. In the leader-follower clustering method, data points are added to a cluster based on a cutoff on their distance from the cluster “leader”. The first leader is chosen by some method, then for each data point, the distance from the leader is evaluated. If it is within the cutoff, it is added to the cluster as a “follower”. If it is not within the cutoff, a new cluster is created, with this data point as a new leader. Subsequent points are tested against each leader, in order.

This method depends on the choice of the first leader, and on the order of the data points (structures). To eliminate the order dependence, the list can be sorted by the fingerprint bit count. The leaders are then chosen in order of the number of bits set. In the Canvas implementation, the first leader is the first structure in the set that you select for clustering or the set after sorting.

```
canvasLC [job-options] {-ifp fpFile|-icsv csvFile [csv-options]} -o outFile
      [options]
```

The job options are listed in Table 5.3. The *csv-options* are the input file options for CSV files listed in Table 5.1 (except -smi). The arguments for this program are described in Table 5.30.

Table 5.30. Arguments for the *canvasLC* command.

Argument	Description
-autoScale	Apply autoscaling to the data in the CSV file. Only valid with -icsv.
-dist <i>cutoff</i>	Distance cutoff for clustering. Default: 0.5.
-group	Group items in the same cluster together. The first item within a cluster is the leader.
-helpMetric	Print types of metrics available with different options and exit.
-icsv <i>csvFile</i>	Input CSV file.
-ifp <i>fpFile</i>	Input fingerprint file.
-metric <i>metric</i>	Metric used to calculate similarity/distance between two molecules. Default is tanimoto if -ifp is used or euclidean if -icsv is used. Allowed values are listed in Table 5.22 by index and name, either of which may be used for <i>metric</i> .
-ocsv <i>outFile</i>	Output CSV file. By default, it contains 3 columns: ID, Cluster, and Leader. ID is the name or row index; the latter is used when the input CSV file has no name column (see -name). Cluster is the index of the cluster into which an observation (input row) falls. Leader is a 0/1 value that marks the leader of each cluster.

Table 5.30. Arguments for the `canvasLC` command. (Continued)

Argument	Description
<code>-odata</code>	Include property data in the output CSV file. For fingerprint input, the non-fingerprint data from <i>fpFile</i> is copied; for CSV input, all fields from <i>csvFile</i> are copied.
<code>-rep {sim leader}</code>	Print out only the representatives in each cluster. If <code>sim</code> is used, the representative has the minimum average distance from other members in the cluster. If <code>leader</code> is used, the representative is the leader in each cluster.
<code>-s</code>	Sort the fingerprints in descending order of the number of “on” bits. Supported metric types with sorted fingerprints are <code>tanimoto</code> , <code>cosine</code> , <code>dice</code> , <code>kulczynski</code> , <code>mcConnaughey</code> , <code>petke</code> , and <code>soergel</code> .

5.4.13 canvasPW2CSV

This program generates a CSV matrix file from a binary pairwise similarity or distance matrix. The command syntax is as follows:

```
canvasPW2CSV -i binaryFile -o csvFile [-norow] [-nocol]
```

The arguments are described in [Table 5.31](#).

Table 5.31. Arguments for the `canvasPW2CSV` command.

Argument	Description
<code>-i <i>binaryFile</i></code>	Input matrix file in binary format.
<code>-o <i>csvFile</i></code>	Output destination file containing CSV data.
<code>-norow</code>	Omit header row that lists all molecule names in output.
<code>-nocol</code>	Omit molecule name as first item of each line in output.

5.5 Model Building and Related Applications

The programs available for statistics, 2D QSAR, neural networks and machine learning are listed in [Table 5.32](#), and are described in the following sections.

Table 5.32. Programs for statistics and 2D QSAR.

Program	Description
canvasMDS	Multidimensional scaling.
canvasMLR	Build and test multiple linear regression models.
canvasMolDescriptors	Calculate molecular descriptors
canvasPCA	Direct principal component generation without intermediate analysis
canvasPCAGen	Principal component generation
canvasPCAProj	Project data along one or more principal components generated by canvasPCAGen
canvasPCAReg	Build and test principal component analysis regression models.
canvasPLS	Build and test partial least square regression models.
canvasKPLS	Build and test kernel-based partial least square regression models.
canvasBayes	Build and test a Bayes model from binary or continuous training data
canvasNnet	Build and test an ensemble model of neural networks
canvasRP	Build and test a recursive partitioning model
canvasSOM	Generate Kohonen self organizing map
canvasSOMBits	

The programs that build models have many options in common. The common input and output arguments are listed in syntax statements as [*io-args*], and have the following syntax:

```
-in inFile [-d delim] [-out outFile] [-plot plotFile]
```

These arguments are described in [Table 5.33](#).

The common build options specify the dependent and independent variable columns, the training and test sets, and the output file containing the model. They are listed in the syntax statements as [*build-options*], and have the following syntax:

```
-y yvar [-omod modelFile]
[-lt trainList | -pt trainFraction [-rand seed]]
[-lx xvarList | -fieldIn list | -fieldOut list]
```

Table 5.33. Common input and output arguments for model-building programs.

Argument	Description
-in <i>inFile</i>	Required. CSV file (delimited text file) containing independent (<i>x</i>) variables and, if building a model, a dependent (<i>y</i>) variable. The first line must contain a unique name for each column, and each name must start with a letter of the alphabet (A-Z or a-z).
-d <i>delim</i>	Optional. Delimiter used to separate values in <i>inFile</i> . Default: comma.
-out <i>outfile</i>	Optional. File for program output. Default: standard output.
-plot <i>plotFile</i>	Optional. CSV file to which observed and calculated <i>y</i> values are written. If using <code>-build</code> with both the training and test sets, the training set rows are written first.

These options are described in [Table 5.34](#).

Table 5.34. Common build options.

Option	Description
-y <i>yvar</i>	Required. Name or index of the dependent variable column. For <code>canvasNnet</code> , multiple dependent variables can be specified for <i>yvar</i> , delimited by spaces. The first column in <i>inFile</i> has an index of 1.
-lt <i>trainList</i>	Training set observations. The second line in the input file is observation 1. <i>trainList</i> is a range specification, as described on page 119 . Observations not in <i>trainList</i> are assigned to the test set. By default, all observations are included in training set.
-pt <i>trainFract</i>	Randomly assign the specified fraction of observations to the training set, with the remainder assigned to the test set. <i>trainFract</i> is a real number between 0 and 1.
-rand <i>seed</i>	Random seed integer for selecting the training set. If omitted, the seed is assigned from the current local time.
-omod <i>modelFile</i>	Write model to a file. This file must be created if the model is to be used later.
-lx <i>xvarList</i>	Independent variable column indices. <i>xvarList</i> is a range specification, as described on page 119 . The first column in the input file is 1. Not available with <code>canvasBayes</code> .
-fieldIn <i>list</i>	Space-delimited list of field names (columns) to include in the model.
-fieldOut <i>list</i>	Space-delimited list of field names to exclude from the model. By default, all fields except the <i>y</i> variable are included when building the model. However, if input data includes non-numerical columns, use one of the above 3 options to specify the numerical columns to include.

The test options are listed in the syntax statements as [*test-options*], and have the following syntax:

```
[-y yvar|-unknown] -imod modelFile
```

These options are described in [Table 5.35](#).

Table 5.35. Test options for regression models

Option	Description
-y <i>yvar</i>	Name or index of the dependent variable column. For <code>canvasNnet</code> , multiple dependent variables can be specified for <i>yvar</i> , delimited by spaces.
-unknown	Dependent variable data are not known.
-imod <i>modelFile</i>	File containing previously created model.

5.5.1 canvasMDS

This program preforms multi-dimensional scaling (metric scaling based on principal coordinate analysis) on the output from `canvasCSVMatrix` or `canvasFPMatrix`, or an external distance matrix in CSV format.

```
canvasMDS {-i dmFile|-icsv dmFile} -o outputFile [options]
```

The input file can be in binary format (-i) or CSV format (-icsv); the output file (-o) is in CSV format.

Table 5.36. Options for the `canvasMDS` command.

Option	Description
-n <i>dim</i>	Number of dimensions to use. Default: 2.
-eigval <i>neig</i>	Print out <i>neig</i> eigenvalues from the principal coordinate analysis.
-d <i>outputDelim</i>	Output delimiter. Default: comma.
-di <i>inputDelim</i>	Delimiter used in the input file if it is in CSV format. Default: comma.
-nameCol <i>col</i>	Column that contains the molecule names in a CSV input file. If the names are in the last column, use <code>last</code> for <i>col</i> . Default: 1.
-noHeader	The input CSV file does not have a header row.

5.5.2 canvasMLR

This program builds and tests multiple linear regression models using data supplied in a CSV (delimited) file. In the build mode, a model is developed on a training set and optionally applied to a test set contained in the same file as the training set. In the test mode, an existing model is applied to a test set contained in its own file. The command syntax is as follows:

```
canvasMLR [job-options] input-args { -build build-options model-options |
  -test test-options }
```

The job options are listed in [Table 5.3](#). The input arguments are described in [Table 5.33](#), the common build options are described in [Table 5.34](#), and the test options are described in [Table 5.35](#). The syntax of the model options is as follows:

```
-single [-noIntercept] |
-best -nx numXvar [-steps n] [-start t1] [-stop t2]
  [-ensemble [-keep m] [-weight]]
```

The model options are described in [Table 5.37](#).

Table 5.37. Build options for the canvasMLR command.

Option	Description
-single	Build a single MLR model from a list of independent variables. One of the options -lx, -fieldIn and -fieldOut is required with this model.
-noIntercept	Suppress y intercept term in the regression. Forces the regression lines to pass through the origin.
-best	Attempt to identify the best model containing a given number of independent variables, using a simulated annealing Monte Carlo technique.
-nx numXvar	Number of independent variables to include in model. Required with -best.
-steps n	Number of Monte Carlo steps. Default: 1000.
-start t1	Initial temperature factor. The actual temperature is the product of t1 and the standard deviation in y. Default: 0.5.
-stop t2	Final temperature factor. Default: 0.05.
-ensemble	Build an algebraically averaged ensemble model from among those with the lowest standard deviation of regression.
-keep m	Number of models to include in ensemble. Default: 5.
-weight	Weight each model by its R-squared value. Default: do not weight.

5.5.3 canvasMolDescriptors

This program calculates molecular descriptors from a set of structures.

```
canvasMolDescriptors -ifmt infile -ofmt outFile inputFileArgs [job-options]
                    [options]
```

The input file is specified with the arguments listed in Table 5.1. The job options are standard Job Control options, listed in Table 5.3. This program can be distributed over multiple processors. The options are described in Table 5.38, and include the output file specifications; the -v3 option in Table 5.2 for SD files is also supported. You can specify as many of the descriptors as you wish, or use -All to calculate all descriptors.

Table 5.38. Options for the canvasMolDescriptors command.

Option	Description
-AlogP	Calculate atomic logP [9].
-method last sum mean	AlogP assignment method. Default is last.
-Custom <i>customFile</i>	Calculate a custom property according to the atom types and values supplied in <i>customFile</i> . The property is the sum of atom values for all atoms in a molecule.
-CustomAvg	Calculate the average atom value for each custom atom type in a molecule.
-CustomCnt	Calculate the count of each custom atom type in a molecule.
-CustomKey	Assign a value for the presence (1) or absence (0) of each custom atom type in a molecule.
-CustomName <i>customPropName</i>	Specify the name of the custom property.
-CustomSum	Calculate the sum of atom values for each custom atom type in a molecule.
-Estate	Calculate electrotopological states [8].
-EstateAvg	Calculate the average of each Estate atom type (sum/count) based on -maxPath and -pow values.
-EstateCnt	Calculate the count of each Estate atom type in a molecule.
-EstateKey	Indicate the presence (1) or absence (0) of each Estate atom type in a molecule.
-EstateSum	Calculate the sum of values for each Estate atom type in a molecule according to the -maxPath and -pow values.

Table 5.38. Options for the `canvasMolDescriptors` command.

Option	Description
<code>-maxPath <i>n</i></code>	Consider only neighboring atoms within <i>n</i> bonds when calculating electrotopological states. Valid only with <code>-Estate</code> or <code>-All</code> . Default: all neighboring atoms are taken into account.
<code>-pow <i>fallOffPower</i></code>	Set how fast the perturbation falls off for electrotopological state calculation. Valid only with <code>-Estate</code> or <code>-All</code> . Default: 2.
<code>-HBA</code>	Count hydrogen bond acceptors.
<code>-HBAfile <i>filename</i></code>	File of SMARTS patterns to overwrite Canvas default definition of hydrogen bond acceptors.
<code>-HBD</code>	Count hydrogen bond donors.
<code>-HBDfile <i>filename</i></code>	File of SMARTS patterns to overwrite Canvas default definition of hydrogen bond donors.
<code>-distinct</code>	Count each hydrogen as a distinct donor. By default a given atom that has one or more donatable hydrogens is counted as one donor.
<code>-MR</code>	Calculate molar refractivity [10].
<code>-MW</code>	Calculate molecular weight.
<code>-Polar</code>	Calculate Miller polarizability [11].
<code>-PSA</code>	Calculate polar surface area.
<code>-RB</code>	Calculate the number of rotatable bonds.
<code>-RBfile <i>filename</i></code>	File of SMARTS patterns to modify or overwrite Canvas default definition of rotatable bonds. See the text below for examples of SMARTS patterns.
<code>-use_ligparse_def</code>	Use the same set of rules for rotatable bonds as in <code>ligparse</code> .
<code>-All</code>	Calculate all descriptors (default).
<code>-helpHBD</code>	Print out Canvas default definition of hydrogen bond donors. You can use this output as a template for your own definitions.
<code>-helpHBA</code>	Print out Canvas default definition of hydrogen bond acceptors. You can use this output as a template for your own definitions.
<code>-helpRB</code>	Print out a rotatable bond template file (with no SMARTS patterns).
<code>-smiles</code>	Include SMILES string as the first column for each molecule in the CSV output file.
<code>-fill</code>	Fill in the line for a molecule that fails to generate descriptors. The line contains no values except the molecule name. This option is useful to create placeholders that preserve positional alignment to external data. Default: skip failed molecules in the output.

Table 5.38. Options for the `canvasMolDescriptors` command.

Option	Description
<code>-n structureRange</code>	Input structures to process. <i>structureRange</i> is a range specification, as defined on page 119 . Default: process all structures.
<code>-file</code>	Interpret <i>structureRange</i> as a file name. Each line in this file should contain a valid row range specification. Not available with <code>-iproj</code> .
<code>-ofmt outputFile</code>	Write results to <i>outputFile</i> in the specified format. Valid formats are <code>mae</code> , <code>sd</code> , and <code>csv</code> . Default: CSV format. If <i>outputFile</i> is omitted with <code>-ocsv</code> , standard output is used. Ignored if <code>-JOB</code> is used: the filenames are assigned as <i>jobname_desc.ext</i> , where <i>ext</i> is <code>csv</code> , <code>maegz</code> , or <code>sdf.gz</code> .
<code>-odata</code>	Copy all data fields in the input SD or Maestro files to the output file. Only available with <code>-isd</code> or <code>-imae</code> .
<code>-uniform</code>	All the molecules in the Maestro or SD input file have the same data fields.
<code>-fieldOnly fields</code>	Only the specified fields (space-separated list of property names) in the SD or Maestro input file are saved to the output file.
<code>-obad badMolFile</code>	Save the molecules that failed to generate descriptors to a file. By default, the list of failed molecules is written to standard output. This option is only available for SMILES, SD or CSV input.

You can create your own SMARTS patterns to define hydrogen-bond donors, acceptors, and rotatable bonds, and supply them in template files to `canvasMolDescriptors`. The format of the template files can be obtained with the `-helpHBD`, `-helpHBA`, and `-helpRB` options. For example, to create a template file for rotatable bonds, use the following command:

```
canvasMolDescriptors -helpRB > myRBtemplate.txt
```

The default definitions are included in the template file for hydrogen-bond donors and acceptors, but not for rotatable bonds. Examples for rotatable bonds are described below.

By default, a rotatable bond in Canvas is defined as a single, non-ring bond, bonded to a non-terminal heavy atom (non-hydrogen). Amide bonds (C–N) and a bond next to a triple bond are excluded because of the high energy barrier.

To include or exclude a user-defined bond type, you must specify both atoms connected to the bond with valid SMARTS. A positive integer value (>0) must follow each SMARTS pattern, separated by a space. Within each file, this integer value must be unique. Only heavy atoms may be specified, and only one SMARTS may be given per line. A line starting with `;` is treated as comment.

For example, the following two lines specifies a secondary amide bond:

```
C(=O)N[!#1] 1
N([!#1])C=O 2
```

If, in addition, you wants to include bonds between OH and SP3 carbon, you can include the following two lines in the “include” section of the template file:

```
C(*) (*) (*) - [OH] 3
[OH] - C(*) (*) * 4
```

If you want to exclude single bonds connecting a carbon attached to a halogen, you can copy the following two lines to the “exclude” section of the file:

```
*-C[I,Br,Cl,F] 5
C([I,Br,Cl,F])-* 6
```

Other SMARTS patterns may be considered in counting rotatable bonds such as acid groups:

```
[OH] - C=O 100
C(=O) - [OH] 200
```

Custom atom types (-Custom) are defined in a file that contain rules for assigning atom types. A property assignment rule has the following format:

```
SMARTS ? prop1 prop2 ...
```

This rule assigns the property *prop1* to the first atom in the SMARTS pattern, *prop2* to the second atom, and so on. *prop1* is a name that can consist of alphanumeric characters (upper and lower case) and underscores, and must start with a letter. The property can be used instead of SMARTS patterns in a rule, as *\$propN*. A type assignment rule has the format:

```
SMARTS > type ; {name} value
```

This rule assigns the atom type *type* to the first atom in the SMARTS pattern. *type* is a positive integer, which indexes the atom types. The name of the atom type can be set by adding the optional *{name}*. The name is used when reporting the per-atom-type properties (Key, Cnt, Sum, Avg). The optional *value* is used to calculate the custom property.

Rules are applied in the order in which they are encountered in the file, so any property assignments must precede their use. The default atom type is 0. Atom types set by one rule can be changed by a subsequent rule.

5.5.4 canvasPCA

This program performs direct principal component generation without intermediate analysis. The command syntax is as follows:

```
canvasPCA [job-options] {-icsv csvFile [csv-options] |-ifp binaryFile}
          -o outputFile -ostat statFile
```

The job options are listed in [Table 5.3](#). The input file can be in CSV format (`-icsv`) or binary fingerprint format (`-ifp`). The output file (`-o`) is in CSV format, and contains the input data projected along all principal components generated. The options are described in [Table 5.39](#).

5.5.5 canvasPCAGen

This program generates and writes out principal components. The command syntax is as follows:

```
canvasPCAGen {-icsv csvFile |-ifp binaryFPFile} -o outputFile [options]
```

The input file can be in CSV format (`-icsv`) or binary fingerprint format (`-ifp`). The output file (`-o`) is in comma-delimited CSV format. The options are the same as for `canvasPCA`, described in [Table 5.39](#), except that `-ostat` is not recognized and `-d` applies only to an input CSV file.

Table 5.39. Options for the canvasPCA and canvasPCAGen commands.

Option	Description
<code>-ostat statFile</code>	Output file with variances and loadings. Not valid for <code>canvasPCAGen</code> .
<code>-n components</code>	Number of principal components to use. Default: 2.
<code>-d delimiter</code>	Input file delimiter if <code>-icsv</code> is used; also output file delimiter for <code>canvasPCA</code> . Default: comma.
<code>-noHeader</code>	The input data has no header row.
<code>-nameCol col</code>	Column index of the names. If names are in the last column, <code>last</code> can be used as the value for <code>col</code> . Default: 1.
<code>-scale</code>	Scales each dimension to unit variance if it is within the threshold.
<code>-t threshold</code>	Threshold for scaling. Only valid with <code>-scale</code> . Default: 0.001.
<code>-mostSig nbits</code>	Keep only the <code>nbits</code> most informative bits across the chosen input set. Default: use all bits. Only valid with <code>-ifp</code> .

5.5.6 canvasPCAProj

This program projects data along one or more principal components generated by canvasPCAGen. The command syntax is as follows:

```
canvasPCAProj -ipc pcFile -idata dataFile -o outputFile [options]
```

The arguments are described in [Table 5.40](#).

Table 5.40. Arguments for the canvasPCAProj command.

Argument	Description
-ipc <i>pcFile</i>	Input data containing principal components generated by canvasPCAGen.
-idata <i>dataFile</i>	Original data file (input to canvasPCAGen) in CSV format.
-o <i>outputFile</i>	Output file for projections, in CSV format.
-load	Project by columns instead of rows to produce PCA loadings.
-n <i>components</i>	Number of principal components (dimensions) desired. Default: 2.
-d <i>delimiter</i>	Delimiter to use in the output file. Default: comma.
-ddata <i>delimiter</i>	Delimiter used in the original data file. Default: comma.
-noHeader	No header row is in the original data file.
-nameCol <i>col</i>	Column index of the names in the original data file. Default: 1.

5.5.7 canvasPCAReg

This program builds and tests principal component analysis regression models using data supplied in a CSV file. In build mode, a model is developed on a training set and optionally applied to a test set contained in the same file as the training set. In test mode, an existing model is applied to a test set contained in its own file. The command syntax is as follows:

```
canvasPCAReg [job-options] io-args {-build build-options model-options |  
-test test-options}
```

The job options are listed in [Table 5.3](#). The input arguments are described in [Table 5.33](#), the common build options are described in [Table 5.34](#), and the test options are described in [Table 5.35](#). The syntax of the model options is as follows:

```
-maxf maxFactors [-autoScaleOff]
```

The model options are described in [Table 5.41](#).

Table 5.41. Model options for the canvasPCAReg command.

Option	Description
-maxf <i>maxFactors</i>	Maximum number of PCA factors to use in building the model. The model contains all factors from 1 to <i>maxFactors</i> .
-autoScaleOff	Do not use auto-scaling when creating the model. Default: use autoscaling.

5.5.8 canvasPLS

This program builds and tests partial least squares regression models using data supplied in a CSV (delimited) file. In build mode, a model is developed on a training set and optionally applied to a test set contained in the same file as the training set. In test mode, an existing model is applied to a test set contained in its own file.

```
canvasPLS [job-options] io-args {-build build-options model-options |  
  -test test-options}
```

The job options are listed in [Table 5.3](#). The input arguments are described in [Table 5.33](#), the common build options are described in [Table 5.34](#), and the test options are described in [Table 5.35](#). The syntax of the model options is as follows:

```
-maxf maxFactors [-autoScaleOff] [-sd sdLimit] [-tmin minTvalue]
```

The model options are described in [Table 5.42](#).

Table 5.42. Model options for the canvasPLS command.

Option	Description
<i>-maxf maxFactors</i>	Maximum number of PLS factors to use in building the model. The model contains all factors from 1 to <i>maxFactors</i> .
<i>-autoScaleOff</i>	Do not use auto-scaling when creating the model. Default: use autoscaling.
<i>-sd sdLimit</i>	If the standard deviation of the regression is less than or equal to <i>sdLimit</i> , stop adding PLS factors during model building. Overrides <i>maxFactors</i> . Default: -1.0.
<i>-tmin minTvalue</i>	Minimum T-value for selecting significant <i>x</i> variables. Must be positive.

5.5.9 canvaskPLS

This program builds and tests kernel-based partial least squares regression models using data supplied in a CSV (delimited) file. In build mode, a model is developed on a training set and optionally applied to a test set contained in the same file as the training set. In test mode, an existing model is applied to a test set contained in its own file.

```
canvasPLS [job-options] io-args {-build build-options model-options |  
  -test test-options}
```

The job options are listed in [Table 5.3](#). The input arguments are described in [Table 5.33](#), the common build options are described in [Table 5.34](#), and the test options are described in [Table 5.35](#). The syntax of the model options is as follows:

```
-maxf maxFactors [-sigma sigma]
```

The model options are described in [Table 5.43](#).

Table 5.43. Model options for the canvaskPLS command.

Option	Description
<i>-maxf maxFactors</i>	Maximum number of KPLS factors to use in building the model. The model contains all factors from 1 to <i>maxFactors</i> .
<i>-sigma sigma</i>	Set the linearity parameter. Larger values mean more linear. Allowed range is from 1 to 100. Default: 30.0.

5.5.10 canvasBayes

This program builds and tests a Bayes model from binary or continuous training data to predict the probability of a molecule at each activity level [14]. The syntax of the command is as follows:

```
canvasBayes [job-options] { -in CSVFile [-d delim] [-prefix string]
  [-name nameCol] | -infp FPBinaryFile }
  [-out outFile] [-plot plotFile [-bin]]
  { -build build-options model-options | -test test-options }
```

The job options are listed in Table 5.3. The input arguments are described in Table 5.44, the common build options are described in Table 5.34, and the test options are described in Table 5.35. The syntax of the model options is:

```
[-category] [-c cutoffsList] [-ebin n] [-s coefficient]
[-KL cutoff] [-KLPos cutoff] [-KLNeg cutoff] [-F cutoff] [-binary] [-nobinary]
[-noise value]
```

The model options are described in Table 5.45.

Table 5.44. Input and output arguments for the canvasBayes command.

Argument	Description
-in <i>CSVFile</i>	CSV file containing independent (<i>x</i>) variables and, if building a model, a dependent (<i>y</i>) variable. The first line must contain a unique name for each column, and each name must start with a letter of the alphabet (A-Z or a-z).
-d <i>delim</i>	Optional. Delimiter used to separate values in <i>inFile</i> . Default: comma.
-prefix <i>string</i>	Prefix for binary columns in the header of the CSV input file. The default prefix of a binary column is BIT.
-name <i>nameCol</i>	Name column in the CSV file. Default: first column.
-infp <i>FPFile</i>	Binary file generated by canvasFPGen, containing fingerprint bits and optionally other molecular property data as independent variables and, if building a model, a dependent variable.
-out <i>outfile</i>	File for program output. Default: standard output.
-plot <i>plotFile</i>	CSV file to which the observed and calculated categories that the dependent variable belongs to are written. If using -build with both the training and test sets, the training set rows are written first.
-bin	If specified, only the indexes, instead of the full range, of the observed and calculated categories are written to <i>plotFile</i> .

Table 5.45. Model options for the `canvasBayes` command.

Option	Description
-category	Chosen activity field is categorical, with no implicit scale or order.
-c <i>cutoffsList</i>	Cutoffs for the numerical activities. Default: use all distinct values.
-ebin <i>n</i>	Divide the training data into <i>n</i> equal sized bins based on the values of dependent variable. <i>n</i> must be at least 2.
-s <i>coefficient</i>	Smoothing coefficient, typically (0-1). Default: 1e-07.
-KL <i>cutoff</i>	Kullback-Leibler distance cutoff used during training. Binary features with a significance less than or equal to <i>cutoff</i> are excluded from the model. Applied to both positively and negatively correlated features. Default: 0.1.
-KLPos <i>cutoff</i>	Kullback-Leibler distance cutoff used during training. Applies to positively correlated features only. Can be used with -KLNeg.
-KLNeg <i>cutoff</i>	Kullback-Leibler distance cutoff used during training. Applies to negatively correlated features only. Can be used with -KLpos.
-F <i>cutoff</i>	Fraction of bits to keep during training, as ranked by Kullback-Leibler significance. Remaining bits are excluded. Default: 1.
-binary	Use only binary data to build a model.
-nobinary	Do not use binary data to build a model.
-noise <i>value</i>	Add random noise with standard deviation of <i>value</i> to all activities. Categorical fields are scrambled with a probability provided by <i>value</i> . New values are chosen by their relative occurrences.

5.5.11 canvasNnet

This program builds and tests an ensemble model of neural networks with data supplied in a CSV (delimited) file. The network has three layers, one input layer, one output layer, and one hidden layer. Networks are trained using a BFGS algorithm. In build mode, a model is developed from a training set, with built-in cross-validation. Optionally, this model can be applied to a test set in the same file as the training set. In test mode an existing model is applied to a test set in its own file. The command syntax is as follows:

```
canvasNnet [job-options] io-args {-build build-options model-options|-test test-options}
```

The job options are listed in [Table 5.3](#). The input and output arguments are described in [Table 5.33](#), the common build options are described in [Table 5.34](#), the model options are described in [Table 5.46](#), and the test options are described in [Table 5.35](#).

Table 5.46. Model options for the canvasNnet command.

Option	Description
-cvp <i>cvPercent</i>	Percentage of the training set used for cross-validation (randomly selected). Default value is 0.1 (10%).
-cycle <i>nCycle</i>	Training cycles for each network. Default is 200 cycles.
-nnet <i>nNetwork</i>	The number of networks to train. Default is 20.
-ensemble <i>n</i>	The number of best networks to form the ensemble model. Default is 5.

5.5.12 canvasRP

This program builds and tests recursive partitioning trees with data supplied in a CSV file. The command syntax is as follows:

```
canvasRP inFile [-d delimiter] [outFile] [job-options] {-build build-options
  model-options}|-test test-options}
```

The input file is a CSV file that contains x and y variables for the training and test sets. The default delimiter is a comma, but you can specify it with the `-d` option. Use `-d " "` for space and `-d "\\t"` for tab. Consecutive spaces are treated as a single delimiter. The output file name is optional: the default is standard output.

The common build options are described in [Table 5.34](#), the model options are described in [Table 5.47](#), and the test options are described in [Table 5.48](#). The job options are described in [Table 5.3](#).

The syntax of the model options is:

```
[-ensemble n [-sample fraction] [-pavg] [-hist]] [-ut] [-r minCorr]
[[-randattr perc] [-uniquattr] ] | [-diffattr]
[-split gini|gain] [-avgsplit ] [-leaf minObs] [-tree]
[-category] | [-c <cutoffsList>] [-nocheck]
```

Table 5.47. Model options for the `canvasRP` command.

Option	Description
<code>-exclude</code>	Treat <i>list</i> in the <code>-lx</code> parameter as the variables to exclude.
<code>-ensemble n</code>	Build an ensemble model of n trees, where each tree is created from a random sample of the training set.
<code>-sample fraction</code>	Training set sample fraction. Default: 0.8.
<code>-pavg</code>	Assign the class of a given observation based on the average probabilities from all leaf nodes to which it is assigned. By default, the probabilities are assigned from the number of votes for each class.
<code>-hist</code>	Show histogram of multiple recursive partitioning tree votes, including the breakdown for each category (TruePositives, FalsePositives, TrueNegatives, FalseNegatives)
<code>-ut</code>	Build each tree with a different root attribute.
<code>-r minCorr</code>	Minimum $x:y$ Pearson correlation coefficient evaluated over the training set. Any x variable not meeting this threshold is ignored.
<code>-randattr perc</code>	For each tree, select attributes from a random subset with the specified percentage of the total attributes.

Table 5.47. Model options for the *canvasRP* command. (Continued)

Option	Description
-uniqattr	If the attributes are selected randomly (-randattr), this argument assures that trees do not share possible attributes. If <i>perc</i> is larger than 1/(number of trees), then this parameter defaults to -diffattr.
-diffattr	Specify that an equal number of attributes (1/(number of trees)) should be selected for each tree.
-split gini gain	Splitting quality measure: gini Gini impurity (default) gain Information gain
-avgsplit	If a range of values yields equivalent splitting quality, use the average.
-leaf <i>minObs</i>	Minimum number of observations in any leaf node. Default: 5.
-tree	Write out trees in human-readable format.
-category	Flag to indicate chosen activity field is categorical, with no implicit scale or order. (Default)
-c <i>cutoffsList</i>	Cutoffs for the numerical activities. Default: use all distinct values.
-nocheck	Omit checks on reasonable assortments among categories.

The syntax of the test options is:

-imod *modelFile* [-unknown] [-hist] [-accdetail]

Table 5.48. Test options for the *canvasRP* command.

Option	Description
-imod <i>modelFile</i>	Input model file. The variable names in this file must be present in <i>inFile</i> .
-unknown	Dependent y variable is not present in <i>inFile</i> .
-hist	Show histogram of multiple recursive partitioning tree votes, including the breakdown for each category (TruePositives, FalsePositives, TrueNegatives, FalseNegatives)
-accdetail	Show accuracy of each tree on test set.

5.5.13 canvassOM and canvassOMBits

These two programs create a Kohonen self-organizing map (SOM) from data in a CSV file or a binary fingerprint file. `canvassOM` creates the map from scaled (real) values, whereas `canvassOMBits` creates the map from binary fingerprint data. Both programs have the same syntax, which is as follows:

```
canvassOM|canvassOMBits [job-options] {-icsv csvFile|-ifp fpFile}  
-outmap outputMap [options]
```

The job options are listed in [Table 5.3](#). The input can come from a CSV file (`-icsv`) or a binary fingerprint file (`-ifp`). For `canvassOM`, the data should be generated with `canvasFPGen` using the `-scaling` option. For both programs, it is recommended to use `-min` and `-max` with `canvasFPGen` to limit the number of bits, though `canvassOMBits` can handle much larger numbers of bits than `canvassOM`. The output map is written to the file specified with `-outmap`. The arguments are described in [Table 5.49](#).

Table 5.49. Arguments for the `canvasSOM` command.

Option	Description
<code>-inmap <i>inputMap</i></code>	Existing SOM map on which to place new data. Overrides <code>-topology</code> , <code>-xdim</code> , and <code>-ydim</code> . can optionally lock or tether any/all cells.
<code>-outdata <i>file</i></code>	Name of file where cell membership info is stored. Default: stdout.
<code>-outrun <i>file</i></code>	Name of file where run information is written. Default: stdout.
<code>-outstats <i>file</i></code>	name of file where class statistics is stored.
<code>-cycles <i>N</i></code>	Maximum number of training cycles. Default: 100.
<code>-update <i>N</i></code>	Status update interval, in cycles. Default: 1.
<code>-topology <i>type</i></code>	Topology. Allowed values are <code>rect</code> (rectangular) and <code>hexa</code> (hexagonal). Overridden if an existing map is read with <code>-inmap</code> . Default: <code>rect</code> .
<code>-xdim <i>N</i></code> <code>-ydim <i>N</i></code>	Lattice 2D dimensions. Default: 10 x 10. Only valid with <code>-topology</code> .
<code>-hcube <i>bits</i></code>	Use hypercube lattice with 2^{bits} cells. Not valid with <code>-topology</code> , and implies <code>-wrap</code> .
<code>-kmeans <i>k</i></code>	Use kmeans lattice with <i>k</i> cells. <i>k</i> should be an integer no smaller than 2. Not valid with <code>-topology</code> .
<code>-wrap</code>	Make lattice topology periodic. Default: nonperiodic.
<code>-noshuffle</code>	Do not randomize compound order. Default: randomize order.
<code>-stop</code>	Terminate training whenever the rms quantization error rises. Default: continue training.
<code>-stopWhen <i>rmsqe</i></code>	Terminate training if the specified rms quantization error is reached.
<code>-decay <i>type</i></code>	Training decay. Allowed values: <code>linear</code> , <code>power</code> , <code>inverse</code> , <code>expo</code> (exponential). Default: <code>expo</code> .
<code>-seed <i>int</i></code>	Use the specified value to seed random number generator.
<code>-idcol <i>N</i></code>	Zero-based index of identifier column. Default: 0.
<code>-catcol <i>N</i></code>	Zero-based index of category column. Default: none.
<code>-blocksize <i>N</i></code>	Number of records to load into memory at a time. Default: 1000.
<code>-fast</code>	Run faster by updating weights selectively. Recommended for large data sets.
<code>-round</code>	Force SOM weights to 0 or 1 after each training cycle. Valid only for <code>canvasSOMBits</code> .
<code>-maxdist <i>float</i></code>	Assign no cell for items whose nearest cell exceeds this distance. Default: no limit.

5.6 Utilities

In addition to the programs described above, various utilities are available for job management, format conversion, and various other tasks. These utilities are listed in [Table 5.50](#).

Table 5.50. Canvas utilities.

Utility	Description
canvas_app	Run jobs that are set up by canvasJob.
canvasConvert	Convert between molecular file formats.
canvasJob	Set up and clean up jobs associated with a Canvas project.
canvasProjectDB	Create and update a Canvas project database.
canvasSDMerge	Merge CSV data with an existing SD file.
canvasSearch	Search a list of target molecules against a set of queries.
canvasMCS	Find the maximum common substructure (MCS) in a set of molecules.
canvasScaffold	Decompose structures into ring-containing scaffolds and sort the scaffolds.

5.6.1 canvas_app

This utility runs a job according to the type and parameters specified in the input file. The syntax is as follows:

```
canvas_app [-exec|-inc] jobName
```

jobName is the name of the job to be run. The input file *jobName.inp* must exist; this file is created by canvasJob. The `-inc` option can be used to incorporate results into the Canvas project from a run that has completed already. The `-exec` option can be used to run a job without incorporating results into the Canvas project.

5.6.2 canvasConvert

This utility converts between SD, Maestro, Canvas compact, CSV and SMILES molecular structure file formats.

```
canvasConvert -ifmt inFile [input-options]
               {-ofmt outFile|-div N [-prefix basename]} [output-options]
```

The available *fmt* format options for both input and output are:

smi	SMILES format
sd	SD format
mae	Maestro format
csv	CSV format

The input file arguments (including options) are the common arguments listed in [Table 5.1](#); the output file arguments (including options) are the common arguments listed in [Table 5.2](#). The remaining output options are described in [Table 5.51](#).

Table 5.51. Output file options for the `canvasConvert` command.

Option	Description
-div <i>N</i>	Divide the original file into <i>N</i> parts, each containing approximately the same number of structures. The file is of the same type as the input file. The header line in an input CSV file, if it exists, is written to each of the divided files. Not valid with <code>-ofmt</code> options.
-prefix <i>basename</i>	When dividing a file with <code>-div</code> , name the files <i>basename_ptK.ext</i> , where <i>K</i> is an integer from 1 to <i>N</i> . If not specified, <i>basename</i> is taken from the input file name (without the path).
-n <i>structureRange</i> [-file]	List of input structures to write to the output file. <i>structureRange</i> is a range specification, as described on page 119 . If <code>-file</code> is used, <i>structureRange</i> is a text file containing the row specifications. Default: process all structures.
-uniform	Each structure block in a Maestro or SD input file contains the same set of properties.
-u	Output to SMILES or CSV file as canonical (unique) SMILES.
-k	Output to SMILES or CSV as kekulized SMILES of aromatic atoms and bonds.
-id	Append molecule name after the SMILES string in the SMILES output file, separated by a space.
-obad <i>badMolFile</i>	Save the molecules that failed to convert in the specified file in the original format. Only available with <code>-ismi</code> , <code>-isd</code> and <code>-icsv</code> . Default: print to stdout.
-compress	Compress output files, including the file of failures <i>badMolFile</i> . The extension <code>.gz</code> is added automatically. Compression can also be requested by specifying the file with the appropriate extension.
-2D	Convert output coordinates to 2D.
-allH	Force output of all hydrogens.
-noH	Do not include hydrogens in output structures.
-silent	Suppress console output.

5.6.3 canvasJob

This program does setup and cleanup for jobs associated with a Canvas project. The command syntax is as follows:

```
canvasJob -setup jobName -run program -proj projName [-helpColumns]
        -arg "arguments" [options]
canvasJob -cleanup jobName
canvasJob -help | -helpSetup
```

The arguments (and options) are described in [Table 5.52](#).

Table 5.52. Arguments and options for the canvasJob command.

Argument	Description
-arg " <i>arguments</i> "	Command line options for <i>program</i> enclosed in double quotes. Omit options that refer to an input or output file for <i>program</i> (e.g., -isd <i>sdFile</i>), or specific rows and columns to used, since those options are deduced from information provided to canvasJob.
-cleanup <i>jobName</i>	Clean up after <code>canvas_app</code> has run the specified job. May not be used in combination with any other options.
-cols <i>list</i>	The indices of the property columns to be used by <i>program</i> . <i>list</i> is a range specification, as defined on page 119 . Use -helpColumns to see the mapping of column indices to column names. Do not include the structure column, since canvasJob determines whether it is needed based on the specified Canvas program. This option is required whenever <i>program</i> operates on properties (such as canvasBayes, canvasMLR), but it should not be used for programs that operate only on structure (such as canvasFPGen, canvasMolDescriptors).
-help	Print usage message and exit.
-helpColumns	Print the correspondence between column indices and column names in the Canvas project <i>projName</i> and exit. No other options are required.
-helpSetup	Print information on setting up a job, including the list of programs that can be used and how to set up the -arg option.
-model <i>prevJob</i> -model <i>modelFile</i>	Use the model created during a previous run with the same project, or specify the file name of a previously built model. The file extension must be one of: canvasPLS .pls canvasMLR .mlr canvasPCAReg .pca canvasNnet .nnet canvasBayes .bayes canvasSOM .som canvasSOMBits .som

Table 5.52. Arguments and options for the `canvasJob` command. (Continued)

Argument	Description
<code>-proj projName</code>	Canvas project name, including absolute path.
<code>-rows list [-file]</code>	The project rows to process. <i>list</i> is a range specification, as defined on page 119 . All rows are processed by default. If <code>-file</code> is used, <i>list</i> is the name of a binary file, written by the Canvas GUI with View > Export Row IDs, that contains the row selection.
<code>-rows2 list [-file2]</code>	Project rows to process for the second input file when setting up a job for <code>canvasDBCS</code> or <code>canvasFPMatrix</code> (the file specified by <code>-ifp2</code>). If omitted, no second input file is used. <i>list</i> is a range specification, as defined on page 119 . If <code>-file2</code> is used, <i>list</i> is the name of a binary file, written by the Canvas GUI with View > Export Row IDs, that contains the row selection.
<code>-run program</code>	Canvas program to run. Valid programs are: <code>canvasBayes</code> <code>canvasPCA</code> <code>canvasFPGen</code> <code>canvasPCAReg</code> <code>canvasMLR</code> <code>canvasPLS</code> <code>canvasMolDescriptors</code> <code>canvasSOM</code> <code>canvasNnet</code> <code>canvasSOMBits</code>
<code>-setup jobName</code>	Set up a Canvas job. The input file <code>jobName.inp</code> is created with options for <code>canvas_app</code> , which actually runs the job.
<code>-y list</code>	Dependent variable columns. <i>list</i> is a range specification, as defined on page 119 . Use only for programs that build models, such as <code>canvasBayes</code> , <code>canvasMLR</code> .

5.6.4 canvasProjectDB

This program can be used to create and update a Canvas project database. It can be run as a regular foreground process, or as a single-CPU job on any host that has access to the project directory.

```
canvasProjectDB [job-options] -proj projName.cnv
                [import-options | export-options]
```

The project name, specified by `-proj`, must include the full path, and end in `.cnv`. The Job Control options are the standard options listed in [Table 5.3](#). The import options include the common input file options described in [Table 5.1](#), except that `-ifmt` and `-fieldAsName` are extended as described below. These and the other import options have the following syntax:

```
-ifmt source { -new | -append | -replace | -update prop1, prop2 [-dup] }
[-helpColumns] [-skipDupStruct] [-merge] [-uniform]
[-fieldAsName field] [-rows rowRange] [-file] [-noIndex]
```

These options are described in [Table 5.53](#).

Table 5.53. Import options for the canvasProjectDB command

Option	Description
<code>-ifmt <i>source</i></code>	Source of structures or properties to be imported. The supported formats are: <code>-imae</code> Maestro file <code>-isd</code> SD file <code>-icsv</code> SMILES, name and properties <code>-ismi</code> Space or tab-separated SMILES and name <code>-ifp</code> Binary fingerprints file generated by Canvas (<code>.fp</code>) <code>-iproj</code> Another Canvas project (<code>.cnv</code>). Must be given as absolute path. Include absolute path in <i>source</i> to avoid file copy when running as a job. Valid extensions are listed on page 119 .
<code>-helpColumns</code>	Print the correspondence between column indices and column names in the Canvas project and exit. Not valid with <code>-new</code> .
<code>-skipDupStruct</code>	Skip duplicate structures in <i>source</i> file during import. Only the first of the duplicates is imported.
<code>-merge</code>	Specify if all fields in a CSV file are to be merged into a project. No name or SMILES field is required in the CSV file. Only valid with <code>-update</code> .
<code>-uniform</code>	Each structure in <i>source</i> contains the same set of properties. Valid only with <code>-isd</code> or <code>-imae</code> .
<code>-new</code>	Create a new Canvas project, then import.
<code>-append</code>	Append to an existing project.

Table 5.53. Import options for the `canvasProjectDB` command (Continued)

Option	Description
<code>-replace</code>	Delete all project records, then import.
<code>-update</code> <code>prop1, prop2</code>	<p>Use the specified property mapping to add new records or update properties for existing records.</p> <p><i>prop1</i> is the name (<code>-imae</code>, <code>-isd</code>, <code>-ifp</code>, <code>-iproj</code>) or column index (<code>-icsv</code>, <code>-ismi</code>, <code>-iproj</code>) of a string or integer property in <i>source</i>. <i>prop1</i> must be 2 with <code>-ismi</code>, must be at least 2 with <code>-iproj</code>. Omit <i>prop1</i> if you wish to use the default name of the molecule or fingerprint in <i>source</i>. If a column name starts with a number, use escape plus single quote around the name. e.g. use '6' for a column named 6.</p> <p><i>prop2</i> is the name of a string or integer property in the existing Canvas project database. If a column name starts with a number, use escape and single quote—see above. Omit <i>prop2</i> if you wish to use the ID that appears in the Structure column of the Canvas spreadsheet. You must use "<code>_uid_</code>" for <i>prop2</i> if you wish to use the unique row id in the Canvas project.</p> <p>If the value of <i>prop1</i> is unique with respect to all <i>prop2</i> values, a new database record is created to hold the imported structure and its properties. If, however, the value of <i>prop1</i> matches the value of <i>prop2</i> for an existing database record, the imported properties are merged into the existing database record, and the structure is not changed.</p>
<code>-dup</code>	Perform the update even if the project contains records with duplicate <i>prop2</i> values. If this flag is not used, <code>canvasProjectDB</code> terminates when duplicate <i>prop2</i> values exist. Valid only with <code>-update</code> .
<code>-fieldAsName field</code>	<p>Get structure names from <i>field</i>. This is the ID that appears in the Structure column of the Canvas spreadsheet. <i>field</i> must be the name (<code>-imae</code>, <code>-isd</code>, <code>-iproj</code>) or column index (<code>-icsv</code>, <code>-ismi</code>) of a string or integer property in <i>source</i>. <i>field</i> is ignored with <code>-ismi</code> since structure names always come from the second column when it is present. If this option is not used, structure names are assigned as follows:</p> <ul style="list-style-type: none"> <code>-imae</code> <code>s_m_title</code> property <code>-isd</code> First line of each CT block <code>-icsv</code> Sequential integers starting with 1 <code>-ismi</code> Second column, or sequential integers starting with 1 <code>-iproj</code> ID that appears in the Structure column
<code>-rows rowRange</code>	Rows to import. <i>rowRange</i> is a range specification, as defined on page 119 . All rows are imported by default.
<code>-file</code>	Interpret <i>rowRange</i> as a file name. Each line in the specified file should contain a valid range specification.
<code>-noIndex</code>	Do not generate an index for substructure matching when adding records to the project.

The export options include the `-v3` output file option given in [Table 5.2](#), and the remaining arguments have the following syntax:

```
-ofmt outputFile [-compress] [-u] [-helpColumns] [-uidAsName]
  [-rows rowRange [-file]] [-cols colRange] [-view viewName]
```

These options are described in [Table 5.54](#).

Table 5.54. Export options for the `canvasProjectDB` command.

Option	Description
<code>-ofmt outputFile</code>	Destination of exported structures or properties. Supported formats: <code>-omae</code> Maestro file <code>-osd</code> SD file <code>-ocsv</code> Comma-separated SMILES, name and properties (<code>.csv</code>) <code>-osmi</code> Space-separated SMILES and name (<code>.smi</code>) <code>-ofp</code> Canvas fingerprint file (<code>.fp</code>). Must choose a single fingerprint column and any number of other non-fingerprint columns.
<code>-compress</code>	Compress output file. If specified, <code>.gz</code> is automatically appended to <i>outputFile</i> . Compression can also be requested by using the appropriate extension on <i>outputFile</i> . Not valid with <code>-ofp</code> .
<code>-u</code>	Write out canonical SMILES strings. Valid only with <code>-osmi</code> and <code>-ocsv</code> .
<code>-uidAsName</code>	Use the unique row ID in the Canvas project as the molecule name.
<code>-rows rowRange</code>	Rows to export. <i>rowRange</i> is a range specification, as defined on page 119 . Default: export all rows.
<code>-file</code>	Interpret <i>rowRange</i> as a file name. Each line in the specified file should contain a valid range specification.
<code>-cols colRange</code>	Columns to export. <i>colRange</i> is a range specification, as defined on page 119 . Column 1 (the structure) is always exported with <code>-omae</code> , <code>-osd</code> , and <code>-osmi</code> . This option must be used when exporting a fingerprint column. Default: export all columns.
<code>-helpColumns</code>	Print the correspondence between column indices and column names in the Canvas project and exit. No other options are required.
<code>-view viewName</code>	Export rows and columns from the saved view <i>viewName</i> .

5.6.5 canvasSDMerge

This utility merges CSV data with an existing SD file into a new SD file, preserving the original SD file order. The command syntax is

```
canvasSDMerge -isd sdfile -icsv csvfile [-key1 sdKey] [-key2 csvKey]
  [-not1 excludeList|-only1 includeList] [-not2 excludeList|-only2 includeList]
  [-o outputSDFile] [-keep1|-keep2|-keepboth]
```

The CSV input file must be comma-delimited. The options are described in [Table 5.55](#).

Table 5.55. Options for the *canvasSDMerge* command.

Option	Description
-key1 <i>sdKey</i>	Name of the field in the SD file to use as the key for merging. Default: molecule title.
-key2 <i>csvKey</i>	Name of the field in the CSV file to use as the key for merging. Default: first column.
-not1 <i>excludeList</i>	Fields in the SD file to exclude from the merge, space-delimited. Default: do not exclude any fields.
-only1 <i>includeList</i>	Fields in the SD file to include in the merge, space-delimited. Default: include all fields.
-not2 <i>excludeList</i>	Fields in the CSV file to exclude from the merge, space-delimited. Default: do not exclude any fields.
-only2 <i>includeList</i>	Fields in the SD file to include in the merge, space-delimited. Default: include all fields.
-o <i>outputSDFile</i>	Output file for results after merging. Default: standard output.
-keep1	Retain value from SD file when field names collide.
-keepboth	Retain value from both files when field names collide. This is the default.
-keep2	Retain value from CSV file when field names collide.

5.6.6 canvasSearch

This program searches a list of target molecules against a set of queries, composed of either molecules or partial structures. This program can also be used to filter a target list based on either standard REOS rules or a user-defined file containing SMILES queries and minimum and maximum number of times a query should be matched. Filtering and searching can be performed separately or sequentially (filtering first). The command syntax is as follows:

```
canvasSearch -ifmt structureFile [job-options] [input-options]
  [-ofmt matchFile] [-ofmt2 failFile] [-no2DCoord] [-v3]
  [-n selection] [-qfmt queryFile [-require n|-qmol mdlFile]] [-exact]
  [-filter [-reos] [-file ruleFile [-d delimiter]] [-maxVio n] ]
  [-index indexFile|-newIndex indexFile] [-noIndex] [-helpREOS]
  [-matchCount countFile -qmap queryMapFile -prefix queryPrefix
  [-showAll] ]
```

The job options are standard Job Control options, listed in [Table 5.3](#). This program can be distributed over multiple processors. The input file options are listed in [Table 5.1](#). The remaining arguments are described in [Table 5.56](#)

Table 5.56. Arguments for the *canvasSearch* command.

Argument	Description
-ifmt <i>targetFile</i>	Structure file containing the target molecules. The supported formats are: -imae Maestro file -isd SD file -ismi Space or tab-separated SMILES and name -iproj Canvas project. Must be given as absolute path. Valid extensions are listed on page 119 .
-n <i>selection</i>	Selected molecules in <i>targetFile</i> to search. <i>selection</i> is a range specification, as defined on page 119 . Default: search all molecules.
-index <i>indexFile</i>	Use the specified previously-generated index file of all the molecules in <i>targetFile</i> in the search. A matching fingerprint is generated for each query. Not valid with -iproj.
-newIndex <i>indexFile</i>	Generate index file of both the target molecules and the queries before search. Not valid with -iproj. The saved index file of the target molecules can be used later with the -index option.
-noIndex	Do not use any index for search, even if present. Supersedes -index and -newIndex.
-filter	Filter the target file based on the maximum and optionally the minimum number of counts of a given set of patterns.

Table 5.56. Arguments for the `canvasSearch` command. (Continued)

Argument	Description
<code>-reos</code>	Use Rapid Elimination Of Swill, a set of rules to identify lead-like molecules. Only valid with <code>-filter</code> .
<code>-file ruleFile</code>	Rule file to use. Each line must contain one SMARTS/SMILES string, followed by the minimum and maximum number of allowed counts, and an optional comment surrounded by double quotes. Only valid with <code>-filter</code> .
<code>-d delimiter</code>	Delimiter used to separate each field in <code>ruleFile</code> . Default: tab character '\t'.
<code>-maxVio n</code>	Maximum number of violations allowed for the rules. Only valid with <code>-filter</code> . Default: 0.
<code>-helpREOS</code>	Print out REOS patterns, each followed by the minimum and maximum number of counts. Tab is used as the delimiter in each line.
<code>-qfmt queryFile</code>	Query file containing list of queries, which can be whole molecules or fragments. The supported formats are: <code>-qmae</code> Maestro file <code>-qmol</code> MDL Mol file (single query only) <code>-qsd</code> SD file <code>-qsmi</code> SMILES or SMARTS patterns, one per line.
<code>-require n</code>	Require the target molecules to match at least <code>n</code> queries. This option should only be used with <code>-qfmt queryFile</code> , except for MDL Mol query files. Default: match all queries.
<code>-exact</code>	Require an exact match for each query. Default is to match by substructure.
<code>-ofmt matchFile</code>	Target molecules that passed the filter (if <code>-filter</code> is used) and matched all or the required number of queries. The supported formats are: <code>smi</code> SMILES file. Input SMILES strings are written for SMILES input, Canvas-generated SMILES strings are written otherwise. <code>mae</code> Maestro file. <code>sd</code> SD file. Valid extensions are listed on page 119 .
<code>-ofmt2 matchFile</code>	Target molecules that failed to match. The supported formats are: <code>smi</code> SMILES file. Input SMILES strings are written for SMILES input, Canvas-generated SMILES strings are written otherwise. <code>mae</code> Maestro file. Input must be in Maestro format. <code>sd</code> SD file. Valid extensions are listed on page 119 .
<code>-no2DCoord</code>	Do not generate coordinates in the output SD or Maestro file if SMILES was used as input. Default: generate 2D coordinates
<code>-v3</code>	Write output SD files in MDL version 3 format.

Table 5.56. Arguments for the `canvasSearch` command. (Continued)

Argument	Description
<code>-matchCount</code> <code>countFile</code>	Calculate the number of matches to each query or filtering pattern, 0 for no match. Counts are saved in <code>countFile</code> , which is in CSV format. If <code>-filter</code> and <code>-qfmt</code> are both used, only queries in the latter are listed. By default, only targets that passed the filter (if <code>-filter</code> is used) and matched all or the required number of queries are printed out to <code>countFile</code> . Use <code>-showAll</code> to include all targets.
<code>-prefix</code> <code>queryPrefix</code>	Specify the heading prefix in the match count CSV file. Each query in <code>countFile</code> is represented by the following format: <code>queryPrefix: :queryn.queryPrefix</code> can be a search name, such as <code>my_search1</code> . Only valid with <code>-matchCount</code> .
<code>-qmap</code> <code>queryMapFile</code>	This file provides the mapping between the <code>queryPrefix: :queryn</code> headings in the <code>countFile</code> and the actual SMARTS/SMILES query patterns. Only valid with <code>-matchCount</code> .
<code>-showAll</code>	Write counts for all targets to <code>countFile</code> . Only valid with <code>-matchCount</code> .

5.6.7 canvasMCS

This utility finds the maximum common substructure (MCS) among a given set of molecules. The command syntax is as follows:

```
canvasMCS -ifmt inputFile [input-options] -ofmt outputFile [output-options] [options]
```

The input file arguments (and options) are the common arguments listed in Table 5.1, and in addition input from Canvas projects is supported with `-iproj`. The valid output formats are `mae`, `sd`, and `csv`, as described in Section 5.1 on page 119. The common output arguments listed in Table 5.2 are accepted. Other options are described in Table 5.57.

Table 5.57. Options for the `canvasMCS` command.

Option	Description
<code>-min minMatch</code>	Minimum number of molecules that must match the MCS. If <code>minMatch</code> exceeds the number of input molecules it is interpreted as requiring all to match. Default: match all molecules.
<code>-max maxMatch</code>	Maximum number of molecules that must match the MCS. If <code>maxMatch</code> exceeds the number of input molecules it is interpreted as requiring all to match. When <code>maxMatch</code> is different from <code>minMatch</code> , a series of solutions spanning this range is produced. Default is all.
<code>-stop size</code>	Stop processing when the MCS atom and bond count falls below this threshold.

Table 5.57. Options for the `canvasMCS` command. (Continued)

Option	Description
<code>-nodetail</code>	Omit output of atom and bond number lists to Maestro and SD files. These are never added to CSV files. Default: include lists.
<code>-n structRange</code> <code>[-file]</code>	Input structures to process. <i>structRange</i> is a range specification, as defined on page 119 . If <code>-file</code> is specified, <i>structRange</i> is a file name that contains the range specification. If <code>-file</code> is used with <code>-iproj</code> , the file must be the binary set file written from the Canvas panel. Default: process all structures.
<code>-rs numMols</code>	Process only a random subset of the input molecules. If <i>numMols</i> is omitted, it is set to the square root of the total.
<code>-limit numMols</code>	Maximum number of molecules to process. Submitting larger sets may exceed available memory. Default: 2000.
<code>-timeout time</code>	Stop the job if the total time taken exceeds <i>time</i> seconds.
<code>-showall</code>	Output all equivalents for each MCS.
<code>-exclusive</code>	If an input molecule matches more than one MCS, report only the match to the largest MCS. Cannot be used with <code>-showall</code> .
<code>-sortname</code>	Sort output on molecule name. Default: keep input order.
<code>-opw filename</code>	Output MCS for all pairs of structures in a text file, one pair per line in the format: <i>title1</i> , <i>SMILES1</i> , <i>title2</i> , <i>SMILES2</i> , <i>MCSof1to2</i> , <i>MCSof2to1</i>
<code>-addring</code>	In output MCS SMARTS patterns, mark each atom as cyclic (R) or acyclic (R0). Default: do not mark atoms.
<code>-addh</code>	In output MCS SMARTS patterns, include hydrogen counts for each atom. Default: do not include hydrogen counts.
<code>-nox</code>	In output MCS SMARTS patterns, suppress addition of a connectivity qualification [nX3] for pyrrolic nitrogens. The qualification is added by default.
<code>-nobreakring</code>	Do not consider partial rings as part of MCS.
<code>-nobreakaring</code>	Do not consider partial aromatic rings as part of MCS.
<code>-allH</code>	Consider hydrogens as explicit atoms. Sets <code>-addh</code> .
<code>-atomtype scheme</code>	Atom typing scheme. Must be an integer value between 1 and 12 or C. The schemes are described in Table 5.7 on page 125 . Default: 7.

5.6.8 canvasScaffold

This utility decomposes structures into a set of scaffolds (scaffold detection), and searches a set of structures to determine which structures contain each scaffold in a set of scaffolds (scaffold matching).

Scaffold detection consists of identifying all unique scaffolds within a set of structures. The largest scaffold in a structure is obtained by stripping off all terminal side-chains with the exception of exocyclic and exolinker double bonds such as $>C=O$. For example, quinone contains two exocyclic double bonds and benzophenone contains a single exolinker double bond. These $C=O$ bonds would be retained as part of the quinone ring system and part of the benzophenone linker, respectively.

A given scaffold is broken into smaller subscaffolds by removing linkers between ring systems, but not by splitting fused ring systems. So, for example, benzophenone would be split into two benzene rings, but naphthalene would not.

Generation of subscaffolds is exhaustive, so if R1, R2, and R3 are ring systems, the scaffold R1–R2–R3 would produce the subscaffolds R1–R2, R2–R3, R1, R2, and R3. Scaffold detection finds all such scaffolds, and stores them in a list that is sorted by decreasing scaffold size. The size of a scaffold is based on the number of ring systems, followed by the number of linkers, followed by the number of heavy atoms, followed by molecular weight, followed by canonical SMILES. Sorting in this manner guarantees that the subscaffolds of a given scaffold will always appear after the scaffold itself.

Scaffold matching consists of taking a set of structures and a set of scaffolds, and determining which scaffolds are contained in each structure. The scaffolds may or may not come from the structures being matched against them. Since matching produces a list of scaffolds for each structure, it is possible to assign the structures into clusters based on the scaffolds they contain.

The command syntax is:

```
canvasScaffold inFile {-detect|-match|-full} [options]
```

where *inFile* is the input structure file, in Maestro, SD, CSV (SMILES and properties), or SMILES format. Maestro and SD files can be compressed. CSV files must contain a column header line and SMILES must appear in the first column. SMILES files have no header, and each SMILES string may be followed by a tab and a structure name. The three run modes are

- detect Detect scaffolds among the input structures.
- match Match the scaffolds in *jobname_scaffolds.dat* against the input structures.
- full Perform scaffold detection and matching on the input structures.

The options are given in [Table 5.58](#).

Table 5.58. Options for the `canvasScaffold` command.

Option	Description
<code>-exocyc file</code>	Extend exocyclic definitions to include the SMARTS patterns in <i>file</i> . For example, use <code>[CH3]a</code> to retain methyls attached to aromatic rings. Not valid with <code>-match</code> .
<code>-exolink file</code>	Extend exolinker definitions to include the SMARTS patterns in <i>file</i> . For example, use <code>O=C([O-,OH])([CD3;CR0])</code> to retain carboxylates attached to a tertiary carbon linker atom. Not valid with <code>-match</code> .
<code>-iscaff file</code>	Use scaffolds in <i>file</i> rather than <code>jobname_scaffolds.dat</code> . This file must have been created by a previous <code>canvasScaffold</code> job run with <code>-detect</code> or <code>-full</code> . Valid only with <code>-match</code> .
<code>-embed</code>	Allow embedded ring system and linker matches when detecting or matching. See text below for a detailed description with an example.
<code>-title prop</code>	Use an alternate property for structure names. Must be the name of a string or integer property in <i>inFile</i> . Not valid with <code>-detect</code> or with SMILES (<code>.smi</code>) input.
<code>-props list</code>	Comma-separated list of properties from <i>inFile</i> to include in the output file <code>jobname_match.csv</code> . By default, only the structure name is included. Not valid with <code>-detect</code> or with SMILES input.
<code>-include list</code>	List of scaffold IDs to consider when matching. See Section 5.1 on page 119 for list syntax. These are the only IDs reported in <code>jobname_match.csv</code> . If <i>list</i> is the name of a file, it is assumed that the file contains the desired IDs, one per line. Valid only with <code>-match</code> .
<code>-exclude list</code>	List of scaffold IDs that should be excluded when matching. See Section 5.1 on page 119 for list syntax. Excluded scaffold IDs never appear in <code>jobname_match.csv</code> . If a structure contains an excluded scaffold, that structure is not reported as a match to any subscaffolds of the excluded scaffold. If <i>list</i> is the name of a file, it is assumed that the file contains the list of IDs, one per line. Valid only with <code>-match</code> .
<code>-fuzzy</code>	Use fuzzy scaffold matching. A given structure is considered to match a scaffold if it contains a substructure with the same generic framework as the scaffold, treating all atoms and bonds as equivalent. Not valid with <code>-detect</code> .
<code>-arom</code>	When doing fuzzy matching, distinguish aromatic and non-aromatic atoms. Only valid with <code>-fuzzy</code> .
<code>-bonds</code>	When doing fuzzy matching, distinguish single, double, triple, and aromatic bonds. Only valid with <code>-fuzzy</code> .

Results of scaffold detection are written to *jobname_detect.csv*, which contains a canonical SMILES string for each scaffold, a unique scaffold ID, and quoted, space-delimited lists of the source input structure numbers, superscaffold IDs, and subscaffold IDs. For example,

```
SMILES, ID, Source, Super, Sub
n1cccc1Cc2cccc2,1,"1","","3 4"
c1cccc1Cc2cccc2,2,"2","","4"
c1ccncc1,3,"1","1",""
c1cccc1,4,"1 2","1 2",""
```

A separate file, *jobname_scaffolds.dat*, is created with all the information required to reconstruct the original scaffolds and the relationships among them.

Results of scaffold matching are written to *jobname_match.csv*. Each structure that matches at least one scaffold is written to *jobname_match.csv*, with canonical SMILES, structure name, an index that corresponds to the position of the structure in the original input file, and the IDs of the scaffolds matched by the structure. For example,

```
SMILES, Name, Index, IDs
O=C(O)c1ccc(cc1)Cc2ccccc2,"CHEM0192587",1,"1 3 4"
O=C(O)c1ccc(cc1)Cc2ccccc2,"CHEM0204591",2,"2 4"
```

Here, CHEM0192587 matches scaffolds 1, 3 and 4, while CHEM0204591 matches scaffolds 2 and 4.

The use of *-embed* needs some explanation. If a mixture of imidazoles and benzimidazoles is being analyzed and *-embed* is used with *-detect* or *-full*, embedded benzimidazole-imidazole relationships would be reported in *jobname_detect.csv*. When an embedded relationship is reported, the letter *e* is appended to the applicable source structure numbers and scaffold IDs. For example,

```
SMILES, ID, Source, Super, Sub
c1cccc(c12)nc[nH]2,1,"1","","2e"
c1[nH]cnc1,2,"1e","1e",""
```

Here, the Source value *1e* on line 3 indicates that scaffold 2 is an embedded scaffold of input structure 1. The Sub value *2e* on line 2 indicates that scaffold 2 is an embedded subscaffold of scaffold 1, and the Super value *1e* on line 3 indicates the reciprocal relationship.

If *-embed* is used with *-match* or *-full*, imidazole would be listed in *jobname_match.csv* as an embedded match to any benzimidazole-containing structure. For example,

```
SMILES, Name, Index, IDs
c1cccc(c12)nc([nH]2)C,"2-methyl-benzimidazole",1,"1 2e"
```

Here, "1 2e" indicates an ordinary match to benzimidazole (scaffold ID = 1), and an embedded match to imidazole (scaffold ID = 2).

By default, no relationship is reported between imidazole and benzimidazole. Note also that if scaffold detection is performed on structures that contain only benzimidazole ring systems, imidazole would not be found as a scaffold, even with `-embed`.

Scaffold detection scales linearly with the number of structures analyzed, and quadratically with the number of unique scaffolds contained in those structures. The quadratic part comes from the identification of superscaffold-subcaffold relationships. The worst case scenario is when there are a large number of unrelated structures, because many of the structures will have a unique backbone, and you will get a new scaffold for each of those backbones.

5.7 Scripting with Canvas

In addition to the command-line utilities, Canvas has a Python API that can be used to write Python scripts. The API is described in the document *Canvas Python API*. Canvas utilities are also available for use with KNIME, along with many other Schrödinger products and utilities.

References

1. Leach, A. R.; Gillet V. J. *An Introduction to Cheminformatics*. Springer: Dordrecht, 2007.
2. Gobbi, A.; Lee, M.-L. *J. Chem. Inf. Comput. Sci.* **2003**, *43*, 317.
3. Lloyd, S. P. *IEEE Transactions on Information Theory* **1982**, *28*, 129.
4. Kelley, L. A.; Gardner, S. P.; Sutcliffe, M. J. An automated approach for clustering an ensemble of NMR-derived protein structures into conformationally-related subfamilies. *Protein Eng.* **1996**, *9*, 1063.
5. Sastry, M.; Lowrie, J. F.; Dixon, S. L.; Sherman, W. Large-Scale Systematic Analysis of 2D Fingerprint Methods and Parameters to Improve Virtual Screening Enrichments. *J. Chem. Inf. Model.* **2010**, *50*, 771.
6. Duan, J.; Dixon, S.L.; Lowrie, J.F; Sherman, W. Analysis and comparison of 2D fingerprints: Insights into database screening performance using eight fingerprint methods. *J. Molec. Graph. Model.* **2010**, *29*, 157.
7. Rosipal R.; Trejo L.J. Kernel Partial Least Squares Regression in Reproducing Kernel Hilbert Space. *J. Machine Learning Res.* **2001**, *2*, 97.
8. Kier, L. B.; Hall, L. H. *Molecular Structure Description: The Electrotopological State*; Academic Press: San Diego, 1999.
9. Ghose, A. K.; Viswanadhan, V. N.; Wendoloski, J. J. Prediction of Hydrophobic (Lipophilic) Properties of Small Organic Molecules Using Fragmental Methods: An Analysis of ALOGP and CLOGP Methods. *J. Phys. Chem A* **1998**, *102*, 3762.
10. Viswanadhan, V. N.; Ghose, A. K.; Revankar, G. R.; Robins, R. K. Atomic Physico-chemical Parameters for Three Dimensional Structure Directed Quantitative Structure-Activity Relationships. 4. Additional Parameters for Hydrophobic and Dispersive Interactions and Their Application for an Automated Superposition of Certain Naturally Occurring Nucleoside Antibiotics. *J. Chem. Inf. Comput. Sci.*, **1989**, *29*, 163.
11. Miller K. J. *J. Am. Chem. Soc.* **1990**, *112*, 8533.
12. Hall, L. H.; Kellogg, G. E.; Haney D. N. Molconn-Z Software Package for Molecular Topology Analysis. User's Guide. Version 4.12, 2008. <http://www.edusoft-lc.com/molconn/manuals/400/>

References

13. Ertl, P.; Rohde, B.; Selzer, P. Fast Calculation of Molecular Polar Surface Area as a Sum of Fragment-Based Contributions and Its Application to the Prediction of Drug Transport Properties. *J. Med. Chem.* **2000**, *43*, 3714.
14. Klon, A. E.; Lowrie, J. F.; Diller, D. J. Improved Naive Bayesian Modeling of Numerical Data for Absorption, Distribution, Metabolism and Excretion (ADME) Property Prediction. *J. Chem. Inf. Model.* **2006**, *46*, 1945.
15. Baell, J. R.; Holloway, G. A. New substructure filters for removal of pan assay interference compounds (PAINS) from screening libraries and for their exclusion in bioassays. *J. Med. Chem.* **2010**, *53*, 2719.

Getting Help

Schrödinger software is distributed with documentation in PDF format. The manuals are also available from the Schrödinger [Support Center](#). If the documentation is not installed in `$(SCHRODINGER)/docs` on a computer that you have access to, you should install it or ask your system administrator to install it.

For help installing and setting up licenses for Schrödinger software and installing documentation, see the *Installation Guide*. For information on running jobs, see the *Job Control Guide*.

Canvas has a help system that contains information about the Canvas interface.

- For information about a panel, click the Help button in the panel. The help topic is displayed in the Help panel.
- For other information in the help set, open the default help topic by choosing Help → Online Help in the main window. You can use the search facility to find topics.
- To search the help, use the search box. You can use quotation marks to search for an exact phrase, and you can use “and”, “or”, “not”, and parentheses to construct complex search expressions. The results are displayed in the display area, and you can click links to display the help topics.

The Help menu provides access to the knowledge base on the external web site, at <http://www.schrodinger.com/kb>, which you can search for information about Canvas. It also provides access to information about known issues with Canvas, at <http://www.schrodinger.com/known-issues>. To diagnose problems with licensing or remote hosts, you can open the Diagnostics panel from the Help menu.

If you have questions that are not answered from any of the above sources, contact Schrödinger using the information below.

Web: <http://www.schrodinger.com/supportcenter>
E-mail: help@schrodinger.com
Mail: Schrödinger, 101 SW Main Street, Suite 1300, Portland, OR 97204
Phone: +1 503 299-1150 (USA, 9am – 5pm Pacific Time)
+49 621 438-55173 (Europe, 9am – 5pm Central European Time)
Fax: +1 503 299-4532 (USA, Portland office)
FTP: <ftp://ftp.schrodinger.com>

Generally, using the web form is best because you can send machine output and upload files, if necessary. Please include the following information:

- All relevant user input and machine output
- Canvas purchaser (company, research institution, or individual)
- Primary Canvas user
- Installation, licensing, and machine information as described below.

To gather the required machine, licensing, and installation information to send to technical support:

1. Open the Diagnostics panel.
 - **Canvas:** Help → Diagnostics
 - **Windows:** Start → All Programs → Schrodinger-2014 → Diagnostics
 - **Mac:** Applications → Schrodinger2014 → Diagnostics
 - **Linux/Command line:** `$SCHRODINGER/diagnostics`

2. When the diagnostics have run, click Technical Support.

A dialog box opens, with instructions. You can highlight and copy the name of the file.

3. Upload the file specified in the dialog box to the web form.

Numerics

- 2D structures
 - converting to 3D 110
 - editing 25
 - importing by name..... 15
- 3D Minimization dialog box..... 111
- 3D Pharmacophore Fingerprints - Advanced
 - Options dialog box..... 72
- 3D Pharmacophore Fingerprints dialog box 71

A

- Apply Logic to View dialog box..... 33
- Apply to Master dialog box 33
- archiving projects..... 11
- atom typing schemes..... 125
- Axis settings dialog boxes 52

B

- Bars dialog box 53
- Bayes Classification - Bins dialog box 97
- Bayes Classification dialog box..... 96
- Binary Fingerprints from Structures dialog
 - box 68
- bins
 - assigning for Bayes classification 97, 100, 166
 - inter-feature distance 71
 - overlap 123
- bit set, exporting and importing 17
- bits
 - fraction kept in Bayes classification 166
 - limiting for self-organizing maps 170
 - number in fingerprint..... 70, 124
 - omitting from fingerprint..... 70, 73, 124
 - selecting for fingerprint export 17
 - sorting by count 86, 151
- bivariate statistics 55

C

- Calculator dialog box..... 57
- Canvas main window 6
- cells
 - clearing 28
 - editing 28
 - selecting 22

- Chemistry Filter dialog box 44
- Choose Random Subset dialog box 30
- class
 - assigning values..... 37
 - changing color 37
 - creating 38
 - editing values..... 38
 - moving rows to 38
 - removing rows from 39
- clustering
 - by properties 82
 - linkage method 83, 146
 - methods 82
- clusters
 - dendrogram..... 83
 - exporting..... 84
 - leader 87
 - membership property..... 87
 - viewing 87
- color
 - chart title, background, foreground 52
 - distance matrix 74
 - for class 35, 37, 38, 39
 - histogram axis 53
 - job record..... 62
 - pie charts 54
 - property values in spreadsheet 40
 - scatter plot axes 52
 - scatter plot from model 89
 - scatter plot symbols 48
 - self-organizing map..... 102
 - similarity matrix 74
 - SOM cell borders 103
 - statistics 56
- columns
 - finding 24
 - hiding..... 29
 - moving..... 19
 - resizing 19
 - sorting 19
- conformers for 3D pharmacophore
 - fingerprints..... 71, 72
- conventions, document vii
- coordinates, discarding on import..... 13
- custom view 30
 - See also* views

D

- data
 - clearing from cells 28
 - exporting 16
 - importing 14
- decision tree 99
- dendrogram, viewing 83
- directed sphere exclusion (DISE) 79
- directory
 - installation 1
 - Maestro working 2
 - project 10
- distance matrix
 - from fingerprints 74
 - from properties 74
 - metric 74, 142
- Distance Matrix from Properties dialog box 74, 75
- Diversity-Based Selection dialog box 78
- duplicate structures
 - detecting in project 47
 - removing from project 116
 - removing on import 12, 13

E

- Edit Partition dialog box 39
- Edit Structure dialog box 25
- electrotopological states—*see* Estates
- environment variable
 - PATH 58
 - SCHRODINGER 1
 - SCHRODINGER_CANVAS_MAX_MEM.... 59, 119
- Estates
 - atom typing 124
 - calculating 64, 156
- Export Options dialog box 16
- External Application dialog box 112

F

- Feature Selection dialog box 66
- Feature Selection Viewer panel 67
- file formats
 - for import 11
- Find panel 23

- fingerprints
 - 3D from pharmacophores 70
 - calculating 68, 122
 - exporting 17
 - importing Canvas 11, 13
 - omitting bits 70, 124
 - precision 70, 71, 123
 - reducing number of bits 124
 - sorting by bit count 86, 151
 - types 123
- Frequency Pie Chart panel 51

H

- Hashed Fingerprints - Advanced Options dialog box 69
- heat map
 - distance matrix 74
 - similarity matrix 74
 - spreadsheet 40
 - statistics 56
- Heat Map dialog box 40
- hiding rows and columns 29
- Hierarchical Clustering Dendrogram panel 85
- Hierarchical Clustering dialog box 83
- Histogram panel 50
- histograms 50
 - axis and bar settings 52
- Hole-Filling and Library Optimization - Property Filter dialog box 80
- Hole-Filling and Library Optimization dialog box 80
- hydrogens, showing and hiding 20

I

- Import SMILES and Properties dialog box 14
- Import Structure File dialog box 12

K

- Kernel-Based Partial Least-Squares Regression dialog box 93
- K-Means Clustering dialog box 86
- KPLS Model Visualization panel 94
- KPLS Regression dialog box 93

L

Leader-Follower Clustering dialog box 86
 Library Comparison dialog box 81
 LigPrep..... 111
 linkage method..... 83, 146
 loadings, plotting 105

M

Maestro
 exporting structures to 17
 exporting structures to Canvas..... 11
 master view 30
 Maximum Common Substructure panel 106
 MCS groups 105
 MCS Matches view 107
 memory, for applications 59, 61, 119
 metric
 description 76
 for diverse compound selection 136
 for property distances 74
 for similarity/distance matrices 142
 model
 building 88
 exporting 89
 Molecular Properties - Advanced Options dialog
 box 64
 Molecular Properties dialog box 63
 multiple linear regression
 choosing subset of variables 89, 155
 Multiple Linear Regression - Best Subsets Options
 dialog box 90
 Multiple Linear Regression dialog box..... 90

N

Neural Networks dialog box 98
 New Class dialog box 39
 New Partition from Property dialog box 36
 New Partition from Random Split dialog box .. 38
 New Partition from Selection dialog box..... 35

P

partial least-squares regression
 number of factors..... 91, 92, 163, 164
 t-value 91, 163
 Partial Least-Squares Regression dialog box.... 91

partition

 adding and rearranging classes 38
 assigning training and test set from 88
 creating from property 36
 creating from selection 36
 creating random 37
 Partition Filter dialog box 43
 pharmacophore feature definitions 71
 pie charts 51
 settings 54
 Plot Symbols dialog box 49
 Preferences dialog box 59
 Principal Components Analysis dialog box.... 104
 Principal Components Regression dialog box.. 95
 product installation 191
 projects
 archiving 11
 creating 10, 176
 importing from 11
 incorporating results into 62
 information available in 18
 opening on startup 5
 recently used 10
 updating 176
 properties
 adding to application results view 32
 calculating 63, 156
 cluster membership 87
 clustering by 82
 coloring plot symbols by 48
 displaying heat maps 40
 distance matrix 74
 diverse structure set 79
 exporting to Maestro 17
 filtering by 41
 for heat map 40
 from command-line jobs 32
 import order 11
 importing 14
 missing values 55
 selecting for export 16
 similarity, naming 75
 sizing plot symbols by 48
 sorting by value of 34
 statistics for 55
 updating on import 13
 Property Filter dialog box 41

- Python Shell panel 114
- R**
- Recursive Partitioning - Advanced Options
 dialog box 100
- Recursive Partitioning dialog box..... 99
- REOS rules 180
- rotatable bonds
 counts of 64
 default Canvas definition 158
 sampling for conformer generation 72
- rows
 exporting random subset..... 16
 hiding 29
 random subset..... 30
 resizing 19
 selecting by row number 22
 sorting by property values 34
- S**
- Save Custom View dialog box 31
- Scaffold Decomposition dialog box..... 108
- Scaffold Decomposition for jobname panel ... 109
- Scatter Plot panel 48
- scatter plots 48
 axis settings 52
 from model predictions 89
 of bivariate statistics variables..... 56
- Schrödinger contact information 191
- scores, principal component 105
- Sectors dialog box..... 54
- selection
 applying from custom to master view 32
 by matched text..... 24
 clearing 22
 expanding to entire columns..... 22
 expanding to entire rows..... 22
 from histograms..... 50
 from pie charts 51
 from scatter plots 49
 inverting 22
 shortcut menu 20
 with arrow/tab keys 20
 with mouse 22
- Self-Organizing Map dialog box 101
- Self-Organizing Map Viewer panel 103
- self-organizing maps
 cell border color 103
 lattice parameters..... 102, 171
 limiting number of bits 170
- Shape Screen - Advanced Options dialog box.. 77
- Shape Screen dialog box 77
- similarity matrix
 from fingerprints..... 74
 metric 142
- Similarity/Distance Matrix from Fingerprints
 dialog box 73
- Similarity/Distance Screen dialog box 76
- Sort dialog box..... 34
- sort keys, reordering 35
- Source Structures panel 110
- sphere exclusion..... 79
- spreadsheet, copying to..... 17
- Start Application dialog box 62
- Statistics dialog box 55, 56
- statistics, univariate and bivariate 55
- structural keys 70
- structures
 coloring..... 24
 copying 24
 diverse, selecting 78
 drawing of..... 20
 duplicates on import 12, 13
 editing 25
 excluding from diverse set..... 79
 exporting to Maestro 17
 importing a random subset 12
 MCS membership..... 106
 pasting 24
 renaming..... 25
 shortcut menu 21
 showing and hiding 20
 similarity to reference..... 75
 size..... 20
 using for chemistry filters..... 44
- Substructure Query dialog box 46
- T**
- training set
 manual selection..... 89, 153
 predefined 117
 random..... 88, 153

U

undoing view changes	32
univariate statistics	55

V

views

application results	32
applying custom to master	32
closing	32
clustering results	87

combining	33
custom, definition	30
deleting	32
editing description	31
limitations	31
master, definition	30
opening	31
renaming	31
saving	32
sorting list of	32
undoing changes	32

120 West 45th Street
17th Floor
New York, NY 10036

155 Gibbs St
Suite 430
Rockville, MD 20850-0353

Quatro House
Frimley Road
Camberley GU16 7ER
United Kingdom

101 SW Main Street
Suite 1300
Portland, OR 97204

Dynamostraße 13
D-68165 Mannheim
Germany

8F Pacific Century Place
1-11-1 Marunouchi
Chiyoda-ku, Tokyo 100-6208
Japan

245 First Street
Riverview II, 18th Floor
Cambridge, MA 02142

Zeppelinstraße 73
D-81669 München
Germany

No. 102, 4th Block
3rd Main Road, 3rd Stage
Sharada Colony
Basaveshwaranagar
Bangalore 560079, India

8910 University Center Lane
Suite 270
San Diego, CA 92122

Potsdamer Platz 11
D-10785 Berlin
Germany

SCHRÖDINGER