

# **CSD Pipeline Pilot Component Collection**

## 1 CSD Pipeline Pilot Component Collection

### 1.1 Conditions of Use

## 2 Introduction

### 2.1 Integration Scheme

### 2.2 Requirements

### 2.3 Installation of the CSD Pipeline Pilot Component Collection Package

#### 2.3.1 Installation on Windows

#### 2.3.2 Installation on Linux

#### 2.3.3 Configuration

#### 2.3.4 Uninstallation on Windows

#### 2.3.5 Uninstallation on Linux

### 2.4 About the CSD PP Component Collection

## 3 CSD Python API Components

### 3.1 Manipulators

#### 3.1.1 Parse Citation

### 3.2 Readers

#### 3.2.1 Get CSD Entry Attributes

#### 3.2.2 Get CSD Crystal Attributes

#### 3.2.3 Get CSD Molecule Attributes

#### 3.2.4 Get Molecule Structure

#### 3.2.5 CSD Substructure Search

#### 3.2.6 CSD Similarity Structure Search

#### 3.2.7 CSD Text Numeric Search

#### 3.2.8 CSD Reduced Cell Search

### 3.3 Viewers

#### 3.3.1 Hermes Viewer & Mercury Viewer

#### 3.3.2 Conformer Report Viewer

#### 3.3.3 Virtual Screening Report Viewer

### 3.4 Virtual Screening

- 3.4.1 Perform Virtual Screening
- 3.4.2 Perform Virtual Screening Validation
- 3.4.3 Generate Enrichment Plot & Generate ROC Plot
- 3.4.4 Perform Conformer Generation

### 3.5 Utilities Components

- 3.5.1 Run Python Script
- 3.5.2 Run Virtual Screening
- 3.5.3 Run Virtual Screening Validation
- 3.5.4 Run Generate Conformers
- 3.5.5 Derive Script Path
- 3.5.6 Throw Script Error Message
- 3.5.7 Check Journal Name
- 3.5.8 Validate Journal Name
- 3.5.9 Gather Database Names
- 3.5.10 Join Data from JSON

## 4 CSD Python API Protocols

### 4.1 CSD Searching

- 4.1.1 Search CSD By Structure
- 4.1.2 Search CSD By Text Numeric Fields
- 4.1.3 Search CSD By Reduced Cell
- 4.1.4 Retrieve Entry and Molecule Attributes
- 4.1.5 Combining Hit Sets - AND
- 4.1.6 Combining Hit Sets – NOT
- 4.1.7 Parsing Citation and Synonyms

### 4.2 Python Examples

- 4.2.1 Run Python Script Example
- 4.2.2 Using Derive Script Path
- 4.2.3 Get CSD Python API Version
- 4.2.4 Count Entries per Decade
- 4.2.5 Count Entries per Year

### 4.3 Virtual Screening and Conformer

- 4.3.1 Queries Identified by File Screening Example
- 4.3.2 Queries Identified by Tag Screening Example
- 4.3.3 Screen Validation Using Tagging
- 4.3.4 Generate Enrichment Plot Example
- 4.3.5 Generate ROC Example
- 4.3.6 Generate Conformers for Molecule
- 4.3.7 Mercury Viewer Example
- 4.3.8 Conformer Writer Example

- 4.3.9 View Conformers in Report Viewer
- 4.3.10 Mercury Viewer Example with Grouping
- 4.3.11 Virtual Screening Report Viewer Example
- 4.3.12 Hermes Viewer Example – Structures

# **1 CSD Pipeline Pilot Component Collection**

Copyright © 2024

Registered Charity No 800579

## **1.1 Conditions of Use**

The Cambridge Structural Database Portfolio (CSD Portfolio) including, but not limited to, the following: ConQuest, CSD-Editor, Decifer, Mercury, Mogul, IsoStar, CSD Conformer Generator, Hermes, GOLD, SuperStar, the CSD Python API, web accessible CSD tools and services, WebCSD, CSD sketchers, CSD data files, CSD data updates, the CSD database, sub-files derived from the foregoing data files, documentation and command procedures, test versions of any existing or new program, code, tool, data files, sub-files, documentation or command procedures which may be available from time to time (each individually a Component) encompasses database and copyright works belonging to the Cambridge Crystallographic Data Centre (CCDC) and its licensors and all rights are protected.

Any use of a Component of the CSD Portfolio, is permitted solely in accordance with a valid Licence of Access Agreement or Products Licence and Support Agreement and all Components included are proprietary. When a Component is supplied independently of the CSD Portfolio its use is subject to the conditions of the separate licence. All persons accessing the CSD Portfolio or its Components should make themselves aware of the conditions contained in the Licence of Access Agreement or Products Licence and Support Agreement or the relevant licence.

In particular:

- The CSD Portfolio and its Components are licensed subject to a time limit for use by a specified organisation at a specified location.
- The CSD Portfolio and its Components are to be treated as confidential and may NOT be disclosed or re-distributed in any form, in whole or in part, to any third party.
- Software or data derived from or developed using the CSD Portfolio may not be distributed without prior written approval of the CCDC. Such prior approval is also needed for joint projects between academic and for-profit organisations involving use of the CSD Portfolio.
- The CSD Portfolio and its Components may be used for scientific research, including the design of novel compounds. Results may be published in the scientific literature, but each such publication must include an appropriate citation as indicated in the Schedule to the Licence of Access Agreement or Products Licence and Support Agreement and on the CCDC website.
- No representations, warranties, or liabilities are expressed or implied in the supply of the CSD Portfolio or its Components by CCDC, its servants or agents, except where such exclusion or limitation is prohibited, void or unenforceable under governing law.

Licences may be obtained from:

Cambridge Crystallographic Data Centre

12 Union Road

Cambridge CB2 1EZ, United Kingdom

Web: <http://www.ccdc.cam.ac.uk>

Telephone: +44-1223-336408

Email: [admin@ccdc.cam.ac.uk](mailto:admin@ccdc.cam.ac.uk)

## 2 Introduction

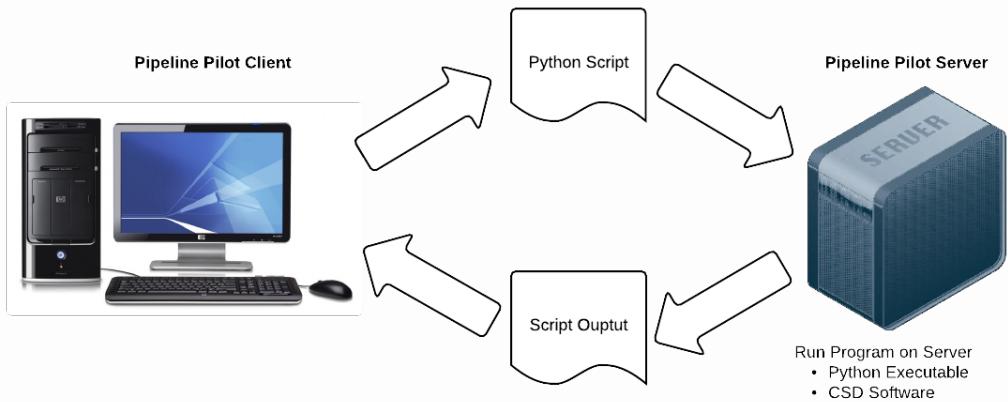
This component collection for Pipeline Pilot allows protocols to integrate functionality from the Cambridge Structural Database (CSD) Python API. This includes capabilities for CSD searching, model validation, conformer generator and virtual screening. There are also components to help in the integration of further Python scripts which may use the API.

Information about the CSD Python API may be found at <https://www.ccdc.cam.ac.uk/solutions/software/csd-python/>. Please do not hesitate to contact [support@ccdc.cam.ac.uk](mailto:support@ccdc.cam.ac.uk) for more information and help.

### 2.1 Integration Scheme

The integration works as follows:

- The components found in the CSD Pipeline Pilot (PP) component collection create the necessary input files from the incoming data, and from the parameters on those components.
- The components find the appropriate Python script in the component collection; various python scripts are supplied.
- The Run Program on Server component runs the Python executable, supplying the script and arguments for that script, which produce the output.



- The components then read the file(s) generated by the script, performing appropriate merges, joins, etc. to produce a stream of data.

## 2.2 Requirements

The Pipeline Pilot server host must have installed:

- Pipeline Pilot server (version 19.1 or later).
- The 2022 release of the CSD, including required licenses.
- Currently supported operating systems are Windows Server 2019 and CentOS 7, although CentOS 8 and RHEL 7 and 8 are also expected to work.
- The minimum specifications for the CSD Python API apply.
- On Linux, the user which the Pipeline Pilot server runs as must have a home directory for the CSD licensing to function.

To use the Mercury Viewer or Hermes Viewer components, Mercury and Hermes must be installed on the client. A free version of Mercury can be downloaded from <http://www.ccdc.cam.ac.uk/Community/csd-community/freemercury/>.

## 2.3 Installation of the CSD Pipeline Pilot Component Collection Package

### 2.3.1 Installation on Windows

All the actions described below should be carried out on the PP server host.

Ensure that all the requirements specified above are met.

If you have a previous version of the package installed, it must be uninstalled first (see below).

Open a PowerShell prompt as Administrator and navigate to the apps folder under the PP server installation. If the PP server has been installed into the default location, this would be C:\Program Files\BIOVIA\PPS.

In the apps folder, inflate the package zip archive. A folder ccdc should then be present.

Activate the CSD Python API environment. Assuming the CSD has been installed into the default location, use the commands:

```
C:\Program  
  Files\CCDC\Python_API_2022\miniconda\shell\condabin\conda-  
  hook.ps1
```

```
conda activate "C:\Program Files\CCDC\Python_API_2022\miniconda"
```

Use the Pkgutil tool to install the package; this registers protocols and components and makes them visible to the Pipeline Pilot client (N.B. a forward slash is required in the package name):

```
...\bin\pkgutil.exe -i ccdc/pythonapi
```

### 2.3.2 Installation on Linux

All the actions described below should be carried out on the PP server host.

Ensure that all the requirements specified above are met.

If you have a previous version of the package installed, it must be uninstalled first (see below).

Open a terminal and navigate to the `apps` directory under the PP server installation. A typical location for the PP server on Linux would be `/opt/BIOVIA/PPS`.

In the `apps` directory, inflate the package zip archive. A directory `ccdc` should then be present.

Activate the CSD Python API environment. If the installation prefix for the CSD Portfolio was `/opt/CCDC`, the appropriate command would be:

```
/opt/CCDC/Python_API_2022/miniconda/bin/activate
```

Set the `CSDHOME` environment variable. If the installation prefix for the CSD Portfolio is as above, the appropriate command would be:

```
export CSDHOME=/opt/CCDC/CSD_2022
```

Set up the environment for the `Pkgutil` tool:

```
.../linux_bin/ppvars.sh
```

Use the `Pkgutil` tool to install the package; this registers protocols and components and makes them visible to the PP client:

```
.../linux_bin/pkgutil -i ccdc/pythonapi
```

**IMPORTANT:** You may see an error message like the following:

- `ccdc/pythonapi` was not installed: Unable to open file `<pps_dir>/apps/ccdc/pythonapi/docs/pipeline_pilot_component_collection.html-tmp` for reading: No such file or directory.

This is a problem with Pipeline Pilot, and while HTML documentation will not have been generated, the package itself should actually have been installed.

### 2.3.3 Configuration

Once the package has been installed, the `PYTHON_HOME` Global Property needs to be set to the path of the CSD Python API directory containing the python executable. This is done using the web-based Administration Portal, accessible from the Server Home Page (see the Help menu of the PP client). Note that the default username is 'scitegicadmin' and the password 'scitegic'.

Under **Admin Pages**, open the **Setup** folder and click on **Global Properties**.

From the Package drop-down, select **CCDC/CSD Pipeline Pilot Collection**.

Select the `PYTHON_HOME` property by clicking on the name and set the value to the appropriate path. If the CSD installation is as above, this would be:

Windows: `C:\Program Files\CCDC\Python_API_2022\miniconda`

Linux: `/opt/CCDC/Python_API_2022/miniconda/bin`

All the other values can be left blank.

### 2.3.4 Uninstallation on Windows

Open a PowerShell prompt as Administrator and navigate to the PPS apps folder.

Use the Pkgutil tool to uninstall the package:

```
..\bin\pkgutil.exe -u cc当地/PythonAPI
```

Delete the cc当地 folder.

### 2.3.5 Uninstallation on Linux

Open a terminal and navigate to the apps directory under the PP server installation.

Set up the environment for the Pkgutil tool:

```
.../linux_bin/ppvars.sh
```

Use the Pkgutil tool to uninstall the package:

```
../linux_bin/pkgutil -u ccdc/pythonapi
```

Remove the ccdc folder.

## 2.4 About the CSD PP Component Collection

This collection was developed in partnership with Finia Consulting.

Finia Consulting was founded in 2013 and is a small software consultancy, specialised in using the Pipeline Pilot platform to develop complex protocols and components including custom components developed using Java and C# and can also provide custom and specialized training for the Pipeline Pilot platform. Finia Consulting have customers in the life sciences, chemicals and academic sectors.

Finia Consulting may be contacted at:

Web: [www.finiasconsulting.com](http://www.finiasconsulting.com)

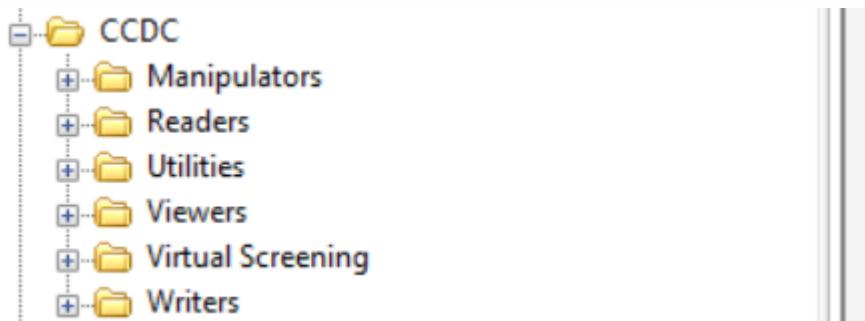
Email: [pcochrane@finiasconsulting.com](mailto:pcochrane@finiasconsulting.com)

Phone: +44 118 981 5993

## 3 CSD Python API Components

The CSD Component Collection consists of a series of components designed for regular usage. These are underpinned by utilities components which can be used to perform lower-level operations. For most day to day usage, it is expected that the higher-level components will be sufficient. A series of protocols are provided with CSD PP Component Collection that show how to combine the different components in a workflow and allow to perform several complex operations such as merged, compared, and processed, according to the logic of the protocol.

The components provided within the CSD PP Component Collection are included in the CCDC folder in the PP Components. The CCDC components are organized in broader categories related to the type of operation performed by the included components such as: Manipulators, Readers, Utilities, Viewers, Virtual screening and Writer. Note that the Virtual Screening and the Writers components are considered as high-level components and are available for CSD-Discovery and CSD Discovery or CSD-Materials users only.

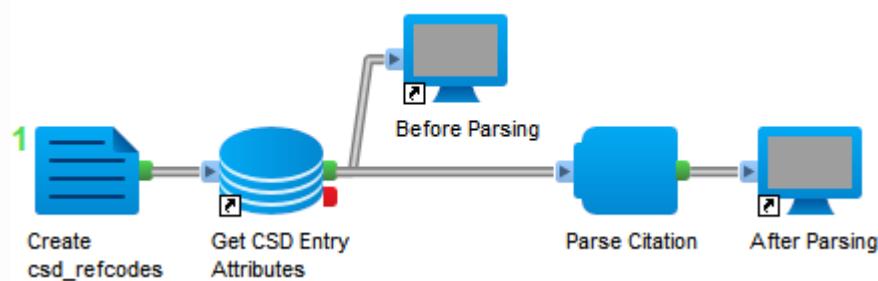


## 3.1 Manipulators

The `manipulator` folder includes components designed to make specific changes to data record properties.

### 3.1.1 Parse Citation

A citation record is retrieved from the CSD as an array property, the values of which represent the various fields of the record (authors, journal etc.). The Parse Citation component takes this property and produces a separate property for each field.



After Parsing						
csd_refcode	authors	journal	volume	year	first_page	doi
AHEPUY	I.D.H.Oswald, W.D.S.Motherwell, S.Parsons, C.R.Pulham	Acta Crystallographica Section E: Structure Reports Online [2001-2014]	58	2002	1290	10.1107/S1600536802018111

Before Parsing	
csd_refcode	publication
AHEPUY	I.D.H.Oswald, W.D.S.Motherwell, S.Parsons, C.R.Pulham Acta Crystallographica Section E: Structure Reports Online [2001-2014] 58 2002 1290 10.1107/S1600536802018111

## 3.2 Readers

There are two type of readers in the CSD Component Collection: those that perform a search on the CSD database (names suffixed 'Search') and those that get attributes for a specific type of object from the CSD (names prefixed 'Get'). The former take a query and return the refcodes of matching database entries; the latter take refcodes as input and add the selected attributes.

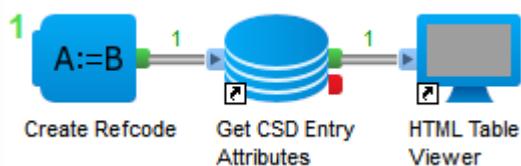
Three types of object are represented, which conceptually hold different types of data relating to a crystallographic experiment. The Entry is the 'top level' object and hold information such as the literature citation and deposition date. The Crystal object holds information such as the unit cell parameters and the 3D coordinates (in CTAB or MOL2 format); note that coordinates sets are available with and without any disordered atoms present. The Molecule object handles derived molecular properties, including the SMILES string where available (it is not computed for disordered structures).

Note that, for convenience, some data are available via multiple objects; an example is the 3D coordinates without disordered atoms, which are available as 'Molecule CTAB' via the Crystal object or as 'CTAB' via the Molecule object.

As the 'Get' components are perhaps simpler, they are described first.

### 3.2.1 Get CSD Entry Attributes

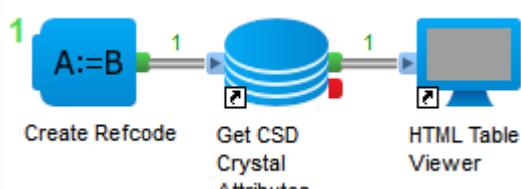
The **Get CSD Entry Attributes** component retrieves Entry attributes for the supplied refcodes; these attributes include the literature citation, deposition date, description of the sample and experimental details.



csd_refcode	bioactivity	color	deposition_date	habit	publication	synonyms
AHEPUY	paracetamol is an analgesic drug	colorless	2003-02-06	plate	I.D.H.Oswald, W.D.S.Motherwell, S.Parsons, C.R.Pulham Acta Crystallographica Section E: Structure Reports Online [2001-2014] 58 2002 1290 10.1107/S1600536802018111	Paracetamol morpholine Acetaminophen morpholine

### 3.2.2 Get CSD Crystal Attributes

The **Get CSD Crystal Attributes** component retrieves Crystal attributes for the supplied refcodes; these attributes include the crystal cell lengths and angles, lattice centring information and the 3D coordinates (with disordered atoms included or excluded).

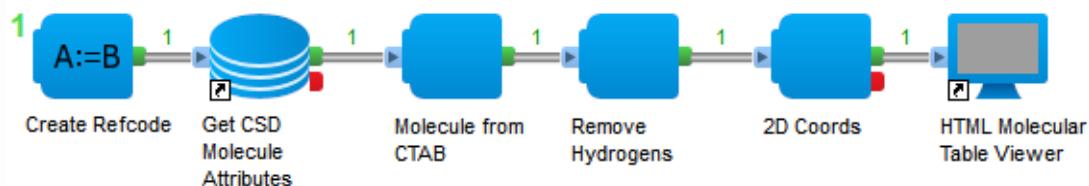


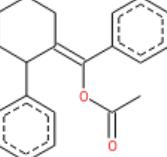
csd_refcode	Cell Angles	Cell Lengths	Cell Volume	Lattice Centring	Packing Coefficient	Void Volume	Z Prime	Z Value
AACMHX10	90 90 90	24.157 16.758 8.5350	3455.2	primitive	0.66106	0	1	8

### 3.2.3 Get CSD Molecule Attributes

The **Get CSD Molecule Attributes** component retrieves Molecule attributes for the supplied refcodes; these attributes include the 3D coordinates in CTAB or MOL2 format (without disordered atoms) and various derived molecular properties, including the SMILES string where available (it is not computed for structure with disordered atoms).

In the example below, a Pipeline Pilot molecule is created from the CTAB and used to generate a depiction via the HTML Molecular Table Viewer. As 2D depictions from 3D coordinates (as stored in the CSD) are generally unsatisfactory, standard Remove Hydrogens and 2D Coords components are included before the Viewer component.



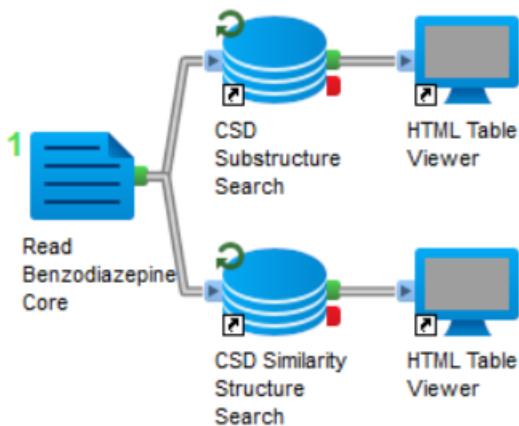
Molecule	csd_refcode	formal_charge	formula	heavy_atoms	molecular_weight	smiles
	AACMHX10	0	C21 H22 O2	23	306.40	CC(=O)OC(=C1CCCCC1c1ccccc1)c1ccccc1

### 3.2.4 Get Molecule Structure

The **Get Molecule Structure** component is a convenience version of the **Get CSD Molecule Attributes** component, with only the CTAB retrieved and converted into a molecular object.

### 3.2.5 CSD Substructure Search

The **CSD Substructure Search** component performs a substructure or exact match search on the CSD. Incoming molecules are treated as queries (mol2 and sdf format are supported). The records for which CSD compounds are found, is output with the CSD refcode, as well as the data for the incoming query structure for which this was a hit. This can be useful when passing in multiple queries.



The example above shows how the **CSD Substructure Search** and CSD Similarity Structure Search components can be combined to read a mol file and perform a substructure and a similarity search in CSD, retrieving the results in two distinct HTML tables.

### 3.2.6 CSD Similarity Structure Search

The **CSD Similarity Structure Search** component performs a similarity search on the CSD. Incoming molecules are treated as queries (mol2 and sdf format are supported). The records for which CSD compounds are found are output with the CSD refcode, along with the data from the incoming record for which this refcode was a hit. This can be useful when passing in multiple queries.

Similarity searches take into account the similarity threshold which all hit structures must exceed. The similarity threshold and the search filters can be edited in the **Parameters** section of the component.

**Parameters**

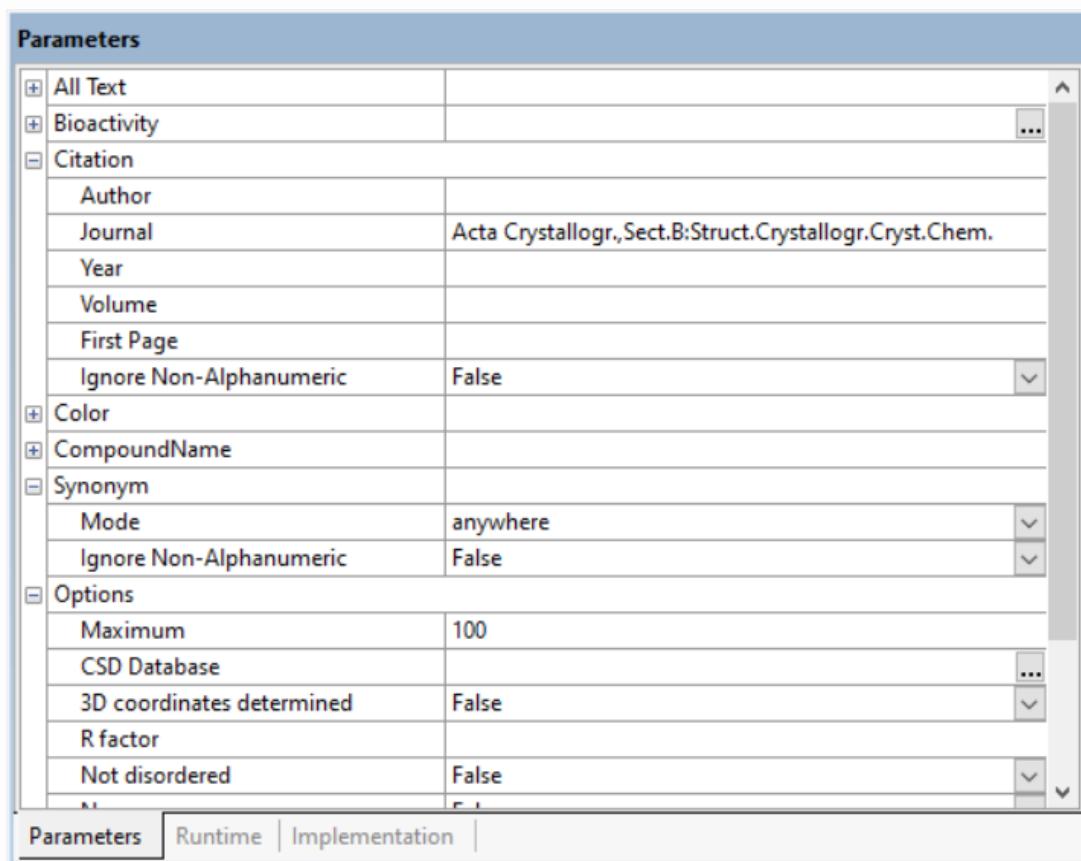
Search Type	Similarity
Options	
Maximum	
CSD Database	...
3D coordinates determined	False
R factor	
Not disordered	False
No errors	False
Not polymeric	False
No ions	False
No powder structures	False
Only	Organics   Organometallics
Similarity Threshold	0.7

Parameters    Runtime    Implementation

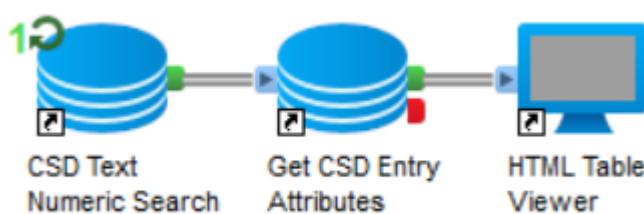
To aid in the interpretation of the results, the component outputs have an additional property which is the similarity value for this compound to the query for which it was found.

### 3.2.7 CSD Text Numeric Search

**CSD Text Numeric Search** component performs a text numeric search against the CSD and produce a stream of CSD refcodes for the hits found by the search. The query is built up from the criteria entered in the component's **Parameters** section.



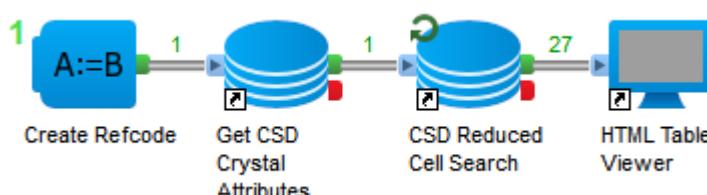
Each element is combined in an **AND** fashion. For example, specifying an author name and a journal will only find an entry with the specified author is in the specified journal.



In the example above, the **CSD Text Numeric Search** component is used to search the CSD for all crystals reported in the *Acta Crystallogr., Sect.B:Struct.Crystallogr.Cryst.Chem.* journal, limiting the search to 100 records by setting the **Maximum** parameter in the **Options** section of the **Parameters** section of the component.

### 3.2.8 CSD Reduced Cell Search

The **CSD Reduced Cell Search** component performs a reduced cell search of the CSD. Reduced cell searches can be carried out in two ways.

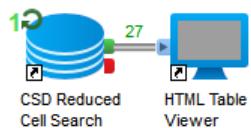


The default is to use data record properties for the cell lengths and angles, as returned by the **Get CSD Crystal Attributes** component are used:

Reduced-cell search results	
csd_refcode	database
AACMHX10	as541be_ASER
BIGRIT	as541be_ASER
BOLCUA01	as541be_ASER
BOLCUA11	as541be_ASER
FIDMIR	as541be_ASER
GEBYIV	as541be_ASER
JENMOE	as541be_ASER

Cell Lengths	
Lengths String	
a	24.157
b	16.758
c	8.535
Cell Angles	
Angles String	
alpha	90
beta	90
gamma	90

The alternative is to enter the a, b, c and alpha, beta, gamma values separately in the appropriate parameters:



With the refcodes found by this search, the chemical name is retrieved and passed to the **Get CSD Entry Attributes** component.

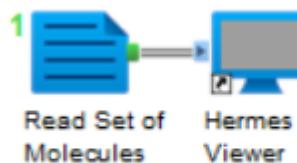
## 3.3 Viewers

A viewer is a component that displays information or results on your client. Viewers are frequently used as the final component in a pipeline however, they can be also used to view intermediate results.

In addition to the Generic Viewer provided by Pipeline Pilot, the CSD PP Component Collection includes four viewer components that are useful when using the CSD components.

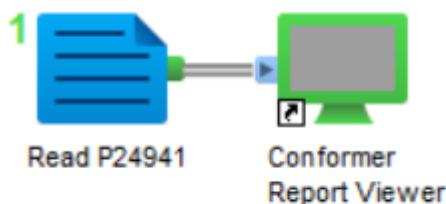
### 3.3.1 Hermes Viewer & Mercury Viewer

**Hermes Viewer & Mercury Viewer** components allow to view the incoming stream of data in Hermes and Mercury, respectively. To aid in identifying the structures loaded, the name of the MOL2 file which will be passed to Hermes or Mercury can be supplied using the **Dataset Name** option of the component's parameters or can be provided as a reader component as showed in the example below.



### 3.3.2 Conformer Report Viewer

The **Conformer Report Viewer** component is available for users with a CSD-Discovery or CSD-Materials licence. This component generates conformers for the incoming file of molecule(s) (SDF and MOL2 formats are supported) and produces a report summarizing the process. This means that the component will return a HTML summary of the settings used, a summary of the conformer generation results and links to the file(s) produced. By default, the summary file and the conformers outputs are generated.



From the **Parameters** section of the **Conformer Report Viewer** component it is possible to change both the **Conformer Options** such as **Max Number of Conformers** and the **Output Options** of the component.

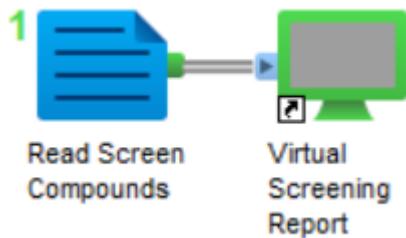
**Parameters**

<b>Conformer Options</b>	
Max Number of Conformers	25
Number of Threads	1
Maximum Unusual Torsions	2
Superimpose	False
<b>Output Options</b>	
Output What	Molecules   Summary
Split Output	False
Report Title	Conformer Generation Report
Reporting Options	Show Plot in Browser

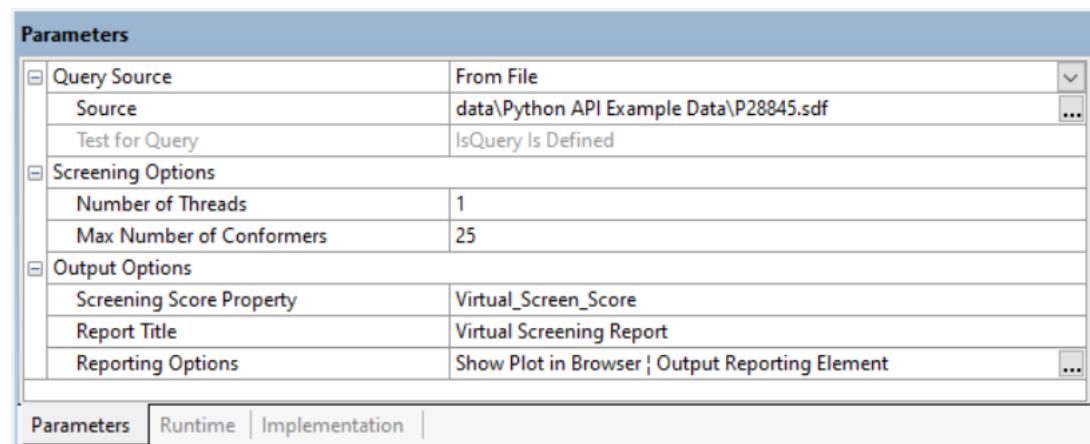
Parameters    Runtime    Implementation

### 3.3.3 Virtual Screening Report Viewer

**Virtual Screening Report Viewer** component is available only for users with a CSD-Discovery licence. This component produces an HTLM report from the virtual screening process. This component receives a stream of incoming molecules and screens against a query set to generate a virtual screening score.



The query set may be supplied as file, specifying the **Source** in the **Parameters** section of the component. The query set file can be supplied either as a MOL2 or SD file.



The **Virtual Screening Report Viewer** component may also receive a stream of query records which are tagged to differentiate them from the screening set. This is achieved by setting the **Query Source** parameter to **From Tag** and declaring the PilotScript which differentiates the query records from the screening records.

The resulting HTML report contains the settings used for the virtual screening, plus links to the files used (screening and query), as well as the results file produced.

The **Virtual Screening Report Viewer** component can be used as part of a larger report by using the option to output reporting elements setting in the **Reporting Options** parameter of the component to include **Output Reporting Element**.

## 3.4 Virtual Screening

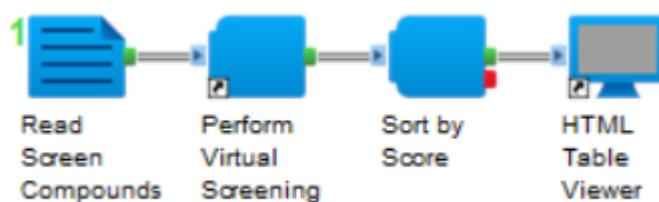
The Virtual screening components includes different components useful to perform ligand based virtual screening but also components that allow to assess the validity of the methodology.

Note that the following components are only available for users with a CSD-Discovery licence, while the **Perform Conformer Generation** component is available for users with CSD-Discovery and/or CSD-Materials licence:

- **Perform Virtual Screening**
- **Perform Virtual Screening Validation**
- **Generate Enrichment Plot**
- **Generate ROC Plot**

### 3.4.1 Perform Virtual Screening

This component receives a stream of incoming records and applies a set of query molecules to that screening set to generate a virtual screening score. The query set may be supplied either as a MOL2 or SD file, specified using the **Source** parameter in the component, or, it may receive a stream of query records which are tagged to differentiate them from the screening set. This is achieved by setting the **Query Source** parameter to From Tag and declaring the PilotScript which differentiates the query records from the screening records.

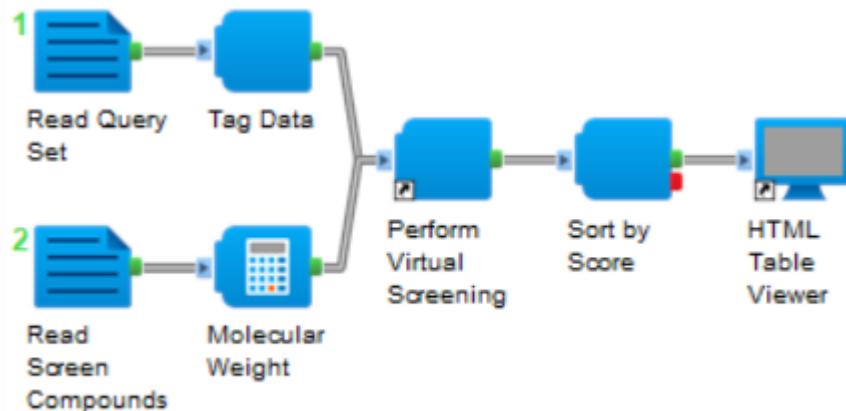


**Parameters**

Query Source	From File	...
Source	data\Python API Example Data\P28845.sdf	...
Test for Query	IsQuery Is Defined	
Screening Options		
Number of Threads	1	
Max Number of Conformers	25	
Output Options		
Screening Score Property	Virtual_Screen_Score	
Ranked Score File	...	

Parameters    Runtime | Implementation

The above example shows the component being used in a manner where the query file already exists and is specified in the **Source** parameter.



**Parameters**

Query Source	From Tag	...
Source		
Test for Query	IsQuery Is Defined	
Screening Options		
Number of Threads	1	
Max Number of Conformers	25	
Output Options		
Screening Score Property	Virtual_Screen_Score	
Ranked Score File	...	

Parameters    Runtime | Implementation

However, as the above example demonstrates, the stream of records can arrive at the component with records tagged appropriately. The component then internally divides the records according to the PilotScript found in **Test for Query** to produce a query set and screening set.

Whether acting on a pre-existing query source file or dynamically generating the query set by tag, the component has the same control over the underlying Python script. These are exposed in the **Screening Options** group parameters.

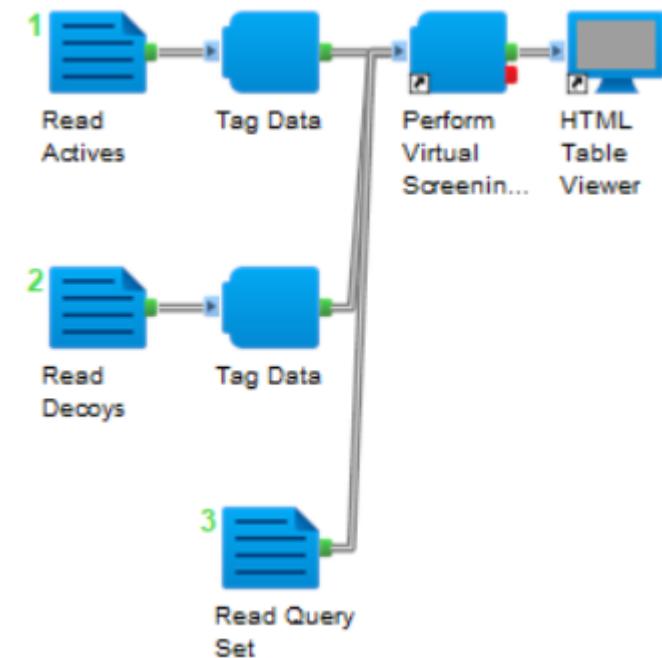
Finally, the data which is output by the component has a new virtual screening score which must be stored in a property. The component allows the user to set the name of this property using **Screening Score Property** parameter.

### 3.4.2 Perform Virtual Screening Validation

The **Perform Virtual Screening Validation** component takes a stream of query records and compares these to known actives and known decoys. This is used to assess the selectivity of the screen by producing a stream of the actives and decoys with scores. The scores are sorted from lowest to highest, where the lowest being the best and the occurrences of genuine actives vs decoys are tagged to help identify them.

As with the **Perform Virtual Screening** component, this component can work with actives and decoys in files specified in the appropriate source parameter, or it can identify the actives and/or decoys in-line.

The data produced has two new properties; one for the active tag (a value of 1 in this property implies the record related to the active set) and another for the virtual screening score. Therefore, the **Output Options** group has parameters to control the naming of these properties; **Active Tag** and **Screening Score Property** respectively.



**Parameters**

<b>Input Options</b>	
<b>Active Source</b>	From Tag
Source	
<b>Test for Active</b>	IsActive Is Defined
<b>Decoy Source</b>	From Tag
Source	
<b>Test for Decoy</b>	IsDecoy Is Defined
<b>Screening Validation Options</b>	
<b>Number of Threads</b>	1
<b>Max Number of Conformers</b>	25
<b>Output Options</b>	
<b>Active Tag</b>	Is_Active
<b>Screening Score Property</b>	Virtual_Screen_Score

Parameters    Runtime    Implementation

The above example shows the use of the **Perform Virtual Screening** Validation component where both the actives and decoys are identified by tag - both **Active Source** and **Decoy Source** are set to From Tag, with the **Test for Active** being the PilotScript IsActive Is Defined, whilst the **Test for Decoy** parameter expressing the PilotScript IsDecoy Is Defined. Anything which neither tag is declared is assumed to be a query record.

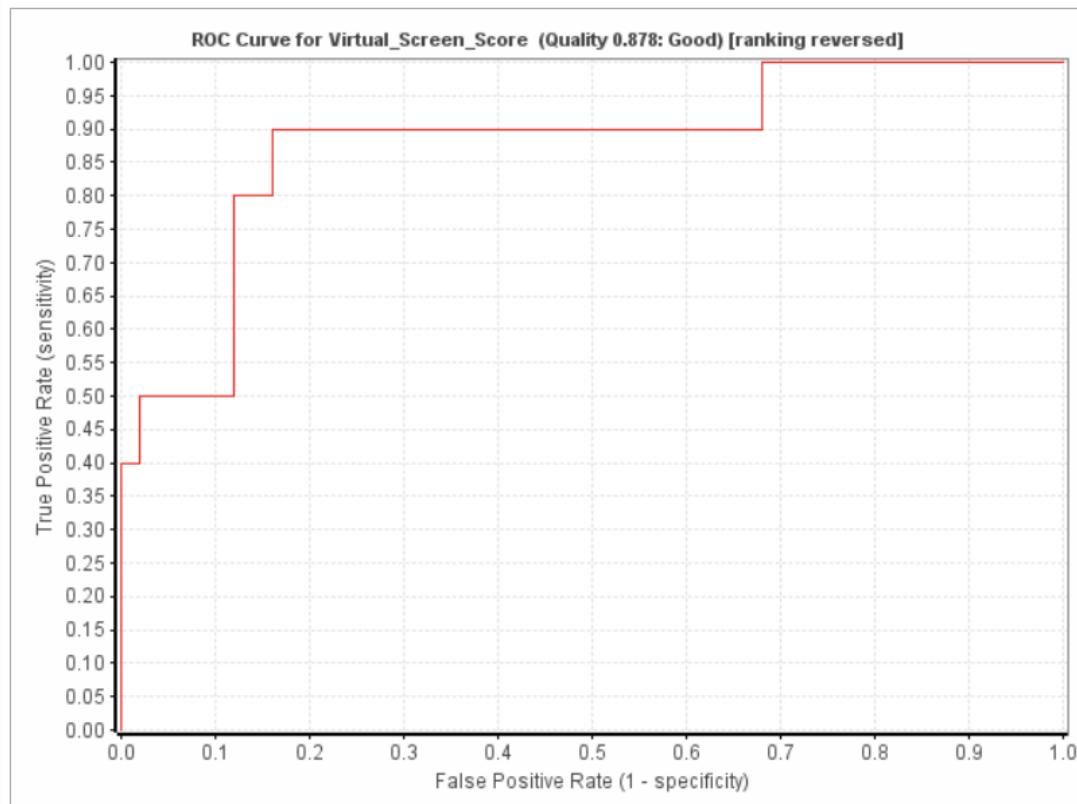
As with the previous component, the script arguments are expressed as options in the **Screening Validation Options** group parameters.

Finally, the **Output Options** include the names to be given to the **Active Tag** (Is\_Active in the example above) and to the **Screening Score Property**.

### 3.4.3 Generate Enrichment Plot & Generate ROC Plot

The **Perform Virtual Screening** Validation component produces “raw” data, data which is more normally presented to the user in a visual form. Therefore, the **Generate Enrichment Plot** component and Generate ROC Plot component exist to present that data in a chart. To that extent, both components are simple extensions of the **Perform Virtual Screening** Validation component, exposing almost identical parameters to the user.

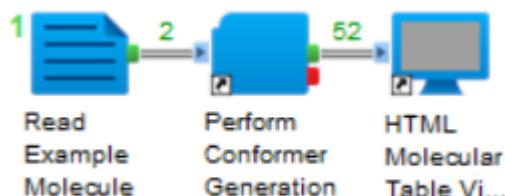
We can then use the **Generate ROC Plot** component to generate the ROC plot for the virtual screen study, that together with Generate Enrichment Plot components provides information about the selectivity of your model.



The main difference between the basic validation component and Generate Enrichment Plot and **Generate ROC Plot** components is that these components can either produce a chart in the browser or produce reporting component which can be incorporated as part of a larger report.

### 3.4.4 Perform Conformer Generation

The **Perform Conformer Generation** component is only available for user with CSD-Discovery or CSD-Materials licence. The Perform Conformer Generation component takes a stream of molecular records and generates multiple conformations for each structure. The component outputs a stream of data where each record represents a conformer for an individual. The number of conformers is limited by the **Max Number of Conformers** in the **Conformer Options** parameter.



In the above example, two records are read in, and for each record, up to twenty-five conformers are generated (default value). Note that if the conformer generation script cannot find the maximum number of conformers, then it will output the number it finds. So, with a higher limit on the maximum number, it's possible to see different numbers of records per starting record. This conformer generation is also found under the virtual screening processes. As with the virtual screening components, this component will ensure that any data on the incoming data is re-attached to the output data.

## 3.5 Utilities Components

There are four internal utilities components that can be used to perform lower-level operations. These include:

- **Run Python Script.**
- **Run Virtual Screening.**
- **Run Virtual Screening Validation.**
- **Run Conformer Generation.**

The first component is capable of running any Python script, whilst the second runs the virtual screening script, and the third runs the virtual screening validation script. The final component generates conformers. These are designed to mirror the execution of these scripts on the command line.

### 3.5.1 Run Python Script

The **Run Python Script** is used internally, either directly or indirectly, in all components and protocols found in the CSD PP Component Collection package. The package contains a single global variable (e.g. @/ccdc/pythonapi/python\_home) which declares the directory in which the python executable is found or installed on the server.

This component is intended for simple usage, and to test that the package is configured correctly. The **Run Python Script** component uses the @/ccdc/pythonapi/python\_home global variable, combined with either a script file, or the content of a script which is written out as a file, together with the arguments to perform whatever task the script is designed to do.

The **Run Python Script** component produces potentially three properties, standard out for the script, standard error for the script and the command executed. The final option is useful for debugging as it can reveal when the script was invoked incorrectly.

The script to be executed can be specified by defining the entire script in the **Script** parameter declared in the **Python** group parameters. The script is entered here as any python script would be written. Internally, the component writes this script out to a file before executing it. Alternatively, the script can be specified as being in a file, located by the **Script File** parameter in the same group. Note that either one or the other of these two parameters must be supplied.

As it is often the case that a script may take a significant amount of time to run, the component exposes a **CustomMessageToClient** parameter to set the message to display whilst the script is running.

The **Arguments** group parameter is an array group which allows the creation of multiple arguments to be supplied to the python script. The arguments take on a pairing of **Switch** and **Value** - where switch is supplied first, followed by the value. Note that group arrays of this type natively support reordering if the order of arguments is important.

Python	
Script	print "Hello World!";
Script File	
CustomMessageToClient	
Arguments	
Output	
Stdout Property Name	script_out
Stderr Property Name	script_err
Command	False

The above example shows a very simple Python script, supplied as a script in the **Script** parameter. This will produce the very simple **Hello World!** message in the property specified in the **Stdout Property Name** script\_out.

### 3.5.2 Run Virtual Screening

This component is a low-level wrapper over the Python script which performs the screening operation. The **Run Virtual Screening** component performs the screening operation found in the Perform Virtual Screening component.

Whereas **Perform Virtual Screening** component handles data and generally ensures the quality of the output, this component simply produces the raw output from the underlying script. For example, this means that the data produced will have lost and additional properties found in screening file. As such, this component can be thought of as being equivalent to the execution of the script from the command line.

Parameters	
Screening Set	...
Query Set	...
Screening Options	
Number of Threads	1
Max Number of Conformers	25
Output Options	
Screening Score Property	Virtual_Screen_Score
Ranked Score File	...

The **Run Virtual Screening** component expects to have a file containing the screening records (**Screening Set** parameter) and another containing the query records (**Query Set** parameter), and exposes parameters which map to arguments for the script. It performs basic validation on these - checking that files are a supported format (by checking for file extensions of .sdf or .mol2), if the files exist, and checking that the content of these files matches the extension (e.g. the .sdf file is a genuine SD file).

In terms of data handling, the component will join the scores produced by the MOL2 file containing the structures screened.

### 3.5.3 Run Virtual Screening Validation

This component is a low-level wrapper over the python script which performs the screening validation operation. The **Run Virtual Screening Validation** component performs the virtual screening validation operation found in the **Perform Virtual Screening Validation** component. As with the **Run Virtual Screening** component, it is designed to simply verify the input for the script (by checking file formats, if they are found), run the script, then join the data produced (MOL2 and CSV files) together. The **Run Virtual Screening Validation** component expects to have a file containing

the query records (**Query Set** parameter), the active records (**Active Set** parameter) and another containing the decoy records (**Decoys Set** parameter).

### 3.5.4 Run Generate Conformers

This component takes a file (either MOL2 or SDF) and generates conformers for each record found in that file. The number of conformers is limited by the **Max Number of Conformers** parameter. The component can work in multiple ways, according to the parameters.

The most basic operation, the default one, is to simply stream out the result of the conformer generation for the input file. The python script generates two outputs. The first is an SD file containing the molecules and the second file is simply a summary which details how many conformations were found for each molecule, and information about the conformers:

In addition to this, the component can generate a single file (by default) or single files containing all the conformers for a given molecule. In this case, the file name reflects the name of the molecule to which the conformers belong. To generate the files, the **Output Directory** is specified. The files will be written to this location, and the molecule can be merged (**Split Output** parameter set to **False**) or split (**Split Output** parameter set to **True**).

### 3.5.5 Derive Script Path

The CSD PP Component Collection package has a script folder which contains the python scripts used to execute on the CSD. To avoid having to reference this folder every time, there is a the Derive Script Path component, which knows where this folder is and will gather the path to the script named in its **Script Filename** parameter.

This component will do more than this though. It may be that different instance of CSD require a different script, due to the CSD Python API changes. Therefore, it will examine the CSD instance

and determine the best script to use. To do this, scripts must have the version for which they are appropriate in the filename: `script.`

`1.2.3.py`

Rather than have to second guess the version of the file to use, or know the versions available, the script can be referred to as:

`script.py` however, it will point to `script.1.2.3.py`

This component will examine all versions of that script, either in the Python API package (if **Script Folder** parameter is blank) or the folder specified in the **Script Folder** parameter and finds the most recent version which is appropriate for this script. For example, if a CSD 0.8.0 version of the script is found, then an earlier version of the script may be used if 0.8.0 is not found (e.g. CSD 0.7.0).

If a single version of the script is not found, this is assumed to mean that `script.py` is appropriate for all versions of CSD and will therefore be used.

### **3.5.6 Throw Script Error Message**

The **Throw Script Error Message** component is a helper component to parse script errors. In particular, this component attempts to find Runtime errors which may relate to licensing and displays those in a better formatted manner - where the Runtime error is positioned at the top of the message to make it clearer.

### **3.5.7 Check Journal Name**

The **Check Journal Name** component will check a journal title against the CSD. This utility component performs a text numeric search on the journal name provided in the **Journal Name** parameter.

### **3.5.8 Validate Journal Name**

The **Validate Journal Name** component checks the validity of a journal title against the CSD. This utility component performs a text numeric search on the journal and returns a **True** or **False** value.

### 3.5.9 Gather Database Names

The **Gather Database Names** component gather names list and split them into array. It checks for each value for being not empty and additionally that the associated file exists. If all the conditions are satisfied than the name goes forward to the return value.

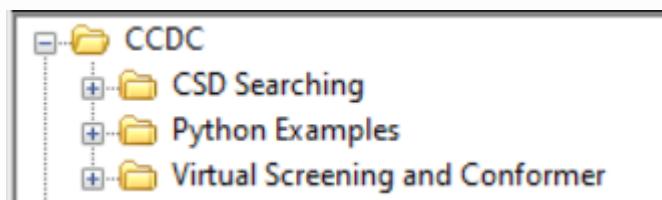
### 3.5.10 Join Data from JSON

The **Join Data from JSON component** joins a JSON-encoded stream of data that can be specified in the **Source** parameter. The source can be a file, a network resource, a data record property or a global property. The name of the parameter to join can be specified in the **JoinUsing** parameter.

## 4 CSD Python API Protocols

The CSD PP Component Collection package is supplied with examples designed to demonstrate the use of the different components. No examples are supplied to demonstrate the use of the utilities components, though they are fully tested and their usage can be determined from the examples.

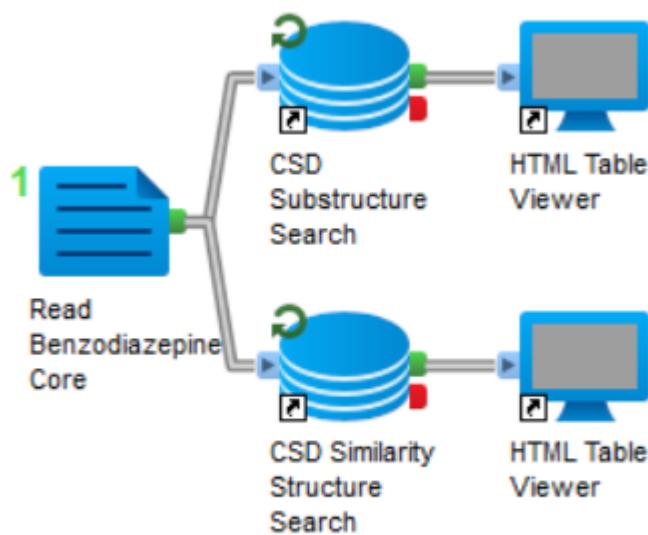
The CSD Python API Protocols are listed under the **CCDC** folder in the **Protocols** tab. The provided protocols are organised in groups based on their purpose: **CSD Searching**, **Python Examples** and **Virtual Screening and Conformer**.



### 4.1 CSD Searching

The **CSD Searching** protocol folder includes seven different ways to search and to combine searches in the CSD using the CSD Python API.

#### 4.1.1 Search CSD By Structure



This protocol uses the benzodiazepine core structure, shipped as part of the Chemistry collection (data\Queries\BenzodiazepineCore.mol), and searches the CSD for structures which contain this as a substructure (top) or as a similar compound (bottom). The query file is specified in the **Source** parameter of the Read Benzodiazepine Core component. The filters of the search are specified in the **Options** parameter of the CSD Substructure Search and **CSD Similarity Structure Search** components. The two searches outputs the results in two HTML pages.

The Substructures HTML page will provide a table containing as a minimum the refcode(s) for hit structures.

Display records  to  of 142

<<First <Previous [Next](#) [Last](#)>

Substructures		
csd_refcode	database	Name
ANIXAX	as540be_ASER	BenzodiazepineCore
AZEGIX	as540be_ASER	BenzodiazepineCore
AZEGOD	as540be_ASER	BenzodiazepineCore
AZEGUJ	as540be_ASER	BenzodiazepineCore
BATDIL	as540be_ASER	BenzodiazepineCore
BAYCUZ	as540be_ASER	BenzodiazepineCore
BAZCLH	as540be_ASER	BenzodiazepineCore
BCHBZP	as540be_ASER	BenzodiazepineCore
BCHBZP01	as540be_ASER	BenzodiazepineCore
BEDZPN10	as540be_ASER	BenzodiazepineCore
BIXBOZ10	as540be_ASER	BenzodiazepineCore
BIZSAE	as540be_ASER	BenzodiazepineCore
BOFDEF	as540be_ASER	BenzodiazepineCore
BOMMUL	as540be_ASER	BenzodiazepineCore
BORXEL	as540be_ASER	BenzodiazepineCore
BUVLEK	as540be_ASER	BenzodiazepineCore
BZPCUC	as540be_ASER	BenzodiazepineCore
CAGWUC	as540be_ASER	BenzodiazepineCore
CASKIT	as540be_ASER	BenzodiazepineCore
CHABZN	as540be_ASER	BenzodiazepineCore
CICSEM	as540be_ASER	BenzodiazepineCore
CLDZPA	as540be_ASER	BenzodiazepineCore
CLDZPA01	as540be_ASER	BenzodiazepineCore
CLDZPB	as540be_ASER	BenzodiazepineCore
COWHIG	as540be_ASER	BenzodiazepineCore

Display records  to  of 142

<<First <Previous [Next](#) [Last](#)>

While similarity search, in addition, take into account the similarity threshold - which all hit structures must exceed. The similarity threshold is specified in the **Similarity Threshold** parameter. To aid in the interpretation of the results, the CSD Similarity Structure Search component outputs an additional property which is the similarity value for this compound to the query for which it was found.

Display records  to  of 42

<<First <Previous [Next](#) [Last](#)>>

Similar			
csd_refcode	database	similarity	Name
DCDAZP	as540be_ASER	0.986187845304	BenzodiazepineCore
BODSOC	as540be_ASER	0.87675070028	BenzodiazepineCore
DOZKEI	as540be_ASER	0.832167832168	BenzodiazepineCore
DIZPAM10	as540be_ASER	0.822580645161	BenzodiazepineCore
DIZPAM11	as540be_ASER	0.822580645161	BenzodiazepineCore
ZZZBNV	as540be_ASER	0.822580645161	BenzodiazepineCore
LUTGIR	as540be_ASER	0.81880733945	BenzodiazepineCore
BAZCLH	as540be_ASER	0.816933638444	BenzodiazepineCore
RUMWAZ	as540be_ASER	0.813211845103	BenzodiazepineCore
VERZEX	as540be_ASER	0.78982300885	BenzodiazepineCore
CERBEG	as540be_ASER	0.780193236715	BenzodiazepineCore
CERBEG01	as540be_ASER	0.780193236715	BenzodiazepineCore
BORXEL	as540be_ASER	0.774403470716	BenzodiazepineCore
SUXFUM	as540be_ASER	0.774403470716	BenzodiazepineCore
SUXFOG	as540be_ASER	0.771058315335	BenzodiazepineCore
FULWUE	as540be_ASER	0.762820512821	BenzodiazepineCore
LUTGAJ	as540be_ASER	0.757961783439	BenzodiazepineCore
NHPBZO	as540be_ASER	0.757961783439	BenzodiazepineCore
QUGGUU	as540be_ASER	0.757961783439	BenzodiazepineCore
QUGHAB	as540be_ASER	0.757961783439	BenzodiazepineCore
FLDAZP	as540be_ASER	0.753164556962	BenzodiazepineCore
OXAPAM10	as540be_ASER	0.742203742204	BenzodiazepineCore
VICFIW	as540be_ASER	0.734567901235	BenzodiazepineCore
ZZZAUS10	as540be_ASER	0.734567901235	BenzodiazepineCore
ZZZAUS20	as540be_ASER	0.734567901235	BenzodiazepineCore

Display records  to  of 42

<<First <Previous [Next](#) [Last](#)>>

#### 4.1.2 Search CSD By Text Numeric Fields



This protocol shows how the CSD can be searched using text and numeric fields. In this particular case, two text and numeric searches are performed.

The first **CSD Text Numeric Search** component looks for all crystals reported in *Acta Crystallogr., Sect.B:Struct.Crystallogr.Cryst.Chem* as specified in the **Citations** parameter. This is limited to 100 records by setting the **Maximum** parameter.

First 100 for Acta Crystallogr.,Sect.B:Struct.Crystallogr.Cryst.Chem.		
csd_refcode	publications	
AADAMC	(Citation authors uX K. Chacko, R.Zand, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u29, year u1973, first_page u2681, doi u10.1107/S056774087307363.)	as540be_ASER
AAMTPX	(Citation authors uX M. V. Vlachos, V. S. Sheldrick, J. Karolak-Wojciechowska, M. Mikalajczyk, B. Zemnicki, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u35, year u1979, first_page u2339, doi u10.1107/S0567740873092529.)	as540be_ASER
AANHOX	(Citation authors uX M. R. Cipolla, F. Lalić, T. Tanecki, P. A. Temussi, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u37, year u1981, first_page u1130, doi u10.1107/S0567740881005256.)	as540be_ASER
AAZTHP	(Citation authors uX B. Kojic-Prodic, V. Rogic, Z. Ruzic-Toros, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1976, first_page u1833, doi u10.1107/S0567740876006481.)	as540be_ASER
ABACUH10	(Citation authors uX R. Ahmed, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u31, year u1975, first_page u1001, doi u10.1107/S0567740875003291.)	as540be_ASER
ABCLUH10	(Citation authors uX C. B. Clegg, C. Lazarus, S. Orman, A. Nakagawa, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u35, year u1977, first_page u2314, doi u10.1107/S0567740877008341.)	as540be_ASER
ABCNCX	(Citation authors uX S. K. Bhattacharya, R. Chacko, R. Zand, R. Water, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u35, year u1975, first_page u395, doi u10.1107/S0567740875003291.)	as540be_ASER
ABCOCX	(Citation authors uX K. K. Chacko, S. K. Bhattacharya, R. Zand, R. Water, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u34, year u1978, first_page u147, doi u10.1107/S0567740878002496.)	as540be_ASER
ABDEC010	(Citation authors uX R. Ramachandra, K. Vijayan, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u29, year u1973, first_page u1945, doi u10.1107/S0567740873005807.)	as540be_ASER
ABFEC010	(Citation authors uX G. C. F. Egan, R. H. Hedges, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u29, year u1973, first_page u1253, doi u10.1107/S0567740873005801.)	as540be_ASER
ABSPON	(Citation authors uX G. C. F. Egan, R. H. Hedges, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u29, year u1973, first_page u1273, doi u10.1107/S0567740873005802.)	as540be_ASER
ABINOR01	(Citation authors uX S. Takagi, S. Norden, G. A. Jeffrey, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u35, year u1977, first_page u981, doi u10.1107/S056774087905379.)	as540be_ASER
ABINOS01	(Citation authors uX S. Takagi, G. A. Jeffrey, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u33, year u1977, first_page u3033, doi u10.1107/S056774087701218.)	as540be_ASER
ABINOS01	(Citation authors uX G. A. Jeffrey, A. Robbins, R. K. McMullan, S. Takagi, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u30, year u1980, first_page u373, doi u10.1107/S0567740880003299.)	as540be_ASER
ABINOS01	(Citation authors uX G. A. Jeffrey, A. Robbins, R. K. McMullan, S. Takagi, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u30, year u1980, first_page u2683, doi u10.1107/S0567740880003299.)	as540be_ASER
ABINOS01	(Citation authors uX S. Takagi, G. A. Jeffrey, R. K. McMullan, S. Takagi, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1981, first_page u2899, doi u10.1107/S0567740881009529.)	as540be_ASER
ABINOS01	(Citation authors uX G. L. Dore, R. C. Sorenson, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1982, first_page u2667, doi u10.1107/S0567740872006484.)	as540be_ASER
ABORCR10	(Citation authors uX F. Taylor Junior, E. A. H. Griffith, E. L. Amme, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1976, first_page u653, doi u10.1107/S0567740876003656.)	as540be_ASER
ABPZOL0	(Citation authors uX J. Lapeyre, A. Escande, J. F. Lapeyre, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1972, first_page u3316, doi u10.1107/S0567740872007897.)	as540be_ASER
ABRPH0	(Citation authors uX H. H. Sutherland, T. H. Hoy, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u25, year u1969, first_page u2386, doi u10.1107/S0567740869005759.)	as540be_ASER
ABRUC010	(Citation authors uX J. Lapeyre, A. Escande, J. F. Lapeyre, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1976, first_page u2074, doi u10.1107/S0567740876004843.)	as540be_ASER
ABSCIC	(Citation authors uX P. Srinivasan, J. V. Jayaram, R. Grimes, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u32, year u1976, first_page u2391, doi u10.1107/S0567740876007747.)	as540be_ASER
ABSFON	(Citation authors uX K. Katai, B. Chaudhuri, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u24, year u1968, first_page u1645, doi u10.1107/S056774088004796.)	as540be_ASER
ABITBR	(Citation authors uX J. Galley, J. P. Declercq, M. van Maarssebeek, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u34, year u1978, first_page u974, doi u10.1107/S056774088004550.)	as540be_ASER
ABTNBA	(Citation authors uX I. Bar, J. Bernstein, journal *Journal/Acta Crystallographica Section B Struct Crystallogr Cryst Chem. [1968-1982]), volume u37, year u1981, first_page u569, doi u10.1107/S0567740880103952.)	as540be_ASER

The second **CSD Text Numeric Search** component also looks in that journal, but now also limits the hits returned to only those where the colour of the crystal is “orange”. This filter is specified in the **Colour** parameter of the second **CSD Text Numeric Search** component. Note here that the maximum is still set, so we still get 100 hits. However, this second search, we only get crystals that are orange, or contain the word orange (e.g. red-orange).

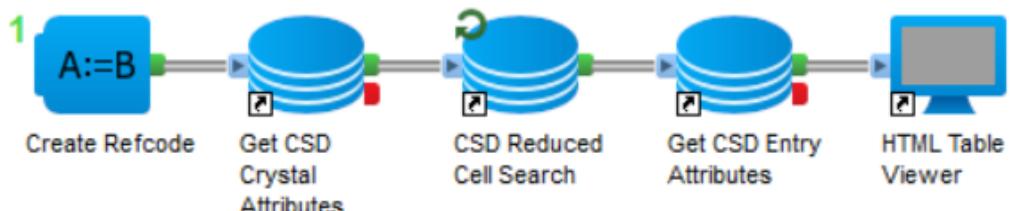
Display records  to  of 100  <<First <Previous [Next](#) [Last](#)>

Orange Crystals in Acta Crystallogr., Sect.B: Struct.Crystallogr.Cryst.Chem.

csd_refcode	color	database
ABTNBA	orange	as540be_ASER
ACDXUS	light orange	as540be_ASER
ACPSMO	orange	as540be_ASER
ACRMSA	orange	as540be_ASER
ACRMSB	red-orange	as540be_ASER
AMALCO10	orange	as540be_ASER
AMAZBRR	orange	as540be_ASER
AMCQUN	orange	as540be_ASER
AMGZCO	dark orange-red	as540be_ASER
AZBNOS	red-orange	as540be_ASER
AZDETO10	orange	as540be_ASER
BADVIL10	yellow-orange	as540be_ASER
BAFXEL	red-orange	as540be_ASER
BAFXUB	orange-red	as540be_ASER
BALACF	orange-red	as540be_ASER
BATCAA01	orange-yellow	as540be_ASER
BAZBIN	orange-red	as540be_ASER
BDMIUP	brown orange	as540be_ASER
BECVEK	red-orange	as540be_ASER
BEJRUD	orange-red	as540be_ASER
BEJSUE	yellow-orange	as540be_ASER
BEPNAL	red-orange	as540be_ASER
BERMIU	yellow-orange	as540be_ASER
BEWKAP	dark orange	as540be_ASER
BEYNEY	pale orange	as540be_ASER

Display records  to  of 100  <<First <Previous [Next](#) [Last](#)>

### 4.1.3 Search CSD By Reduced Cell



This example starts with a Refcode that is specified in the **Expression** parameter of the Create Refcode component. For this Refcode, crystal structure details are retrieved, cell lengths, cell angles, and lattice centring information that are specified in the **Attributes < CSD RefCode Property** parameter of the Get CSD Crystal Attributes component. These details are then used in a reduced cell search.

Reduced cell searches can be carried out in two ways.

In this example, the cell lengths and angles returned by the crystal attributes are used. These are strings of the form:

CellLengths (a=8.4708, b=10.0492, c=14.0363)

CellAngles (alpha=86.016, beta=79.914, gamma=71.818)

The alternative is to enter the a, b, c, and alpha, beta, gamma values separately in the appropriate parameters of the CSD Reduced Cell Search component.

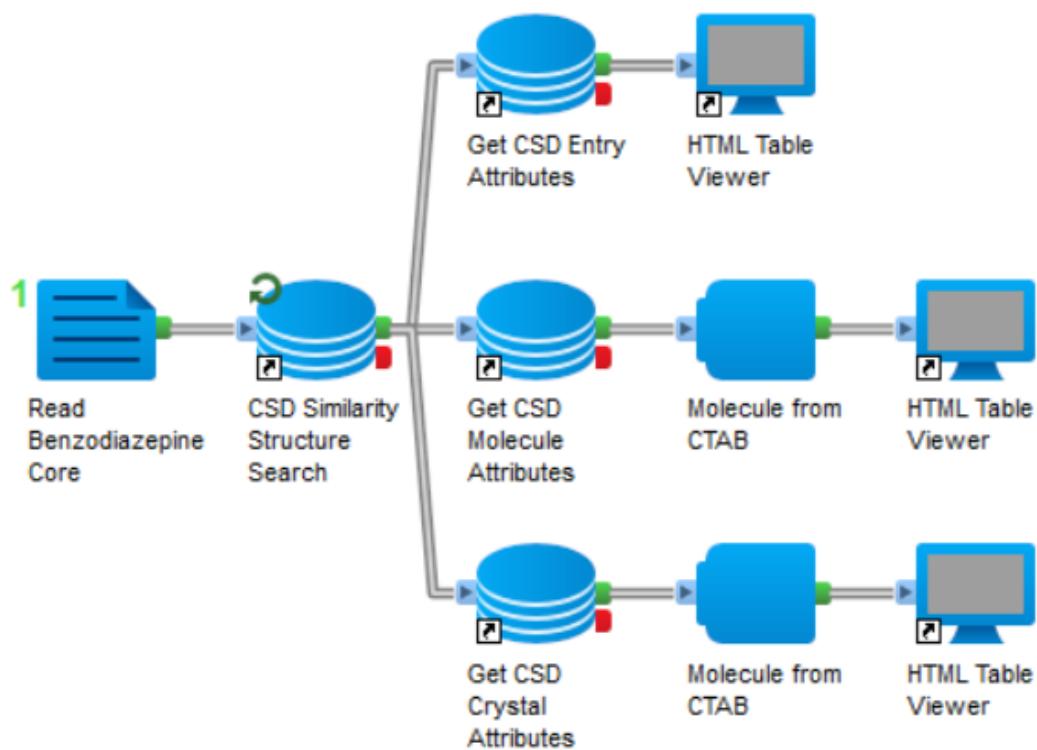
With the refcodes found by this search, the chemical name is retrieved, and the results reported in a HTML table.

Display records  to  of 27  <<First <Previous [Next](#) [Last](#)>>

csd_refcode	chemical_name	database
AACMHX10	-Acetoxy-2-anti-diphenylmethylene-cyclohexane	as540be_ASER
BIGRIT	bis(2-Pyrrolidino)-tetramethyl-di-indium(iii)	as540be_ASER
BOLCUA01	Mesitil	as540be_ASER
BOLCUA11	2,2',4,4',6,6'-Hexamethylbenzil	as540be_ASER
FIDMIR	bis(3,4,7,8-tetrahydro-2H,6H-pyrimido[1,2-a]pyrimidinyl)bis(chloro)diboron	as540be_ASER
GEBYIV	bis((2-Dimethyl phosphonato-P,O)-dicarbonyl-(trifluoroacetato-O,O')-(trimethylphosphite-P)-molybdenum)	as540be_ASER
JENMOE	bis(5-Cyclopentadienyl)-(2-6,6-1,2,3,4-tetramethyl-1,4-dibora-2,5-cyclohexadienyl)-di-nickel	as540be_ASER
MONZAR	4-(4,5-Diphenyl-1H-imidazol-2-yl)benzonitrile	as540be_ASER
PCYPCR	(2,2)Paracyclophane bis-O-18-crown-6 ether	as540be_ASER
PCYPCR01	4,5,15,16-bis(18-Crown-6)(2,2)paracyclophane	as540be_ASER
PCYPCR10	4,5,15,16-bis(18-Crown-6)(2,2)paracyclophane	as540be_ASER
QOBPUT	(R,R)-(6-1-(Pyrrolidinyl)-2-(1-(dimethylamino)ethyl)benzene)-tricarbonyl-chromium	as540be_ASER
RABNIT	catena-[tetrakis(benzimidazole)-zinc(ii) tris(2-oxalato)-di-zinc(ii)]	as540be_ASER
REDJUF	catena-tris((2-5-Cyclopentadienyl)-(5-cyclopentadienyl)-lead) toluene solvate	as540be_ASER
TEJFUM	(S,Z)-5-chloro-2-fluoro-4-((tetrahydro-2H-pyran-3-yl)methyl)amino)-N-(thiazol-2(3H)-ylidene)benzenesulfonamide	as540be_ASER
VEZSOI	Chloro-(6-hexamethylbenzene)-trimethylphosphine-(5-trimethylsiloxy-5,5-diphenylpent-1,3-diyne)-ruthenium	as540be_ASER
VOXDIV	(R,S)-3,3,6,6-Tetramethyl-1,4-diphenyl-2,5,7-trioxabicyclo(2.2.1)heptane	as540be_ASER
WIRNOC	(2,6-bis(3,5-bis(trifluoromethyl)phenyl)-4-phenylphosphinine)-(5-pentamethylcyclopentadienyl)-chloro-iridium	as540be_ASER
WOPKER	(6-2,2,2,2-Buta-1,3-diyne-1,4-diy)- (2-hydrido)-heptadecacarbonyl-(5-cyclopentadienyl)-penta-ruthenium-tungsten	as540be_ASER
WOPKOB	(6-2,2,2,2-Buta-1,3-diyne-1,4-diy)- (2-hydrido)-heptadecacarbonyl-(5-cyclopentadienyl)-iron-tetra-ruthenium-tungsten	as540be_ASER
WOPKOB01	(6-2,2,2,2-Buta-1,3-diyne-1,4-diy)- (2-hydrido)-heptadecacarbonyl-(5-cyclopentadienyl)-iron-tetra-ruthenium-tungsten	as540be_ASER
WOPLAO	(6-2,2,2,2-Buta-1,3-diyne-1,4-diy)- (2-hydrido)-heptadecacarbonyl-(5-cyclopentadienyl)-di-iron-tri-ruthenium-tungsten	as540be_ASER
WOPLAO01	(6-2,2,2,2-Buta-1,3-diyne-1,4-diy)- (2-hydrido)-heptadecacarbonyl-(5-cyclopentadienyl)-di-iron-tri-ruthenium-tungsten	as540be_ASER
WOTREC	(2-4,6-Diacetylresorcinolato)-bis(1,10-phenanthroline)-di-copper(ii) bis(tetrafluoroborate)	as540be_ASER
XATPOX	N-Methyl-N-(4-methyl-6-phenyl-2H-pyran-2-ylidene)-N-phenylammonium iodide	as540be_ASER

Display records  to  of 27  <<First <Previous [Next](#) [Last](#)>>

#### 4.1.4 Retrieve Entry and Molecule Attributes



This protocol uses the benzodiazepine core structure, shipped as part of the Chemistry collection (data\Queries\BenzodiazepineCore.mol), and searches the CSD for structures which contain this as a similar compound. The query file is specified in the **Source** parameter of the Read Benzodiazepine Core component.

The **CSD Similarity Search** components return refcodes only. To gather details about that refcode, the **Get CSD Entry Attributes**, the Get CSD Molecule Attributes and the **Get CSD Crystal Attributes** components can be used.

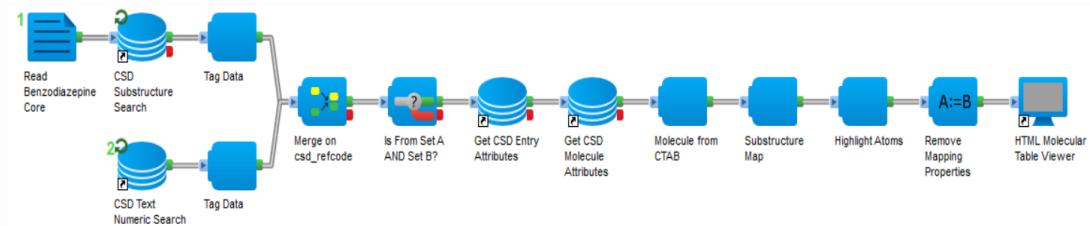
A CSD entry contains attributes that are beyond the concepts of chemistry and crystallography. An example of such an attribute would be the publication details of a CSD entry.

The CSD molecule represents the chemistry and chemical attributes associated with an entry such as the structure, the molecular weight, etc. The structure is found in the following forms: SMILES, CTAB and MOL2.

In some cases, a molecule may not have a canonical SMILES representation (e.g when the structure has unknown atoms or bonds). In these cases, the value will be None. An example would be AJABIX01, the SMILES string property is removed to avoid confusion.

This protocol returns three HTML reports, one for the Entry Attributes, one for the Molecule Attributes and one for the Crystal Attributes of the CSD entries similar to the query.

#### 4.1.5 Combining Hit Sets - AND



This protocol demonstrates how various queries can be combined, by performing them separately and then applying logic to the list of refcodes produced.

In this case, we perform a substructure search for benzodiazepine core (shipped as part of the Chemistry collection (data\Queries\BenzodiazepineCore.mol), which leads to one stream of CSD refcodes. Then, we tag this stream such that each record receives an arbitrary property "A" - just so we can identify these refcodes downstream- defined in the **TagName** parameter.

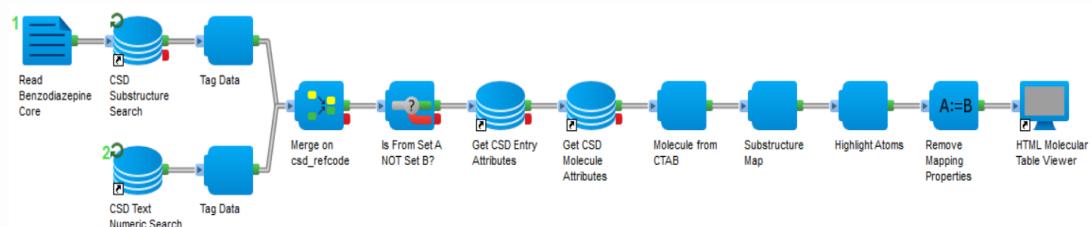
We then perform a search by citation for all crystals in Acta Crystallographica,Section B:Struct.Crystallogr.Cryst.Chem, using the appropriate abbreviation in the **Citation** parameter. This produces a second stream of refcodes, which then we tag with an arbitrary property of "B".

We merge these streams using the refcodes. A filter is then applied to the resultant list of refcodes; if both the "A" and "B" properties are found, then this must have come from both the CSD substructure search and journal search and therefore this passes the filter. Otherwise, the record is passed to the fail port.

With the passing data, we collect both the publication information (specified in the **Attributes** parameter of the Get CSD Entry Attributes component) and the structure information, to demonstrate that the journal is indeed as we had asked for, and the structure contains a benzodiazepine core. The structure is highlighted using the original query to emphasise this point and results reported in an HTML page.

Molecule	csd_refcode	publications	database	Name
	BAYCUZ	(Citation)(authors=u'H Miyamae, A Obata, H Kawazura', journal=JournalActa Crystallographica Section B Struct Crystallogr Cryst Chem [1968-1982]], volume=u'38', year=u'1982', first_page=u'272', doi=u'10.1107/S0567740882002593.')	as540be_ASER	BenzodiazepineCore
	BEDZPN10	(Citation)(authors=u'Z Ruzic-Toros, B Kojic-Prodic, N Bresljan-Pahor, G Nardin, L Randaccio', journal=JournalActa Crystallographica Section B Struct Crystallogr Cryst Chem [1968-1982]], volume=u'38', year=u'1982', first_page=u'272', doi=u'10.1107/S0567740882002593.')	as540be_ASER	BenzodiazepineCore

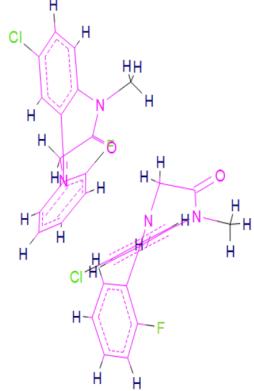
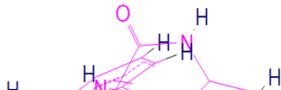
## 4.1.6 Combining Hit Sets – NOT



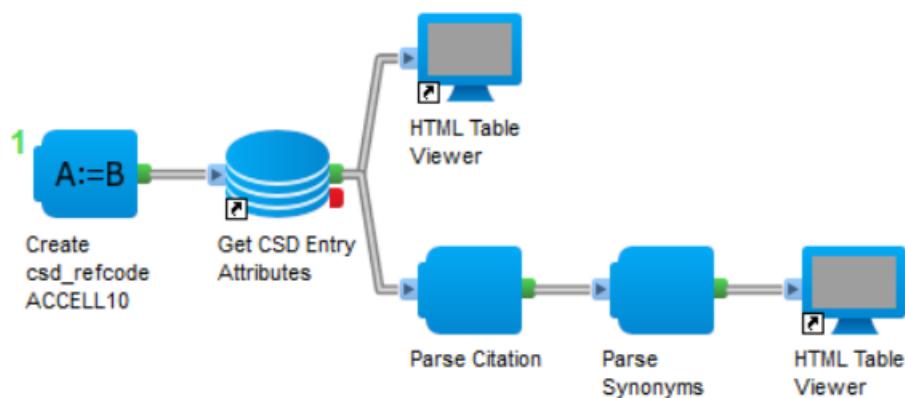
This protocol takes the same premise as for Combining Hit Sets - AND, and queries for both a structure and a journal.

This time the stream of refcodes produced is combined such that the data is passed only if the structure search found it. This means that if the journal search found the data, then the record is failed. This filter is specified in the **Expression** parameter of the Is From Set A Note Set B? component.

As before, the stream data produced are tagged. This time the logic is to pass if tag "A" is defined and tag "B" is not defined. The appropriate attributes (such as structure and journal) are retrieved, and the substructure is highlighted in the structures before reporting.

Display records 1 to 100 of 133 <input type="button" value="Update"/> <<First <Previous <a href="#">Next</a> >Last>		publications	database	Name
Molecule	csd_refcode			
	ANIXAX	(Citation/authors=u'L. E. Brigger, R. Hendrickx, L. Kloo, J. Rosdahl, P. H. Svensson', journal=u'Journal(ChemMedChem)', volume=u'6', year=2011, first_page=u'60', doi=u'10.1002/cmde.201000405.')	as540be_ASER	BenzodiazepineCore
				

#### 4.1.7 Parsing Citation and Synonyms



This example shows how the data returned by the Python API can be further enhanced.

In this case, we pass a refcode (ACCELL10) and take the publication and synonyms in the citations. These are specified in the **Attribute** parameter of the component. Then we produce separate properties

for each field in the citation (publications and synonyms in this case). We take the synonyms and convert them into arrays. In both cases, the encoding used by python is stripped away also.

Finally, we convert the Unicode found in the synonyms into their appropriate characters (e.g. \u03b2 becomes a "beta" character **β**).

## 4.2 Python Examples

This group of protocols shows a few examples of how to use the CSD Python API.

### 4.2.1 Run Python Script Example



This protocol takes a simple example script shipped as part of the CSD PP Component Collection (data\Example Scripts\example.py) and displays it to the user before executing it. The script is simple, taking a set of numeric inputs and summing them together. The script is specified in the **Source** parameter of the Read Script component.

```
Example Python Script.txt - Notepad
File Edit Format View Help
Text
import argparse

parser = argparse.ArgumentParser(description='Process some integers.')
parser.add_argument('integers', metavar='N', type=int, nargs='+',
                    help='an integer for the accumulator')
parser.add_argument('--sum', dest='accumulate', action='store_const',
                    const=sum, default=max,
                    help='sum the integers (default: find the max)')

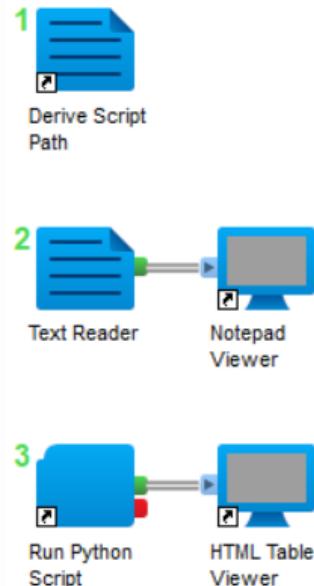
args = parser.parse_args()
print(args.accumulate(args.integers))

-----
```

Switching the final argument to `--max` will return the max value instead.

It demonstrates how the Run Python component can be configured to point to any Python script and execute it.

#### 4.2.2 Using Derive Script Path



The CSD PP Component Collection package has a script folder which contains the python scripts used to execute on the CSD. To avoid having to reference this folder every time, there is a the Derive

Script Path component, which knows where this folder is and will gather the path to the script named in its **Script Filename** parameter.

In this protocol, the **Script Filename** points to a simple script named `validate_journal.py`. The script is displayed to the user before executing it. The script takes precisely one argument; `-j` for the journal name to test. The arguments for the script are specified in the **Arguments** parameter of the **Run Python Script** component of the protocol.

Parameters	
Python	
Script	
Script File	\$(script_location)
CustomMessageToClient	
Arguments	
Argument 1	
Switch	-j
Value	J.Med.Chem.
Output	
Stdout Property Name	is_valid
Stderr Property Name	ErrorText
Command	False

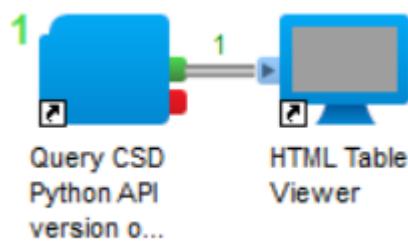
Parameters    Runtime    Implementation

The python script produces two properties; as specified in the **Output** parameters:

**Stdout Property Name** named “`is_valid`” which will contain “`True`” if the journal name is valid, else it will return “`False`” and **Stderr Property Name** named “`ErrorText`”. If there are errors in the execution of the script, these will appear in the “`ErrorText`” property. The results of the script will be reported in an HTML page.

is_valid	ErrorText
True	

#### 4.2.3 Get CSD Python API Version

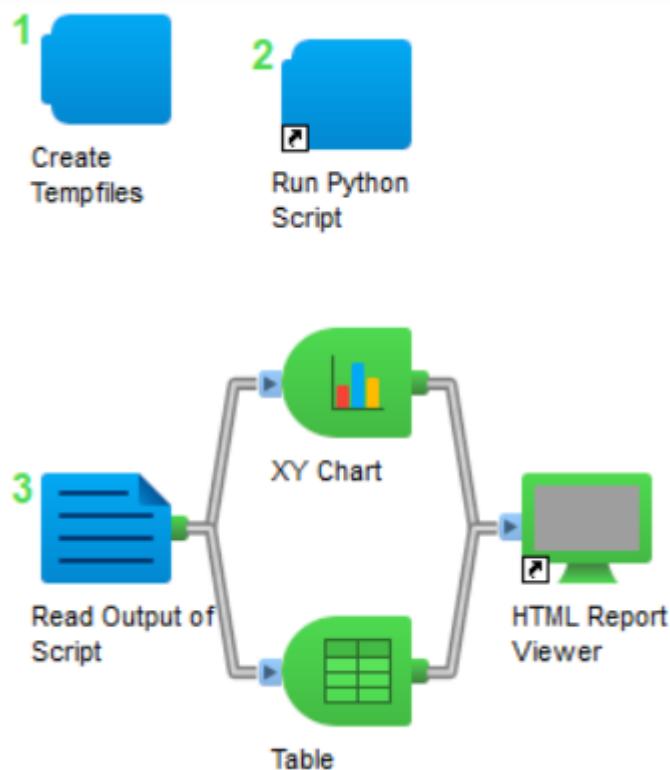


This protocol runs a python script on the server to determine the installed CSD Python API version and the installed python interpreter.

The script is defined in the **Script** parameter of the Query CSD Python API version on server component. The script produces two properties; as specified in the **Output** parameters: **Stdout Property Name** named “script\_out” which will contain the CSD Python API version and the python interpreter and **Stderr Property Name** named “script\_err”. If there are errors in the execution of the script, these will appear in the script\_err property. The results of the script will be reported in an HTML page.

CSD Python API version		
	script_out	script_err
2.2.0	<pre>sys.version_info(major=2, minor=7, micro=15, releaselevel='final', serial=0)</pre>	

#### 4.2.4 Count Entries per Decade



This protocol shows how to implement a more complex script, one which uses the CSD Python API.

It starts by creating a global variable to contain the path to a temporary file(`tmpCSV`). Then the script, specified in the **Script** parameter, is passed this path in an argument, `-o`, so the script may use this path.

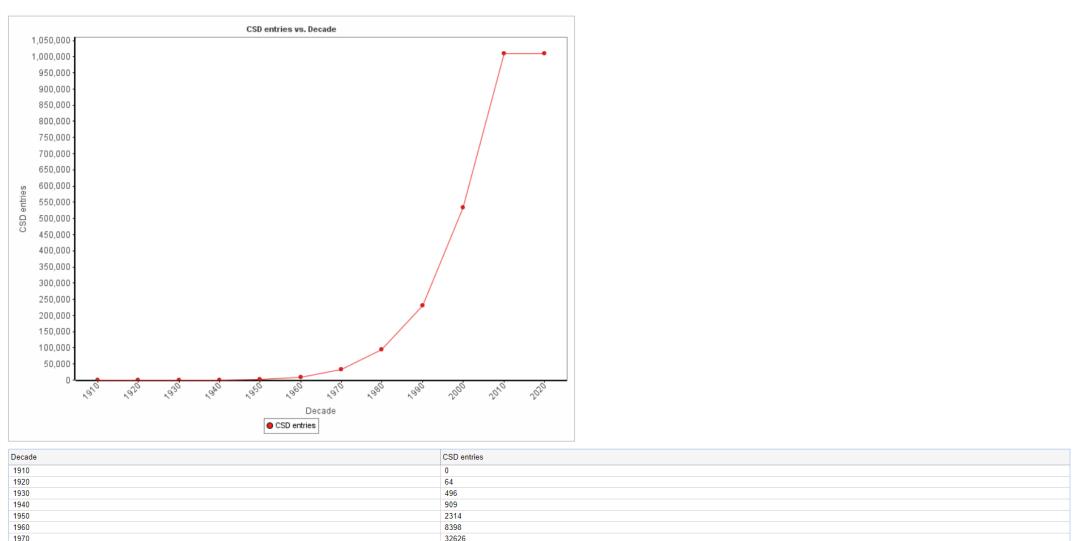
**Parameters**

<b>Python</b>	
Script	'''Extract growth of the CSD over the decades.''' import argparse f...
Script File	
CustomMessageToClient	
<b>Arguments</b>	
Argument 1	
Switch	-o
Value	<code>\$(tmpCSV)</code>
Argument 2	
Switch	
Value	
Argument 3	
<b>Output</b>	
Stdout Property Name	<code>script_out</code>
Stderr Property Name	<code>script_err</code>
Command	True

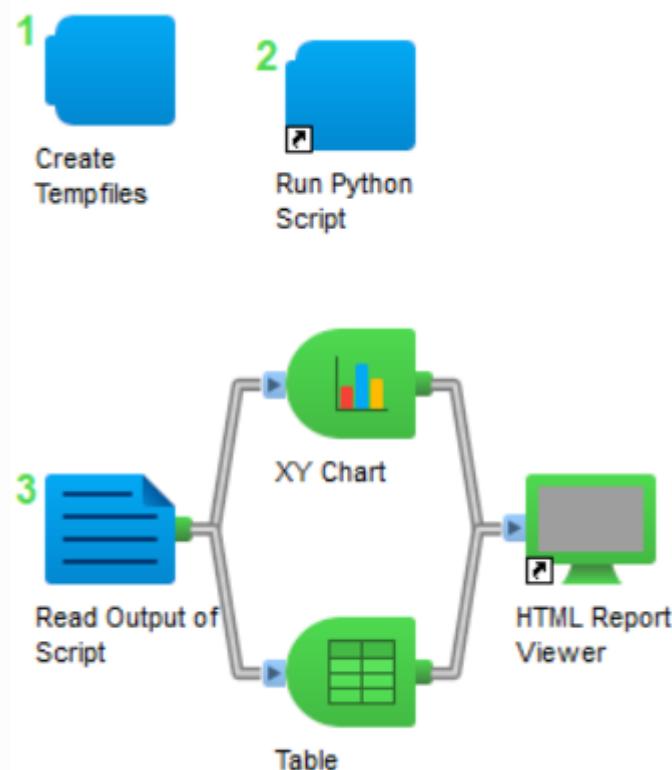
Parameters    Runtime | Implementation |

The script populates the file with the breakdown of entries in CSD by decade of their reporting in the literature.

The protocol then picks this temporary file up and reports on it generating an HTML page with the plot of the CSD entries per decade and a table with these values.

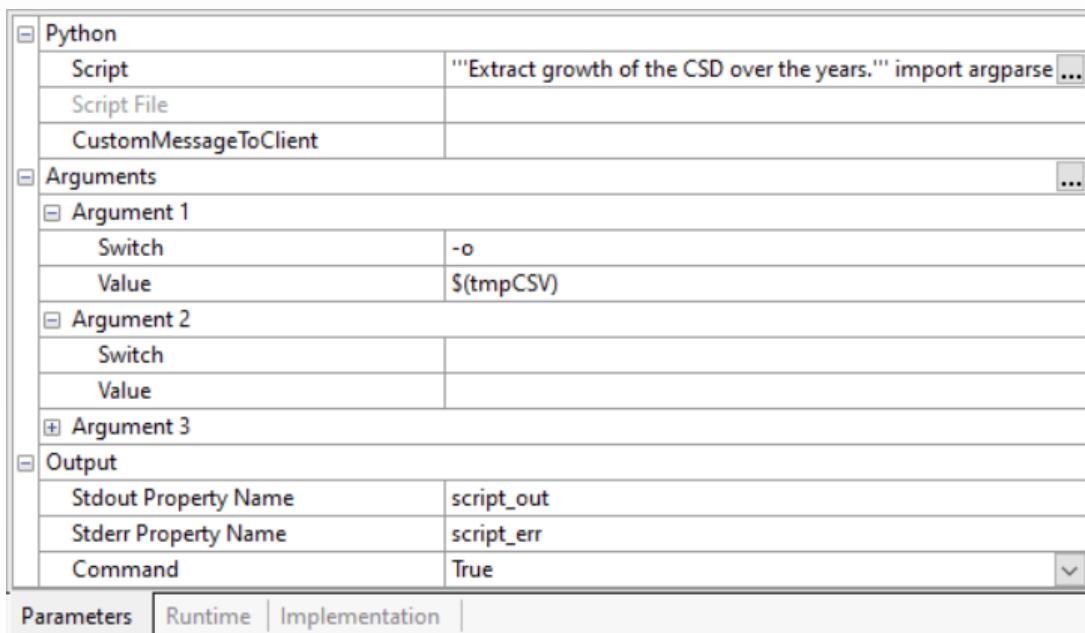


#### 4.2.5 Count Entries per Year



As the Count Entries per Decade protocol, this example shows how to implement a more complex script, one which uses the CSD Python API.

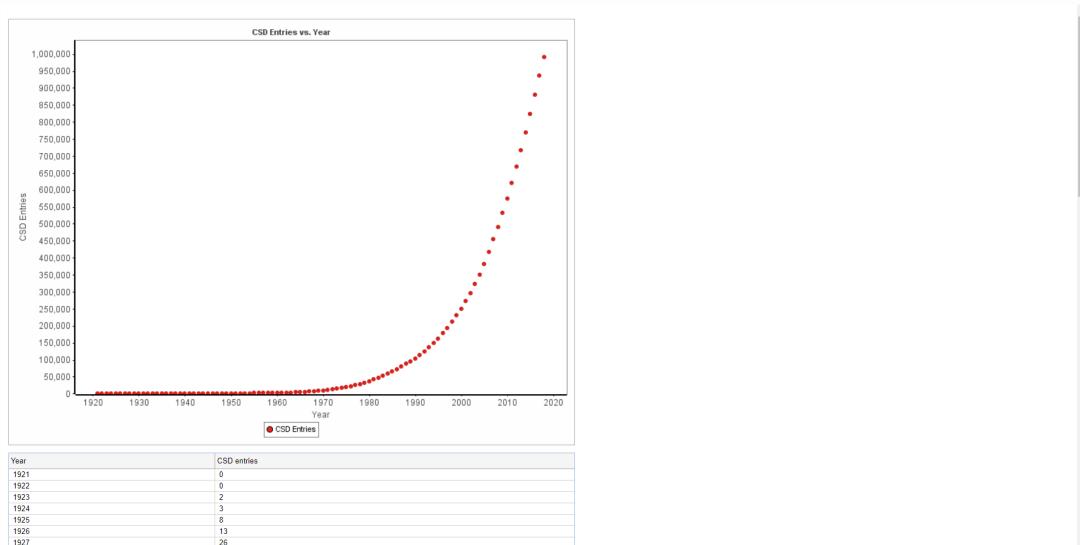
It starts by creating a global variable to contain the path to a temporary file(`tmpCSV`). Then, the script specified in the **Script** parameter passes this path to the temporary file in an argument, `-o`, so the script may use this path.



The screenshot shows the protocol configuration interface for a Python script. The 'Parameters' tab is selected. The 'Script' field contains the code: `'''Extract growth of the CSD over the years.''' import argparse`. The 'Arguments' section contains three arguments: 'Argument 1' with 'Switch' set to `-o` and 'Value' set to `$(tmpCSV)`; 'Argument 2' with 'Switch' set to `-o` and 'Value' set to `$(tmpCSV)`; and 'Argument 3' which is collapsed. The 'Output' section includes 'Stdout Property Name' set to `script_out`, 'Stderr Property Name' set to `script_err`, and 'Command' set to `True`.

The script populates the file with the breakdown of entries in CSD by year of their reporting in the literature.

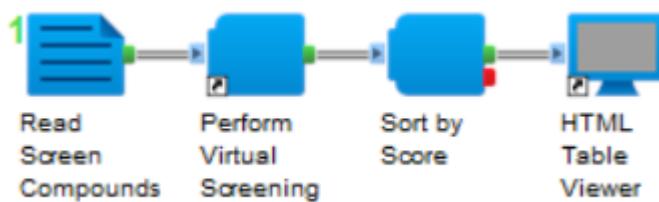
The protocol then picks this temporary file up and reports on it generating an HTML page with the plot of the number of CSD entries per year and a table with these values.



## 4.3 Virtual Screening and Conformer

These protocols are available for user with a CSD-Discovery and/or CSD-Materials licence. The twelve protocols provide different workflows on how to perform virtual screening, virtual screening validation and conformer generation using CSD Python API.

### 4.3.1 Queries Identified by File Screening Example



This protocol takes a stream of incoming molecules, and screens against a query set identified by a file in order to generate a virtual screening score.

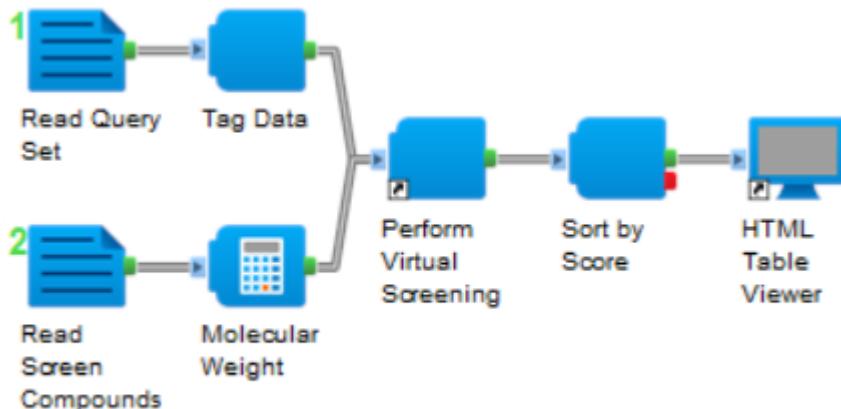
The stream of molecules used in this protocol is provided as part of the CSD PP Component Collection (`data\Python API Example Data\P28845_actives.sdf`).

The query set may be supplied either as a MOL2 or SD file and it is specified in the **Source** parameter of the Perform Virtual Screening component. The query used in this protocol is provided as part of the CSD PP Component Collection (`data\Python API Example Data\P28845.sdf`). Up to twenty-five conformations of each molecule are generated and the calculation is distributed over one thread as specified in the **Screening Options** parameters of the Perform Virtual Screening component.

The virtual screening score is calculated and then the incoming data are sorted by screening score from lowest to highest. The sort criteria are defined by one or more properties in the **Sort By** parameter whose values are used to order the records. The results are then displayed in an HTML page.

Display records	1	to	25	of 164	<a href="#">Update</a>	<<First	<Previous	<a href="#">Next</a>	Last>
Data Image	CHEMBL460962 P10000288	Mo2_MolInfo_Type	Mo2_MolInfo_Charge	Mo2_MolInfo_Comment	Mo2_MolInfo_StatusBits	Mo2_Substructures	Name	Virtual_Screen_Score	
	CHEMBL460962 P10000288	SMALL	USER_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 DONOR_STRONG 5 GROUP 0 4 DONOR_STRONG 0 4 ACCEPTOR_MEDIUM 11 GROUP 0 3 ACCEPTOR_MEDIUM 0 5 DONACC_25 GROUP 0 1 DONACC 0 6 DONOR_MEDIUM 32 GROUP 0 5 DONOR_MEDIUM 0	CHEMBL460962 P10000288	-127.996	
	CHEMBL460962 P10000287	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_MEDIUM 0 GROUP 0 3 ACCEPTOR_MEDIUM 0 4 DONACC_24 GROUP 0 1 DONACC 0 5 DONOR_MEDIUM 31 GROUP 0 5 DONOR_MEDIUM 0	CHEMBL460962 P10000287	-124.099	
	CHEMBL221158 P10000194	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_MEDIUM 10 GROUP 0 3 ACCEPTOR_MEDIUM 0 4 DONOR_MEDIUM 19 GROUP 0 5 DONOR_MEDIUM 0 5 ACCEPTOR_WEAK 28 GROUP 0 6 ACCEPTOR_WEAK 0	CHEMBL221158 P10000194	-114.883	
	CHEMBL510937 P10000220	SMALL	USER_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_STRONG 11 GROUP 0 7 ACCEPTOR_STRONG 0 4 ACCEPTOR_MEDIUM 13 GROUP 0 3 ACCEPTOR_MEDIUM 0	CHEMBL510937 P10000220	-112.336	
	CHEMBL222576 P10000001	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_MEDIUM 10 GROUP 0 3 ACCEPTOR_MEDIUM 0 4 DONOR_MEDIUM 19 GROUP 0 5 DONOR_MEDIUM 0	CHEMBL222576 P10000001	-111.418	
	CHEMBL310617 P10000222					1 NONPOLAR 10 AD 1 DONOR 8 2 NONPOLAR 10 AD 0			

### 4.3.2 Queries Identified by Tag Screening Example



This protocol shows the ability to stream in both query and screening molecules, as long as the query molecules are tagged so they may be identified. The tag name is specified in the **TagName** parameter of the Tag Data component. The query and the stream of molecule are specified in the **Source** parameter of the Read components. For the screening molecules, only the first fifty molecules will be used as specified in the **Maximum** parameter of the Read Screen Compounds component.

**Parameters**

Source	data\Python API Example Data\P28845_actives.sdf	...
Maximum	50	
SourceTag	None	...
Keep Properties		...
UTF-8 Auto Detect	False	...
<b>Additional Options</b>		

Parameters    Runtime

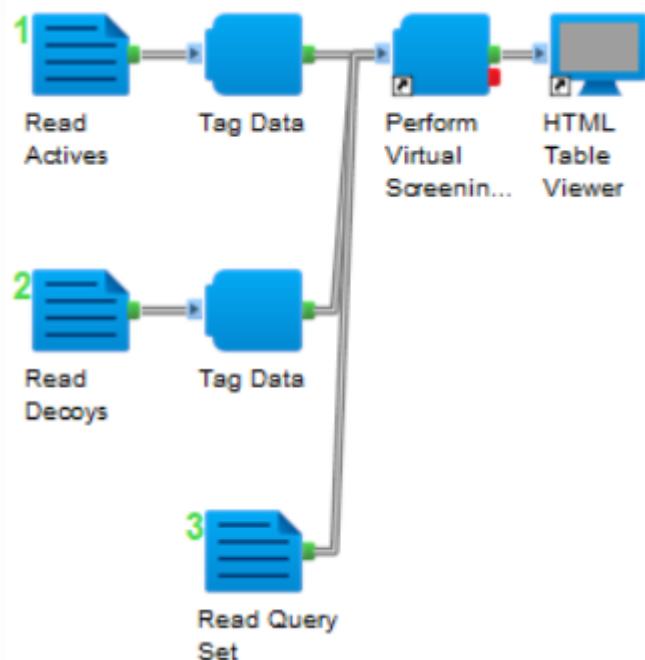
The **Molecular weight** property is added to the screening set to demonstrate that the data defined on the original dataset is passed through this component.

The **Perform Virtual Screening** component then internally divides the records according to the **Test for Query**, "IsQuery is Defined" to produce a query et and screening set.

The data, which is output by the component, has a new virtual screening score which is stored in a property, those are then sorted by screening score from lowest to highest. The results are then displayed in an HTML page.

Display records	1	to	25	of 50	Update	<<First	<Previous	Next>	Last>>
Data Image	Mo2_MolInfo_Name	Mo2_MolInfo_Type	Mo2_MolInfo_Charge	Mo2_MolInfo_Comment	Mo2_MolInfo_StatusBits	Mo2_Substructures	Name	Virtual_Screen_Score	Molecular_Weight
	CHEMBL510937 P10000220	SMALL	USER_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_STRONG 11 GROUP 0 7 ACCEPTOR_STRONG 0 4 ACCEPTOR_MEDIUM 13 GROUP 0 3 ACCEPTOR_MEDIUM 0	CHEMBL510937 P10000220	-113.304	424.51
	CHEMBL401359 P10000034	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 ACCEPTOR_MEDIUM 18 GROUP 0 3 ACCEPTOR_MEDIUM 0 3 DONOR_MEDIUM 10 GROUP 0 5 DONOR_MEDIUM 0 4 unknown 25 GROUP 0 0 unknown 0	CHEMBL401359 P10000034	-111.481	349.43
	CHEMBL222576 P10000001	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_MEDIUM 11 GROUP 0 3 ACCEPTOR_MEDIUM 0 4 ACCEPTOR_MEDIUM 19 GROUP 0 3 ACCEPTOR_MEDIUM 0	CHEMBL222576 P10000001	-110.916	408.32
	CHEMBL464387 P10000020	SMALL	USER_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 unknown 3 GROUP 0 0 unknown 0 3 ACCEPTOR_STRONG 11 GROUP 0 7 ACCEPTOR_STRONG 0 4 ACCEPTOR_MEDIUM 13 GROUP 0 3 ACCEPTOR_MEDIUM 0	CHEMBL464387 P10000020	-110.745	440.51
	CHEMBL249227 P10000190	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 ACCEPTOR_MEDIUM 6 GROUP 0 3 ACCEPTOR_MEDIUM 0 3 unknown 28 GROUP 0 0 unknown 0	CHEMBL249227 P10000190	-110.042	382.48
	CHEMBL399455 P10000029	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 0 NONPOLAR 0 2 ACCEPTOR_MEDIUM 6 GROUP 0 3 ACCEPTOR_MEDIUM 0 3 unknown 28 GROUP 0 0 unknown 0	CHEMBL399455 P10000029	-109.539	364.44

### 4.3.3 Screen Validation Using Tagging



This protocol takes a stream of actives (tagged), decoys (tagged), and query molecules, and performs a virtual screening validation. The output has a new property of “Is\_Active” (1 or 0 depending upon whether the data was considered active), “Virtual\_Screen\_Score” (the score for the model).

**Parameters**

<b>Input Options</b>	
<input type="checkbox"/> Active Source	From Tag
Source	
Test for Active	IsActive Is Defined
<input type="checkbox"/> Decoy Source	From Tag
Source	
Test for Decoy	IsDecoy Is Defined
<b>Screening Validation Options</b>	
Number of Threads	1
Max Number of Conformers	25
<b>Output Options</b>	
Active Tag	Is_Active
Screening Score Property	Virtual_Screen_Score

Parameters    Runtime    Implementation

The data are then reported in an HTML page.

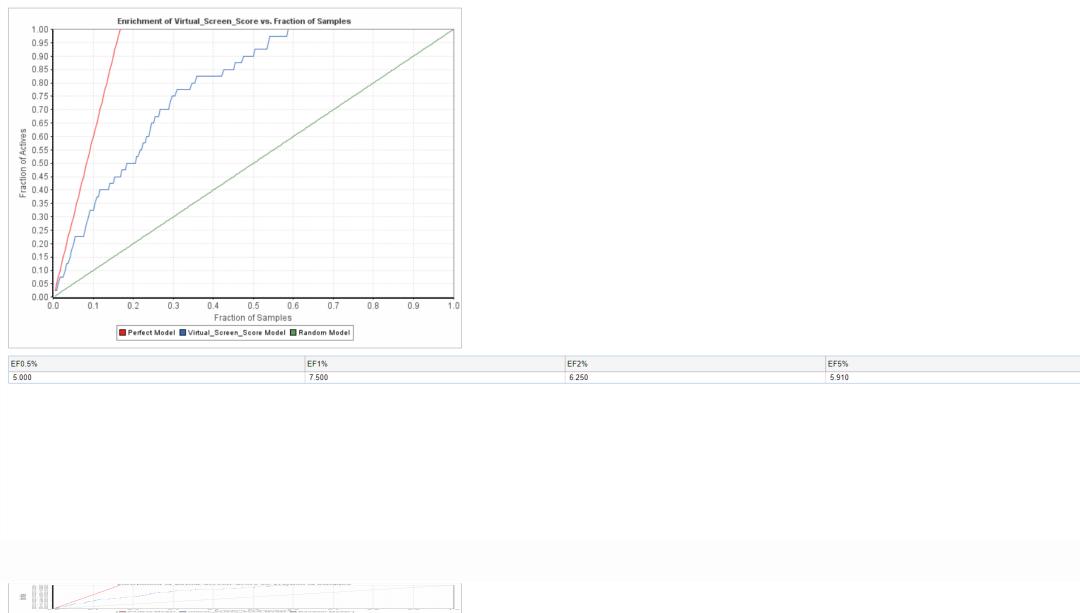
Display records 1 to 25 of 60 <input type="button" value="Update"/>	<<First	<Previous	<a href="#">Next</a>	<a href="#">Last&gt;</a>							
Data Image	Mo2_MolInfo_Name	Mo2_MolInfo_Type	Mo2_MolInfo_Charge	Mo2_MolInfo_Comment	Mo2_MolInfo_StatusBits	Mo2_Substructures	Name	Virtual_Screen_Score	Is_Active	IsActive	IsDecoy
	CHEMBL401359 P1000034	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 1 8 NONPOLAR 2 2 ACCEPTOR, MEDIUM 10 GROUP 0 3 ACCEPTOR, MEDIUM 0 3 DONOR, MEDIUM 24 GROUP 0 5 DONOR, MEDIUM 0 4 unknown 25 GROUP 0 unknown 0	CHEMBL401359 P1000034	-109.630	1	true	
	CHEMBL249227 P10000190	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 1 8 NONPOLAR 0 2 ACCEPTOR, MEDIUM 8 GROUP 0 3 ACCEPTOR, MEDIUM 0 3 unknown 28 GROUP 0 unknown 0	CHEMBL249227 P10000190	-109.585	1	true	
	CHEMBL251403 P10000264	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 1 8 NONPOLAR 2 2 DONOR, MEDIUM 7 GROUP 0 5 DONOR, MEDIUM 0 3 unknown 8 GROUP 0 unknown 0 4 unknown 11 GROUP 0 2 ACCEPTOR, MEDIUM 0 3 ACCEPTOR, MEDIUM 0 5 ACCEPTOR, WEAK 26 GROUP 0 6 ACCEPTOR, WEAK 0	CHEMBL251403 P10000264	-108.088	1	true	
	CHEMBL509257 P10000033	SMALL	USER_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 unknown 3 GROUP 0 unknown 0 3 unknown 10 GROUP 0 1 ACCEPTOR, STRONG 0 4 ACCEPTOR, MEDIUM 13 GROUP 0 3 ACCEPTOR, MEDIUM 0 5 DONOR, MEDIUM 27 GROUP 0 5 DONOR, MEDIUM 0	CHEMBL509257 P10000033	-105.236	1	true	
	CT2478338 P112759218	SMALL	NO_CHARGES	Generated from the CSD	****	1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 ACCEPTOR, MEDIUM 5 GROUP 0 1 ACCEPTOR, MEDIUM 0 3 DONOR, WEAK 6 GROUP 0 2 DONOR, WEAK 0 4 DONACC 22 GROUP 0 1 DONACC 0 5 unknown 26 GROUP 0 unknown 0	CT2478338 P112759218	-104.937	0	true	
	CT2478338 P112759					1 NONPOLAR 1 GROUP 0 8 NONPOLAR 0 2 ACCEPTOR, MEDIUM 5 GROUP 0 1 ACCEPTOR, MEDIUM 0 3 DONOR, WEAK 6 GROUP 0 2 DONOR, WEAK 0 4 DONACC 22 GROUP 0 1 DONACC 0 5 unknown 26 GROUP 0 unknown 0					

This data can be plotted in ROC or enrichment plots to reveal the selectivity of the model. There are separate components which offer those capabilities.

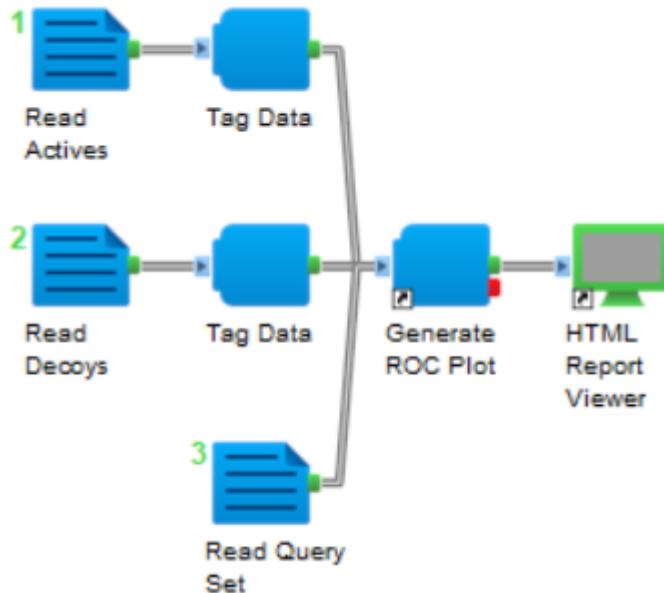
#### 4.3.4 Generate Enrichment Plot Example



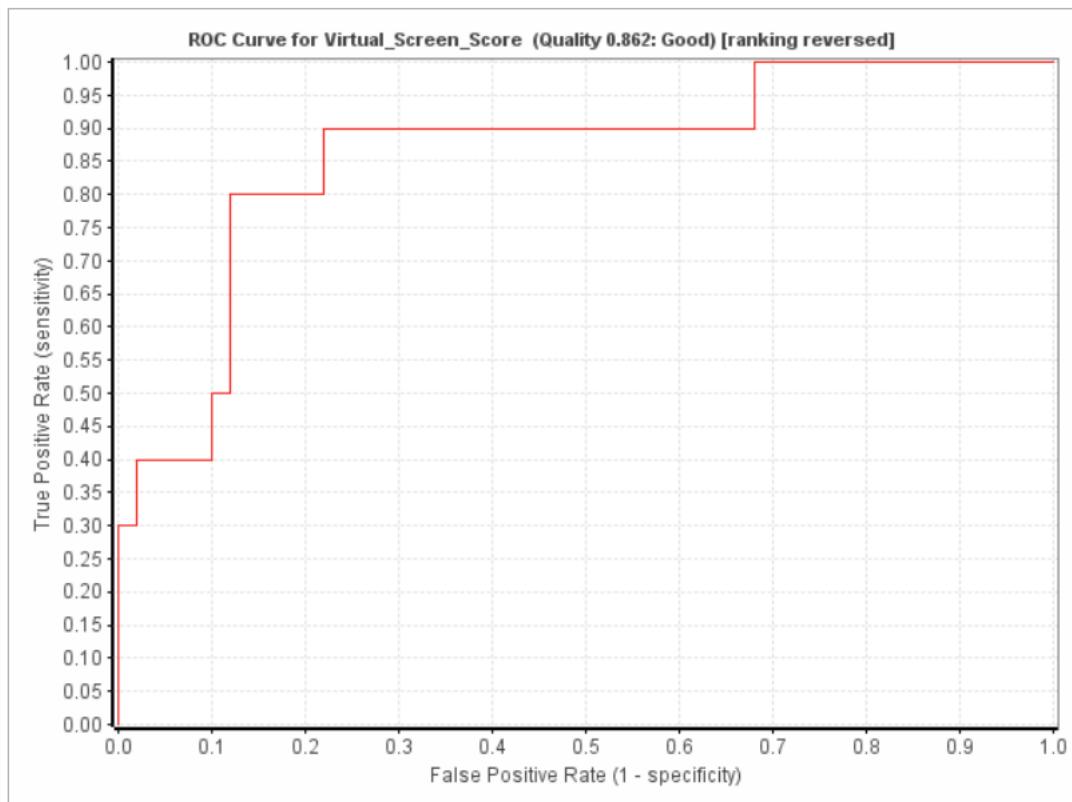
As for the Screen Validation Using Tagging protocol, this example takes a stream of actives (tagged), decoys (tagged), and query molecules, and performs a virtual screening validation however, it uses the data produced to generate an enrichment plot that is useful to help in the understanding of the selectivity of the model. The enrichments at 0.5%, 1%, 2%, and 5% is calculated and an HTML report is returned.



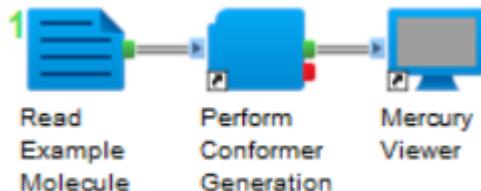
#### 4.3.5 Generate ROC Example



This protocol takes a stream of active molecule (tagged), decoys (tagged), and query molecules, performs a virtual screening validation and, using the data produced, generates a ROC plot. This is useful to help in the understanding of the overall efficiency of the model to separate active ligands from inactive molecules. Note that, in this example, the **Maximum** parameter is set to ten for actives and fifty for decoys.

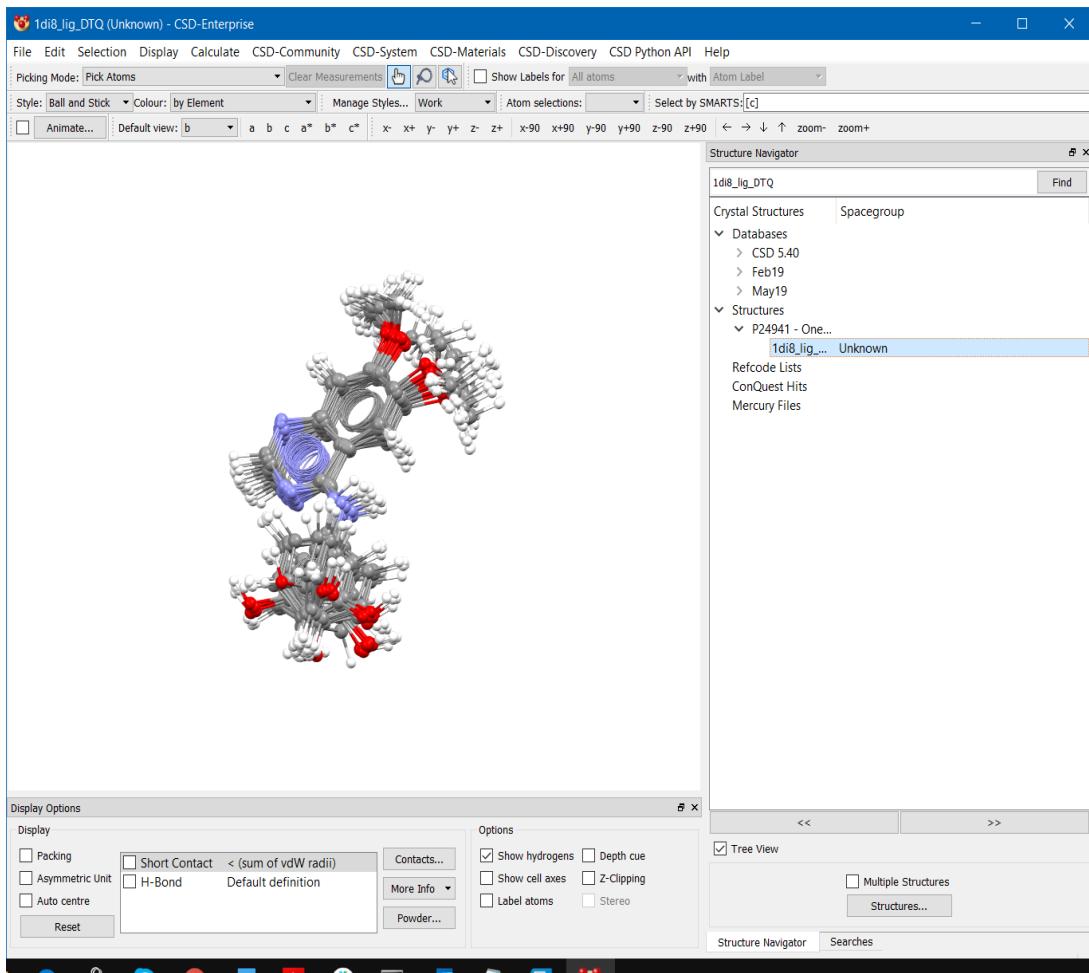


#### 4.3.6 Generate Conformers for Molecule

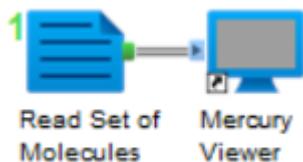


This protocol takes a molecule and generates up to twenty-five conformers for that molecule, using conformer generation in the CSD Python API.

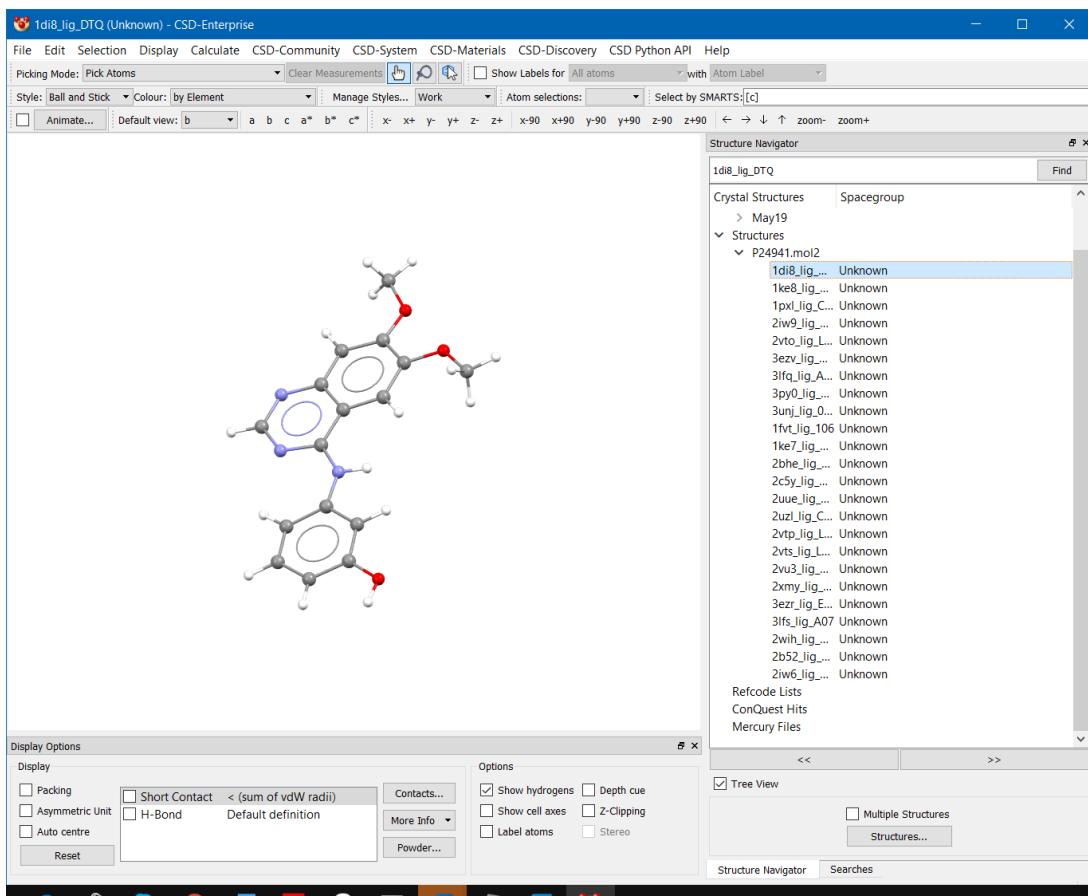
The query molecule used in this protocol is provided as part of the CSD PP Component Collection ([data\Python API Example Data\P24941sdf](#)). The generated conformers are superimposed on one another as specified by switching the **Superimpose** parameter in the **Conformer Options** parameters to “True”. The overlaid results are then viewed in Mercury.



### 4.3.7 Mercury Viewer Example



This simple example demonstrates how molecular records can be piped into Mercury to be viewed. As the structures often contain 3D coordinates, it makes sense that they should be viewed in a way that makes sense of this. The set molecules used in this protocol is provided as part of the CSD PP Component Collection (data\Python API Example Data\P24941sdf).



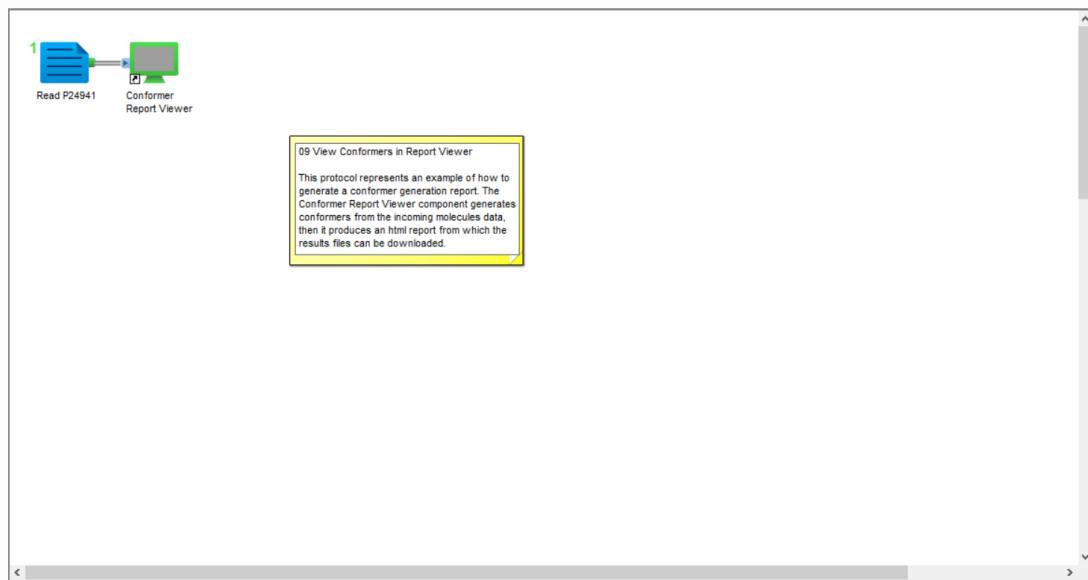
### 4.3.8 Conformer Writer Example

This protocol generates up to twenty-five conformers for the first three molecules provided by the incoming record file and split the molecular output. This is achieved by setting the **Split output** parameter of the Conformer Writer component to “True”. The set molecules used in this protocol is provided as part of the CSD PP Component Collection (data\Python API Example Data\P24941sdf)

In addition to the conformers, the summary of the conformer generation is also provided as specified in **Output Options** parameter of the Conformer Writer component. The conformer summary file is then displayed in an HTML report page.

Conformer Summary					
Molecule name	max_log_probability	min_log_probability	n_conf_gen_clust	n_rotamers_with_no_observations	rotamers_with_no_observations
1d8_lig.DTO	-2.437	-20.366	25	0	0
1ke8_lig.LSA	-5.279	-24.779	25	0	0
1pxl_lig.CK4	-7.092	-13.163	25	0	0

## 4.3.9 View Conformers in Report Viewer



This protocol represents an example of how to generate a conformer generation report. In this example, the **Conformer Report Viewer** component generates up to twenty-five conformers for the first three molecules of the incoming data, then it produces an HTML report from which the results files can be downloaded. The set molecules used in this protocol is provided as part of the CSD PP Component Collection (data\Python API Example Data\P24941sdf). The result files included in the report are specified in the **Output What** of the **Conformer Report Viewer** component. In the example here, the molecules and the summary of the conformer generation process are generated.

[Conformer Generation Report](#)

**Settings**

Names		Conformer Settings		Values	
Max Number of Conformers	25	Number of Threads	1	Max Conformal Torsions	2
Superimpose	False				

**Summary**

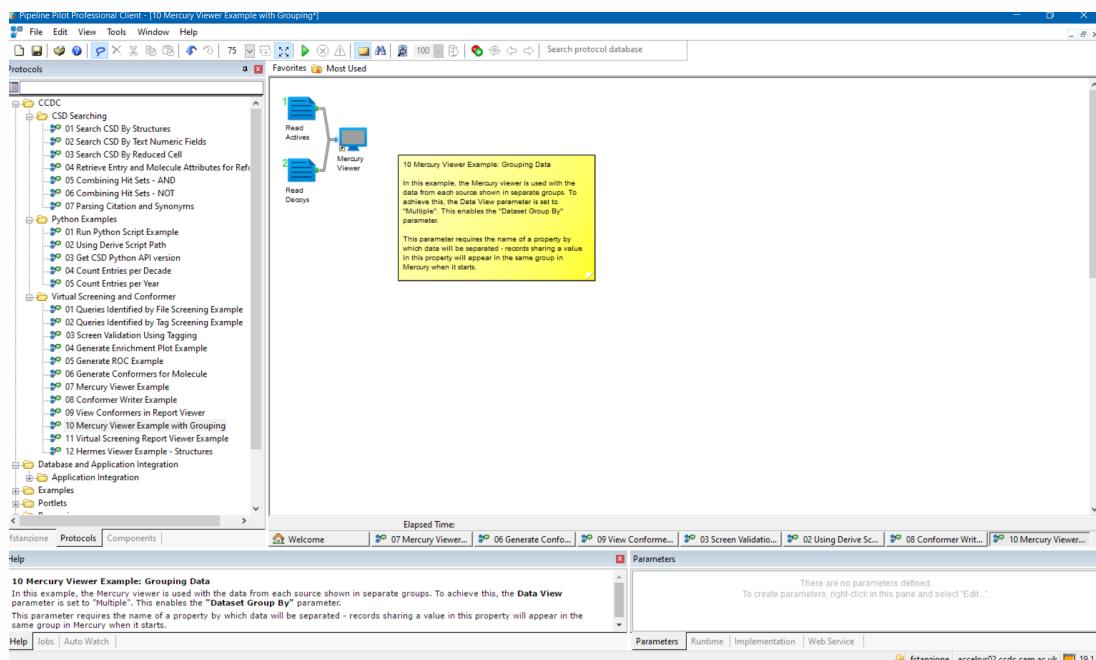
Summary Table					
Molecule name	max_mg_probability	min_mg_probability	n_conf_gen_clust	n_rotamers_with_no_observations	rotamers_with_no_observations
1cd_3g_DTQ	-2.417	-20.566	25	0	0
1kev_3g_L54	-6.279	-24.779	25	0	0
1pol_3g_CK4	-7.092	-13.163	25	0	0

**Downloads**

The following file(s) were created.

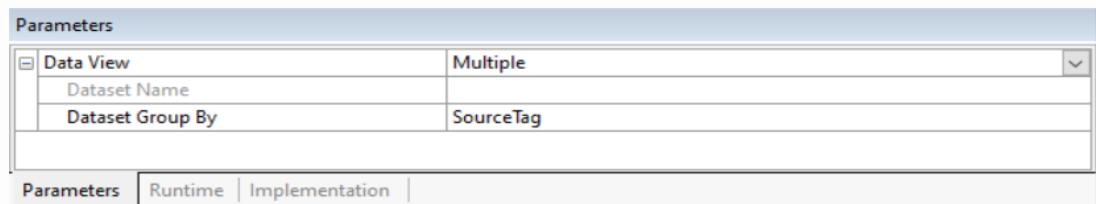
[conformers.sdf](#)

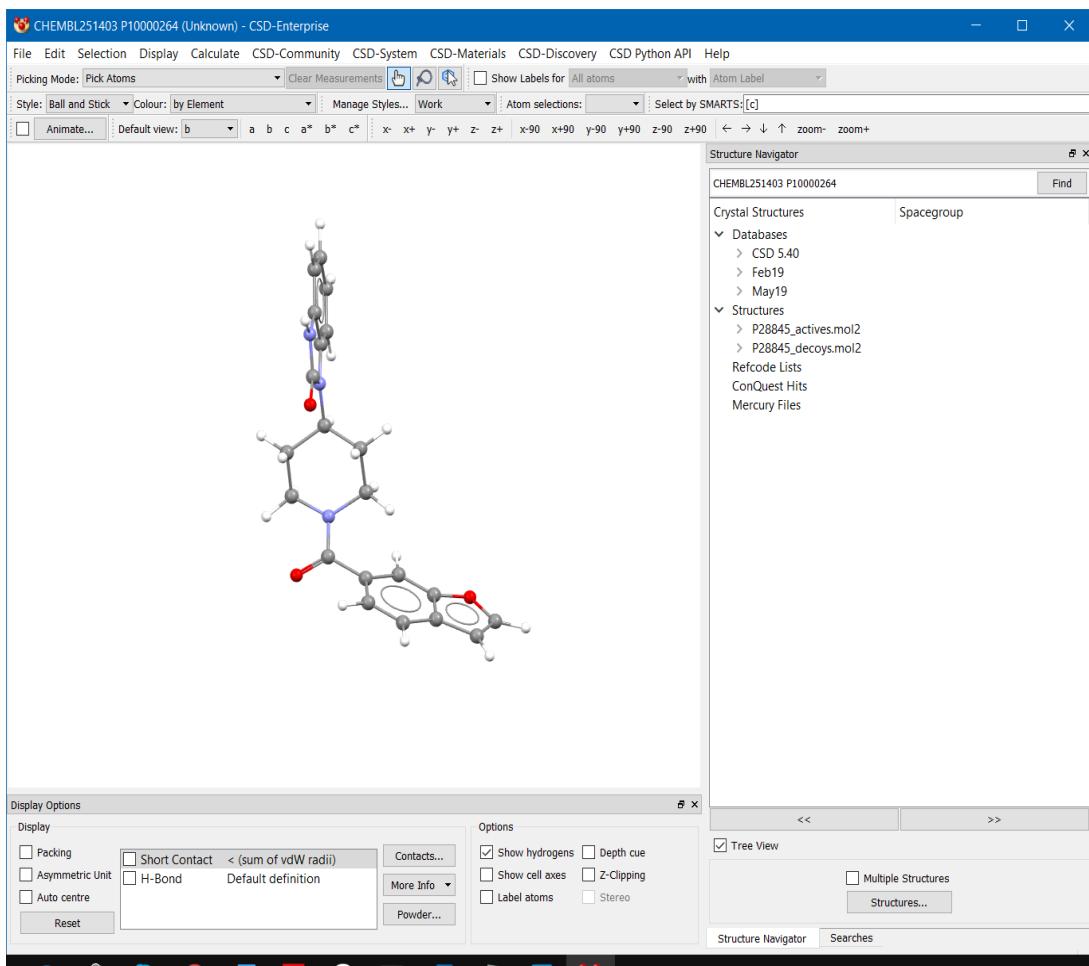
## 4.3.10 Mercury Viewer Example with Grouping



In this example, the Mercury viewer is used to visualise the data from each source specified in Read Actives and Read Decoys components but in separate groups in the Mercury Structure Navigator.

To achieve this, the **Data View** parameter in the Mercury Viewer component is set to "Multiple". This enables the **Dataset Group By** parameter. This parameter requires the name of a property by which data will be separated. In this example the **Dataset Group by** parameter is set to "SourceTag" that is defined as "Filename" in the Read Actives and Read Decoys components.





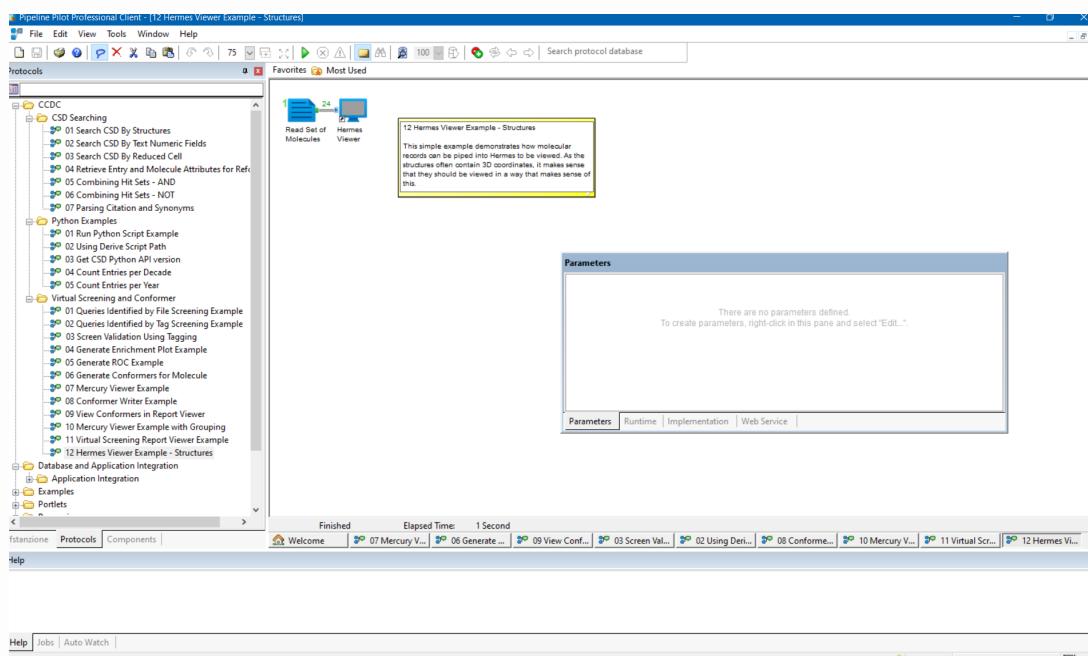
### 4.3.11 Virtual Screening Report Viewer Example

The screenshot shows the Pipeline Pilot Professional Client interface. The title bar reads "Pipeline Pilot Professional Client - [1] Virtual Screening Report Viewer Example". The menu bar includes File, Edit, View, Tools, Window, and Help. The toolbar includes standard icons for file operations. The left sidebar shows a tree view of protocols, including "CCDC", "Python Examples", "Virtual Screening and Conformer", "Database and Application Integration", "Examples", and "Portlets". The main workspace displays a flowchart for "11 Virtual Screening Report Viewer Example", which consists of a "Read Screen Compounds" node connected to a "Virtual Screening Report" node. A tooltip for this flowchart states: "11 Virtual Screening Report Viewer Example Takes a stream of incoming molecules, and screens against a query set identified by file." Below the flowchart, a "Parameters" dialog is open, showing a message: "There are no parameters defined. To create parameters, right-click in this pane and select 'Edit...'". The bottom status bar includes "Elapsed Time", "Welcome", "07 Mercury Vie...", "06 Generate Co...", "09 View Confor...", "03 Screen Valida...", "02 Using Devic...", "08 Conformer...", "10 Mercury Vie...", "11 Virtual Scree...", "Help", "Jobs | Auto Watch", and "Parameters Runtime Implementation Web Service".

This protocol represents an example of how to generate a virtual screening report. In the example here, the Virtual Screening Report component takes a stream of the first ten incoming molecules provided (the file is available in `data\Python API Example Data\P28845_actives.sdf`), and screens against a query set identified by file (provided in `data\Python API Example Data\P28845.sdf`), then it produces an HTML report from which the Results File together with the Screening File, the Query File and the Ranked Scores File can be downloaded.



### 4.3.12 Hermes Viewer Example – Structures



This simple example demonstrates how molecular records can be piped into Hermes to be viewed. In this example the molecules provided by the Read Set of Molecules component are displayed in Hermes.

