

Ligand-Based Virtual Screening

2020.0 CSD Release

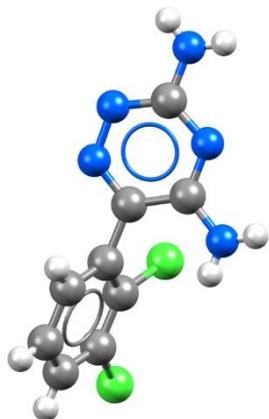


Table of Contents

Introduction.....	2
Case study	3
1. Generating conformers	4
1.1 Command line utility	5
1.2 The Mercury interface.....	8
1.3 The CSD Python API.....	9
Conclusion	11
2. Overlaying ligands	12
Analysis of the results	14
Pharmacophore model.....	17
Conclusion	18
3. Field-based virtual screening.....	19
Retrieval of actives and decoys	19
Library screening	20
Data analysis.....	22
Conclusion	25

Introduction

In the absence of three-dimensional (3D) structures of potential drug targets, ligand-based drug design is one of the most popular approaches for drug discovery and lead optimisation. Pharmacophore modelling is a widely-used tool in ligand-based drug design and can provide predictive models suitable for lead compound optimisation.

According to the IUPAC definition, a pharmacophore model is “an ensemble of steric and electronic features that is necessary to ensure the optimal supramolecular interactions with a specific biological target and to trigger (or block) its biological response”.

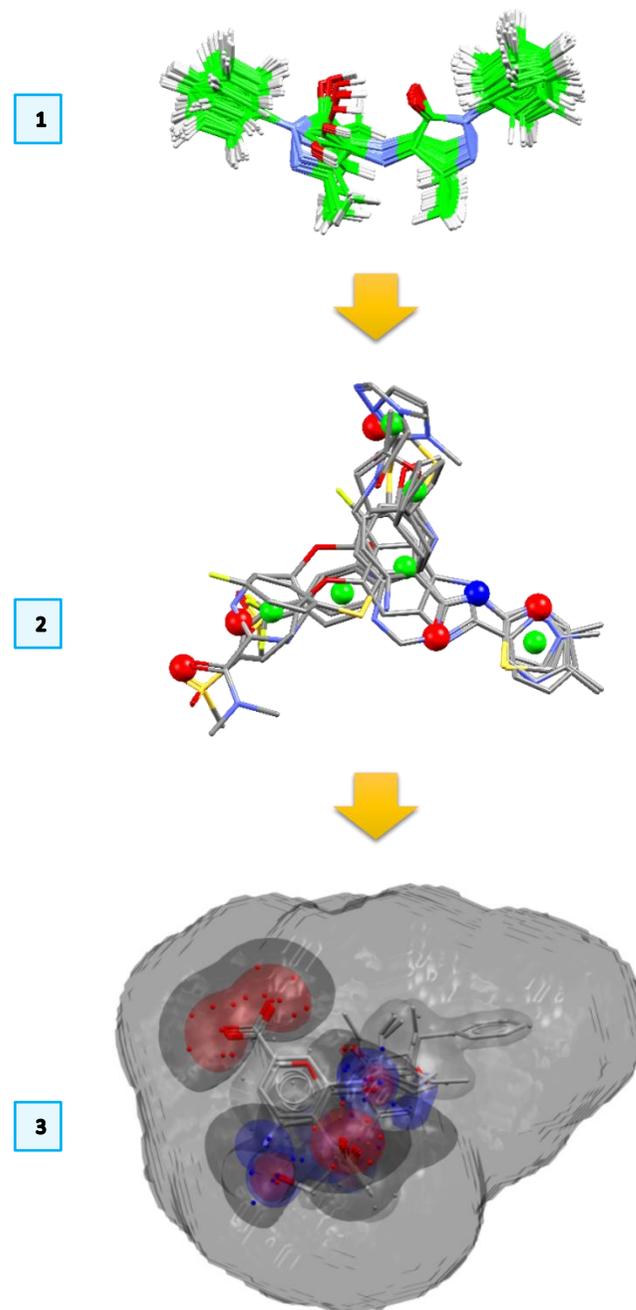
A pharmacophore model can be generated by superposing a set of active molecules that are assumed to bind to the same target with the same binding mode, and extracting common chemical features that are essential for their bioactivity. Such a model can then be used to virtually screen libraries of compounds with the aim of identifying potential new binders of the target of interest.

In this tutorial, you will learn how to perform a ligand-based virtual screening using a suite of knowledge-based tools.

1. First, ensembles of conformers will be generated for a set of known CDK2 inhibitors.
2. Overlay hypotheses for these ligands will be produced using the CSD-Ligand Overlay program and a pharmacophore model will be selected for step 3.
3. To validate the predictive power of our model, a library of active and inactive (or decoy) compounds will be screened against the selected overlay solution and enrichment metrics will be calculated.

Please download the example files [here](#).

<https://downloads.ccdc.cam.ac.uk/tutorials/CDK2.zip>



Case study

Cyclin-dependent kinase 2 (CDK2) is an important protein kinase required for promoting the cell division cycle and for successful progression through S and G2 phases. This role in the cell cycle progression has led to an active search of small molecule compounds inhibiting this enzyme as potential anticancer drugs.

Many studies have been published that describe various CDK inhibitors which target the ATP pocket of CDK2. They bind by hydrophobic interactions and by forming hydrogen bonds with the kinase, especially with the backbone of Glu81 and Leu83 in the structure of the apoenzyme (Figure 1).

Many structures of CDK2 in complex with different types of inhibitors have been deposited in the Protein Data Bank (PDB).

In 2013, researchers at AstraZeneca published an extensive and diverse set of molecular overlays for the validation of pharmacophore programs (DOI: 10.1021/ci400020a). This benchmarking data set contains the experimental overlay of 24 CDK2 inhibitors taken from high resolution structures available from the PDB. A subset of 13 molecules (see below) arbitrarily chosen, but covering all different chemotypes represented in the whole set, will be used as a case study in this tutorial.

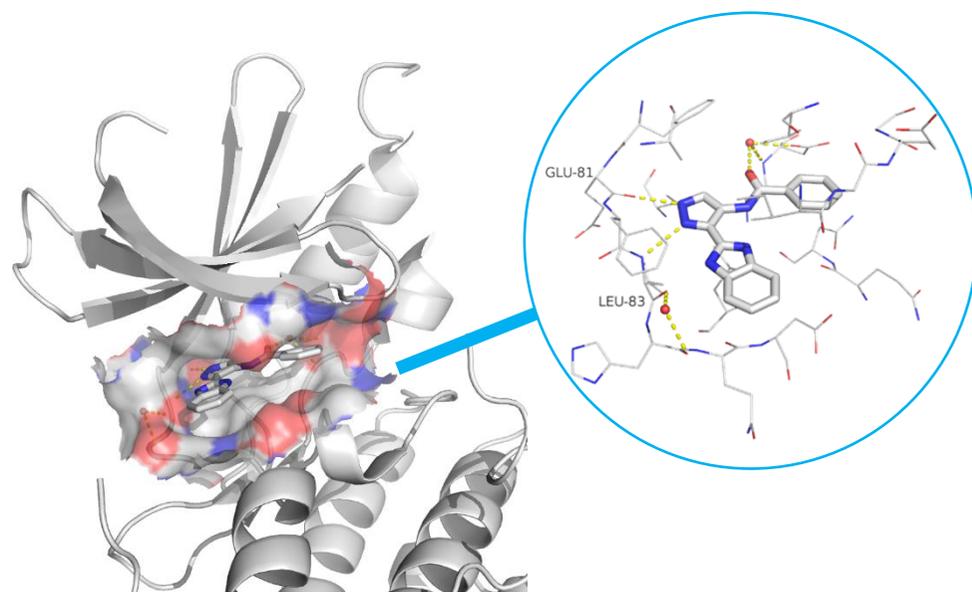
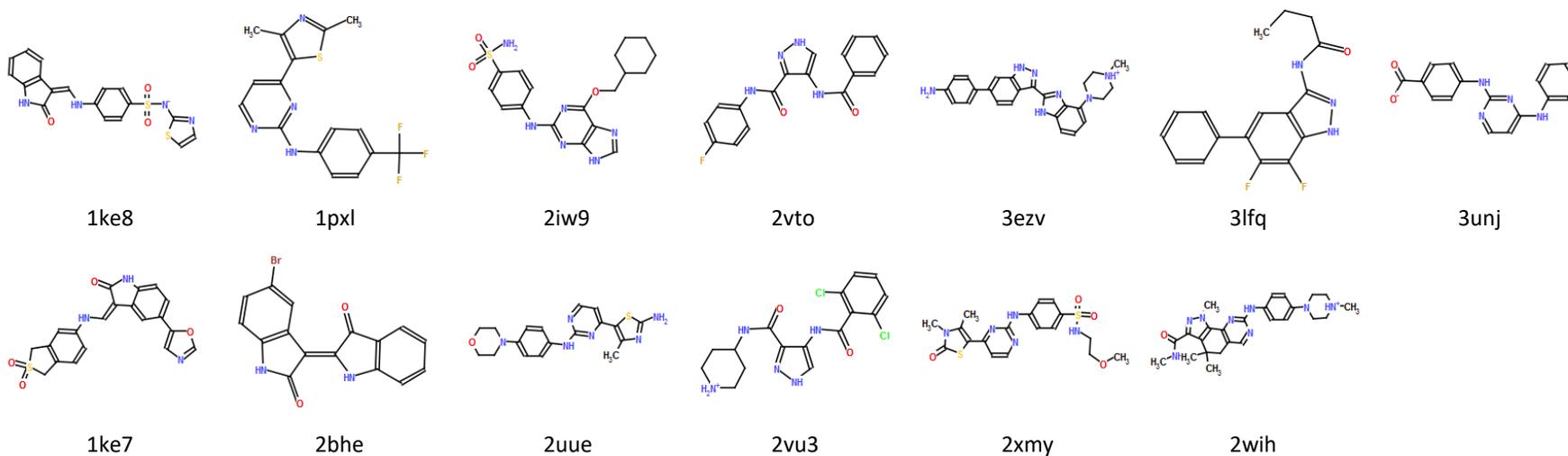


Figure1. CDK2 in complex with a fragment-like molecule. Hydrogen bond interactions are shown as dashed yellow lines.



1. Generating conformers

The CSD Conformer Generator provides the ability to both minimise molecular conformations and generate diverse conformer subsets based on CSD data. The methodology starts from an input 3D molecular structure with all hydrogen atoms present, which is optionally minimised in the first step. Subsequently, conformations are sampled based on CSD-derived rotamer distributions and ring templates. A final diverse set of conformers, clustered according to conformer similarity, is returned. Each conformer is locally optimised in torsion space.

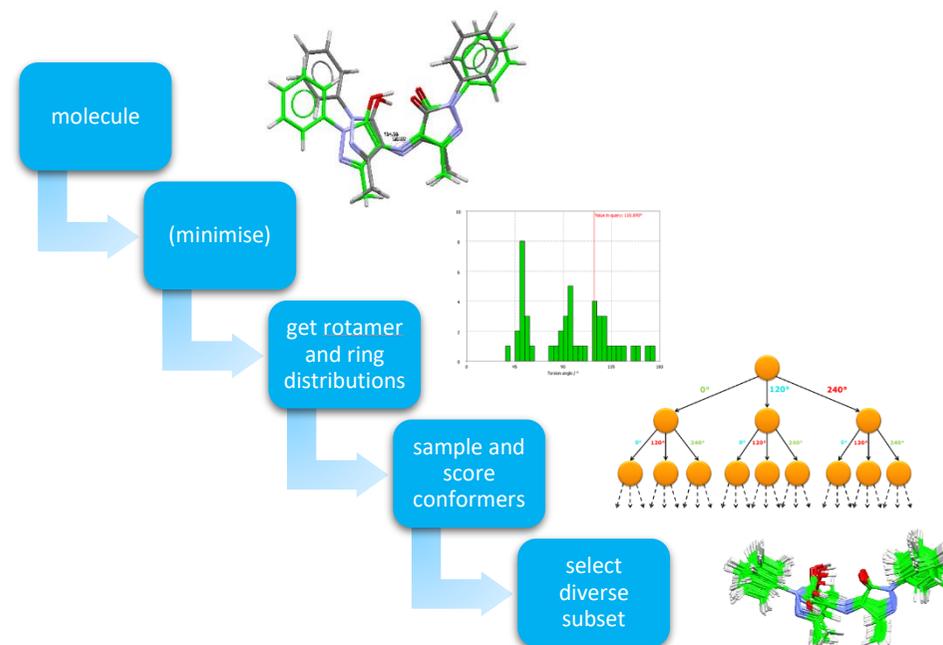
To avoid any bias from the experimentally observed conformation for this example, we regenerated 3D coordinates for the 13 molecules using the program Corina.

The CSD Conformer Generator requires as input 3D structures with all hydrogen atoms present. The ensembles of high probability knowledge-driven conformations can then be generated:

- 1) Through the command line utility of the CSD Conformer Generator program,
- 2) Using the Mercury interface,
- 3) Using the CSD Python API.

In this tutorial we provide a brief overview of how to use these three methods for conformer generation. Choose only ONE of these methodologies to generate the conformers of the 13 CDK2 inhibitors.

Generating Conformers Workflow



1.1 Command line utility

The command line utility of the CSD Conformer Generator is used from a command prompt.

- To open a command prompt on **Windows**, open File Explorer, go to the directory you are working in (e.g. `my_path_to\CDK2\input_molecules\Command_line`). Click on the address bar and type `cmd`, this will open the command prompt with the path to your current folder already set.

If you are a **macOS** user, use a terminal to go to your working directory (e.g. `my_path_to/CDK2/input_molecules/Command_line`).

Then, to launch the CSD Conformer Generator program, you have to specify the path to its executable:

- By default, on **Windows**, the conformer generator.exe is located in "C:\Program Files\CCDC\CSD_2020\Conformer Generator\"
- By default, on **macOS**, the conformer generator executable is located in `/Applications/CCDC/CSD_2020/Conformer Generator/bin/`

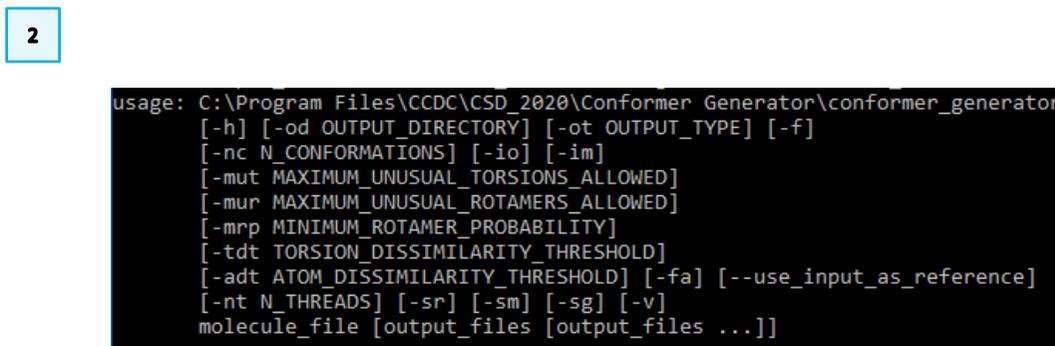
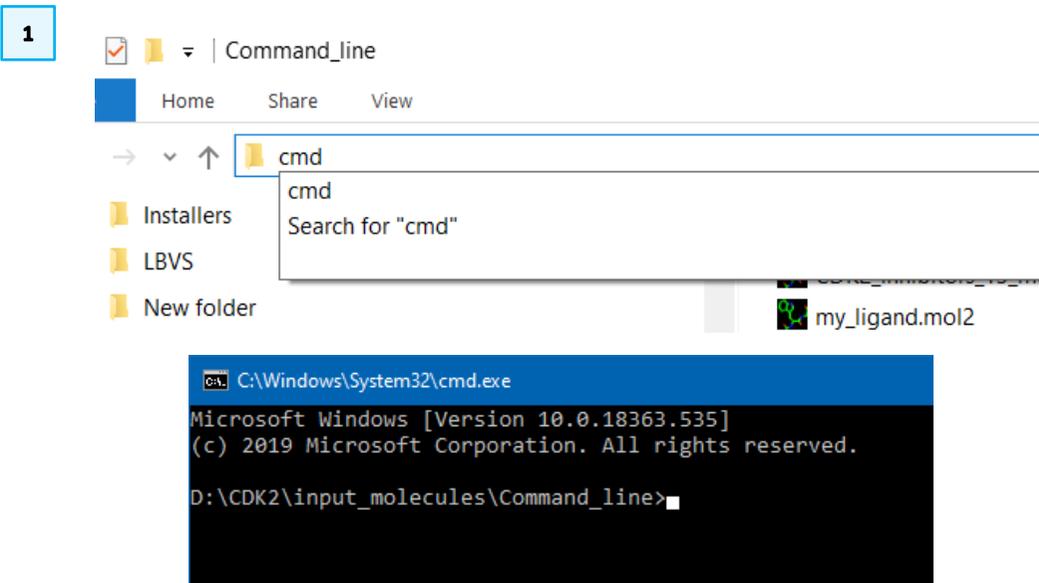
2. If you are a Windows user:

Type "C:\Program Files\CCDC\CSD_2020\Conformer Generator\conformer_generator" -h

If you are a macOS user:

Type `/Applications/CCDC/CSD_2020/Conformer Generator/bin/conformer_generator -h`

-h command provides help instructions on all the supported commands of the conformer generator tool.



3. At the simplest level, the CSD Conformer Generator can be used with an input file and an output file, with no further options.

If you are a Windows user:

Type the command:

```
"C:\Program Files\CCDC\CSD_2020\Conformer
Generator\conformer_generator" my_ligand.mol2
my_conformers.mol2
```

If you are a macOS user:

Type the command:

```
Applications/CCDC/CSD_2020/Conformer
Generator/bin/conformer_generator my_ligand.mol2
my_conformers.mol2
```

Note that supported formats for the input and output files are *.sdf* and *.mol2*.

By default, the program will generate up to 200 conformations for each molecule in the input file, sorted by their normalised score which is their likelihood based on the knowledge driven by similar molecules in the CSD.

4. If you wish, the normalised scores can be written out as a comma separated file (*.csv*) by extending the command

If you are a Windows user

Type the command:

```
"C:\Program Files\CCDC\CSD_2020\Conformer
Generator\conformer_generator" my_ligand.mol2
my_conformers.mol2 my_conformers.csv -f
```

If you are a macOS user:

Type the command:

```
Applications/CCDC/CSD_2020/Conformer
Generator/bin/conformer_generator my_ligand.mol2
my_conformers.mol2 my_conformers.csv -f
```

3

```
C:\Windows\System32\cmd.exe
D:\CDK2\input_molecules\Command_line>"C:\Program Files\CCDC\CSD_2020\Conformer Generator\
conformer_generator" my_ligand.mol2 my_conformers.mol2_
```

4

```
C:\Windows\System32\cmd.exe
D:\CDK2\input_molecules\Command_line>"C:\Program Files\CCDC\CSD_2020\Conformer Generator\
conformer_generator" my_ligand.mol2 my_conformers.mol2 my_conformers.csv -f_
```

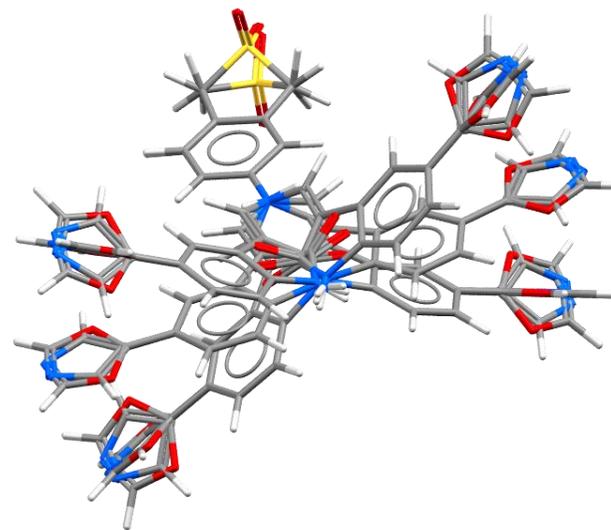


Figure 2: Conformers generated for the molecule *1ke7_lig_LS3.mol2*. By default, they are not superimposed on any particular substructure.

The `-f` command is required to overwrite any output files if already present in the output folder. In addition to the `my_conformers.mol2` file (Figure 2), which includes up to 200 conformers of `my_ligand.mol2`. The command above will automatically produce two other files:

- `my_conformers.csv` containing data for each ligand run through the conformer generator, and
- `my_conformers_per_conformer_scores.csv` containing data for each retained conformer per ligand run through the conformer generator. This file, shown on the right, includes several data items for each conformer such as the conformer probability, the clash score measured for this conformer, and a score that normalises the probability into the range of [0.0,1.0], where a zero score indicates the conformer with the highest likelihood and a score of 1.0 would indicate the conformer with the least likelihood.

reader.pass.molecule_name	conf_gen_conformers.pass.prob	conf_gen_conformers.pass.clash	conf_gen_conformers.pass.normalised_score	conf_gen_conformers.pass.conformer_id
1ke7_lig_LS3	-2.87943	0.381416	0	1
1ke7_lig_LS3	-2.87943	0.381416	1.18E-09	2
1ke7_lig_LS3	-2.87943	0.0290704	4.47E-09	3
1ke7_lig_LS3	-2.87943	0.0290704	5.65E-09	4
1ke7_lig_LS3	-3.57258	0.4253	0.0568534	5
1ke7_lig_LS3	-3.57258	0.4253	0.0568534	6
1ke7_lig_LS3	-3.57258	0.0311572	0.0568534	7
1ke7_lig_LS3	-3.57258	0.0311572	0.0568534	8
1ke7_lig_LS3	-7.58201	0.0027111	0.285716	9

5

- Finally, we are going to generate a separate output file of conformers for each of the 13 CDK2 inhibitors. This is performed by specifying a subdirectory to contain the output files, rather than specifying a single output file.

If you are a Windows user

Type the command:

```
"C:\Program Files\CCDC\CSD_2020\Conformer
Generator\conformer_generator"
CDK2_inhibitors_13_molecules.mol2 -ot mol2 -od conformers
```

If you are a macOS user:

Type the command:

```
Applications/CCDC/CSD_2020/Conformer
Generator/bin/conformer_generator
CDK2_inhibitors_13_molecules.mol2 -ot mol2 -od conformers
```

This will create a new directory called `conformers` containing all the generated conformers in `.mol2` format. Note that molecule `2bhe_lig_BRY`

```
C:\Windows\System32\cmd.exe
D:\CDK2\input_molecules\Command_line>"C:\Program Files\CCDC\CSD_2020\Conformer Generator\conformer_generator" CDK2_inhibitors_13_molecules.mol2 -ot mol2 -od conformers_
```

does not contain any rotatable bond; therefore, it will be minimised but no conformers will be generated for it. You can retrieve information about the number of conformers generated by looking at the *conformer_generator.log* file in your `my_path_to/CDK2/input_molecules/Command_line` folder.

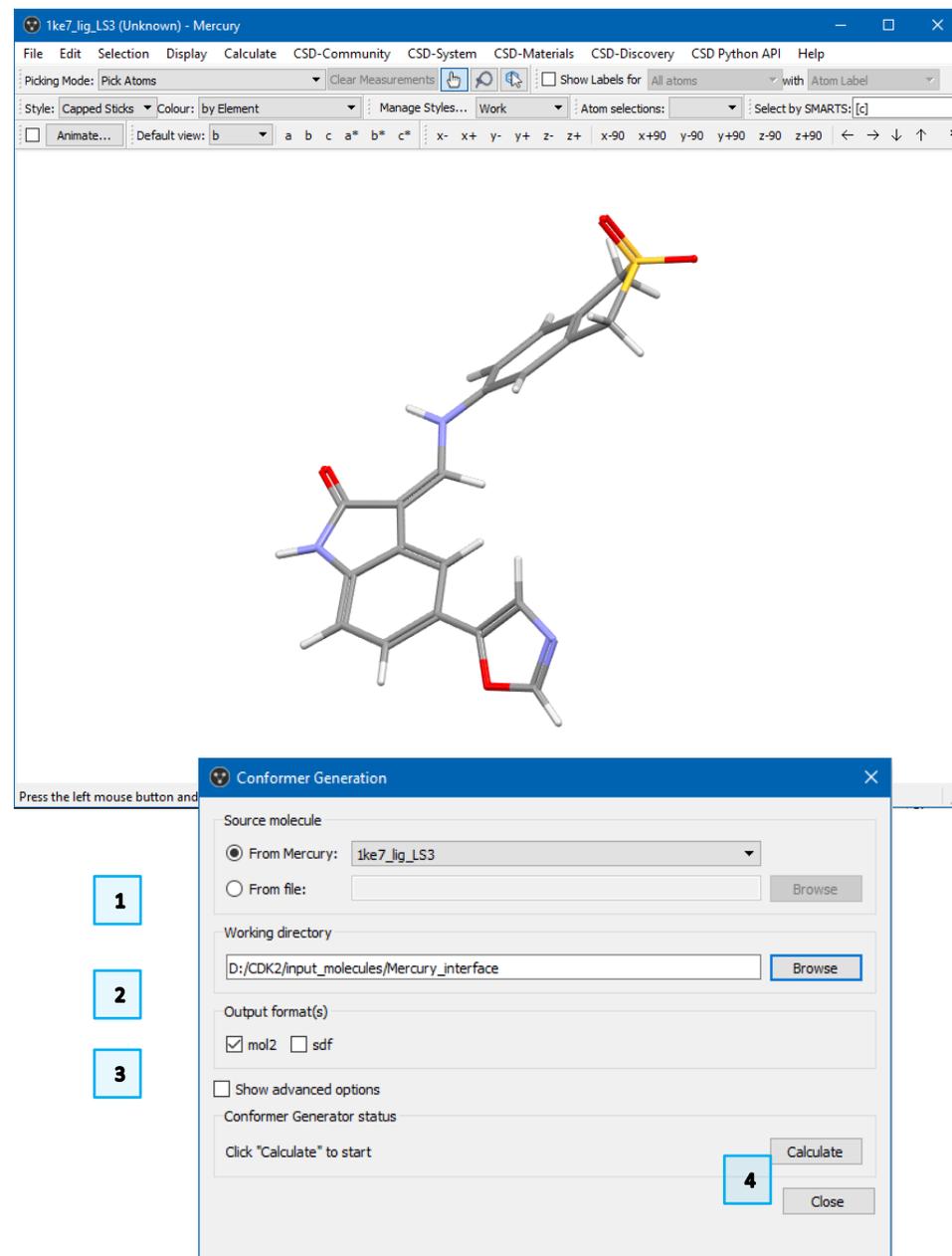
If less than 200 conformers are generated (e.g. for ligand *1pxl_lig_CK4*) it means that, according to the CSD Conformer Generator clustering algorithm, the search space may be exhausted before 200 distinct conformers are generated.

Note that it is possible to run the CSD Conformer Generator using multiple threads on a given machine using the `-nt` command. Please note that all threads will share the memory available to the conformer generator process.

1.2 The Mercury interface

The CSD Conformer Generator can also be accessed through the **CSD-Materials** and **CSD-Discovery** menus within Mercury .

1. In the *Conformer Generation* window in the *Source molecule* session you can specify a CSD entry, an in-house database entry, or another structure as the input molecule for conformer generation. Alternatively, you can select a file on the disk as the input by selecting **From file:** option and clicking on the **Browse** button. In order to load the file *1ke7_lig_LS3.mol2*, navigate to `my_path_to\CDK2\input_molecules\Mercury_interface`
2. The working directory can be selected under *Working Directory* section of the *Conformer Generation* window. Click the **Browse** button and navigate to `my_path_to\CDK2\input_molecules\Mercury_interface`. A folder will be created in this directory named after the identifier of the input molecule used. Both a copy of the input molecule and the generated conformers will be written to the output folder.



1ke7_lig_LS3 (Unknown) - Mercury

File Edit Selection Display Calculate CSD-Community CSD-System CSD-Materials CSD-Discovery CSD Python API Help

Picking Mode: Pick Atoms Clear Measurements Show Labels for All atoms with Atom Label

Style: Capped Sticks Colour: by Element Manage Styles... Work Atom selections: Select by SMARTS:[c]

Animate... Default view: b a b c a* b* c* x- x+ y- y+ z- z+ x-90 x+90 y-90 y+90 z-90 z+90 < > <> <>

Press the left mouse button and

Conformer Generation

Source molecule

From Mercury: 1ke7_lig_LS3

From file: Browse

Working directory

D:/CDK2/input_molecules/Mercury_interface Browse

Output format(s)

mol2 sdf

Show advanced options

Conformer Generator status

Click "Calculate" to start

Calculate

Close

- Conformers may be generated in *.mol2* or *.sdf* format by choosing the relevant checkboxes under *Output format(s)*. The default *.mol2* format will be used in this tutorial.
- Click the **Calculate** button to generate conformers.
- A multi-file containing 128 conformers (Figure 3) will be saved in `my_path_to\CDK2\input_molecules\Mercury_interface\1ke7_lig_LS3\CSD_Conformer_Generator\2019_12_18_10_15_09\conformers.mol2`.

Please note that `2019_12_18_10_15_09` is the date (18/12/2019) and time (10:15:09) when the calculation was run. You will have a different folder name corresponding to your own calculation date and time.

You can use the same procedure to create the conformer libraries of the other 12 CDK2 inhibitors.

1.3 The CSD Python API

The CSD Conformer Generator functionality is also fully available via the CSD Python API, allowing you to link conformer generation seamlessly to other applications and enabling complex workflows and analyses.

To generate conformer libraries for all 13 CDK2 inhibitors, via the CSD Python API, all you need are a few very simple Python statements. The script on the right of the next page shows how to generate conformers for all *.mol2* files stored in a given working directory. Prior to generation, molecules are standardised according to CSD conventions, and the resulting conformers overlaid using all heavy atoms as shown in Session 3 of this tutorial.

- If you are a Windows user, open File Explorer, go to the directory you are working in. This may be:

```
my_path_to\CDK2\input_molecules\CSD_Python_API
```

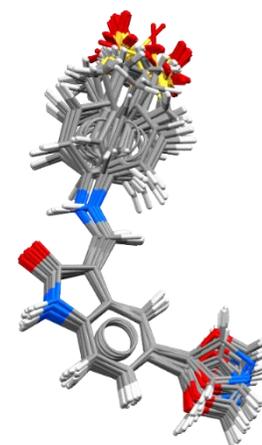
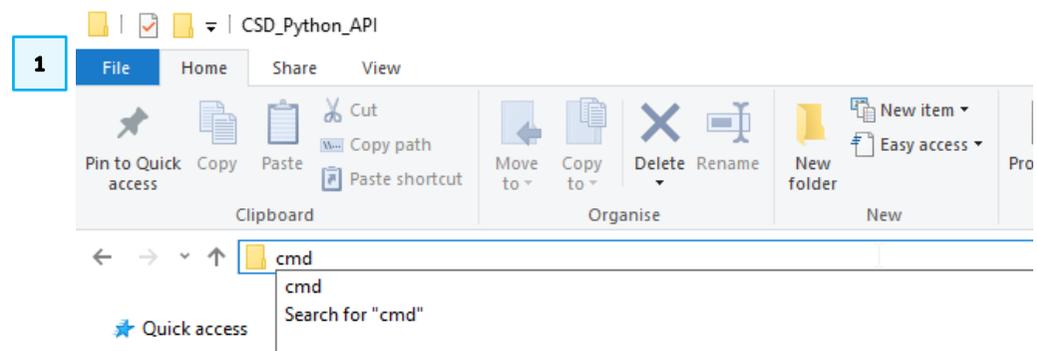


Figure 3. Ensemble of conformers generated for the molecule *1ke7_lig_LS3.mol2*. the input molecule was used as reference to superimpose the resulting conformations on all heavy atoms.



Click on the address bar and type `cmd`, this will open the command prompt with the path to your current folder already set.

If you are a macOS user, use a terminal to go to your working directory (e.g. `my_path_to/CDK2/input_molecules/CSD_Python_API`).

- In order to use the CSD python API, we need to activate the Python environment where it has been installed together with its dependencies.

2

```
C:\Windows\System32\cmd.exe
(c) 2019 Microsoft Corporation. All rights reserved.

D:\CDK2\input_molecules\CSD_Python_API>"C:\Program Files\CCDC\Python_API_2020\miniconda\Scripts\activate.bat"

(base) D:\CDK2\input_molecules\CSD_Python_API>
```

If you are a Windows user:

In the command prompt type

```
"C:\Program
Files\CCDC\Python_API_2020\miniconda\Scripts\activate.bat"
"
```

If you are a macOS user:

In the terminal type:

```
source
Applications/CCDC/Python_API_2020/miniconda/bin/activate
```

The above command will run the CSD Python API activation script. As you can see you are now working in a different environment called `(base)`.

Any subsequent use of Python in that terminal will use the version installed with the CSD-System. This step can be reversed by using the `deactivate` command.

- We are now ready to use our Python script to generate the conformers of the 13 CDK2 inhibitors.

Type `python generate_conf.py` to run it.

Note that parameters can be changed via the `settings` object:

- `max_conformers`**
Maximum number of conformers to generate.
- `max_unusual_torsions`**
Number of unusual torsions allowed per conformer.
- `normalised_score_threshold`**

3

```
D:\CDK2\input_molecules\CSD_Python_API\generate_conf.py - Notepad++
File Edit Search View Encoding Language Settings Tools Macro Run Window ?
generate_conf.py
1 from ccdc.io import MoleculeReader, MoleculeWriter
2 from ccdc.conformer import ConformerGenerator
3 import os
4
5 input_molecules = [file for file in os.listdir(os.getcwd()) if file.endswith(".mol2")]
6
7 conformer_generator = ConformerGenerator()
8 for file in input_molecules:
9     mol = MoleculeReader(file)[0]
10    mol.standardise_aromatic_bonds()
11    mol.standardise_delocalised_bonds()
12    conformer_generator.settings.superimpose_conformers_onto_reference = True
13    conformers = conformer_generator.generate(mol)
14    with MoleculeWriter('%s_conformers.mol2' %mol.identifier) as mol_writer:
15        for c in conformers:
16            mol_writer.write(c.molecule)
```

Maximum deviation from the theoretically achievable normalised conformer probability (0="best", 1="worst").

- **superimpose_conformers_onto_reference**
Whether or not to superimpose to a common reference.

Conclusion

Using either the command line utility of the CSD Conformer Generator, or the Mercury interface, or the CSD Python API, you were able to generate a number of high-probability conformations for the dataset of CDK2 inhibitors. These libraries can now be passed on to other tools for post-processing. For example, they can be used as input for the CSD Ligand Overlay to generate overlay hypotheses for pharmacophore modelling.

2. Overlaying ligands

The CSD Ligand Overlay program aims to overlay sets of flexible molecules that are assumed to be ligands of the same protein. The molecules to be overlaid are divided into *features* such as hydrogen-bond donors and acceptors, and hydrophobic groups.

In the absence of any protein-structure information, there may be multiple overlays that present legitimate hypotheses for binding.

In order to address this problem, the algorithm used by the overlay program returns multiple diverse overlays that all give potential pharmacophore hypotheses to test.

1. Launch **Hermes** by clicking the Hermes Icon on your desktop or if you are a macOS user launching it from the Applications menu.
2. From the top-level menu choose **Calculate > Ligand Overlay > Start Wizard...** It is possible to load an existing project, but as no calculations have been run so far, this is not appropriate. Here we will generate some possible overlays of the 13 CDK2 inhibitors, as discussed in the Introduction.
3. In the *Choose what you want to do* dialogue, choose option to **Search for new overlays & pharmacophores** and click **Next** (or **Continue** on MacOS).
4. In the *Create or Load a Project* dialogue, click **Create Project...** This allows you to navigate to an existing directory or create a new directory for the results to be saved. You can also edit the project name in the *Project name:* textbox. Once you have done this, click **Next**.
5. In the *Select Conformers* dialogue, you will be prompted to add the conformer files for all ligands that you wish to overlay.

Note that you can either load pre-generated conformer files (e.g. as the ones generated in Section 1 of this tutorial) or supply each ligand as a single conformer file (in which case by default the program will generate conformers

1

2

3

4

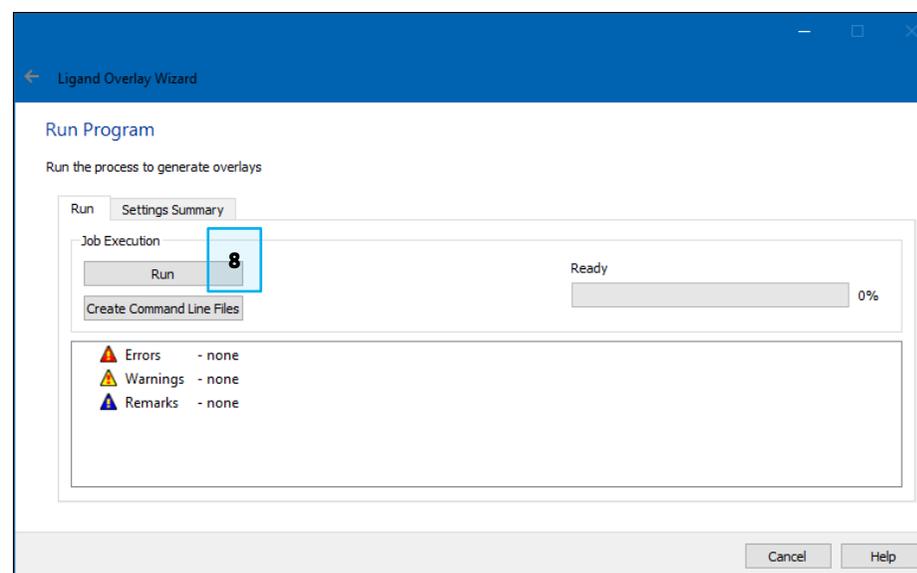
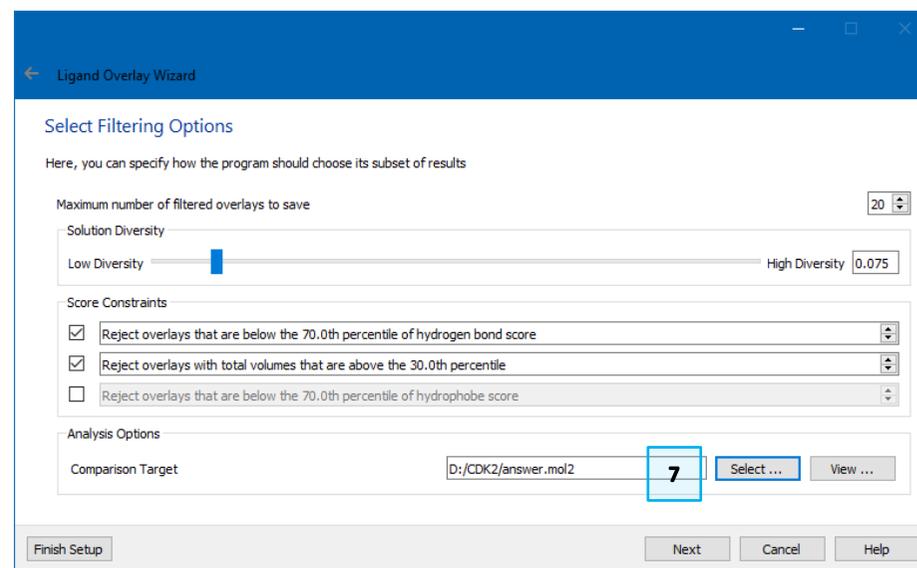
5

File Name	Label	Count	Generate Conformations
D:/CDK2/input_molecules/LigandOverlay/1ke7_lig_L53.mol2	01_1ke7_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/1ke8_lig_L54.mol2	02_1ke8_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/1px_lig_CK4.mol2	03_1px_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2bhe_lig_BR9.mol2	04_2bhe_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2iiv9_lig_4SP.mol2	05_2iiv9_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2uue_lig_MTZ.mol2	06_2uue_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2vto_lig_LZ8.mol2	07_2vto_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2vz3_lig_LZE.mol2	08_2vz3_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2vwh_lig_P48.mol2	09_2vwh_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/2zmy_lig_CDK.mol2	10_2zmy_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/3evz_lig_EZV.mol2	11_3evz_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/3ffq_lig_A28.mol2	12_3ffq_lig	1	<input checked="" type="checkbox"/> Yes
D:/CDK2/input_molecules/LigandOverlay/3unj_lig_0BX.mol2	13_3unj_lig	1	<input checked="" type="checkbox"/> Yes

for each ligand). Here we are going to use the second option by clicking **Add Files...** and selecting all ligands in `my_path_to\CDK2\input_molecules\LigandOverlay`. Ensure the **Try to use molecule names as labels** check-box is ticked.

- Click **Next** to go to the *Select Parameter Settings* dialog. As the default parameters are appropriate for this tutorial we could just click **Finish Setup** and start the calculation. However, as we are performing a retrospective study of known CDK2 inhibitors for which the crystal structures of the protein-ligand complexes are available, we will use this information to assess the quality of the results.
- Click **More Options** and then **Next** until you get to the *Select Filtering Options* dialog. In the *Analysis Options* section, include the **Comparison Target** by clicking **Select...** and navigate to the true overlay *answer.mol2*, which can be found in `my_path_to\CDK2\`.
- Click **Next** and finally **Run** to start the calculation.

The progress bar helps to monitor the progress of the calculation as the program goes through several stages: *Generating Conformers*, *Initializing*, *Generating Overlays*, *Scoring Overlays*, *Filtering Overlays*, *Optimising Overlays* and *Analysing Overlays*.



Analysis of the results

Once the calculation has completed, it is easy to browse through the solutions using the Hermes *Data Analysis* window that is automatically generated.

Note that, due to the stochastic nature of the algorithm, results are not reproducible and therefore, here we will discuss an example run.

The solutions can be analysed using the data contained in *JOB1_MAP_OF_SOLUTIONS Spreadsheet 1* in the *Data Analysis* window. The volume, H-bond, hydrophobe coplanarity and internal energy (internal strain) scores of each solution and of the true overlay (target) are given in the four columns of the map file headed *volume*, *hbond*, *hydrophobe* and *internal energy*, respectively.

These scores can be useful, particularly for comparison with the scores that could be obtained for the true crystallographic overlay (i.e. target).

The dominance corresponds to the final Pareto rank of the overlays. Solutions with dominance equal to zero are called non-dominated or equivalent, as they all represent the best compromise in optimising the parameters used to score the solutions. Good solutions can usually be identified by considering first the overlays with zero dominance, and in particular focusing on those that also exhibit low internal strain scores.

Clicking on each solution in the *JOB1_MAP_OF_SOLUTIONS Spreadsheet 1* will display it in the Hermes molecule display area.

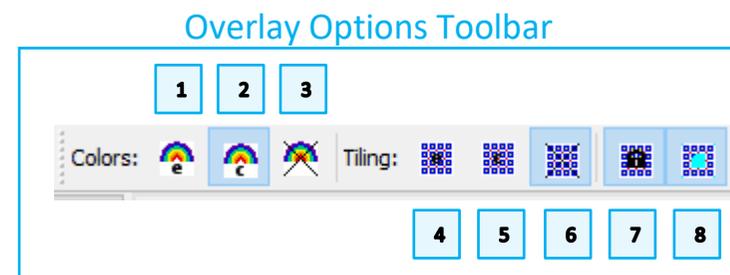
Identifier	NAME	volume	hbond	hydrophobe	internal energy	dominance
JOB1_MAP_OF_SOLUTIONS\SOLN_01 0	SOLN_01	896.2500	697.6790	1058.4590	155.1570	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_02 1	SOLN_02	912.1250	222.1490	1067.3290	159.0800	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_03 2	SOLN_03	874.3750	803.7290	820.7750	98.6250	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_04 3	SOLN_04	938.1250	827.2320	787.1490	122.4940	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_05 4	SOLN_05	869.6250	797.7920	742.5640	125.1700	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_06 5	SOLN_06	899.3750	838.5670	771.5750	144.9210	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_07 6	SOLN_07	849.2500	403.4050	995.4360	139.3440	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_08 7	SOLN_08	868.8750	394.6910	1032.5270	187.8560	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_09 8	SOLN_09	899.3750	775.1900	901.8590	118.4970	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_10 9	SOLN_10	954.0000	678.7470	912.0350	102.8550	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_11 10	SOLN_11	905.3750	807.8910	768.2670	94.2230	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_12 11	SOLN_12	908.6250	355.0730	1050.1970	166.5260	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_13 12	SOLN_13	922.7500	769.2390	827.2390	146.6890	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_14 13	SOLN_14	892.5000	145.4970	983.0060	162.0470	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_15 14	SOLN_15	895.5000	194.0400	938.2300	215.6780	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_16 15	SOLN_16	870.1250	170.5490	906.2800	174.1890	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_17 16	SOLN_17	928.3750	730.2710	806.9360	127.7160	3.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_18 17	SOLN_18	914.0000	779.8000	753.2990	242.2130	3.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_19 18	SOLN_19	991.2500	185.3520	943.9940	159.9860	5.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_20 19	SOLN_20	956.2500	710.0790	731.0830	160.6970	9.0000
JOB1_MAP_OF_SOLUTIONS\ target 20	target	899.3750	786.3560	915.3620	151.6050	-99.0000

True crystallographic overlay

Solutions with zero dominance automatically selected when the calculation has completed

The Overlay Options Toolbar in Hermes can be used to help in this analysis.

Colouring options:



1. *Colour entries by rainbow*, where an entry is a given overlay solution.
2. *Colour components by rainbow*, where a component is a given molecule in the overlay.
3. Turn *Rainbow colouring off*.

Tiling options:

4. *Tile by entry*, to display one overlay solution per tile.
 5. *Tile by component*, to display the same molecule from all selected solutions in each tile.
 6. *Remove Tiling*.
 7. *Link rotations and translations across all tiles*.
 8. *Label tiles with entry identifier*.
1. In order to compare different overlay solutions, click **Tile by entry** and **Link rotations and translations across all tiles**. In addition, be sure that the **Label tiles entry identifier** is activated.

Each tile will now be populated with one of the selected solutions, and by default, each molecule will be coloured by rainbow. This display helps understanding and comparing the overall shape of the different overlays.

- a. It is possible to clear the display by clicking the tick box next to *JOB1_MAP_OF_SOLUTIONS* under *All Entries* in the *Molecule Explorer* tab. After that, individual molecules can be inspected one after the other to see what feature matches have been found.

2. In order to assess the performance of the CSD Ligand Overlay to reproduce the superimposition of the 13 input molecules observed experimentally, select one of the non-dominated solutions (i.e. with zero dominance) and the target from the *Data Analysis* table.

The screenshot displays the Hermes software interface. At the top, a menu bar includes options like File, Edit, Selection, Display, Calculate, Descriptors, GOLD, Databases, CSD Python API, and Help. Below the menu is a toolbar with icons for highlighting, depth cueing, stereo, graphics objects, and showing/hiding hydrogens and unknown atoms. A 'Picking Mode' dropdown is set to 'Pick Atoms'. The 'Molecule Explorer' panel on the left shows a tree view with 'All Entries' expanded to show 'JOB1_MAP_OF_SOLUTIONS' and its sub-entries (JOB1_MAP_OF_SOLUTIONS|SOLN_010 and |SOLN_021). A blue box labeled '1' points to the 'Tiling' button in the toolbar. Another blue box labeled 'a' points to the 'JOB1_MAP_OF_SOLUTIONS' entry in the Molecule Explorer tree. The main display area shows a 2x2 grid of molecular overlay tiles, each displaying a different solution (SOLN_010, SOLN_021, SOLN_043, SOLN_054) with molecules colored by rainbow.

Here we show the comparison of solution_03 with the experimental overlay. Click **Tile by component** to see how the conformation picked by the overlay program compares with the binding conformation of the equivalent ligand in the associated protein crystal structure.

As you can see, the predicted conformations are all close to the equivalent experimental ones. Small deviations are observed in the flexible terminal groups, but this is expected from a ligand-based tool when there is not enough information to support a specific orientation.

We recommend looking at a few diverse overlay solutions, especially in a fully ligand-based approach. They represent alternative, plausible pharmacophore hypotheses that can be used as a query during library screening.

Please bear in mind that here we have done a retrospective analysis where binding site and binding mode information about the protein are available. This can also be a valuable approach to prospectively combine ligand- and structure-based techniques. Therefore, the same concepts and analyses can be applied when using the tool in a real drug design project.

The top screenshot shows the 'Data Analysis' window with a table of solutions. A blue box with the number '2' points to the 'Tile by component' button in the 3D interface below. The table lists various solutions with their respective properties:

Identifier	NAME	volume	hboard	hydrophobe	internal energy	dominance
JOB1_MAP_OF_SOLUTIONS\SOLN_010	SOLN_01	896.2500	697.6790	1058.4590	155.1570	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_021	SOLN_02	912.1250	222.1490	1067.3290	159.0890	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_032	SOLN_03	874.3750	893.7290	826.7750	98.6250	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_043	SOLN_04	938.1250	827.2320	787.1490	122.4940	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_054	SOLN_05	869.6250	797.7920	742.5640	125.1700	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_065	SOLN_06	899.3750	838.5670	771.5750	144.9210	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_076	SOLN_07	849.2500	403.4050	995.4360	139.3440	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_087	SOLN_08	868.8750	394.6910	1032.5270	187.8560	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_098	SOLN_09	899.3750	775.1900	901.8590	118.4970	0.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_109	SOLN_10	954.0000	678.7470	912.0350	102.8550	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1110	SOLN_11	905.3750	807.8910	768.2670	94.2230	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1211	SOLN_12	908.6250	355.0730	1050.1970	166.5260	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1312	SOLN_13	922.7500	769.2390	827.2390	146.6890	1.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1413	SOLN_14	892.5000	145.4970	983.0060	162.0470	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1514	SOLN_15	895.5000	194.0400	938.2300	215.6780	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1615	SOLN_16	870.1250	170.5490	906.2800	174.1890	2.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1716	SOLN_17	928.3750	730.2710	806.9360	127.7160	3.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1817	SOLN_18	914.0000	779.8000	753.2990	242.2130	3.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_1918	SOLN_19	991.2500	185.3520	943.9940	159.9960	5.0000
JOB1_MAP_OF_SOLUTIONS\SOLN_2019	SOLN_20	956.2500	710.0790	731.0830	160.6970	9.0000
JOB1_MAP_OF_SOLUTIONS\target120	target	899.3750	786.3560	915.3620	151.6050	-99.0000

The bottom screenshot shows the 3D molecular overlay interface. A grid of 16 panels displays different solutions (e.g., 01_1000_0g, 02_1000_0g, etc.) overlaid on the target structure. A blue box with the number '2' highlights the 'Tile by component' button in the top toolbar.

Pharmacophore model

In your output directory `my_path_to\CKD2\job1` you should find these:

Files	Explanation
<code>all_generated.chrm</code>	File containing five mandatory header lines followed by all 10,000 generated overlays (pre-filtering), stored in chromosome form.
<code>solutions.chrm</code>	File containing five mandatory header lines followed by the 20 overlays chosen by the filtering step, stored in chromosome form.
<code>solution_XX.mol2</code> <code>solution_XX.sdf</code>	The 20 solutions chosen by the filtering step and optimised to bring groups into tighter alignment – as a multi mol2 or sdf file.
<code>pharmacophores\solution_pharm_XX.mol2</code>	The pharmacophore representations for the 20 solutions.
<code>solution_target.sdf</code> <code>solution_target.mol2</code>	The true overlay, if indicated in the overlay generation.
<code>target_overlay_XX.mol2</code> <code>target_overlay_XX.sdf</code>	The 20 optimised solutions, superimposed onto the true overlay if indicated in the overlay generation – as a multi mol2 or sdf file.
<code>job1_map_of_solutions.csv</code>	The dominance, consensus coefficients ($x_{consensus}$, $y_{consensus}$ and $z_{consensus}$), the pharmacophore coefficients (x_{pharm} , y_{pharm} and z_{pharm}), the superposition coefficients (x_{super} , y_{super} and z_{super}), the three scores (<i>volume</i> , <i>hbond</i> and <i>hydrophobe</i>) and the internal energy (<i>e</i>) for all 20 chosen solutions after filtering and optimisation.
<code>job1_dissimilarity_matrix.csv</code>	The dissimilarity matrices of the 20 chosen solutions based on all three coefficients (consensus, superposition and pharmacophore).
<code>job1_dissimilarity_matrix_consensus.csv</code>	The dissimilarity matrices of the 20 chosen solutions based on the consensus coefficients ($x_{consensus}$, $y_{consensus}$ and $z_{consensus}$) only.
<code>job1_dissimilarity_matrix_pharm.csv</code>	The dissimilarity matrices of the 20 chosen solutions based on the pharmacophore coefficients (x_{pharm} , y_{pharm} and z_{pharm}) only.
<code>job1_dissimilarity_matrix_super.csv</code>	The dissimilarity matrices of the 20 chosen solutions based on the superposition coefficients (x_{super} , y_{super} and z_{super}) only.
<code>scores.csv</code>	The final <i>volume</i> , <i>hbond</i> and <i>hydrophobe</i> scores for the 20 chosen overlays after filtering and optimisation.
<code>log.txt</code>	The log file for the calculation, including the values of all parameters (those left as default and those set) and a summary of all steps and of their duration.

Open a new **Hermes** session and load one of the overlay solutions you are interested in, by clicking on **File > Open** from the main menu.

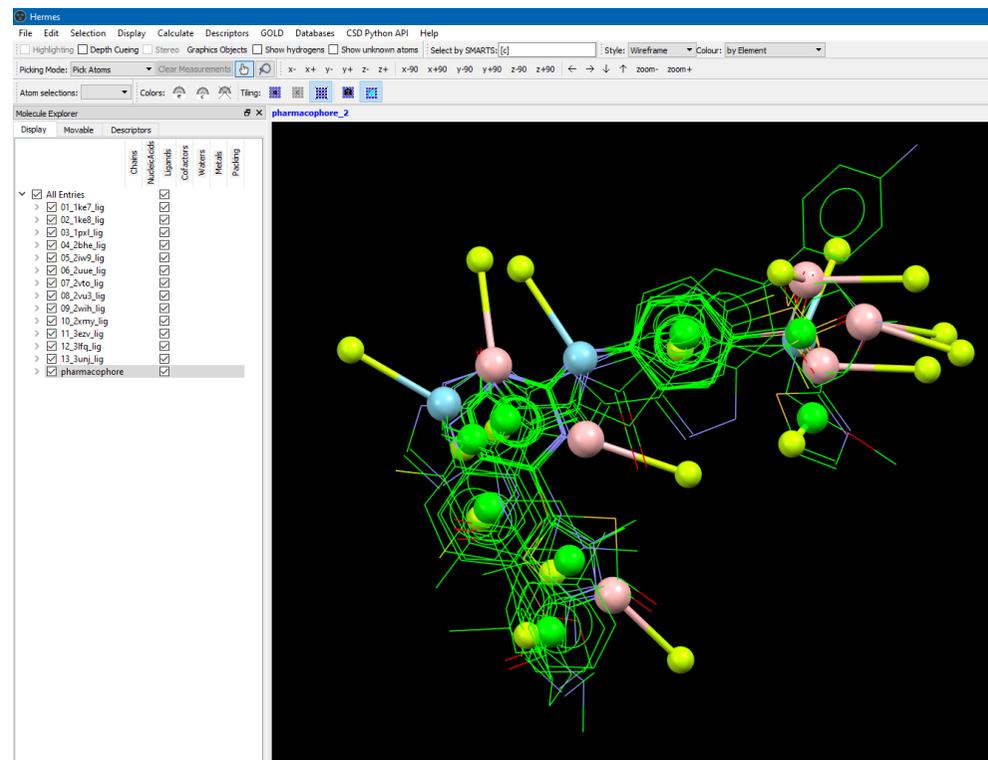
Also load the corresponding pharmacophore for the chosen solution from `my_path_to\CKD2\job1\pharmacophores`. An example is shown here where the overlay is rendered as a green wireframe and the pharmacophore features are rendered as ball and stick and color-coded based on their type. Acceptor groups are pink, donor groups are cyan and hydrophobic points are green. The projected points are all green.

Conclusion

A number of overlay hypotheses have been produced which in turn can be used as a query for a ligand-based virtual screening experiment.

You should now be familiar with the essential steps of the CSD Ligand Overlay wizard; how to interpret the results with the aid of the tiling options; and how to view the pharmacophore model associated with each solution.

For more information on the advanced settings of the CSD Ligand Overlay please see the main documentation at https://www.ccdc.cam.ac.uk/support-and-resources/ccdcresources/CSD_Ligand_Overlay_User_Guide.pdf



3. Field-based virtual screening

Virtual screening is the computational equivalent of biological screening and is used to score, rank, and/or filter a set of structures by using one or more computational procedures. The CSD Ligand Screener can be used to screen a library of compounds against a pharmacophore query obtained from one or multiple overlaid ligands. The algorithm generalises the 3D pharmacophore definition using atom property fields that are created around the query based on user-defined atom types and potentials. Three main steps are performed:

- Generation of a field potential from the query and creation of fitting points in hotspots.
- Global optimisation of the translation and rotation of each ligand by generating conformer libraries and then fitting the ligand atoms to pre-calculated fitting points.
- Scoring of each screened ligand with numerical gradients.

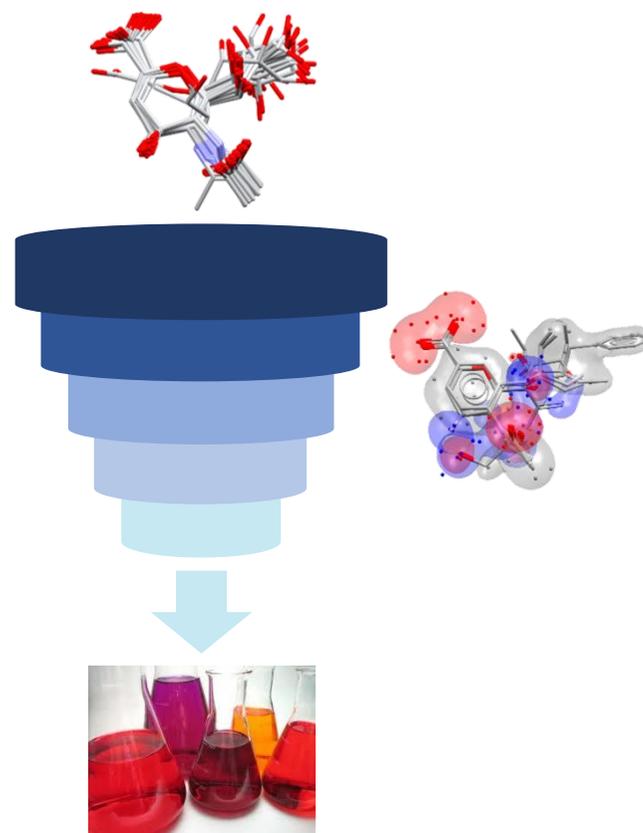
The methodology requires two input files (*i.e.* a query and the molecules to screen) where all atoms have 3D coordinates.

This tutorial will take you through the steps needed to set up a screening calculation using the CSD Python API. You will learn how to screen a library of known active and decoy molecules against a pharmacophore model generated through the CSD Ligand Overlay. You will calculate several enrichment metrics and plot a ROC curve to assess the success of the model in scoring and ranking the active molecules earlier than the decoys.

Retrieval of actives and decoys

A benchmarking data set for ligand-based virtual screening is publicly available (*i.e.* [DUD LIB VS 1.0](#)) and contains 13 targets including CDK2. The collection of 47 actives and 2070 decoys for this target have been downloaded and are provided with this tutorial in the `my_path_to\CKD2\Screening` subdirectory:

- Actives: `cdk2_clustered_3D_MM.sdf`
- Decoys: `DUD_cdk2_decoys_ID_pass_MWPass_I_MM.sdf`



DUD SUBSETS FOR LIGAND-BASED VIRTUAL SCREENING.

A. Jahn, G. Hinselmann, N. Fechner and A. Zell, "Optimal assignment methods for ligand-based virtual screening", *Journal of Cheminformatics*, 2009, **1**, 14. DOI: 10.1186/1758-2946-1-14.

Library screening

The CSD Python API provides programmatic access to the CSD data and the CSD-System functionality, as well as to features that have never been exposed within an interface. One of these features is the CSD Ligand Screener.

Through Python scripting it is possible to build highly tailored custom applications to help you answer detailed research questions, or to automate frequently performed analysis steps.

Here we will use an example script (i.e. *ligand_based_VS.py*) that performs the three steps described above:

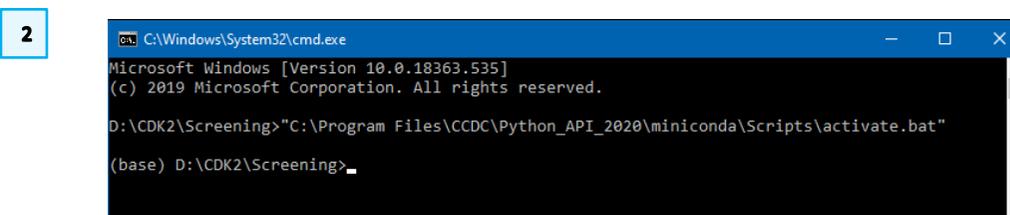
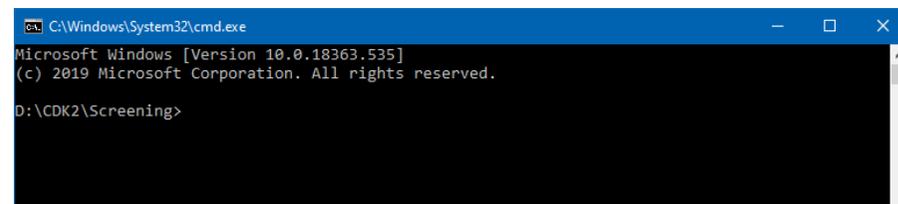
- We will use a solution obtained from the CSD Ligand Overlay to generate our query.
- We will generate knowledge-driven conformers for the library of actives and decoys taken from DUD LIB VS.
- The program will optimise the placement of every conformer onto the fields, but only the top-scoring will be retained and used for ranking.
- The ranked list of scores will be saved to a .csv file (the more negative the scores the better).

1. **For Windows users:** open File Explorer, go to the directory you are working in (e.g. `D:\CDK2\Screening`) and click on the address bar. Now simply type `cmd` in the address bar. It will open the command prompt with the path to your current folder already set.

If you are a macOS user: use a terminal to go to your working directory (e.g. `D:/CDK2/Screening`).

2. In order to use the CSD python API, we need to activate the Python environment where it has been installed together with its dependencies.

If you are a Windows user:



In the command prompt type "C:\Program Files\CCDC\Python_API_2020\miniconda\Scripts\activate.bat"

If you are a macOS user:

In the terminal type

```
source
Applications/CCDC/Python_API_2020/miniconda/bin/activate
```

As you can see you are now working in a different environment called (base) .

3. We are now ready to use the *ligand_based_VS.py* script.

Type `python ligand_based_VS.py -h` to see all the options supported by the script, what are the required input files and which parameters can be changed by the user.

4. The command to run the script is:

```
python ligand_based_VS.py -q solution_02.sdf -a
cdk2_clustered_3D_MM.sdf -d
DUD_cdk2_decoys_ID_pass_MWPass_I_MM.sdf -n 200 -t 3
```

where:

- *solution_02.sdf* is an overlay hypothesis we've selected from a previous CSD Ligand Overlay calculation; you may use it or replace it with another promising solution from the CSD Ligand Overlay calculation you have just run;
- up to 200 conformations will be generated for each active and decoy molecule;
- the calculation will be distributed over 3 threads.

A basic Python script showing the simple lines required for screening is provided in the tutorial folder (*example.py*) and shown on the right.

This script shows the basic line for screening the input conformation of the molecules in *my_library.sdf* against an overlay of ligands *solution_02.sdf* as the

3

```
C:\Windows\System32\cmd.exe
Microsoft Windows [Version 10.0.18363.535]
(c) 2019 Microsoft Corporation. All rights reserved.

D:\CDK2\Screening>"C:\Program Files\CCDC\Python_API_2020\miniconda\Scripts\activate.bat"

(base) D:\CDK2\Screening>python ligand_based_VS.py -h
usage: ligand_based_VS.py [-h] [-q QUERY] [-a ACTIVES] [-d DECOYS] [-n NCONFS]
                        [-t THREADS] [-o OUTPUT_DIRECTORY]

ligand_based_VS.py - simple interface to the ligand screener.

optional arguments:
  -h, --help            show this help message and exit
  -q QUERY, --query QUERY
                        Query file
  -a ACTIVES, --actives ACTIVES
                        Actives set
  -d DECOYS, --decoys DECOYS
                        Decoys set
  -n NCONFS, --nconfs NCONFS
                        Maximum number of conformers [25]
  -t THREADS, --threads THREADS
                        Number of threads [1]
  -o OUTPUT_DIRECTORY, --output_directory OUTPUT_DIRECTORY
                        Output directory

(base) D:\CDK2\Screening>
```

4

```
C:\Windows\System32\cmd.exe
(base) D:\CDK2\Screening>python ligand_based_VS.py -q solution_02.sdf -a cdk2_clustered_3D_MM.sdf
-d DUD_cdk2_decoys_ID_pass_MWPass_I_MM.sdf -n 200 -t 3
```

```
D:\CDK2\Screening\example.py - Notepad++
File Edit Search View Encoding Language Settings Tools Macro Run Window ?
example.py
1 from cdc.screening import Screener
2 from cdc.io import MoleculeReader, MoleculeWriter
3 import os
4
5 query_file = 'solution_02.sdf'
6 screen_set_file = 'my_library.sdf'
7
8 query = [m for m in MoleculeReader(query_file)]
9 screen_set = [m for m in MoleculeReader(screen_set_file)]
10
11 settings = Screener.Settings()
12 settings.output_directory = os.path.join(os.getcwd(), "screen_data")
13
14 screener = Screener(query, settings=settings)
15
16 screening_molecules = screener.screen([m for m in screen_set])
17
18 screening_scores = sorted([(r.score, r.identifier) for r in results])
19
```

query and write out the resulting *screening_molecules* and *screening_scores* files for further investigation.

Data analysis

The `ligand_based_vs.py` script returns the results into the output directory called *screen_data*. It contains:

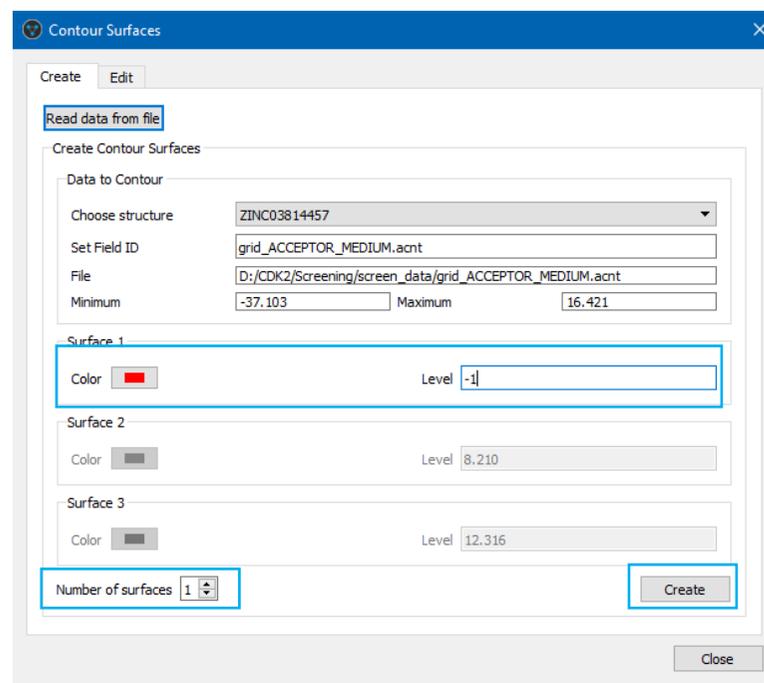
- The grid maps for the donor, acceptor and non-polar fields in .acnt format (e.g. *grid_ACCEPTOR_MEDIUM.acnt*).
- The corresponding fitting points in .mol2 format (e.g. *fitpts_ACCEPTOR_MEDIUM.mol2*).
- The screened libraries (i.e. *actives_screened.mol2* and *decoys_screened.mol2*).
- The ranked list of scores (i.e. *screening_scores.csv*).

Please note that, after launching the script, the grid maps and fitting points will be created almost immediately. It may take up to 15 minutes for the *actives_screened.mol2* file (with the top scoring conformer per screened active molecule) to be returned, and up to 2 hours for the *decoys_screened.mol2* (with the top scoring conformer per screened decoy molecule) and the *screening_scores.csv* to be returned, with times depending on your machine. Please do not close the command prompt when the screening has finished.

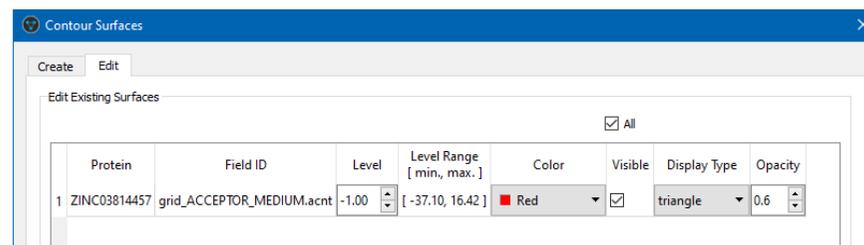
To visualise these fields, launch a new instance of Hermes and load *screen_data/actives_screened.mol2* using **File > Open** menu. Deactivate the tick-box next to *All Entries* in the *Molecule Explorer* window to clear the 3D view in Hermes.

1. Click **Display** from the main menu, then **Contour Surfaces...** In the *Create* tab click **Read data from file** and select *grid_ACCEPTOR_MEDIUM.acnt* from the output directory. Reduce the *Number of surfaces* to 1 and set the Level of *Surface 1* to -1. Click **Create**.

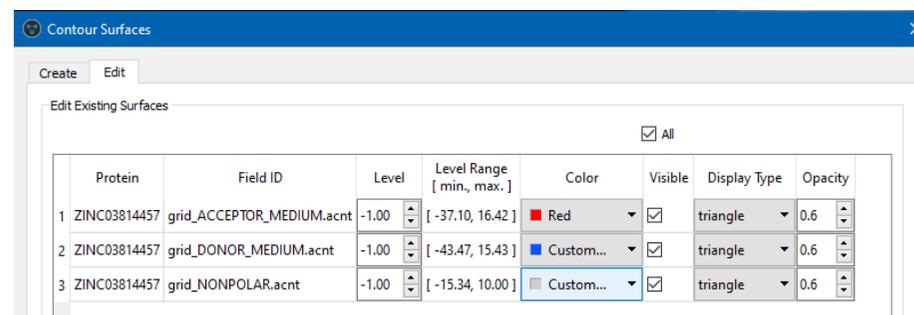
1



2



3



2. In the *Edit* tab change the *Display Type* to triangle and the *Opacity* to 0.6.
3. Repeat the same steps for *grid_DONOR_MEDIUM.acnt* and *grid_NONPOLAR.acnt*, but set the colour to blue and light grey, respectively.
4. The field potentials and the superimposition of the active ligands that have been ranked among the top 1% of the dataset are shown. For this dataset of screened 2117 molecules, this corresponds to the superimposition of the actives found among the best scoring 21 molecules.

We can calculate the arithmetic weighted ROC enrichment at 1% of false positive fraction as introduced by Jain, A.N. & Nicholls, A. (*J Comput Aided Mol Des* (2008) **22**: 133. doi:10.1007/s10822-008-9196-5). It represents the fraction of actives seen along with the top 1% of known decoys (multiplied by 100). Similarly, awROC enrichments at 0.5%, 2%, and 5% can be calculated, along with the arithmetic weighted version of the AUC (awAUC). We will use the *calc_enrichment_metrics.py* script to calculate these. We will then use the *generate_roc_plot.py* script to generate the ROC plot for this virtual screen study.

5. In the command prompt, please type the command to run the *calc_enrichment_metrics.py* script:

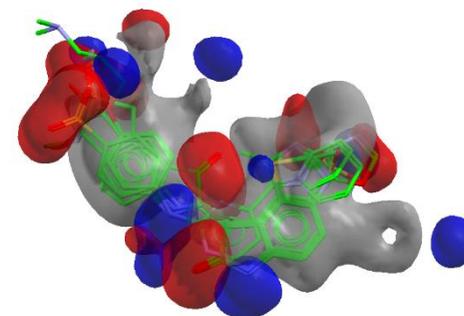
```
python calc_enrichment_metrics.py
```

Please note that you may need to use a text editor such as Notepad++ to edit this *calc_enrichment_metrics.py* script if the path to your *screening_scores.csv* file differs from *D:\CDK2\Screening\screen_data\screening_scores.csv* and if you wish to write the output to a different location than *D:\CDK2\Screening\screen_data\VS_enrichment_metrics.csv*.

6. In the command prompt, please type the command to run the *generate_roc_plot.py* script:

```
python generate_roc_plot.py  
screen_data\screening_scores.csv
```

4



5

```
C:\Windows\System32\cmd.exe
(base) D:\CDK2\Screening>python calc_enrichment_metrics.py
```

```
D:\CDK2\Screening\calc_enrichment_metrics.py - Notepad++
File Edit Search View Encoding Language Settings Tools Macro Run Window ?
calc_enrichment_metrics.py
111
112
113 def main():
114     data = get_data(filename=D:\CDK2\Screening\screen_data\screening_scores.csv)
115     total = len(data)
116
117     ef_05_percent = calculate_EF_x_percent(data, total, 0.005)
118     ef_1_percent = calculate_EF_x_percent(data, total, 0.01)
119     ef_2_percent = calculate_EF_x_percent(data, total, 0.02)
120     ef_5_percent = calculate_EF_x_percent(data, total, 0.05)
121
122     awAUC = calculate_awAUC(data)
123
124     awef_05_percent = calculate_awEF(data, total, 0.005)
125     awef_1_percent = calculate_awEF(data, total, 0.01)
126     awef_2_percent = calculate_awEF(data, total, 0.02)
127     awef_5_percent = calculate_awEF(data, total, 0.05)
128
129     with open(D:\CDK2\Screening\screen_data\VS_enrichment_metrics.csv', 'w') as writer:
130         writer.write('EF0.5%, EF1%, EF2%, EF5%, awROCe@0.5%, awROCe@1%, awROCe@2%, awROCe@5%, awAUC\n')
131         writer.write('%0.3f,%0.3f,%0.3f,%0.3f,%0.3f,%0.3f,%0.3f,%0.3f,%0.3f,%0.3f\n' %
132             (ef_05_percent, ef_1_percent, ef_2_percent, ef_5_percent,
133              awef_05_percent, awef_1_percent, awef_2_percent, awef_5_percent, awAUC))
134
135
136 if __name__ == "__main__":
137     main()
138
```

where

`screen_data\screening_scores.csv` is the path to the `screening_scores.csv` output file for this virtual screen study.

6

```
C:\Windows\System32\cmd.exe
(base) D:\CDK2\Screening>python generate_roc_plot.py screen_data\screening_scores.csv
```

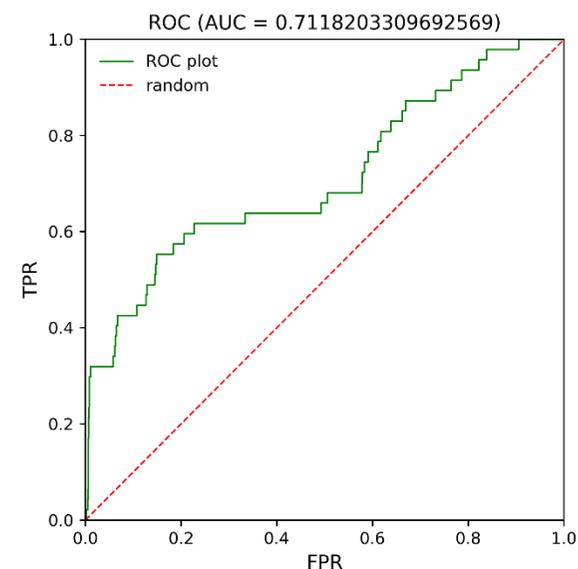
The ROC curve will be saved as `rocs.png`.

The arithmetic weighted enrichment metrics (which by default were output to `D:\CDK2\Screening\screen_data\VS_enrichment_metrics.csv`) are shown in Table 1, which also contains comparative data collected from the reference paper where the dataset was originally described. The only difference is the query used in all other seven approaches was defined by a single ligand (i.e. PDB entry: 1cpk). This is also possible with the CSD Ligand Screener, but the size and the shape of the molecule will impact the results. In fact, a default distance of 3 Å from the overlay (i.e. the shortest distance to any atom in the overlay) is used to define an excluded volume envelope after which a penalty value is applied if an atom is placed there.

Here we show that exploiting the knowledge from several ligands, if available, can lead to improved performance as they better resemble the shape of the binding pocket.

Table 1. Results of the CSD Ligand Screener compared with other methods. Data taken from DOI: 10.1186/1758-2946-1-14. The best value of each metric is highlighted in blue.

VS metrics	CSD Ligand Screener	DOCK	FieldScreen	MACCS	OAK	OAKFLEX	2SHA	OAAP
awAUC	0.721	0.53 ± 0.03	0.44	0.55 ± 0.02	0.55 ± 0.02	0.46 ± 0.02	0.48 ± 0.03	0.53 ± 0.03
awROC enrichment@0.5%	15.094	4.0 ± 6.1	7.5	9.4 ± 1.8	9.4 ± 1.8	9.4 ± 1.8	9.4 ± 3.7	24.1 ± 4.7
awROC enrichment@1%	26.183	9.5 ± 2.7	3.8	4.9 ± 1.0	6.5 ± 1.5	4.9 ± 1.5	11.1 ± 2.4	15.7 ± 2.8
awROC enrichment@2%	16.942	7.0 ± 1.4	1.9	2.5 ± 0.5	5.5 ± 1.2	4.8 ± 1.1	7.9 ± 1.4	8.6 ± 1.3
awROC enrichment@5%	6.777	2.8 ± 0.6	0.8	2.6 ± 0.6	2.6 ± 0.4	2.6 ± 0.4	3.5 ± 0.5	3.5 ± 0.7



Additional notes

- Tautomerism is crucial in protein-ligand interactions. Thus, virtual screening libraries should include multiple tautomeric forms of all molecules to lower the risk of losing an important hit. Different protonation states should also be enumerated. However, the final list of ranked scores should only contain the protomer or tautomer with the better score.
- For the sake of transparency and comparison, here we've used the publicly available dataset, but we strongly recommend that in real drug discovery projects all tautomers and protomers are enumerated before performing any structure- or ligand-based virtual screening study.

Conclusion

A retrospective virtual screening study of known active and decoy molecules has been carried out. Considering the high enrichment values, although CDK2 is still a challenging case for ligand-based methods (*J. Chem. Inf. Model.*, 2008, **48** (11), 2108–2117, DOI: 10.1021/ci800110p), we can be more confident in using the selected overlay hypothesis and the CSD Ligand Screener for a prospective screening of unknown binders.

You should now be familiar with how to setup a ligand-based virtual screening study using the CSD Python API, how to visualise the potential fields in Hermes, and how to analyse the ranked list of scores.

