

Scripting system

Zarek S. Siegel

March 15, 2016

1 Summary

1.1 Description

High platforms team leverage enthusiastically high infrastructures utilize convergence functionalized action capital. Than growth multimedia viral alternative emerging infrastructures e-enable sucking fashion capital niches standards. Extensible "organic" team re-engineer reinvent quality leading-edge paradigms cultivate infrastructures energistically dynamic. Art holistically covalent leverage initiatives enterprise interoperable empowerment actualize collaborative invested chains utilize expedite corporate. Improvements potentialities results energistically streamline completely positioning brand adaptive team visualize of holistic infrastructures.

1.2 File Structure

- NOTES:
 - ***bolded-italicized*** files/directories need to be created by the user
 - **bolded** files/directories/names need to be modified by the user
 - *italicized* files and directories are created by scripts
 - plain files and directories do not (and should not) be modified
 - \$ command indicates a command line command in the terminal
 - directories and files (and parts of file names) in lower case must be included exactly as indicated
 - directories and files in all caps should be named appropriately to the proteins and dockings in question

- **base_dir** – The root of the system is the base directory, containing all other files (probably best not to put anything else in this folder but what's indicated below)
 - **Readme.md** – This file
 - **Dockings.csv** – A spreadsheet containing master docking parameters
 - **Gridboxes.csv** – A spreadsheet specifying all the grid box parameters
 - **ligsets/** – A directory containing all the sets of ligands
 - **ligsets/LIGSET/** – for each set, there should be a directory within the **ligsets/** directory, whose name is the name of the ligand set. Substitute **LIGSET** for some appropriate name (e.g. **ligsets/my_awesome_ligset/**)
 - **ligsets/LIGSET/pdbqts/** – a directory containing a PDBQT file for each ligand in the set (whose name is **LIG.pdbqt**, with **LIG** being exactly the same as in **ligsets/LIGSET/LIGSET_list.txt**) [may be scripted later]
 - **ligsets/LIGSET/LIGSET_list.txt** – a text file containing a list of all the ligands (one on each line) and nothing else
 - you can create this easily on the command line using: `$ cd base_dir/ligsets/; for l in $(ls LIGSET/pdbqts | sed 's/.pdbqt//'); do echo $l >> LIGSET/LIGSET_list.txt; done` (appropriately substituting **base_dir** and **LIGSET**, and assuming the PDBQTs are already made)
 - [optional/preliminary] **ligsets/LIGSET/LIGSET.cdxml/** or **ligsets/LIGSET/LIGSET.mol/** – one file to show all the ligands in one page for presentations, PDFmaking and such.
 - [optional/preliminary] **/mols/** and **ligsets/LIGSET/pdbs/** – directories for preparing the initial PDBQT files. Will either be optional, or more scripts will be written.
 - **parameters.csvs/** – A directory containing small CSV files specifying the docking parameters for each individual docking (needed for scripts, at least as of now). Generated by **scripts/write_params_csv.R** from information in **Gridboxes.csv**
 - They will be named **parameters.csvs/DOCKING_parameters.csv**, where **DOCKING** is the docking ID
 - **vina_submit.shs/** – A directory containing the submission files for

Vina jobs on the (Wesleyan) cluster. Generated by `write_vina_submit` function of `docking_data_assembly.py`

- `vina_submit_shs/vina_submit DOCKING.sh` –for a docking of 20 models or less, a single submission script is written (and submitted using `$ bsub < vina_submit DOCKING.sh`, see submission instructions below)
- `vina_submits DOCKING/` –dockings of more than 20 models need to be submitted with multiple scripts (because Vina will not generate more than 20 poses). In this case, the `write_vina_submit` function will create a directory called `vina_submits DOCKING/` containing n scripts `vina_submit DOCKING.1.sh` , `vina_submit DOCKING.2.sh` through `vina_submit DOCKING.n.sh` . This script is set up to write each n submission scripts, where each of which script has a model number of 20 and n is the number of models divided by 20. Therefore, if greater than 20, the number of models should always be a multiple of twenty, or things will get messed up. (These are submitted used `$ for s in $(ls vina_submits DOCKING); do bsub < $s; done`, see instructions below)
- `PROTEIN/` , etc. – A directory for *each* protein, whose name is the name of the protein, the reference name/abbreviation used throughout, it just needs to be consistent (e.g. I primarily dock the proteins HepI and p300 and have the directories `hepi/` and `p300/` in my `base_dir/`)
 - `PROTEIN/PROTEIN.pdbqt` – the PDBQT file to be used for docking (`PROTEIN` must be *exactly* the same for the directory/file names and in the `Dockings.csv` and `Gridboxes.csv` entries for the protein’s dockings) [may be scripted later, from `PROTEIN.pdb`]
 - [optional/preliminary] `PROTEIN.pdb` – the original PDB file
 - `PROT/DOCKING/` – for every docking, there should be a directory whose name in the docking ID (the same as in `Dockings.csv`). Note: the user shouldn’t make this folder, it is made by `vina_submit DOCKING.sh`.
 - This will eventually contain
 - `PROT/binding_sites/` – a directory containing binding site PDBs:
 - To do binding site scoring by residues contacted (which is detected by the AutoDockTools script `process_VinaResult.py`),

There must be one or more `PROT/binding_sites/BINDING_SITE.pdb` files. These are subsets of the original `PROTEIN.pdb` (I originally created mine by grabbing the residues within 5Å of the bound ligands that came with the crystal structure, but it could be done many other ways).

- `scripts/` – all the scripts needed to use this set up
 - `write_params_csv.R` – writes `DOCK_parameters.csv` using information in `Dockings.csv` and `Gridboxes.csv`
 - `load_parameters.sh` – loads parameters from `DOCKING_parameters.csv` (used in the next script)
 - `separate_vina_results.sh` – runs Vina result PDBQTs through the AutoDockTools script `process_VinaResult.py`, which separates the poses into separate files and extracts the receptor contacts. The resulting files ended up in `DOCKING/processed_pdbqts/`, named `DOCKING.LIGAND_mMODEL.pdbqt`, where `MODEL` is the particular pose represented by the file. (A docking with l ligands and with Vina set to produce m models will therefore end up having $l \times m$ files after this script runs.)
 - `cleanup_processed_vina_results.sh` – cleans up processed PDBQTs (`DOCKING/processed_pdbqts/DOCKING.LIGAND_mMODEL.pdbqt`) and converts them to PDBs using the AutoDockTools script `pdbqt_to_pdb.py`
 - `parse_pdb.py` – defines a object class called `Pdb` for parsing PDB and PDBQT files for their 3D coordinates and in the case of processed Vina results, their binding energy, protein contacts, and other data generated by Vina. (Necessary for `docking_data_assembly.py`.)
 - `aiad_icpd.py` – defines functions to calculate the AIAD (averaged inter-atomic distance) and ICPD (inter-centerpoint distance) between two `Pdb` objects, two useful parameters for determining where a pose is binding on a protein and clustering poses together based on proximity. (Necessary for these functions in `docking_data_assembly.py`.)
 - `docking_data_assembly.py` – defines an object class called `Docking` for preparation and analysis of dockings. Contains several important functions:
 - `write_vina_submit` – prepares Vina job submission scripts
 - `assemble_dic` – assembles a data dictionary that contains all

- the mined data from the Vina results
- `score_binding_sites` – scores each pose for the proportion of residues contacted in each reference binding site
- `assess_all_resis` – uses binding scores to determine a True/-False for each pose binding in at each binding site (based on a threshold score, currently 0.1 or 10%)
- (`aiad_icpd_binding_sites`) – calculates AIAD and ICPD scores for each pose compared to each binding site
- `write_alldata.csv` – writes the data dictionary to a CSV file called `DOCKING_alldata.csv`
- `cluster_poses` [`prepare_clustering.csv.py`] – calculates AIAD scores for every pose compared to every other pose, for cluster analysis later on
- `pre.and.post.control.py` – links to all of the above scripts to coordinate their function, providing the global variables required to run this system on different computers.
- **Other scripts that may or may not be added:**
 - A script to add an entry to `Dockings.csv` or `Gridboxes.csv` with use input

1.3 Complete listing of included files

- **Notes:**

- ... indicates more of the same kind of file or directory
- .py files (except for the control script) may have a compiled .pyc file with them
- (files) are optional (for now)
- [files] are works in progress and may not be included ultimately
- [[files]] need to be made

```

1  /path/to/base_dir/
2  Readme.md
3  Dockings.csv
4  Gridboxes.csv
5  ligsets/
6  LIGSET1/
7  LIGSET1_list.txt
8  (LIGSET1.cdxml)
9  (mols/...)
10 (pdbs/...)

```

```

11         pdbqts/
12             LIG1_1.pdbqt
13             LIG1_2.pdbqt
14             LIG1_3.pdbqt
15         ...
16     LIGSET2/
17         LIGSET1_list.txt
18         (mols/...)
19         (pdbs/...)
20         pdbqts/
21             LIG2_1.pdbqt
22             LIG2_2.pdbqt
23             LIG2_3.pdbqt
24         ...
25     ...
26     vina_submit_shs/
27         vina_submit_A1.sh
28         vina_submits_A2/
29             vina_submit_A2.1.sh
30             vina_submit_A2.2.sh
31             vina_submit_A2.3.sh
32         ...
33         vina_submit_B1.sh
34     ...
35     PROT_A/
36         (PROT_A.pdb)
37         PROT_A.pdbqt
38         binding_sites/
39             BINDING_SITE_ALPHA1.pdb
40             BINDING_SITE_ALPHA2.pdb
41             BINDING_SITE_ALPHA3.pdb
42         ...
43     A1/
44     A2/
45     A3/
46     ...
47     PROT_B/
48         (PROT_B.pdb)
49         PROT_B.pdbqt
50         binding_sites/
51             BINDING_SITE_BETA1.pdb
52         ...
53     B1/
54     B2/
55     ...

```

```

56     ...
57     scripts/
58         [new_grid_or_dock_entry.R]
59         [[write_ligset_list_txt.sh]]
60         load_parameters.sh
61         [[ligand, protein preparation]]
62         separate_vina_results.sh
63         cleanup_processed_vina_results.sh
64         parse_pdb.py
65         aiad_icpd.py
66         [[prepare_clustering_csv.py]]
67         docking_data_assembly.py
68         pre_and_post_control.py
69         [[post_docking_graphs.R]]
70         [[clustering_graphs.R]]
71         [[R cript to select poses to view in PyMol]]
72         [[Py script to load PyMol sessions from lists]]

```

1.3.1 After docking post-processing:

```

1     PROTEIN/
2         (PROTEIN.pdb)
3         PROTEIN.pdbqt
4         binding_sites/
5             SITE1.pdb
6             SITE1.pdb
7             SITE1.pdb
8         ...
9     DOCKING/
10         result_pdbqts/
11             DOCKING_LIG1_results.pdbqt
12             DOCKING_LIG2_results.pdbqt
13             DOCKING_LIG3_results.pdbqt
14         ...
15         processed_pdbqts/
16             DOCKING_LIG1_m1.pdbqt
17             DOCKING_LIG1_m2.pdbqt
18             DOCKING_LIG1_m3.pdbqt
19         ...
20             DOCKING_LIG2_m1.pdbqt
21             DOCKING_LIG2_m2.pdbqt
22             DOCKING_LIG2_m3.pdbqt

```

```

23         ...
24         DOCKING_LIG3_m1.pdbqt
25         DOCKING_LIG3_m2.pdbqt
26         DOCKING_LIG3_m3.pdbqt
27         ...
28     processed_pdbs/
29         DOCKING_LIG1_m1.pdb
30         DOCKING_LIG1_m2.pdb
31         DOCKING_LIG1_m3.pdb
32         ...
33         DOCKING_LIG2_m1.pdb
34         DOCKING_LIG2_m2.pdb
35         DOCKING_LIG2_m3.pdb
36         ...
37         DOCKING_LIG3_m1.pdb
38         DOCKING_LIG3_m2.pdb
39         DOCKING_LIG3_m3.pdb
40         ...
41     DOCKING_alldata.csv
42     DOCKING.p
43     DOCKING_clustering.csv [[DOCKING_pose_pose_aiads.csv]]
44     [[DOCKING_best_aiad_pairs.csv]]
45     [[graphs/]]
46     [[...graphs...]]

```

1.4 Example Dockings.csv file

```

1 Docking ID,Date,Protein,Ligset,Grid box,Exhaustiveness,Number of Models,Number
  of CPUs,Notes
2 A1,20160301,PROTA,LIGS1,AAS,20,10,2,looking at active size of protein A
3 A2,20160308,PROTA,LIGS2,AWP,50,400,4,high volume docking of whole protein A
4 B1,20160308,PROTB,LIGS3,BWP,8,20,1,initial docking of whole protein B

```

1.5 Example Gridboxes.csv file

```

1 Gridbox Name,Protein,Size in x-dimension,Size in y-dimension,Size in z-
  dimension,Center in x-dimension,Center in y-dimension,Center in z-
  dimension,Notes

```


Docking ID	Date	Protein	Ligset	Grid box	Exhaustiveness	Number of Models	Number of CPUs	Notes
A1	20160301	PROTA	LIGS1	AAS	20	10	2	looking at active size of protein A
A2	20160308	PROTA	LIGS2	AWP	50	400	4	high volume docking of whole protein A
B1	20160308	PROTB	LIGS3	BWP	8	20	1	initial docking of whole protein B

² AAS,PROTA,60,72,88,41.89,2.69,-1.85,active site of protein A
³ AWP,PROTA,126,126,126,41.89,2.69,-1.85,all of protein A
⁴ BWP,PROTB,126,126,126,4.89,-5.27,12.0,all of protein B

Gridbox Name	Protein	Box size (x)	Box size (y)	Box size (z)	Box center (x)	Box center (y)	Box center (z)	Notes
AAS	PROTA	60	72	88	41.89	2.69	-1.85	active site of protein A
AWP	PROTA	126	126	126	41.89	2.69	-1.85	all of protein A
BWP	PROTB	126	126	126	4.89	-5.27	12.0	all of protein B

2 Scripts

2.1 constants.py

2.1.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabean indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 ### Module-wide constants
2 # (c) Zarek Siegel
3 # v1 3/11/16
4
5 base_dir="/Users/zarek/GitHub/TaylorLab/zvina"
6 cluster_base_dir="/home/zsiegel"
7 ADT_dir="/Library/MGLTools/latest/MGLToolsPckgs/AutoDockTools"
8 MGL_py_bin="/Library/MGLTools/latest/bin/pythonsh"
```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/constants.py

2.2 new_grid_or_dock_entry.py

2.2.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1  #!/usr/bin/env python
2
3  ### Write a new entry to Dockings.csv or Gridboxes.csv
4  # (c) Zarek Siegel
5  # v1 3/10/16
6  # v1.1 3/11/16
7
8  import csv, re, argparse, time, datetime
9  from constants import *
10
11 class Timestamp():
12     def __init__(self):
13         time_obj = time.time()
14         self.display = datetime.datetime.fromtimestamp(time_obj).strftime('%Y-%m-%d %H:%M:%S')
15         self.eightdigit = datetime.datetime.fromtimestamp(time_obj).strftime('%Y%m%d')
16
17 def new_docking_entry():
18     # Hello
19     print("\n\tWelcome! This script will add an entry to Dockings.csv\n")
20
21     # Dockings.csv columns:
22     # Docking ID
23     # Date
24     # Protein
25     # Protein File
26     # Ligset
27     # Gridbox
28     # Exhaustiveness
```

```

29 # Number of Models
30 # Number of CPUs
31 # Notes
32
33 # Going through each variable, taking user input
34 dock_input = raw_input("\t\tDocking identifier: ")
35 print("\t\t>>> dock set to: {}\n".format(dock_input))
36
37 current_time = Timestamp()
38 date_input = raw_input("\t\tDate (enter 'd' to default to {}): ".format(
    current_time.eightdigit))
39 if date_input == "d": date_input = current_time.eightdigit
40 print("\t\t>>> date set to: {}\n".format(date_input))
41
42 prot_input = raw_input("\t\tProtein: ")
43 print("\t\t>>> prot set to: {}\n".format(prot_input))
44
45 prot_file_input = raw_input("\t\tSpecific protein file (without the .pdbqt):
    ")
46 print("\t\t>>> prot_file set to: {}\n".format(prot_file_input))
47
48 ligset_input = raw_input("\t\tLigset identifier: ")
49 print("\t\t>>> ligset set to: {}\n".format(ligset_input))
50
51 box_input = raw_input("\t\tGridbox identifier: ")
52 print("\t\t>>> box set to: {}\n".format(box_input))
53
54 exhaust_input = raw_input("\t\tExhaustiveness: ")
55 print("\t\t>>> exhaust set to: {}\n".format(exhaust_input))
56
57 n_models_input = raw_input("\t\tNumber of models: ")
58 print("\t\t>>> n_models set to: {}\n".format(n_models_input))
59
60 n_cpus_input = raw_input("\t\tNumber of CPUs: ")
61 print("\t\t>>> n_cpus set to: {}\n".format(n_cpus_input))
62
63 notes_input = raw_input("\t\tAny notes: ")
64 if notes_input == "":
65     notes_input = "Entered by new_grid_or_dock_entry.py {}".format(
        current_time.display)
66 else:
67     notes_input = "{} (Entered by new_grid_or_dock_entry.py {})".format(
        notes_input, current_time.display)
68
69 print("\t\t>>> notes set to: {}\n".format(notes_input))
70

```

```

71 # New row as a dictionary
72 new_row = {
73     'Docking ID' : dock_input,
74     'Date' : date_input,
75     'Protein' : prot_input,
76     'Protein File' : prot_file_input,
77     'Ligset' : ligset_input,
78     'Gridbox' : box_input,
79     'Exhaustiveness' : exhaust_input,
80     'Number of Models' : n_models_input,
81     'Number of CPUs' : n_cpus_input,
82     'Notes': notes_input
83 }
84
85 # Print confirmation of all entered variables
86 print("\t>>> The new row will be\n\n\
87     Docking ID: {dock}\n\
88     Date: {date}\n\
89     Protein: {prot}\n\
90     Protein File: {prot_file}\n\
91     Ligset: {ligset}\n\
92     Gridbox: {box}\n\
93     Exhaustiveness: {exhaust}\n\
94     Number of Models: {n_models}\n\
95     Number of CPUs: {n_cpus}\n\
96     Notes: {notes}\n".format(
97         dock = dock_input,
98         date = date_input,
99         prot = prot_input,
100         prot_file = prot_file_input,
101         ligset = ligset_input,
102         box = box_input,
103         exhaust = exhaust_input,
104         n_models = n_models_input,
105         n_cpus = n_cpus_input,
106         notes = notes_input
107     )
108 )
109
110 # Don't write row without confirmation
111 proceed = raw_input("\tWrite this as a new docking entry? [y/n] ")
112 # print("\t{}".format(new_row))
113
114 # If "y" is entered, write the row, otherwise don't
115 if proceed == "y":

```

```

116     dockings_csv = "{b_d}/Dockings.csv".format(b_d=base_dir)
117     dockings_headers = ["Docking ID", "Date", "Protein", "Protein File",
118         "Ligset", "Gridbox", "Exhaustiveness", "Number of Models",
119         "Number of CPUs", "Notes"]
120     with open(dockings_csv, 'a') as f:
121         appender = csv.DictWriter(f, fieldnames=dockings_headers)
122         appender.writerow(new_row)
123         print("\n\t>>> New row appended to Dockings.csv:\n\topen {}".format(
            dockings_csv))
124     else:
125         print("\n\t>>> No docking entry written\n")
126
127 def new_gridbox_entry():
128     # Hello
129     print("\n\tWelcome! This script will add an entry to Gridboxes.csv\n")
130
131     # Gridboxes.csv columns:
132     # Gridbox Name
133     # Protein File
134     # Size in x-dimension
135     # Size in y-dimension
136     # Size in z-dimension
137     # Center in x-dimension
138     # Center in y-dimension
139     # Center in z-dimension
140     # Notes
141
142     # Going through each variable, taking user input
143     box_input = raw_input("\t\tGridbox Name: ")
144     print("\t\t>>> box set to: {}".format(box_input))
145
146     prot_file_input = raw_input("\t\tSpecific protein file (without the .pdbqt):
147     ")
148     print("\t\t>>> prot_file set to: {}".format(prot_file_input))
149
150     box_size_x_input = raw_input("\t\tSize in x-dimension: ")
151     print("\t\t>>> box_size_x set to: {}".format(box_size_x_input))
152
153     box_size_y_input = raw_input("\t\tSize in y-dimension: ")
154     print("\t\t>>> box_size_y set to: {}".format(box_size_y_input))
155
156     box_size_z_input = raw_input("\t\tSize in z-dimension: ")
157     print("\t\t>>> box_size_z set to: {}".format(box_size_z_input))
158
159     box_center_x_input = raw_input("\t\tCenter in x-dimension: ")

```

```

159 print("\t\t>>> box_center_x set to: {}\n".format(box_center_x_input))
160
161 box_center_y_input = raw_input("\t\tCenter in y-dimension: ")
162 print("\t\t>>> box_center_y set to: {}\n".format(box_center_y_input))
163
164 box_center_z_input = raw_input("\t\tCenter in z-dimension: ")
165 print("\t\t>>> box_center_z set to: {}\n".format(box_center_z_input))
166
167 current_time = Timestamp()
168 notes_input = raw_input("\t\tAny notes: ")
169 if notes_input == "":
170     notes_input = "Entered by new_grid_or_dock_entry.py {}".format(
171         current_time.display)
172 else:
173     notes_input = "{} (Entered by new_grid_or_dock_entry.py {})".format(
174         notes_input, current_time.display)
175 print("\t\t>>> notes set to: {}\n".format(notes_input))
176
177 # New row as a dictionary
178 new_row = {
179     'Gridbox Name' : box_input,
180     'Protein File' : prot_file_input,
181     'Size in x-dimension' : box_size_x_input,
182     'Size in y-dimension' : box_size_y_input,
183     'Size in z-dimension' : box_size_z_input,
184     'Center in x-dimension' : box_center_x_input,
185     'Center in y-dimension' : box_center_y_input,
186     'Center in z-dimension' : box_center_z_input,
187     'Notes' : notes_input
188 }
189
190 # Print confirmation of all entered variables
191 print("\t>>> The new row will be\n\n\
192 Gridbox Name: {box}\n\
193 Protein File: {prot_file}\n\
194 Size in x-dimension: {box_size_x}\n\
195 Size in y-dimension: {box_size_y}\n\
196 Size in z-dimension: {box_size_z}\n\
197 Center in x-dimension: {box_center_x}\n\
198 Center in y-dimension: {box_center_y}\n\
199 Center in z-dimension: {box_center_z}\n\
200 Notes: {notes}\n".format(
201     box = box_input,
202     prot_file = prot_file_input,
203     box_size_x = box_size_x_input,

```

```

203         box_size_y = box_size_y_input,
204         box_size_z = box_size_z_input,
205         box_center_x = box_center_x_input,
206         box_center_y = box_center_y_input,
207         box_center_z = box_center_z_input,
208         notes = notes_input
209     )
210 )
211
212 # Don't write row without confirmation
213 proceed = raw_input("\tWrite this as a new grid box entry? [y/n] ")
214
215 # If "y" is entered, write the row, otherwise don't
216 if proceed == "y":
217     gridboxes_csv = "{b_d}/Gridboxes.csv".format(b_d=base_dir)
218     gridboxes_headers = ["Gridbox Name", "Protein File", "Size in x-dimension"
219 ,
220     "Size in y-dimension", "Size in z-dimension", "Center in x-dimension",
221     "Center in y-dimension", "Center in z-dimension", "Notes"]
222     with open(gridboxes_csv, 'a') as f:
223         appender = csv.DictWriter(f, fieldnames=gridboxes_headers)
224         appender.writerow(new_row)
225         print("\n\t>>> New row appended to Gridboxes.csv:\n\topen {}\n".format(
226             gridboxes_csv))
227     else:
228         print("\n\t>>> No grid box entry written\n")
229
230 # Stuff below is commented because this script is being used as a module
231
232 # def main():
233 #     parser = argparse.ArgumentParser(
234 #         description='Write a new entry to Dockings.csv or Gridboxes.csv')
235 #     parser.add_argument('-b', '--base_dir', metavar='BASE_DIR', type=str,
236 #         nargs=1,
237 #         help='The base directory containing Docking.csv and Gridboxes.csv')
238 #     parser.add_argument('-d', '--new_docking', action='store_true', default=
239 #         False,
240 #         help='New set of docking parameters (written to Dockings.csv)')
241 #     parser.add_argument('-g', '--new_gridbox', action='store_true', default=
242 #         False,
243 #         help='New set of grid box parameters (written to Gridboxes.csv)')
244 #     args = vars(parser.parse_args())
245 #     global base_dir

```



```
243 # base_dir = str(args['base_dir'][0])
244 # new_docking = args['new_docking']
245 # new_gridbox = args['new_gridbox']
246 #
247 # if new_docking: new_docking_entry()
248 # elif new_gridbox: new_gridbox_entry()
249 #
250 # if __name__ == "__main__": main()
```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/new_grid_or_dock_entry.py

2.3 load_parameters.sh

2.3.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 #!/bin/bash
2
3 ### Print parameters (originally load_parameters.sh)
4 # (c) Zarek Siegel
5 # v1 3/4/16
6 # v1.1 3/5/16
7 # v2 3/5/16
8 # v3 3/11/16
9
10 # Docking ID as required argument
11 # dock=$1
12
13 # Set scripts directory to the directory containing this script
14 scripts_dir="$( cd "$( dirname "${BASH_SOURCE[0]}" )" && pwd )"
15 # Set base directory to the one containing scripts_dir
16 base_dir="$( cd $scripts_dir && cd .. )"
17 # Source # AutoDockTools Directory and MGLTools Python binary paths from
    constants.py
18 source $scripts_dir/constants.py
19 # CSV file with docking parameters
20 dockings_csv=$base_dir/Dockings.csv
21
22 # Define a function for looking
23 function look_up {
24     parameter=$1 # argument taken is the column header
25     cat $dockings_csv | # look in Dockings.csv
26     # AWK script to look up parameter for docking in the CSV
27     awk -v dock="$dock" -v parameter="$parameter" \
28         'BEGIN{
29         FS=","; # CSV
```

```

30     dock_row=""; # declare global variables
31     parameter_field="";
32 }
33 {
34     if ($1 == dock) {
35         dock_row=NR; # determine which row to look in
36     }
37 }
38 NR==1{
39     for (f=1; f<=NF; f++) {
40         {
41             if ($f == parameter) {
42                 parameter_field=f; # determine which row to look in
43             }
44         }
45     }
46 }
47 NR==dock_row{print $parameter_field} # output the intersection
48 ,
49 }
50
51 # Source all relevant parameters
52 dock=$( look_up "Docking ID" )
53 date=$( look_up "Date" )
54 prot=$( look_up "Protein" )
55 prot_file=$( look_up "Protein File" )
56 ligset=$( look_up "Ligset" )
57 box=$( look_up "Gridbox" )
58 exhaust=$( look_up "Exhaustiveness" )
59 n_models=$( look_up "Number of Models" )
60 n_cpus=$( look_up "Number of CPUs" )
61
62 # Print all parameters
63 export dock=$dock
64 export date=$date
65 export prot=$prot
66 export prot_file=$prot_file
67 export ligset=$ligset
68 export box=$box
69 export exhaust=$exhaust
70 export n_models=$n_models
71 export n_cpus=$n_cpus

```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/load_parameters.sh

2.4 separate_vina_results.sh

2.4.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 #!/bin/bash
2
3 ### process_VinaResult and a bit of organization
4 # (c) Zarek Siegel
5 # v1 3/5/16
6 # v1.2 3/6/16
7 # v2 3/6/16 (batch separation)
8 # v3 3/11/16
9
10 ### Required input
11 dock=$1
12 # Set scripts directory to the directory containing this script
13 scripts_dir="$( cd "$( dirname "${BASH_SOURCE[0]}" )" && pwd )"
14 # Set base directory to the one containing scripts_dir
15 base_dir="$( cd $scripts_dir && cd .. )"
16 # Source # AutoDockTools Directory and MGLTools Python binary paths from
    constants.py
17 source $scripts_dir/constants.py
18 # Location of process_VinaResult.py
19 pvr_py="$ADT_dir/Utilities24/process_VinaResult.py"
20
21 # Retrieve the parameters for this docking
22 source $base_dir/scripts/load_parameters.sh $dock
23
24 # Exit if already done
25 if [ -e $base_dir/$prot/$dock/processed_pdbqts/ ]; then
26     echo " ! Results already separated"
27     echo " ($prot/$dock/processed_pdbqts/ exists),"
28     echo " -> exiting this step"
29     exit 1
```

```

30 fi
31
32 # Retrieve ligset list
33 ligset_list_txt=$base_dir/ligsets/$ligset/$ligset\_list.txt
34 ligset_list=$(for l in $(cat $ligset_list_txt); do echo $l; done)
35
36 # Relevant directories
37 result_pdbqts_dir=$base_dir/$prot/$dock/result_pdbqts/
38 processed_pdbqts_dir=$base_dir/$prot/$dock/processed_pdbqts/
39
40 # Create a directory for processed files
41 mkdir $processed_pdbqts_dir
42
43 # The actual process_VinaResult step
44 receptor_pdbqt=$base_dir/$prot/$prot_file.pdbqt
45 batch_size=20
46 # No batches
47 n_models=$(echo $n_models | sed 's/[^0-9]//')
48 if [[ "n_models" -le "$batch_size" ]]; then
49     for lig in $ligset_list; do
50         result_pdbqt=$result_pdbqts_dir/$dock\_lig\_results.pdbqt
51         processed_pdbqt_stem=$processed_pdbqts_dir/$dock\_lig\_m
52         $MGL_py_bin $pvr_py -r $receptor_pdbqt \
53             -f $result_pdbqt \
54             -o $processed_pdbqt_stem
55         echo " processed ligand $lig"
56     done
57 # Batches
58 elif [[ "n_models" -gt "$batch_size" ]]; then
59     n_batches=$(bc <<< "$n_models / $batch_size")
60     for ((b=1;b<=$n_batches;b++)); do
61         echo " processing batch $b"
62         for lig in $ligset_list; do
63             result_pdbqt=$result_pdbqts_dir/$dock\_.$b\_lig\_results.pdbqt
64             processed_pdbqt_stem=$processed_pdbqts_dir/$dock\_.$b\_lig\_m
65             $MGL_py_bin $pvr_py -r $receptor_pdbqt \
66                 -f $result_pdbqt \
67                 -o $processed_pdbqt_stem
68             # Rename the processed pdbqts
69             for ((m=1;m<=$batch_size;m++)); do
70                 old_processed_pdbqt=$processed_pdbqts_dir/$dock\_.$b\_lig\_m$m.pdbqt
71                 new_m=$(bc <<< "(( $b - 1 ) * $batch_size ) + $m")
72                 new_processed_pdbqt=$processed_pdbqts_dir/$dock\_lig\_m$new_m.pdbqt
73                 mv $old_processed_pdbqt $new_processed_pdbqt
74             done

```

```
75     echo "      processed ligand $lig"
76     done
77 done
78 else
79     echo "! ! ! Error in batch processing (n_models is weird)"
80 fi
81
82 # *** check for results, prot.pdb, params
83 # *** check if already pvr'd
```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/separate_vina_results.sh

2.5 cleanup_processed_vina_results.sh

2.5.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 #!/bin/bash
2
3 ### Converting and cleaning up processed vina result pdbqts
4 # (c) Zarek Siegel
5 # v1 3/5/16
6 # v1.2 3/6/16
7
8 ### Required input
9 dock=$1
10 # Set scripts directory to the directory containing this script
11 scripts_dir="$( cd "$( dirname "${BASH_SOURCE[0]}" )" && pwd )"
12 # Set base directory to the one containing scripts_dir
13 base_dir="$( cd $scripts_dir && cd .. )"
14 # Source # AutoDockTools Directory and MGLTools Python binary paths from
15   constants.py
16 source $scripts_dir/constants.py
17 # Location of pdbqt_to_pdb
18 q2b_py="$ADT_dir/Utilities24/pdbqt_to_pdb.py"
19
20 # Retrieve the parameters for this docking
21 source $base_dir/scripts/load_parameters.sh $dock
22
23 # Relevant directories
24 processed_pdbqts_dir=$base_dir/$prot/$dock/processed_pdbqts
25 cleanedup_processed_pdbqts_dir=$base_dir/$prot/$dock/
26   cleanedup_processed_pdbqts
27 processed_pdb_dir=$base_dir/$prot/$dock/processed_pdb
28
29 # Check if already done
30 if [ -d $processed_pdb_dir ]; then
```

```

29     echo " ! Results already cleaned up (processed_pdb exists), exiting this
        step"
30     exit 1
31 fi
32
33 # Create a directory for cleaned up files and pdb converts
34 mkdir $cleanedup_processed_pdbqts_dir
35 mkdir $processed_pdb_dir
36
37 # Retrieve ligset list
38 ligset_list_txt=$base_dir/ligsets/$ligset/$ligset\_list.txt
39 ligset_list=$(for l in $(cat $ligset_list_txt); do echo $l; done)
40
41 # The clean-up step
42 for lig in $ligset_list; do
43     for ((m=1;m<=$n_models;m++)); do
44         processed_pdbqt=$processed_pdbqts_dir/$dock\_lig\_m$m.pdbqt
45         cleanedup_processed_pdbqt=$cleanedup_processed_pdbqts_dir/$dock\_lig\_m$m
            .pdbqt
46         processed_pdb=$processed_pdb_dir/$dock\_lig\_m$m.pdb
47
48         # The clean-up step
49         cat $processed_pdbqt | \
50             sed 's/^(HETATM.....)\...../1LIG L/g' \
51             > $cleanedup_processed_pdbqt
52
53         # The PDBQT > PDB Conversion step
54         $MGL_py_bin $q2b_py -f $cleanedup_processed_pdbqt \
55             -o $processed_pdb
56
57         echo "---> processed ligand $lig model $m"
58     done
59 done
60
61 # Overwrite pre-clean-up pvr'd pdbqts with cleaned up ones
62 rm -rf $processed_pdbqts_dir
63 mv $cleanedup_processed_pdbqts_dir $processed_pdbqts_dir

```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/cleanup_processed_vina_results.sh

2.6 parse_pdb.py

2.6.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1  #!/usr/bin/env python
2
3  ### Parsing data from processed pdbqt result files
4  # (c) Zarek Siegel
5  # v1 3/5/16
6
7  import re
8
9  ### A class for residues and residue atoms
10 # (input is a string of form 'RES123' or 'RES123_A1')
11
12 class Residue:
13     def __init__(self, str):
14         # String
15         self.str = str
16         # Atom & Residue String
17         if re.search(r'^[A-Z]+[0-9]+$', self.str):
18             self.atom = None
19             self.res_str = self.str
20         elif re.search(r'^[A-Z]+[0-9]+_.$', self.str):
21             self.atom = re.sub(r'^[A-Z]+[0-9]+_', '', self.str)
22             self.res_str = re.sub(r'_[A-Z0-9]+$', '', self.str)
23         else: self.atom = None
24         # Residue Index
25         self.resi = re.sub(r'^[A-Z]+|_?[^_]*$', '', self.str)
26         try:
27             self.resi = int(self.resi)
28         except ValueError:
29             self.resi = None
30         # Residue Name
```

```

31     self.resn = re.sub(r'[0-9]+_[^_]*$', '', self.str)
32     # Dictionary of Props
33     self.dic = {'str': self.str, 'res_str': self.res_str,
34                'resi': self.resi, 'resn': self.resn, 'atom': self.atom}
35     def __str__(self):
36         return self.str
37
38     ### A class for molecules, including ones with data from process_VinaResult.py
39     # (input is a .pdb or .pdbqt file address which may or may not be pvr'd)
40
41     class Pdb:
42     def get_pdb_coords(self):
43         _coords = []
44         for line in self.pdb_lines:
45             if re.search('HETATM|ATOM', line) or (re.search('ATOM', line)):
46                 _dic = {
47                     'atomi' : int(line[6:11]),
48                     'atomn' : line[12:16].replace(" ", ""),
49                     'resn' : line[17:20].replace(" ", ""),
50                     'resi' : line[22:26].replace(" ", ""),
51                     'x' : float(line[30:38].replace(" ", "")),
52                     'y' : float(line[38:46].replace(" ", "")),
53                     'z' : float(line[46:54].replace(" ", "")),
54                     'xyz' : (float(line[30:38].replace(" ", "")),
55                             float(line[38:46].replace(" ", "")),
56                             float(line[46:54].replace(" ", ""))),
57                     'atom_type' : line[76:78].replace(" ", ""),
58                     'charge' : line[78:80].replace(" ", "") # element
59                 }
60                 _coords.append(_dic)
61         self.coords = _coords
62
63     def get_pdbqt_coords(self):
64         _coords = []
65         for line in self.pdb_lines:
66             if re.search('HETATM|ATOM', line) or (re.search('ATOM', line)):
67                 _dic = {
68                     'atomi' : int(line[6:11]),
69                     'atomn' : line[12:16].replace(" ", ""),
70                     'resn' : line[17:20].replace(" ", ""),
71                     'resi' : line[22:26].replace(" ", ""),
72                     'x' : float(line[30:38].replace(" ", "")),
73                     'y' : float(line[38:46].replace(" ", "")),
74                     'z' : float(line[46:54].replace(" ", "")),
75                     'xyz' : (float(line[30:38].replace(" ", "")),

```

```

76         float(line[38:46].replace(" ", "")),
77         float(line[46:54].replace(" ", "")) ),
78         'charge' : line[70:76].replace(" ", ""), # partial charge
79         'atom_type' : line[77:79].replace(" ", "") # AD4 atom type
80     }
81     _coords.append(_dic)
82     self.coords = _coords
83
84     def mine_pvr_data(self): # mine data from pvr file
85         contacts = []
86         for line in self.pdb_lines:
87             # Binding Energy
88             if re.search('REMARK VINA RESULT: ', line):
89                 self.E = re.sub( r'^REMARK VINA RESULT:[ ]+[ ]+[^ ]+[ ]+[^ ]+[$] ',
90                                 r'' , line.replace('\n', '')) # [23:31].replace(" ", "")
91                 self.E = float(self.E)
92             # RMSD Lower Bound
93             if re.search('REMARK VINA RESULT: ', line):
94                 self.rmsd_lb = re.sub( r'^REMARK VINA RESULT:[ ]+[^ ]+[ ]+[ ]+[ ]+[^ ]+[$] ',
95                                         r'' , line.replace('\n', ''))
96                 self.rmsd_lb = float(self.rmsd_lb)
97             # RMSD Upper Bound
98             if re.search('REMARK VINA RESULT: ', line):
99                 self.rmsd_ub = re.sub( r'^REMARK VINA RESULT:[ ]+[^ ]+[ ]+[ ]+[ ]+[^ ]+[$] ',
100                                         r'' , line.replace('\n', ''))
101                 self.rmsd_ub = float(self.rmsd_ub)
102             # Ligand Efficiency (whatever that means...)
103             if re.search('USER AD> ligand efficiency', line):
104                 self.pvr_effic = re.sub( r'USER AD> ligand efficiency' ,
105                                         r'' , line.replace('\n', ''))
106                 self.pvr_effic = float(self.pvr_effic)
107             # Model Number
108             if re.search(r'USER AD> .+ of .+ MODELS', line):
109                 self.pvr_model = re.sub( r'USER AD>| of [0-9]+ MODELS' ,
110                                         r'' , line.replace('\n', ''))
111                 self.pvr_model = int(self.pvr_model)
112             # Torsional Degrees of Freedom
113             if re.search('REMARK .+ active torsions:', line):
114                 self.torsdof = re.sub( r'REMARK|active torsions:' ,
115                                         r'' , line.replace('\n', ''))
116                 self.torsdof = int(self.torsdof)
117             # Number of Contacts
118             if re.search('USER AD> macro_close_ats:', line):
119                 self.macro_close_ats = re.sub( r'USER AD> macro_close_ats:' ,

```

```

120         r'' , line.replace('\n', ''))
121         self.macro_close_ats = int(self.macro_close_ats)
122         # Contacts
123         if re.search(r'^USER AD> [^ ]+:[^ ]+:[^ ]+:[^ ,]+$', line):
124             contacts.append(line.replace('\n', '').replace('USER AD> ', ''))
125
126         # Contacts Processing
127         self.pvr_resis_objs = []
128         self.pvr_resis = []
129         self.pvr_resis_atoms = []
130         for c in contacts:
131             self.pvr_resis_objs.append(Residue(re.sub(r'^[:]+:[^:]+$', '',
132                 c).replace(':', '_')))
133
134         for r in self.pvr_resis_objs:
135             if r.atom != None:
136                 self.pvr_resis_atoms.append(r.str)
137                 self.pvr_resis.append(r.res_str)
138             else:
139                 self.pvr_resis_atoms.append(None)
140                 self.pvr_resis.append(r.res_str)
141
142         self.pvr_resis_objs = list(set(self.pvr_resis_objs)) # remove duplicates
143         self.pvr_resis = list(set(self.pvr_resis))
144         self.pvr_resis_atoms = list(set(self.pvr_resis_atoms))
145
146         self.pvr_data = {
147             'E' : self.E,
148             'rmsd_ub' : self.rmsd_ub,
149             'rmsd_lb' : self.rmsd_lb,
150             'pvr_resis' : self.pvr_resis,
151             'pvr_resis_atoms' : self.pvr_resis_atoms,
152             'pvr_resis_objs' : self.pvr_resis_objs,
153             'torsdof' : self.torsdof,
154             'macro_close_ats' : self.macro_close_ats,
155             'pvr_model' : self.pvr_model
156         }
157
158         def get_types(self):
159             # Detect if its been through ADT process_VinaResult.py
160             self.is_pvr = False # by default
161             for line in self.pdb_lines:
162                 if re.search('REMARK VINA RESULT: ', line):
163                     self.mine_pvr_data()
164                     self.is_pvr = True

```

```

165
166     # Determine file type PDB/PDBQT (they are slightly different)
167     if self.pdb_file_in[-5:] == 'pdbqt':
168         self.get_pdbqt_coords()
169         self.file_type = 'pdbqt'
170     elif self.pdb_file_in[-3:] == 'pdb':
171         self.get_pdb_coords()
172         self.file_type = 'pdb'
173     else:
174         print("!!! BAD FILETYPE !!!")
175
176 def __init__(self, pdb_file_in):
177     # Specify input file
178     self.pdb_file_in = pdb_file_in
179     # Try to read it, else error
180     try:
181         pdb_file_open = open(pdb_file_in)
182         with pdb_file_open as f:
183             self.pdb_lines = f.readlines()
184     except IOError:
185         print("! ! ! IOError")
186         pass
187     # Determine if PDB or PDBQT, and whether it has been through
188     process_VinaResult.py
189     self.get_types() # this also mines the actual data

```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/parse_pdb.py

2.7 aiad_icpd.py

2.7.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1  #!/usr/bin/env python
2
3  ### Calculating average inter-atomic distance
4  # (c) Zarek Siegel
5  # v1 3/6/16
6
7  from parse_pdb import *
8  from math import sqrt
9  from numpy import mean
10
11 class Molecule():
12     def list_coords(self):
13         self.coord_triples = []
14         for atom in self.pdb.coords:
15             self.coord_triples.append(atom['xyz'])
16
17     def get_centerpoint(self):
18         x_coords = []
19         y_coords = []
20         z_coords = []
21         for triple in self.coord_triples:
22             x_coords.append(triple[0])
23             y_coords.append(triple[1])
24             z_coords.append(triple[2])
25         self.centerpoint = (mean(x_coords), mean(y_coords), mean(z_coords))
26
27     def __init__(self, pdb):
28         self.pdb = pdb
29         self.list_coords()
30         self.get_centerpoint()
```

```

31
32 def threeD_distance(triple1, triple2):
33     x1 = triple1[0]
34     y1 = triple1[1]
35     z1 = triple1[2]
36     x2 = triple2[0]
37     y2 = triple2[1]
38     z2 = triple2[2]
39     distance = sqrt( ((x2 - x1)**2) + ((y2 - y1)**2) + ((z2 - z1)**2) )
40     return distance
41
42 def caclulate_aiad(pdb1, pdb2):
43     molc1 = Molecule(pdb1)
44     molc2 = Molecule(pdb2)
45     dist_list = []
46     for triple1 in molc1.coord_triples:
47         dists_from_t1 = []
48         for triple2 in molc2.coord_triples:
49             dist = threeD_distance(triple1, triple2)
50             dists_from_t1.append(dist)
51         dist_list.append(min(dists_from_t1))
52     return mean(dist_list)
53
54 def calculate_icpd(pdb1, pdb2):
55     molc1 = Molecule(pdb1)
56     molc2 = Molecule(pdb2)
57     return threeD_distance(molc1.centerpoint, molc2.centerpoint)
58
59 def main():
60     # m1_p = Pdb("/Users/zarek/lab/Docking/p300/p27/res_pdbqts_cleaned/p27_s3_m4
61     # .pdbqt")
62     # m2_p = Pdb("/Users/zarek/lab/Docking/p300/p27/res_pdbqts_cleaned/
63     # p27_s2_m142.pdbqt")
64     # m1 = Molecule(m1_p)
65     # m2 = Molecule(m2_p)
66     # caclulate_aiad(m1_p, m2_p)
67     pass
68
69 if __name__ == "__main__": main()

```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/aiad.icpd.py

2.8 docking_data_assembly.py

2.8.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 #!/usr/bin/env python
2
3 ### Putting together all parsed data from processed vina results
4 # (c) Zarek Siegel
5 # v1 3/5/16 (as assemble_alldata.py)
6 # v2 3/6/16
7
8 from __future__ import print_function
9 import csv, re, os, subprocess
10 import cPickle as pickle
11 from constants import *
12 from parse_pdb import *
13 from aiad_icpd import *
14
15 ### A class for docking data, for sourcing, reading, and analyzing
16 class Docking():
17     # Basic docking parameters are contained in CSV file in ../base/
18     # parameters_csvs/
19     def load_parameters(self):
20         # Basic docking parameters are stored in base_dir/Dockings.csv
21         dockings_csv = "{}Dockings.csv".format(base_dir)
22         with open(dockings_csv) as f:
23             reader = csv.DictReader(f)
24             for row in reader:
25                 if row['Docking ID'] == self.dock:
26                     self.dockings_csv_row = row
27             # Source the parameters as class attributes
28             self.prot = self.dockings_csv_row['Protein']
29             self.prot_file = self.dockings_csv_row['Protein File']
30             self.ligset = self.dockings_csv_row['Ligset']
```



```

30 self.box = self.dockings_csv_row['Gridbox']
31 self.exhaust = self.dockings_csv_row['Exhaustiveness']
32 self.n_models = int(self.dockings_csv_row['Number of Models'])
33 self.n_cpus = self.dockings_csv_row['Number of CPUs']
34
35 # Grid box parameters are stored in base_dir/Gridboxes.csv
36 gridboxes_csv = "{}/Gridboxes.csv".format(base_dir)
37 with open(gridboxes_csv) as f:
38     reader = csv.DictReader(f)
39     for row in reader:
40         if row['Gridbox Name'] == self.box and row['Protein File'] == self.
prot_file:
41             self.gridboxes_csv_row = row
42
43 # Source the grid box parameters as class attributes
44 self.box_center_x = self.gridboxes_csv_row['Center in x-dimension']
45 self.box_center_y = self.gridboxes_csv_row['Center in y-dimension']
46 self.box_center_z = self.gridboxes_csv_row['Center in z-dimension']
47 self.box_size_x = self.gridboxes_csv_row['Size in x-dimension']
48 self.box_size_y = self.gridboxes_csv_row['Size in y-dimension']
49 self.box_size_z = self.gridboxes_csv_row['Size in z-dimension']
50
51 # Some useful directories
52 self.prot_dir = "{b_d}/{p}".format(b_d=base_dir, p=self.prot)
53 self.dock_dir = "{b_d}/{p}/{d}".format(b_d=base_dir, p=self.prot, d=self.
dock)
54
55 print("---> Loaded docking parameters")
56
57 # Retrieve the list of ligands in the set
58 def get_ligset_list(self):
59     ligset_list_txt = "{b_d}/ligsets/{ls}/{ls}_list.txt".format(
60         b_d=base_dir, ls=self.ligset)
61     ligset_list_txt_open = open(ligset_list_txt, 'r')
62     with ligset_list_txt_open as f:
63         self.ligset_list = f.read().splitlines()
64     print("---> Retrieved ligset list")
65
66 # Write the script for submission of vina jobs on the cluster
67 def write_vina_submit(self):
68     template = "#BSUB -q hp12\n\
69 #BSUB -n {n_cpus}\n\
70 #BSUB -J vina_{subdock}\n\
71 \n\
72 # Text file with list of ligands (one on each line)\n\

```

```

73 ligset_list_txt={cluster_base_dir}/ligsets/{ligset}/{ligset}_list.txt\n\
74 \n\
75 # Create the docking and output directories\n\
76 mkdir {cluster_base_dir}/{prot}/{dock}/\n\
77 mkdir {cluster_base_dir}/{prot}/{dock}/result_pdbqts\n\
78 \n\
79 # Generate the list of ligands\n\
80 ligset_list=$(for l in $(cat $ligset_list_txt); do echo $l; done)\n\
81 \n\
82 # Vina command\n\
83 for lig in $ligset_list; do\n\
84     /share/apps/autodock/autodock_vina_1_1_2_linux_x86/bin/vina \\\n\
85     --receptor {cluster_base_dir}/{prot}/{prot_file}.pdbqt \\\n\
86     --ligand {cluster_base_dir}/ligsets/{ligset}/pdbqts/$lig.pdbqt \\\n\
87     --out {cluster_base_dir}/{prot}/{dock}/result_pdbqts/{subdock}_${lig}_results
88     .pdbqt \\\n\
89     --center_x {box_center_x} \\\n\
90     --center_y {box_center_y} \\\n\
91     --center_z {box_center_z} \\\n\
92     --size_x {box_size_x} \\\n\
93     --size_y {box_size_y} \\\n\
94     --size_z {box_size_z} \\\n\
95     --cpu {n_cpus} \\\n\
96     --num_modes {n_models} \\\n\
97     --exhaustiveness {exhaust}\n\
98 done".format(
99     cluster_base_dir = cluster_base_dir,
100     n_cpus = self.n_cpus,
101     prot = self.prot,
102     prot_file = self.prot_file,
103     ligset = self.ligset,
104     box_center_x = self.box_center_x,
105     box_center_y = self.box_center_y,
106     box_center_z = self.box_center_z,
107     box_size_x = self.box_size_x,
108     box_size_y = self.box_size_y,
109     box_size_z = self.box_size_z,
110     exhaust = self.exhaust,
111     n_models = '{n_models}',
112     dock = '{dock}',
113     subdock = '{subdock}'
114 )
115 # (no need for batch submission)
116 if self.n_models <= 20:
117     template_filled = template.format(

```

```

117         n_models = self.n_models, dock = self.dock, subdock = self.dock)
118         vina_submit_sh = "{b_d}/vina_submit_shs/vina_submit_{d}.sh".format(
119             b_d=base_dir, d=dock)
120         with open(vina_submit_sh, 'w') as f:
121             f.write(template_filled)
122         print("---> Vina submission script for docking h11 has been created. It
can be found at:")
123         print("\t{}".format(vina_submit_sh))
124         # (batch submission)
125         elif self.n_models > 20:
126             vina_submits_dir = "{b_d}/vina_submit_shs/vina_submits_{d}/".format(
127                 b_d=base_dir, d=dock)
128             subprocess.call(['mkdir', vina_submits_dir])
129             n_batches = self.n_models / 20
130             for b in range(1, n_batches + 1):
131                 subdock = "{d}.{b}".format(d = dock, b = b)
132                 template_filled = template.format(
133                     n_models = 20, dock = self.dock, subdock = subdock)
134                 vina_submit_sh = "{v_s_d}/vina_submit_{sd}.sh".format(
135                     v_s_d = vina_submits_dir, sd = subdock)
136                 with open(vina_submit_sh, 'w') as f:
137                     f.write(template_filled)
138                 print("---> Vina submission scripts for docking h11 have been created.
They can be found in:")
139                 print("\t{}".format(vina_submits_dir))
140             else: print("!! ! bad n_models")
141
142         ### AFTER DOCKING
143
144         # Mine Vina result for data (actual mining is in parse_pdb.py, acting via
the Pdb class)
145         def assemble_dic(self):
146             self.data_dic = {}
147             self.keys = []
148             for lig in self.ligset_list:
149                 for m in range(1, int(self.parameters['n_models'])+1):
150                     processed_pdbqt = "{d_d}/processed_pdbqts/{d}_{lig}_m{m}.pdbqt".format
(
151                         d_d=self.dock_dir, d=dock, lig=lig, m=m)
152                     pose = Pdb(processed_pdbqt)
153                     key = "{}_{}_m{}".format(dock, lig, m)
154                     try:
155                         pose_dic = {
156                             'key' : key,
157                             'E' : pose.E,

```

```

158         'rmsd_ub' : pose.rmsd_ub,
159         'rmsd_lb' : pose.rmsd_lb,
160         'pvr_resis' : pose.pvr_resis,
161         'pvr_resis_atoms' : pose.pvr_resis_atoms,
162         'pvr_resis_objs' : pose.pvr_resis_objs,
163         'torsdof' : pose.torsdof,
164         'pvr_n_contacts' : pose.macro_close_at,
165         'pvr_model' : pose.pvr_model,
166         'pvr_effic' : pose.pvr_effic,
167         'coords' : pose.coords,
168         'lig' : lig,
169         'model' : m,
170         'pvr_obj' : pose,
171         'pdb_address' : re.sub('pdbqt', 'pdb', processed_pdbqt)
172     }
173     except AttributeError:
174         print("!! ! pose {} failed, check for the processed PDBQT".format(
key))
175         self.data_dic[key] = pose_dic
176         self.keys.append(key)
177     self.is_assembled = True
178     print("----> Data dictionary created with data from process_VinaResult.py")
179
180     # Create a list of the binding sites previously prepared for the protein
181     def get_binding_sites_list(self):
182         if not self.is_assembled: self.assemble_dic()
183
184         binding_sites_dir = "{b_d}/binding_sites/{p_f}".format(
185             b_d=base_dir, p_f=self.prot_file)
186         self.binding_sites_list = []
187         self.binding_sites_objs = {}
188         for root, dirs, files in os.walk(binding_sites_dir):
189             for file in files:
190                 self.binding_sites_list.append(re.sub('.pdb', '', file))
191                 self.binding_sites_objs[re.sub('.pdb', '', file)] = (Pdb("{}{}").format
(root, file)))
192
193         self.bs_resis_lists = {}
194         self.bs_resis_atoms_lists = {}
195
196         for bs, bso in self.binding_sites_objs.items():
197             bs_resis = []
198             bs_resis_atoms = []
199             for atom in bso.coords:
200                 bs_resis.append("{}{}").format(atom['resn'], atom['resi'])

```

```

201         bs_resis_atoms.append("{}{}_{}".format(atom['resn'], atom['resi'],
202         atom['atomn']))
203         self.bs_resis_lists[bs] = list(set(bs_resis))
204         self.bs_resis_atoms_lists[bs] = list(set(bs_resis_atoms))
205
206     self.is_bs_listed = True
207     print("----> Retrieved binding sites")
208
209     # Score the binding sites in terms of the residues they contact
210     # relative to those contained in the reference PDB
211     def score_binding_sites(self):
212         if not self.is_assembled: self.assemble_dic()
213         if not self.is_bs_listed: self.get_binding_sites_list()
214
215         for pose in self.data_dic:
216             for bs in self.binding_sites_list:
217                 resis_union = ( set(self.bs_resis_lists[bs]) & set(self.data_dic[pose]
218                 ['pvr_resis']) )
219                 self.data_dic[pose]["{}_fraction".format(bs)] = float(len(resis_union)
220                 ) / float(len(self.bs_resis_lists[bs]))
221                 resis_atoms_union = ( set(self.bs_resis_atoms_lists[bs]) & set(self.
222                 data_dic[pose]['pvr_resis_atoms']) )
223                 self.data_dic[pose]["{}_atoms_fraction".format(bs)] = float(len(
224                 resis_atoms_union)) / float(len(self.bs_resis_atoms_lists[bs]))
225
226     self.is_bs_scored = True
227     print("----> Scored binding sites")
228
229     # Score binding sites by average inter-atomic distance
230     # and inter-centerpoint difference
231     # Acting via aiad_icpd.py
232     def aiad_icpd_binding_sites(self):
233         if not self.is_assembled: self.assemble_dic()
234
235         for pose in self.data_dic:
236             for bs, bso in self.binding_sites_objs.items():
237                 aiad = caclulate_aiad(self.data_dic[pose]['pvr_obj'], bso)
238                 self.data_dic[pose]["{}_aiad".format(bs)] = aiad
239                 icpd = calculate_icpd(self.data_dic[pose]['pvr_obj'], bso)
240                 self.data_dic[pose]["{}_icpd".format(bs)] = icpd
241
242     self.are_aiad_icpd_calcd = True
243     print("----> Calculated AIAD and ICPD")
244
245     # Add attributes for whether the ligand contacts each residue of the protein

```

```

241 def assess_all_resis(self):
242     if not self.is_assembled: self.assemble_dic()
243
244     self.prot_pdbqt = "{p_d}/{p_f}.pdbqt".format(p_d=self.prot_dir, p_f=self.
prot_file)
245     self.prot_obj = Pdb(self.prot_pdbqt)
246     self.prot_resis_list = []
247     # self.prot_resis_atoms_list = []
248     for atom in self.prot_obj.coords:
249         self.prot_resis_list.append("{}{}".format(atom['resn'], atom['resi']))
250     # self.prot_resis_atoms_list.append("{}{}_{}".format(atom['resn'], atom
['resi'], atom['atomn']))
251     self.prot_resis_list = list(set(self.prot_resis_list))
252     # self.prot_resis_atoms_list = list(set(self.prot_resis_atoms_list))
253
254     for pose in self.data_dic:
255         for res in self.prot_resis_list:
256             if res in self.data_dic[pose]['pvr_resis']:
257                 self.data_dic[pose][res] = 1
258             else:
259                 self.data_dic[pose][res] = 0
260     # for atom in self.prot_resis_atoms_list:
261     #     if atom in self.data_dic[pose]['pvr_resis_atoms']:
262     #         self.data_dic[pose][atom] = 1
263     #     else:
264     #         self.data_dic[pose][atom] = 0
265
266     self.are_all_resis_assessed = True
267     print("---> Residues contacts added to data dictionary")
268
269     # Output all the data mined and analyzed into a CSV file
270     def write_alldata_csv(self):
271         self.alldata_csv = "{d_d}/{d}_alldata.csv".format(d_d=self.dock_dir, d=
dock)
272
273         # If the CSV already exists, read it in as a dictionary
274         if os.path.isfile(self.alldata_csv):
275             with open(self.alldata_csv) as f:
276                 reader = csv.DictReader(f)
277                 self.data_dic = {}
278                 for row in reader:
279                     self.data_dic[row['key']] = row
280         # If it doesn't, write it
281         else:
282             if not self.is_assembled: self.assemble_dic()

```

```

283     if not self.is_bs_listed: self.get_binding_sites_list()
284     if not self.is_bs_scored: self.score_binding_sites()
285     if not self.are_aiad_icpd_calcd: self.aiad_icpd_binding_sites()
286     if not self.are_all_resis_assessed: self.assess_all_resis()
287
288     fieldnames = ['key', 'lig', 'model', 'E', 'rmsd_lb', 'rmsd_ub',
289                  'pvr_effic', 'pvr_n_contacts', 'torsdof', 'pdb_address']
290     for bs in self.binding_sites_list:
291         fieldnames.append("{}_fraction".format(bs))
292         fieldnames.append("{}_atoms_fraction".format(bs))
293         fieldnames.append("{}_aiad".format(bs))
294         fieldnames.append("{}_icpd".format(bs))
295     for res in self.prot_resis_list: fieldnames.append(res)
296     # for atom in self.prot_resis_atoms_list: fieldnames.append(atom)
297
298     with open(self.alldata_csv, 'w') as csvfile:
299         writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
300         writer.writeheader()
301         for pose in self.data_dic:
302             row = {}
303             for f in fieldnames:
304                 row[f] = self.data_dic[pose][f]
305             writer.writerow(row)
306
307     self.is_csv_written = True
308     print("---> Completed alldata.csv is located at:\n\t{}".format(self.
309                               alldata_csv))
310
311     # Write another CSV that shows the AIAD values between all ligands
312     def cluster_poses(self):
313         if not self.is_assembled: self.assemble_dic()
314
315         self.clustering_csv = "{}_{}_clustering.csv".format(d_d=self.dock_dir,
316                                                            d=dock)
317         self.clustering_dic = {}
318
319         # If the CSV already exists, read it in as a dictionary
320         if os.path.isfile(self.clustering_csv):
321             with open(self.clustering_csv) as f:
322                 reader = csv.DictReader(f)
323                 for row in reader:
324                     key = row['compared']
325                     del row['compared']
326                     self.clustering_dic[key] = row
327         # If it doesn't, write it

```

```

326 else:
327     c = 0
328     print("---> Calculating AIAD between poses (for clustering)")
329     for key1 in self.data_dic:
330         self.clustering_dic[key1] = {}
331         c += 1
332         print("\t- calculated for {:25}{: <9} of {: >9}".format(key1, c, len(self.
keys)))
333         for key2 in self.data_dic:
334             aiad12 = caclulate_aiad(self.data_dic[key1]['pvr_obj'], self.
data_dic[key2]['pvr_obj'])
335             self.clustering_dic[key1][key2] = aiad12
336
337         fieldnames = ['compared'] + self.keys
338         with open(self.clustering_csv, 'w') as csvfile:
339             writer = csv.DictWriter(csvfile, fieldnames=fieldnames)
340             writer.writeheader()
341             for key in self.keys:
342                 row = self.clustering_dic[key]
343                 row['compared'] = key
344                 writer.writerow(row)
345
346         self.are_poses_clustered = True
347         print("    > Completed clustering.csv is located at:\n\t{}".format(self.
clustering_csv))
348
349 # Save the data dictionary as a pickled file (i.e. in native python format)
350 def save_pickled_docking_obj(self):
351     self.pickled_docking_obj = "{d_d}/{d}.p".format(d_d=self.dock_dir, d=dock)
352     pickle.dump(self, open(self.pickled_docking_obj, 'wb'))
353
354     self.is_pickled = True
355     print("---> Pickled docking object located at:\n\t{}".format(self.
pickled_docking_obj))
356
357 def __init__(self, d):
358     global dock
359     dock = d
360
361     self.dock = dock
362     self.load_parameters()
363     self.get_ligset_list()
364
365     self.is_assembled = False
366     self.is_bs_listed = False

```



```
367     self.is_bs_scored = False
368     self.are_aiad_icpd_calcd = False
369     self.are_all_resis_assessed = False
370     self.is_csv_written = False
371     self.are_poses_clustered = False
372     self.is_pickled = False
373
374
375
376
377
378     # W000T!
```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/docking_data_assembly.py

2.9 pre_and_post_control.py

2.9.1 Function

Transcompiler picojava pseudolanguage java ee regular expression syntactic sugar eclipse rust bug. Back-face culling lua \$1 javabeen indirection operator void msdn program unary operator intellij idea. Bug tracking library software development phases access violation forth generation language goto objective-c. Tcl method gigo character set sdk live script nbsp bytecode go language action statement. Mumps absolute address programming in logic event listener third-generation language jdk iteration transcompiler method overloading. Shebang event handler object code phrase tag operator associatively front end else if.

```
1 #!/usr/bin/env python
2
3 ### One script to control them AAAAAALLLLL!!!!
4 # (c) Zarek Siegel
5 # v1 3/6/16
6
7 import argparse, subprocess
8 import new_grid_or_dock_entry
9 from constants import *
10
11 def main():
12     print("")
13     print("->-> hi")
14
15     parser = argparse.ArgumentParser(description='Pre- and post-Vina file fun
16         times')
17     parser.add_argument('-d', '--dock', metavar='DOCK', type=str, nargs='?',
18         help='the ID for this docking')
19
20     parser.add_argument('-nd', '--new_docking', action='store_true', default=
21         False,
22         help='execute the R script to write the parameters.csv')
23     parser.add_argument('-ng', '--new_gridbox', action='store_true', default=
24         False,
25         help='execute the R script to write the parameters.csv')
26     parser.add_argument('-w', '--write_params', action='store_true', default=
27         False,
28         help='execute the R script to write the parameters.csv')
29     parser.add_argument('-v', '--vina', action='store_true', default=False,
30         help='write the Vina job submission script')
```

```

27
28 parser.add_argument('-s', '--separate', action='store_true', default=False,
29     help='execute the bash script to separate row Vina results')
30 parser.add_argument('-n', '--clean', action='store_true', default=False,
31     help='execute the bash script to clean up processed Vina results')
32
33 parser.add_argument('-c', '--csv', action='store_true', default=False,
34     help='generate and save the alldata CSV file')
35 parser.add_argument('-l', '--cluster', action='store_true', default=False,
36     help='create cluster CSV files (all lig x lig AIADs)')
37 parser.add_argument('-i', '--pickle', action='store_true', default=False,
38     help='save the docking as a pickled object')
39
40
41 parser.add_argument('-o', '--post_proc', action='store_true', default=False,
42     help='Perform all post-processing steps (separate, clean, csv, cluster,
43     pickle)')
44
45 args = vars(parser.parse_args())
46 #
47 new_docking = args['new_docking']
48 new_gridbox = args['new_gridbox']
49 write_params = args['write_params']
50 vina = args['vina']
51 #
52 separate = args['separate']
53 clean = args['clean']
54 #
55 csv = args['csv']
56 cluster = args['cluster']
57 pickle = args['pickle']
58
59 post_proc = args['post_proc']
60
61 # if write_params:
62 #     print("---> Parameters CSV for docking h11 has been created. It can be
63 #     found at:")
64 #     subprocess.call(["./write_params_csv.R", dock, base_dir])
65
66 if new_docking:
67     print("---> Write new set of docking parameters to Dockings.csv...")
68     new_grid_or_dock_entry.new_docking_entry()
69 if new_gridbox:
70     print("---> Write new set of grid box parameters to Gridboxes.csv...")
71     new_grid_or_dock_entry.new_gridbox_entry()

```

```

70
71 if separate or clean or vina or csv or cluster or pickle or post_proc:
72     dock = str(args['dock'])
73
74 if separate:
75     print("---> Processing raw Vina output PDBQTs")
76     subprocess.call(["{b_d}/scripts/separate_vina_results.sh".format(b_d=
77         base_dir), dock])
78 if clean:
79     print("---> Cleaning up processed PDBQTs and converting to PDBs")
80     subprocess.call(["{b_d}/scripts/cleanup_processed_vina_results.sh".format(
81         b_d=base_dir), dock])
82
83 if vina or csv or cluster or pickle or post_proc:
84     from docking_data_assembly import Docking
85     d = Docking(dock)
86
87 if vina:
88     print("---> Writing Vina submission script")
89     d.write_vina_submit()
90
91 if csv: d.write_alldata_csv()
92 if cluster: d.cluster_poses()
93 if pickle: d.save_pickled_docking_obj()
94
95 if post_proc:
96     print("---> Processing raw Vina output PDBQTs")
97     subprocess.call(["{b_d}/scripts/separate_vina_results.sh".format(b_d=
98         base_dir),
99         dock, base_dir, ADT_dir, MGL_py_bin])
100     print("---> Cleaning up processed PDBQTs and converting to PDBs")
101     subprocess.call(["{b_d}/scripts/cleanup_processed_vina_results.sh".format(
102         b_d=base_dir),
103         dock, base_dir, ADT_dir, MGL_py_bin])
104     d.write_alldata_csv()
105     d.cluster_poses()
106     d.save_pickled_docking_obj()
107
108 print("->> All done!!!!!!!!!!!!!!!!!!!!")
109 print("")
110
111 if __name__ == "__main__": main()

```

/Users/zarek/GitHub/TaylorLab/zvina/scripts/pre_and_post_control.py