

**Supporting Information for:**

**Data Mining the Cambridge Structural Database for Hydrate-Anhydrate  
Pairs with SMILES Strings**

**Jen E. Werner and Jennifer A. Swift**

**Georgetown University, Department of Chemistry, Washington, DC 20057-**

**1227**

**Table of Contents for Scripts**

<b>I. Identifying Hydrates Page Number</b>	
1. Find hydrates with water SMILES strings	1
2. Find solvates and hydrates with no water SMILES strings	3
3. Find forms without smiles strings	6
4. Find forms with metals	8
5. Check for SMILES of solvent molecule in solvates	10
6. Find potential duplicate none hydrate structures with formula match	18
7. Find potential duplicate none hydrate structures with same identifier prefix	24
8. Find potential duplicate SMILES hydrate structures with formula and SMILES string match	26
9. Find duplicate SMILES hydrate structures with packing similarity	37
10. Find potential duplicate none waterless forms with formula match	40
11. Find potential duplicate none waterless forms with same identifier prefix	44
12. Find potential duplicate SMILES waterless forms with formula match	46
13. Find potential duplicate SMILES waterless forms with SMILES string match	49
14. Find duplicate SMILES waterless forms with packing similarity	60
15. Find duplicates that fail packing similarity check	63
16. Remove overlapping duplicate structures	68
17. Find stoichiometrically distinct duplicate structures	72
18. Find distinct identifier prefix duplicate structures	76
19. Check chirality of duplicates with distinct identifier prefixes	78
20. Remove duplicates with distinct chirality from list of duplicates	108
<b>II. Identifying Hydrate-Anhydrate Pairs</b>	
21. Find hydrates with waterless forms using SMILES string method	111
22. Find hydrates with waterless forms using None method	126
23. Find hydrate-anhydrate pairs with distinct chirality	130
24. Create the three main classes	159
25. Find stoichiometrically distinct hydrate-anhydrate pairs	163
26. Determine numbers for flowchart	167

27. List of hydrate-anhydrate pairs	175
-------------------------------------	-----

III. Analyzing Paired and Unpaired Hydrate and Anhydrous Forms	
28. Determine water molecule stoichiometry	195
29. Find pairs with 1 and 2+ components and 2000 unpaired subsets	199
30. Determine crystal system symmetry for hydrate-anhydrate pairs	203
31. Determine crystal system symmetry for different classes	206
32. Find hydrate-anhydrate pairs determined at the same temperature	209
33. Determine the packing fraction of hydrate-anhydrate pairs at the same temperature	212

```

1 #1 Find hydrates with water SMILES strings
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7
8 #Defines restrictions for structures that are pulled from the CSD
9 #There are no limitations on the number of structures to pull, the R-factor of a
10 #structure, or the disorder of a structure
11 #Structures with reported errors, organometallic or polymeric entities, or 2D
12 #coordinates are not permitted
13 class Runner(argparse.ArgumentParser):
14     def __init__(self):
15         super(self.__class__, self).__init__(description=__doc__)
16         self.add_argument(
17             '-i', '--input', default='CSD',
18             help='input database [CSD]')
19         self.add_argument(
20             '-o', '--output', default='entries_with_water_smiles_string.gcd',
21             help='output file [entries_with_water_smiles_string.gcd]')
22         self.add_argument(
23             '-m', '--maximum', default=0, type=int,
24             help='Maximum number of structures to find [all]')
25         self.add_argument(
26             '-R', '--r_factor', default=100.0, type=float,
27             help='Maximum acceptable R-factor [100.0]')
28         self.add_argument(
29             '-E', '--errors', default=True, action='store_false',
30             help='Whether structures with errors are acceptable [No]')
31         self.add_argument(
32             '-M', '--organometallic', default=True, action='store_false',
33             help='Whether organometallic structures are acceptable [No]')
34         self.add_argument(
35             '-P', '--polymeric', default=True, action='store_false',
36             help='Whether polymeric structures are acceptable [No]')
37         self.add_argument(
38             '-T', '--two_d', default=True, action='store_false',
39             help='Whether 2d structures are acceptable [No]')
40
41     args = self.parse_args()
42
43
44     self.args = args
45     self.settings = search.Search.Settings()
46     self.settings.max_hit_structures = self.args.maximum
47     self.settings.max_r_factor = self.args.r_factor
48     self.settings.no_errors = self.args.errors
49     self.settings.only_organic = self.args.organometallic
50     self.settings.not_polymeric = self.args.polymeric
51     self.settings.has_3d_coordinates = self.args.two_d
52     #Any structures containing the following elements were not pulled from the CSD
53     self.settings.must_not_have_elements = ['Be', 'Mg', 'Ca', 'Sr', 'Ba', 'Ra',
54     'B', 'Si', 'As', 'Se', 'Te', 'At',
55                                         'He', 'Ne', 'Ar', 'Kr', 'Xe', 'Rn']
56
57     #The entry reader object is used to generate the SMILES string for each structure
58     def run(self):
59         if self.args.input == 'CSD':

```

```

65     reader = io.EntryReader()
66 else:
67     database = self.args.input
68     reader = io.EntryReader(database, format='identifiers')
69
70 count = 0
71 total = 0
72
73 with io.EntryWriter(self.args.output) as hydrates_writer:
74     with io.EntryWriter("entries_with_no_water_smiles_string.gcd") as writer:
75         for i, entry in enumerate(reader):
76             #The number of structures with water SMILES strings found so far is
77             #printed for every 10,000 structures analyzed
78             if i and i % 10000 == 0:
79                 print 'Analysed %d structures from %d so far...' % (count, i)
80
81             if not self.settings.test(entry):
82                 continue
83             total += 1
84
85             try:
86                 molecule = entry.molecule
87             except RuntimeError:
88                 continue
89
90             #The SMILES string for water is searched for
91             #Structures without a water SMILES string but that meet all other
92             #search criteria form the OTHER list
93             if any(component.smiles == 'O' for component in molecule.components):
94                 hydrates_writer.write(molecule)
95                 count += 1
96             else:
97                 writer.write(molecule)
98
99             #A final count of structures with water SMILES strings are printed
100            print 'Found %d hydrates from %d valid structures' % (count, total)
101
102 if __name__ == '__main__':
103     # This runs the script
104     r = Runner()
105     r.run()

```

```

1 #2 Find solvates and hydrates with no water SMILES strings
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #Input is no longer the CSD but the output files of the previous script
10 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\entries_with_water_smiles_string.txt"
11 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\entries_with_no_water_smiles_string.txt"
12
13 class Runner(argparse.ArgumentParser):
14
15     def __init__(self):
16         super(self.__class__, self).__init__(description=__doc__)
17         self.add_argument(
18             '-i', '--input', default=filepath1,
19             help='input database filepath')
20
21         self.add_argument(
22             '-o', '--output', default='solvated_hydrates.gcd',
23             help='output file [solvated_hydrates.gcd]')
24
25         self.add_argument(
26             '-m', '--maximum', default=0, type=int,
27             help='Maximum number of structures to find [all]')
28
29
30     args = self.parse_args()
31     self.args = args
32     self.settings = search.Search.Settings()
33     self.settings.max_hit_structures = self.args.maximum
34
35     def run(self):
36
37         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
38         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
39
40         with io.EntryWriter(self.args.output) as solvate_writer:
41             with io.EntryWriter("hydrates.gcd") as writer1:
42                 with io.EntryWriter("solvated_waterless_forms.gcd") as writer2:
43                     with io.EntryWriter("waterless_forms.gcd") as writer3:
44
45                         count1 = 0
46                         total1 = 0
47                         #The title of each entry with a water SMILES string is searched for
48                         #the word "solvate"
49                         #Solvated hydrates are separated from non-solvated hydrates based
50                         #on the presence of "solvate" in the chemical name
51                         for a in range(len(entry_reader1)):
52                             entry1 = entry_reader1[a]
53                             if entry1.chemical_name != None:
54                                 title1 = entry1.chemical_name
55                                 if 'solvate' in title1:
56                                     solvate_writer.write(entry_reader1[a])
57                                     count1 += 1
58                                 else:
59                                     writer1.write(entry_reader1[a])
60
61                                 elif entry1[a].identifier == 'KIPWOY' or entry1[a].identifier
62                                 == 'KIPWUE':
63                                     solvate_writer.write(entry_reader1[a])
64
65                                     total1 += 1
66
67             print 'Found %d solvated structures from %d hydrate structures' %

```

```

(count1, total1)

63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106

(count1, total1)

count2 = 0
count3 = 0
total2 = 0
#The title of each entry without a water SMILES string is searched
for "hydrate" and "deuterium oxide"
#The formula of each entry without a water SMILES string is
searched for "H2 O1" and "D2 O1"
#Once again structures with "solvate" in their chemical name are
separated from those without "solvate"
for b in range(len(entry_reader2)):
    entry2 = entry_reader2[b]
    if entry2.chemical_name != None:
        title2 = entry2.chemical_name
        pieces2 = entry2.formula.split(',')
        if 'deuterium oxide' in title2 and any(x == 'D2 O1' for x
in pieces2) and 'solvate' in title2 or 'deuterium oxide' in
title2 and any('(D2 O1)' in x for x in pieces2) and
'solvate' in title2:
            solvate_writer.write(entry_reader2[b])
            count3 += 1
        elif 'deuterium oxide' in title2 and any(x == 'D2 O1' for x
in pieces2) or 'deuterium oxide' in title2 and any('(D2
O1)' in x for x in pieces2):
            writer1.write(entry_reader2[b])
        elif 'hydrate' in title2 and any(x == 'H2 O1' for x in
pieces2) and 'solvate' in title2 or 'hydrate' in title2 and
any('(H2 O1)' in x for x in pieces2) and 'solvate' in title2:
            solvate_writer.write(entry_reader2[b])
            count3 += 1
        elif 'hydrate' in title2 and any(x == 'H2 O1' for x in
pieces2) or 'hydrate' in title2 and any('(H2 O1)' in x for
x in pieces2):
            writer1.write(entry_reader2[b])
#Structures that are hydrates but are not indicated as such
by both "H2 O1" or "D2 O1" in the formula and "hydrate" in
the chemical name are accounted for here
#FIFYUQ and QIMROV have "H2 O1" in the formula but no
"hydrate" in the chemical name
elif entry_reader2[b].identifier == 'FIFYUQ' or
entry2[b].entry_reader2[b].identifier == 'QIMROV':
    solvate_writer.write(entry_reader2[b])
#KIJFAN and LOYXAA have "hydrate" in the chemical name but
no "H2 O1" or "D2 O1" in the formula
elif entry_reader2[b].identifier == 'KIJFAN' or
entry_reader2[b].identifier == 'LOYXAA':
    writer1.write(entry_reader2[b])
elif 'solvate' in title2:
    writer2.write(entry_reader2[b])
    count2 += 1
else:
    writer3.write(entry_reader2[b])

else:
    print entry_reader2[a].identifier

total2 += 1
print 'Found %d solvated hydrate structures from %d waterless
structures' % (count3, total1)
print 'Found %d solvated structures from %d waterless structures' %
(total2, total2)

if __name__ == '__main__':
    r = Runner()

```



```

1 #3 Find forms without smiles strings
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\solvated_hydrates.txt"
10 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\solvated_waterless_forms.txt"
11 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\hydrates.txt"
12 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\waterless_forms.txt"
13
14 class Runner(argparse.ArgumentParser):
15
16     def __init__(self):
17         super(self.__class__, self).__init__(description=__doc__)
18         self.add_argument(
19             '-i', '--input', default=filepath1,
20             help='input database filepath')
21     )
22     self.add_argument(
23         '-o', '--output', default='smiles_solvated_hydrates.gcd',
24             help='output file [smiles_solvated_hydrates.gcd]')
25     )
26     self.add_argument(
27         '-m', '--maximum', default=0, type=int,
28             help='Maximum number of structures to find [all]')
29     )
30
31     args = self.parse_args()
32
33     self.args = args
34     self.settings = search.Search.Settings()
35     self.settings.max_hit_structures = self.args.maximum
36
37     def run(self):
38
39         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
40         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
41         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
42         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
43
44         total = 0
45         count = 0
46
47         with io.EntryWriter(self.args.output) as SMILES_writer:
48             with io.EntryWriter("smiles_solvated_waterless_forms.gcd") as writer1:
49                 with io.EntryWriter("smiles_hydrates.gcd") as writer2:
50                     with io.EntryWriter("smiles_waterless_forms.gcd") as writer3:
51                         with io.EntryWriter("None_solvated_hydrates.gcd") as writer4:
52                             with io.EntryWriter("None_solvated_waterless_forms.gcd") as writer5:
53                                 with io.EntryWriter("None_hydrates.gcd") as writer6:
54                                     with io.EntryWriter("None_waterless_forms.gcd") as writer7:
55
56                                         loop = 0
57                                         entry1 = entry_reader1
58                                         while loop < 4:
59                                             #The entry reader is used to generate an
59                                             #entry SMILES string
59                                             #Any SMILES strings that return as None are

```

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

```
91 if __name__ == '__main__':
92     # This runs the script
93     r = Runner()
94     r.run()
```

separated as hydrate and anhydrate structures without SMILES strings  
#None will be returned for both absent SMILES strings and incomplete SMILES strings (at least one molecule SMILES string is present)

```
for a in range(len(entry1)):
    if loop == 0:
        if entry1[a].molecule.smiles == None:
            writer4.write(entry1[a])
        else:
            SMILES_writer.write(entry1[a])
    if loop == 1:
        if entry1[a].molecule.smiles == None:
            writer5.write(entry1[a])
        else:
            writer1.write(entry1[a])
    if loop == 2:
        if entry1[a].molecule.smiles == None:
            writer6.write(entry1[a])
        else:
            writer2.write(entry1[a])
    if loop == 3:
        if entry1[a].molecule.smiles == None:
            writer7.write(entry1[a])
        else:
            writer3.write(entry1[a])
loop += 1
if loop == 1:
    entry1 = entry_reader2
elif loop == 2:
    entry1 = entry_reader3
elif loop == 3:
    entry1 = entry_reader4
```







```

hash('petrol'):'CCCCCC', hash('n-alkane'):'CCCCC',
hash('dichloromethanal'):'ClC(Cl)=O',
hash('1-fluoro-4-methylbenzene'):'Cc1ccc(F)cc1',
hash('1,2,3-trimethoxybenzene'):'COc1cccc(OC)c1OC',
hash('1,3,5-trimethoxybenzene'):'COc1cc(OC)cc(OC)c1',
hash('chloroform'):'C1C(Cl)Cl', hash('methanol'):'CO',
hash('hexane'):'CCCCCC', hash('ethanol'):'CCO',
hash('acetonitrile'):'CC#N', hash('toluene'):'Cc1ccccc1',
hash('tetrahydrofuran'):'C1CCOC1', hash('acetone'):'CC(C)=O',
hash('dichloromethane'):'ClCCl', hash('1,2-dichloroethane'):'ClCCl',
hash('dimethylsulfoxide'):'CS(C)=O', hash('1,4-dioxane'):'C1COCCO1',
hash('1,2-dimethoxyethane'):'COCCOC', hash('o-xylene'):'Cc1cccc1c',
hash('propanol'):'CCO', hash('dioxane'):'C1COCCO1',
hash('n-pentane'):'CCCCCC', hash('4-cyanopyridine'):'N#Cc1ccncc1',
hash("4,4'-bipyridine"):'c1cc(ccn1)clccncc1', hash('isopropanol'):'CC(C)O',
hash('1-butanol'):'CCCCO', hash('nitromethane'):'CN(=O)=O',
hash('n-hexane'):'CCCCCC', hash('dimethylformamide'):'CN(C)C=O',
hash('N,N-dimethylformamide'):'CN(C)C=O', hash('cyclohexane'):'C1CCCCC1',
hash('hydroquinone'):'Oc1ccc(O)cc1', hash('benzene'):'clcccc1',
hash('octane'):'CCCCCC', hash('cyclohexanone'):'O=C1CCCCC1',
hash('isobutanol'):'CC(C)CO', hash('p-xylene'):'Cc1ccc(C)cc1',
hash('pentane'):'CCCCCC', hash('pyridine'):'clccncc1',
hash('t-butanol'):'CC(C)(C)O', hash('m-xylene'):'Cc1cccc(C)c1',
hash('phenol'):'Oc1ccccc1', hash('o-cresol'):'Cc1ccccc1O',
hash('p-cresol'):'Cc1ccc(O)cc1',
hash('2,4,6-trimethylpyridine'):'Cc1cc(C)nc(C)c1',
hash('1-methylpyrrolidin-2-one'):'CN1CCCC1=O',
hash('1,1,2,2,3,3,4,4,5,5,6,6,7,7,8,8-hexadecafluoro-1,8-diiodooctane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)I',
hash('trifluoroethanol'):'OCC(F)(F)F', hash('dichloroethane'):'ClCCl',
hash('1,1,2,2-tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('2-propanol'):'CC(C)O', hash('chlorobenzene'):'Cc1cccc1',
hash('naphthalene'):'c1ccc2ccccc2c1', hash('propan-2-ol'):'CC(C)O',
hash('benzonitrile'):'N#Cc1ccccc1', hash('nitrobenzene'):'O=N(=O)c1ccccc1',
hash('dodecane'):'CCCCCCCCCCCC',
hash('1,2-dichlorobenzene'):'Clc1ccccc1Cl', hash('n-heptane'):'CCCCCC',
hash('methanesulfonate'):'CS(=O)(=O)[O-]',
hash('3-methylbutan-1-ol'):'CC(C)CCO',
hash('1,1,2-trichloroethane'):'ClCC(Cl)Cl', hash('n-butane'):'CCCC',
hash('dichlorine'):'ClCl', hash('o-dichlorobenzene'):'Clc1ccccc1Cl',
hash('N,N-dimethylacetamide'):'CN(C)C(C)=O',
hash('4-methylaniline'):'Cc1ccc(N)cc1', hash('morpholine'):'C1COCCN1',
hash('1,3-dimethylimidazolidin-2-one'):'CN1CCN(C)C1=O',
hash('1,3-dimethoxybenzene'):'COc1cccc(OC)c1', hash('n-propanol'):'CCCO',
hash('mesitylene'):'Cc1cc(C)cc(C)c1',
hash('acetylacetone'):'CC(=O)CC(C)=O', hash('diethylether'):'CCOCC',
hash('hydrazine'):'NN', hash('dimethylamine'):'CNC',
hash('1,2-difluorobenzene'):'Fclccccc1F', hash('cycloheptane'):'C1CCCCC1',
hash('formamide'):'NC=O', hash('2-butanol'):'CCC(C)O',
hash('ethylacetate'):'CCOC(C)=O', hash('butan-1-ol'):'CCCCO',
hash('bromomethane'):'CBr', hash('piperazine'):'C1CNCCN1',
hash('n-nonane'):'CCCCCCCCC',
hash('1,2,4-trimethoxybenzene'):'COc1ccc(OC)c(OC)c1',
hash('1,2-diaminoethane'):'NCCN', hash('n-butanol'):'CCCCO',
hash('N-methylformamide'):'CNC=O', hash('pyrrolidine'):'C1CCNCl',
hash('actone'):'CC(C)=O', hash('n-octane'):'CCCCCC',
hash('('S)-2-methyl-1-butanol'):'CCC(C)CO',
hash('N,N-diethylformamide'):'CCN(CC)C=O',
hash('tetracyanoethylene'):'N#CC(C#N)=C(C#N)C#N',
hash('fluorobenzene'):'Fclccccc1', hash('heptane'):'CCCCCC',
hash('tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('methylcyclohexane'):'CC1CCCCC1',
hash('tetrachloromethane'):'ClC(Cl)(Cl)Cl', hash('1-propanol'):'CCCO',
hash('acetamide'):'CC(N)=O', hash('propane'):'CCC',
hash('diethylamine'):'CCNCC',
hash('hexafluorophosphate'):'F[P-](F)(F)(F)(F)F',
hash('hexafluoropropan-2-ol'):'OC(C(F)(F)F)C(F)(F)F',

```

```

hash('1-chloro-4-methylbenzene'):'Cc1ccc(Cl)cc1',
hash('acetaldehyde'):'CC=O', hash('propiononitrile'):'CCC#N',
hash('bromobenzene'):'Brclcccc1'}
48 deuterated_dictionary = {hash('dideutero-dichloromethane'):'ClCCl',
hash('perdeutero-toluene'):'Cclcccc1', hash('octadeutero furan'):'C1CCOC1',
hash('deutero-ethanol'):'CCO', hash('deuterochloroform'):'C1C(Cl)Cl',
hash('deutero-chloroform'):'C1C(Cl)Cl', hash('deuterobenzene'):'c1cccc1',
hash('deuteromethanol'):'CO', hash('hexadeutero-benzene'):'c1cccc1',
hash('perdeuteroacetone'):'CC(C)=O',
hash('dideuterodichloromethane'):'ClCCl', hash('perdeuteromethanol'):'CO',
hash('hexadeutero benzene'):'c1cccc1'}
49 acids_dictionary = {hash('trifluoroacetic'):'OC(=O)C(F)(F)F',
hash('acetic'):'CC(O)=O', hash('formic'):'OC=O',
hash('sulfuric'):'OS(O)(=O)=O', hash('maleic'):'OC(=O)C=CC(O)=O',
hash('benzoic'):'OC(=O)c1cccc1', hash('succinic'):'OC(=O)CCC(O)=O'}
50 ethers_dictionary = {'t-butyl methyl ether':'COC(C)(C)C',
hash('di-isopropyl'):'CC(C)OC(C)C', hash('dipropyl'):'CCCOCCCC',
hash('diethyl'):'CCOCC', hash('diisopropyl'):'CC(C)OC(C)C',
hash('isopropyl'):'CC(C)OC(C)C', hash('t-butyl'):'COC(C)(C)C',
hash('acetic'):'CCOC(C)=O'}
51 acetates_dictionary = {hash('ethyl'):'CCOC(C)=O',
hash('amyl'):'CCCCOC(C)=O', hash('methyl'):'COC(C)=O',
hash('butyl'):'CCCCOC(C)=O', hash('isobutyl'):'CC(C)COC(C)=O'}
52 glycol_dictionary = {hash('ethylene'):'OCOCO'}
53 oxalates_dictionary = {hash('dipropyl'):'CCCO(=O)C(=O)OCOC'}
54 acid_neutral_dictionary = {hash('oxalic'):'OC(=O)C(O)=O',
hash('nitric'):'ON(=O)=O', hash('trans-but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('fumaric'):'OC(=O)C=CC(O)=O', hash('phosphonic'):'OP(O)=O',
hash('perchloric'):'O[Cl](=O)(=O)=O', hash('phosphoric'):'OP(O)(O)=O',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(O)=O',
hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(O)=O)c(cc1C(O)=O)C(O)=O',
hash('trifluoroacetic'):'OC(=O)C(F)(F)F'}
55 acid_charged_dictionary = {hash('oxalic'):'OC(=O)C(=O)[O-]',
hash('nitric'):'O=N(=O)[O-]', hash('trans-but-2-enedioic'):'OC(=O)C=CC(=O)[O-]',
hash('but-2-enedioic'):'OC(=O)C=CC(=O)[O-]', hash('fumaric'):'OC(=O)C=CC(=O)[O-]', hash('phosphonic'):'OP(=O)[O-]', hash('perchloric'):'O=[Cl](=O)(=O)[O-]', hash('phosphoric'):'OP(O)(=O)[O-]', hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(=O)[O-]', hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(=O)[O-])c(cc1C(=O)[O-])C(O)=O', hash('trifluoroacetatic'):'FC(F)(F)C(=O)[O-]'}
56
57 partial_name = [hash('ketone'), hash('propanoate'), hash('carbonate'),
hash('formate'), hash('iodine'), hash('tetrachloride'), hash('sulfoxide'),
hash('dichloride'), hash('chloride'), hash('benzoate'), hash('disulfide'),
hash('oxide')]
58 full_solvate_name_dictionary = {'methyl isobutyl ketone':'CC(C)CC(C)=O',
'deuterium oxide':'O', 'ethyl propanoate':'CCOC(=O)CC', 'dimethyl
carbonate':'COC(=O)OC', 'ethyl formate':'CCOC=O', '12]cycloparaphenylene
iodine':'II', 'carbon tetrachloride':'C1C(Cl)(Cl)Cl', 'dimethyl
sulfoxide':'CS(C)=O', 'methylene dichloride':'ClCCl', 'methylene
chloride':'ClCCl', 'ethyl benzoate':'CCOC(=O)c1cccc1', 'carbon
disulfide':'S=C=S'}
59
60 #These lists keep track of the solvent names where all the evaluated
entries contain the corresponding solvent SMILES string
61 complete_smiles_partial_name = [hash('ether'), hash('peroxide'),
hash('oxide'), hash('acid'), hash('glycol'), hash('oxalate'),
hash('oxide')), hash('sulfoxide'), hash('alcohol'), hash('anhydride'),
hash('acetate'), hash('disulfide'), hash('dioxide'), hash('triamide'),
hash('fluoride'), hash('dulfoxide'), hash('sulfide'),
hash('hydroperoxide'), hash('propionate'), hash('acrylate'),
hash('ather'), hash('chloride'), hash('benzoate'), hash('triacetate'),
hash('propanoate'), hash('acid')), hash("2,2'-oxydiacetate"),
hash('ester'), hash('2-methylbutanoate')]

```

```

complete_smiles_full_solvate_name_dictionary = {'trideutero-oxonium
deuterium oxide':'[D+](OH2)[OH2]', 'hydrogen peroxide':'OO', 'propionic
acid':'CCC(O)=O', '(S)-propylene glycol': 'CC(O)CO', 'acrylic
acid':'OC(=O)C=C', 'diethyl oxalate': 'CCOC(=O)C(=O)OCC', 'dimethyl
oxalate': 'COC(=O)C(=O)OC', 'di(deuterium oxide)': 'O', 'butyric
acid': 'CCCC(O)=O', 'dimethyl ether': 'COC', 'perdeuterodimethyl
sulfoxide': 'CS(C)=O', 'picric acid': 'Oc1c(cc(cc1N(=O)=O)N(=O)=O)N(=O)=O',
'propanoic acid': 'CCC(O)=O', '(trifluoromethyl)benzyl
alcohol': 'OC(c1ccccc1)C(F)(F)F', 't-butyl alcohol': 'CC(C)(C)O', 'acetic
anhydride': 'CC(=O)OC(C)=O', 't-butyl acetate': 'CC(=O)OC(C)(C)C',
'chloroacetic acid': 'OC(=O)CCl', 'isoamyl acetate': 'CC(C)CCOC(C)=O',
't-butylcarboxylic acid': 'CC(C)(C)C(O)=O', 'benzyl alcohol': 'Oc1cccc1',
'ethanedioic acid': 'OC(=O)C(O)=O', 'diphenyl ether': 'O(c1ccccc1)c1ccccc1',
"(2,7-diethylbenzo[1,2-b:6,5-b']bisthiene-4,5-diy)C60fullerene
disulfide": 'S=C=S', 'n-butyl acetate': 'CCCCOC(C)=O', 'carbon
dioxide': 'O=C=O',
'2-amino-3-(2-hydroxy-2-phenylethyl)pyrazolo[1,5-a]pyrimidin-5(4H)-one
acetate': 'CC(O)=O', 'pentanedioic acid': 'OC(=O)CCCC(O)=O', 'dimethyl
sulfoxide': 'CS(C)=O', 'isopropyl alcohol': 'CC(C)O', 'isoamyl
alcohol': 'CC(C)CCO', 'bis(trideuteromethyl) sulfoxide': 'CS(C)=O',
'trifluoromethanesulfonic acid': 'OS(=O)(=O)C(F)(F)F',
'2-amino-6-(4-chlorophenyl)-5-((3-chlorophenyl)diazenyl)nicotinate
acetate': 'CC(O)=O', 'pivalic acid': 'CC(C)(C)C(O)=O', '2,6-dihydroxybenzoic
acid': 'OC(=O)c1c(O)cccc1O', '1-hydroxynaphthalene-2-carboxylic
acid': 'OC(=O)c1ccc2cccc2c1O', 'decanedioic acid': 'OC(=O)CCCCCCCC(O)=O',
'hydroxy(phenyl)acetic acid': 'OC(O)=O)c1ccccc1', 'sulfur
dioxide': 'O=S=O', 'salicylic acid': 'OC(=O)c1ccccc1O', 'butanoic
acid': 'CCCC(O)=O', 'phthalic acid': 'OC(=O)c1ccccc1C(O)=O', '2-phenylacetic
acid': 'OC(=O)Cc1ccccc1', 'hexamethylphosphoric acid
triamide': 'CN(C)P(=O)(N(C)C)N(C)C', '2,5-dihydroxybenzoic
acid': 'OC(=O)c1cc(O)cccc1O', 'squaric acid': 'OC1=C(O)C(=O)C1=O',
'hyponitrous acid': 'ON=NO', '2,2,2-trifluoroacetic acid': 'OC(=O)C(F)(F)F',
'trichloroacetic acid': 'OC(=O)C(Cl)(Cl)Cl', 'sulfamic
acid': '[NH3+]S(=O)(=O)[O-]', 't-butyl(methyl) ether': 'COC(C)(C)C',
'4-aminobenzoic acid': 'Nc1ccc(cc1)C(O)=O', 'hydrogen fluoride': 'F',
'methanoic acid': 'OC=O', '2,3,5-triodobenzoic
acid': 'OC(=O)c1cc(I)cc(I)c1I', 'phenyl sulfoxide': 'CS(=O)c1ccccc1',
'3,4-dimethylphenyl sulfoxide': 'Cc1ccc(cc1)S(C)=O', 'o-tolyl
sulfoxide': 'Cc1cccc1S(C)=O', 'tartaric acid': 'OC(C(O)C(O)=O)C(O)=O',
'dimethyl dulfoxide': 'CS(C)=O', 'propylene sulfide': 'CC1CS1', 'isopropyl
acetate': 'CC(C)OC(C)=O', 'carbonic acid': 'OC(O)=O', '4-methoxybenzoic
acid': 'COc1ccc(cc1)C(O)=O', '1,4-phenylenediacetic
acid': 'OC(=O)Cc1ccc(CC(O)=O)cc1', '(1,4-dioxan-2-yl
hydroperoxide)': 'OOC1COCCO1', 'triethylamine oxide': 'CCN(=O)(CC)CC',
'diethyl 2,2'-oxydiacetate': 'CCOC(=O)COCC(=O)OCC', 'hexadeuterodimethyl
sulfoxide': 'CS(C)=O', 'methyl propionate': 'CCC(=O)OC', 'methyl
acrylate': 'COC(=O)C=C', 'bis(2-methoxyethyl) ether': 'COCCOCCOC',
'tris(pyrrolidino)phosphine oxide': 'O=P(N1CCCC1)(N1CCCC1)N1CCCC1',
'pyromellitic anhydride': 'O=C1OC(=O)c2cc3C(=O)OC(=O)c3cc12', 'diethyl
ether': 'CCOCC', 't-butyl chloride': 'CC(C)(C)Cl', 'p-nitrobenzoic
acid': 'OC(=O)c1ccc(cc1)N(=O)=O', 'methyl benzoate': 'COC(=O)c1ccccc1',
't-butylbenzoic acid': 'CC(C)(C)c1ccc(cc1)C(O)=O', 'aminosulfonic
acid': '[NH3+]S(=O)(=O)[O-]', 'ethanoic acid': 'CC(O)=O', 'perdeuterodimethyl
sulfoxide': 'CS(C)=O', 'propane-1,2,3-triyl
triacetate': 'CC(=O)OCC(COC(C)=O)OC(C)=O', 'hydrochloric acid': 'Cl',
'isobutyric acid': 'CC(C)C(O)=O', 'methyl propanoate': 'CCC(=O)OC', 'n-butyl
alcohol': 'CCCCO', '2-hydroxybenzoic acid': 'OC(=O)c1ccccc1O',
'orthophosphoric acid': 'OP(O)(O)=O', 'hexanedioic acid': 'OC(=O)CCCC(O)=O',
'pimelic acid': 'OC(=O)CCCCC(O)=O', 'octanedioic
acid': 'OC(=O)CCCCCCC(O)=O', 'bis(trichloroacetic
acid)': 'OC(=O)C(Cl)(Cl)Cl', 'propylene glycol': 'CC(O)CO', 'di-imino-oxalic
acid diethyl ester': 'CCOC(=N)C(=N)OCC', 'ethyl
2-methylbutanoate': 'CCOC(=O)C(C)CC'}

```

```

complete_smiles_dictionary =
{hash('1,1,2,2,3,3,4,4,5,5,6,6-dodecafluoro-1,6-diiodohexane'): 'FC(F)(I)C(F)(
F)C(F)(F)C(F)(F)C(F)(F)I',
hash('2,2,2-trifluoroethanol'): 'OCC(F)(F)F',

```

```

hash('1,1,2,2-tetrachloro-1,2-dideuteroethane'):'ClC(Cl)C(Cl)Cl',
hash("2,2'-bipyrimidine"):'c1cnc(nc1)c1nccn1',
hash('2-methyl-2,4-pentanediol'):'CC(O)CC(C)(C)O',
hash('1,1,1,3,3,3-hexafluoropropan-2-ol'):'OC(C(F)(F)F)C(F)(F)F',
hash('pyrazine'):'clcnccn1', hash('2-methyl-3-butanol'):'CC(C)C(C)O',
hash('2-pentanol'):'CCCC(C)O', hash('2-hexanol'):'CCCCCC(C)O',
hash('2-heptanol'):'CCCCCC(C)O', hash('ethane-1,2-diol'):'OCCO',
hash('N,N-diethylethanamine'):'CCN(CC)CC',
hash('4-ethylpyridine'):'CCC1ccncc1', hash('2-methylpyridine'):'Cc1ccccn1',
hash('acidN,N-dimethylformamide'):'CN(C)C=O',
hash('dimethylsulfone'):'CS(C)(=O)=O',
hash('2-methylpentane-2,4-diol'):'CC(O)CC(C)(C)O',
hash('7-methyl-2-oxo-3,6,8-trioxadecane'):'CCOC(C)OCCOC(C)=O',
hash('4-nitrophenol'):'Oc1ccc(cc1)N(=O)=O',
hash('1,10-phenanthroline'):'c1cnc2c(c1)ccc1cccnc21',
hash('1-methyl-2-pyrrolidinone'):'CN1CCCC1=O',
hash('S)-butane-1,2-diol'):'CCC(O)CO',
hash('S)-propane-1,2-diol'):'CC(O)CO',
hash('perdeuterodimethylsulfoxide'):'CS(C)=O',
hash('S)-butan-2-ol'):'CCC(C)O', hash('2-methylpropanol'):'CC(C)(C)O',
hash('trans-1,4-cyclohexanediol'):'OC1CCC(O)CC1',
hash('quinoline'):'c1ccc2ncccc2c1', hash('1-undecanol'):'CCCCCCCCCCC',
hash('acetylene'):'C#C', hash('1-naphthol'):'Oc1cccc2cccc12',
hash('2-naphthol'):'Oc1ccc2cccc2c1',
hash('2-methoxy-2-methylpropane'):'CO(C(C)C)C',
hash('hexane-1,6-diol'):'OCCCCCC',
hash('1,3,5-trimethylbenzene'):'Cc1cc(C)cc(C)c1',
hash("N,N'-dimethylformamide"):'CN(C)C=O',
hash("N,N'-dimethylacetamide"):'CN(C)C(C)=O',
hash('2-ethylhexanol'):'CCCCC(CC)CO', hash('octan-1-ol'):'CCCCCCCO',
hash('4-amino-1,2,4-triazole'):'NN1C=NN=C1',
hash('perdeutero-dimethylsulfoxide'):'CS(C)=O', hash('methanamine'):'CN',
hash('chloroacetonitrile'):'ClCC#N', hash('malonate'):'CCOC(=O)CC(=O)OCC',
hash('succinate'):'CCOC(=O)CCC(=O)OCC', hash('n-butan-1-ol'):'CCCCO',
hash('ethane-1,2-diol'):'OCCO', hash('4-methylpiperidine'):'CC1CCNCC1',
hash('2-butanone'):'CCC(C)=O', hash('ethyleneglycol'):'OCCO',
hash('diiodine'):'II', hash('acetotnitrile'):'CC#N',
hash('nonane'):'CCCCCCCCC',
hash('bis(phthalate)-N,N-dimethylformamide'):'CN(C)C=O',
hash('butylamine'):'CCCCN', hash('n-decanol'):'CCCCCCCCCCC',
hash('n-pentanol'):'CCCCCO', hash('m-cresol'):'Cc1cccc(O)c1',
hash('pentan-2-ol'):'CCCC(C)O',
hash('1-(4-fluorophenyl)ethanol'):'CC(O)c1ccc(F)cc1',
hash('2-ethoxyethanol'):'CCOCOC', hash('2-propoxyethanol'):'CCCOCCO',
hash('hexanol'):'CCCCCCO', hash('pentanol'):'CCCCCCO',
hash('1,6-diamino-hexane'):'NCCCCCN', hash('phenylmethanol'):'OCc1cccc1',
hash("5,5'-bipyrimidine"):'c1ncc(cn1)c1cncn1',
hash('6-nitrobenzimidazole'):'O=N(=O)c1ccc2NC=Nc2c1',
hash('bromo(trichloro)methane'):'ClC(Cl)(Cl)Br',
hash('3,4-Diacetyl-15,21-dioxatetraacyclo[23.4.0.02,7.06,11]nonacosa-1(29),2,4,6,8,10,25,27-octaene-14,22-dione'):'CC(=O)c1cc2c3CCC(=O)OCCCCOC(=O)CCc4ccccc4c(c2ccc3)c1C(C)=O', hash('butanol-1,2-diol'):'CCC(O)CO',
hash('ethylbenzene'):'CCc1cccc1', hash('1-pentanol'):'CCCCCCO',
hash('undecanol'):'CCCCCCCCCCC', hash('iodomethane'):'CI',
hash('1,1,2,2,3,3,4,4,5,5,6,6,7,7,8,8,9,9,10,10-icosafaluoro-1,10-diiododecane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)I',
hash('1,1,2,2,3,3,4,4-octafluoro-1,4-diiodobutane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)I', hash('1,2-fluorobenzene'):'Fc1cccc1F',
hash('ethanedihydrazide'):'NNC(=O)C(=O)NN', hash('bromoform'):'BrC(Br)Br',
hash('p-bromophenol'):'Oc1ccc(Br)cc1',
hash('N-methylacetamide'):'CNC(C)=O', hash('n-propylamine'):'CCCN',
hash('4-methylbenzenesulfonamide'):'Cc1ccc(cc1)S(N)(=O)=O',
hash('methoxyethanol'):'CCOCO', hash('tetrabromide'):'BrC(Br)(Br)Br',
hash('N-methylpyrrolidone'):'CN1CCCC1=O', hash('ethoxyethane'):'CCOCC',
hash('pentan-3-one'):'CCC(=O)CC', hash('ethane'):'CC',
hash('methylcyclopentane'):'CC1CCCC1',

```

```

hash('hexafluorobenzene'):'Fc1c(F)c(F)c(F)c(F)c1F',
hash('1,4-butanediol'):'OCCCO',
hash('1-chloro-2-methylbenzene'):'Cc1cccc1Cl',
hash('dimethyl-formamide'):'CN(C)C=O', hash('butan-2-one'):'CCC(C)=O',
hash('aniline'):'Nc1ccccc1', hash('1,4-dichlorobutane'):'ClCCCCl',
hash('1,4-difluorobutane'):'FCCCCF',
hash('1,2,4-trichlorobenzene'):'Clc1ccc(Cl)c(Cl)c1',
hash('iodoethane'):'CCI', hash('hexadeutero-dimethylsulfoxide'):'CS(C)=O',
hash('hydroxyacetonitrile'):'OCC#N', hash('butan-2-amine'):'CCC(C)N',
hash('({S})-butan-2-amine'):'CCC(C)N', hash('({R})-butan-2-amine'):'CCC(C)N',
hash('dibromomethane'):'BrCBr', hash('1,2-dibromobenzene'):'Brclcccc1Br',
hash('perdeuterodimethylsulphoxide'):'CS(C)=O',
hash('2-methylprop-1-ene'):'CC(C)=C',
hash('1-bromonaphthalene'):'Brclcccc2cccc12',
hash('n-decane'):'CCCCCC',
hash('2-methylpropane-2-peroxol'):'CC(C)(C)OO',
hash('2-phenylpropane-2-peroxol'):'CC(C)(OO)c1cccc1',
hash('2-isopropoxypropane'):'CC(C)OC(C)C',
hash('1,2-dibromoethene'):'BrC=CBr', hash('1,2-dibromoethane'):'BrCCBr',
hash('2-methylcyclohexanone'):'CC1CCCCC1=O', hash('isopentane'):'CCC(C)C',
hash('1,3,5-trinitrobenzene'):'O=N(=O)c1cc(cc(c1)N(=O)=O)N(=O)=O',
hash('2,4-pentanedione'):'CC(O)=CC(C)=O',
hash('4-methylpyridine'):'Cc1ccncc1', hash('acetonitrle'):'CC#N',
hash('butane-1,4-diol'):'OCCCO',
hash('1,3-diphenylacetone'):'O=C(Cc1cccc1)Cc1cccc1',
hash('2-ethylpyridine'):'Cc1ccccn1', hash('1,2,4-triazole'):'N1C=NC=N1',
hash('1-methylpyrrolidine'):'CN1CCCC1',
hash('trichloroethane'):'ClCC(Cl)Cl', hash('piperidine'):'C1CCNCC1',
hash('but-2-yne-1,4-diol'):'OCC#CCO', hash('methanedithione'):'S=C=S',
hash('2-t-butyl-4-methylphenol'):'Cc1ccc(O)c(c1)C(C)(C)C',
hash('aniline'):'Nc1ccccc1', hash('(methylsulfinyl)methane'):'CS(C)=O',
hash('1,1,3,3-tetramethylurea'):'CN(C)C(=O)N(C)C',
hash('cyclooctane'):'C1CCCCCCC1',
hash('1,3-dimethyltetrahydropyrimidin-2(1H)-one'):'CN1CCN(C)C1=O',
hash('2,2-difluoroethanol'):'OCC(F)F',
hash('trideuteroacetonitrile'):'CC#N',
hash('1-methyl-3-(methylsulfinyl)benzene'):'Cc1cccc(c1)S(C)=O',
hash('1-(ethylsulfinyl)-4-methylbenzene'):'CCS(=O)c1ccc(C)cc1',
hash('dimethylaminobenzene'):'CN(C)c1cccc1', hash('ammonia'):'N',
hash('perdeuterobenzene'):'c1cccc1',
hash('perdeuterochloroform'):'ClC(Cl)Cl', hash('benzylamine'):'NCc1cccc1',
hash('cyclopentanone'):'O=C1CCCC1',
hash('2-nitroaniline'):'Nc1ccccc1N(=O)=O',
hash('1-butoxybutane'):'CCCCOCCCC', hash('propanenitrile'):'CCC#N',
hash('hexamethylbenzene'):'Cc1c(C)c(C)c(C)c(C)c1C',
hash('1,2,3-trimethylbenzene'):'Cc1cccc(C)c1C',
hash('1,2,4-trimethylbenzene'):'Cc1ccc(C)c(C)c1',
hash('1,4-dimethylbenzene'):'Cc1ccc(C)cc1',
hash("N,N,N',N'-tetramethylmethane-1,2-diamine"):'CN(C)CCN(C)C',
hash('ethylenediamine'):'NCCN', hash('p-hydroquinone'):'Oc1ccc(O)cc1',
hash('N,N-Dimethylformamide'):'CN(C)C=O', hash('urea'):'NC(N)=O',
hash('butane-2,3-diol'):'CC(O)C(C)O',
hash('benzene-1,4-diol'):'Oc1ccc(O)cc1',
hash('furan-2-carbaldehyde'):'O=CC1=CC=CO1',
hash('2,6-dimethylpyridine'):'Cc1cccc(C)n1',
hash('perdeuterobenzene'):'c1cccc1',
hash('hexafluoroisopropanol'):'OC(C(F)(F)F)C(F)(F)F',
hash('3,6-dioxaoctane'):'CCOCOCOC',
hash('4-methoxytoluene'):'COc1ccc(C)cc1',
hash('1,4-dichlorobenzene'):'Clc1ccc(Cl)cc1',
hash('N-ethylananine'):'CCNCC', hash('iso-butanol'):'CC(C)CO',
hash('4-hydroxy-4-methylpentan-2-one'):'CC(=O)CC(C)(C)O',
hash('2-methylpropan-1-ol'):'CC(C)CO', hash('propane-1,2-diol'):'CC(O)CO',
hash('iodoform'):'IC(I)I', hash('1-phenylethanol'):'CC(O)c1cccc1',
hash('benzylamine'):'NCc1ccccc1',
hash('1,2-bis(4-pyridyl)ethene'):'c1cc(ccn1)C=Cc1ccncc1',
hash('1,4-dihydroxybenzene'):'Oc1ccc(O)cc1',

```

```

hash('ethoxycarbonylisothiocyanate'):'CCOC(=O)N=C=S',
hash('tribromo(nitro)methane'):'BrC(Br)(Br)N(=O)=O',
hash('2,2,2-tribromooacetamide'):'NC(=O)C(Br)(Br)Br',
hash('tetrahydrafuran'):'C1CCOC1',
hash('N-methylbenzamide'):'CNC(=O)c1ccccc1',
hash('phenylacetylene'):'C#Cc1cccc1',
hash('isonicotinamide'):'NC(=O)c1ccncc1',
hash('benzotriazole'):'N1N=Nc2ccccc12', hash('dichlormethane'):'ClCCl',
hash('N-methyl-2-pyrrolidone'):'CN1CCCC1=O', hash('lactone'):'O=C1CCCO1',
hash('diiodomethane'):'ICI',
hash('4-hydroxy-5-methylpentan-2-one'):'CC(=O)CC(C)(C)O',
hash('dichlorofluoromethane'):'FC(Cl)Cl',
hash('1,4-cyclohexanedione'):'O=C1CCC(=O)CC1',
hash('hexadeutero-acetone'):'CC(C)=O',
hash('1,2-bis(methoxy)benzene'):'COc1ccccc1OC',
hash('hexafluoro-sulfane'):'FS(F)(F)(F)F', hash('methylamine'):'CN',
hash('2-n-butyl-2-t-butyl-4,6-diphenyl-1H-1,3,5-triazine'):'CCCCC1(NC(=NC(=N1)c1ccccc1)c1ccccc1)C(C)(C)C', hash('hexanedinitrile'):'N#CCCCCC#N',
hash('2,3,11,12-dibenzo-1,4,10,13,16-hexaoxocyclo-octadeca-2,11-diene'):'C1COc2ccccc2OCCOCCOc2ccccc2OCCO1', hash('cyclopentanol'):'OC1CCCC1',
hash('methylsulfinyl)benzene'):'CS(=O)c1ccccc1',
hash('2-methylpropane'):'CC(C)C',
hash('pentafluorobenzaldehyde'):'Fc1c(F)c(F)c(C=O)c(F)c1F',
hash('ethanaol'):'CCO', hash('1-pentene'):'CCCC=C',
hash('1-chlorobutane'):'CCCC1', hash('triodomethane'):'IC(I)I',
hash('tetramethylethane-1,2-diamine'):'CN(C)CCN(C)C',
hash('3-ethyltoluene'):'CCc1cccc(C)c1', hash('neopentane'):'CC(C)(C)C',
hash('butane'):'CCCC', hash('1,2-diaminoethenol'):'NC=C(N)O',
hash('dimethylisosorbide'):'COC1COC2C(COC12)OC',
hash('5-chloropentanenitrile'):'ClCCCC#N',
hash('propanedinitrile'):'N#CCC#N', hash('o-toluidine'):'Cc1ccccc1N',
hash('p-toluidine'):'Cc1ccc(N)c1',
hash('hexachloroethane'):'C1C(Cl)(Cl)C(Cl)(Cl)Cl',
hash('1,4-dibromobutane'):'BrCCCCBr', hash('cyclohexenone'):'O=C1CCCC=C1',
hash('tetrahydro-2H-pyran'):'C1CCOCC1',
hash('1-(4-chlorophenyl)ethanol'):'CC(O)c1ccc(Cl)cc1',
hash('1-phenylethanone'):'CC(=O)c1ccccc1',
hash('4-ethyltoluene'):'CCc1ccc(C)cc1', hash('tetrahydropyran'):'C1CCOCC1',
hash('benzoquinone'):'O=C1C=CC(=O)C=C1',
hash('terephthalaldehyde'):'O=Cc1ccc(C=O)cc1', hash('acetonitrile'):'CC#N',
hash('2methylpyridine'):'Cc1cccn1', hash('3-chloropyridine'):'Clc1ccnc1',
hash('1,1-dichloroethane'):'CC(Cl)Cl',
hash('1-(4-methylphenyl)ethanol'):'CC(O)c1ccc(C)cc1',
hash('3-ethylpentan-3-ol'):'CCC(O)(CC)CC',
hash('S-2-methylbutan-1-ol'):'CCC(C)CO',
hash('1,3,5-trifluoro-2,4,6-triiodobenzene'):'Fc1c(I)c(F)c(I)c(F)c1I',
hash('N,N,4-trimethylbenzene-1-sulfonamide'):'CN(C)S(=O)(=O)c1ccc(C)cc1',
hash('iso-propanol'):'CC(C)O', hash('deuterodimethylsulfoxide'):'CS(C)=O',
hash('dibenzylamine'):'C(NCc1ccccc1)c1ccccc1',
hash('hexamethylphosphortriamide'):'CN(C)P(=O)(N(C)C)N(C)C',
hash('2,2'-oxydiethanol'):'OCCOCO',
hash('trichlorofluoromethane'):'FC(Cl)(Cl)Cl',
hash('dichlorodifluoromethane'):'FC(F)(Cl)Cl',
hash('chlorodifluoromethane'):'FC(F)Cl',
hash('1,1,1-trifluoroethane'):'CC(F)(F)F',
hash('2,2-dichloropropane'):'CC(C)(Cl)Cl',
hash('2,4,6-trimethylphenol'):'Cc1cc(C)c(O)c(C)c1',
hash('1,1,1,3,3,3-hexafluoroisopropanol'):'OC(C(F)(F)F)C(F)(F)F',
hash('cyclo-octylamine'):'NC1CCCCCCC1', hash('aetonitrile'):'CC#N',
hash('propan-2-one'):'CC(C)=O', hash('biphenyl'):'c1ccc(cc1)c1ccccc1',
hash('18-crown-6'):'C1COCCOCCOCCOCOCOCCO1',
hash('2,5-dihydrofuran'):'C1OCC=C1',
hash('carbazole'):'N1c2ccccc2c2ccccc12', hash('indole'):'N1C=C2CCCCC12',
hash('N,N-dimethylformamide'):'CN(C)C=O', hash('di-iodomethane'):'ICI',
hash('2,6-dichlorophenol'):'Oc1c(Cl)cccc1Cl',
hash('1-(morpholin-4-yl)ethanone'):'CC(=O)N1CCOCC1',
hash('4-methylbenzaldehyde'):'Cc1ccc(C=O)cc1',

```

```
248         entry1 = entry_reader2
249
250
251
252
253
254 if __name__ == '__main__':
255     # This runs the script
256     r = Runner()
257     r.run()
258
```

```

1 #6 Find potential duplicate none hydrate structures with formula match
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 from decimal import Decimal
9
10 #All the hydrate structures without entry SMILES strings will be compared to each other
11 #and hydrates with entry SMILES strings
12 #This script finds all the hydrate structures with the same entry formula as the
13 #hydrate structures without entry SMILES strings
14 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\hydrates_None_without_metals.txt"
15 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\hydrates_SMILES_without_metals.txt"
16
17 class Runner(argparse.ArgumentParser):
18
19     def __init__(self):
20         super(self.__class__, self).__init__(description=__doc__)
21         self.add_argument(
22             '-i', '--input', default=filepath1,
23             help='input database filepath')
24         self.add_argument(
25             '-o', '--output', default='hydrates_Non_all_potential_duplicates.gcd',
26             help='output file [hydrates_Non_all_potential_duplicates.gcd]')
27         self.add_argument(
28             '-m', '--maximum', default=0, type=int,
29             help='Maximum number of structures to find [all]')
30
31     args = self.parse_args()
32
33     self.args = args
34     self.settings = search.Search.Settings()
35     self.settings.max_hit_structures = self.args.maximum
36
37     def run(self):
38
39         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
40         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
41
42         with io.EntryWriter(self.args.output) as new_lists_writer:
43
44             #Two lists of formulas are created; one for hydrate structures with entry
45             #SMILES strings and one for hydrate structures without entry SMILES strings
46             None_formulas = []
47             formulas = []
48             positions = []
49             numbers = []
50             for a in range(len(entry_reader1)):
51                 component_formulas = entry_reader1[a].formula.split(',')
52                 if entry_reader1[a].identifier == 'LOYXAA':
53                     component_formulas.append('H2 O1')
54                 if entry_reader1[a].identifier != 'KIJFAN':
55                     component_formulas.append('x(H2 O1)')
56                 if all('x' not in t for t in component_formulas) == True and all('n('
57                 not in s for s in component_formulas) == True:
58                     factors = []
59                     for z in range(len(component_formulas)):
60                         if '(' and ')' in component_formulas[z]:
61                             pin = component_formulas[z].index('(')
62                             if component_formulas[z][:pin].isdigit() == True:
63                                 factors.append(int(component_formulas[z][:pin]))
64                             else:

```

```

64 factors.append(Decimal(component_formulas[z][:pin]))
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
factors.append(Decimal(component_formulas[z][:pin]))
else:
    factors.append(1)
normalized_factors = []
for y in range(len(factors)):
    normalized_factors.append(str(factors[y]/min(factors)))
if factors != normalized_factors:
    for x in range(len(normalized_factors)):
        if normalized_factors[x] == '1':
            if '(' and ')' in component_formulas[x]:
                start = component_formulas[x].index('(')
                end = component_formulas[x].index(')')
                component_formulas[x] =
                component_formulas[x][start+1:end]
        else:
            if '(' and ')' in component_formulas[x]:
                start = component_formulas[x].index('(')
                component_formulas[x] = normalized_factors[x] +
                component_formulas[x][start]
            else:
                component_formulas[x] = normalized_factors[x] + "("
                + component_formulas[x] + ")"
pre_formulas = []
#The stoichiometry of each molecule in the entry was excluded when
comparing the formulas
#The stoichiometry of structures that are considered duplicates was
checked later on
#Structures with different stoichiometry are considered distinct entries
for b in range(len(component_formulas)):
    if '(' and ')' in component_formulas[b]:
        start = component_formulas[b].index('(')
        end = component_formulas[b].index(')')
        letter_formula = component_formulas[b][start+1:end]
        elements = letter_formula.split(' ')
    #The elements of each molecular formula was screened to change
    deuterium to hydrogen
    #Structures with the same packing that exist in non-deuterated and
    deuterated forms are considered duplicate entries
else:
    elements = component_formulas[b].split(' ')
deuteriums = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'D']
if deuteriums != []:
    if len(deuteriums) != 1:
        print entry_reader1[a].identifier
        hydrogens = [i for i in elements if re.sub(r'[0-9]+', '', i) ==
        'H']
    if hydrogens == []:
        spot = elements.index(deuteriums[0])
        deuteriums[0] = deuteriums[0].replace('D', 'H')
        elements[spot] = deuteriums[0]
    else:
        #In cases where there are hydrogen and deuterium atoms in a
        molecule, the stoichiometry for each is used to determine
        the stoichiometry of the resulting hydrogen only formula
        if len(hydrogens) != 1:
            print entry_reader1[a].identifier
            spot = elements.index(deuteriums[0])
            spot2 = elements.index(hydrogens[0])
            stoichiometry1 = re.sub(r'[A-Z]+', '', deuteriums[0])
            stoichiometry2 = re.sub(r'[A-Z]+', '', hydrogens[0])
            total = int(stoichiometry1) + int(stoichiometry2)
            hydrogens[0] = hydrogens[0].replace(stoichiometry2,
            str(total))
            elements[spot2] = hydrogens[0]
            elements.remove(elements[spot])
        if '(' and ')' in component_formulas[b]:

```

```

120
121         letter_formula = ' '.join(elements)
122         component_formulas[b] =
123             component_formulas[b].replace(component_formulas[b][start+1:end], letter_formula)
124     else:
125         component_formulas[b] = ' '.join(elements)
126     pre_formulas.append(hash(component_formulas[b]))
127
128     pre_formulas = list(set(pre_formulas))
129     pre_formulas = sorted(pre_formulas)
130     None_formulas.append(pre_formulas)
131     formulas.append(pre_formulas)
132     #The position of each structure in the input file was recorded to
133     #update the list of formulas later on in the script
134     positions.append(a)
135     #The spacegroup number of each structure was also recorded in a
136     #separate list
137     #Structures that could not generate the spacegroup number were listed
138     #separately with their number determined manually in Mercury
139     if hash(entry_reader1[a].identifier) == hash('KECYBU15'):
140         numbers.append(87)
141     else:
142
143         numbers.append(entry_reader1[a].crystal.spacegroup_number_and_setting
144             [0])
145
146     #The same steps are taken for hydrate structures with entry SMILES strings
147     for d in range(len(entry_reader2)):
148         component_formulas = entry_reader2[d].formula.split(',')
149         if entry_reader2[d].identifier == 'LOYXAA':
150             component_formulas.append('H2 O1')
151         if entry_reader2[d].identifier != 'KIJFAN':
152             component_formulas.append('x(H2 O1)')
153         if all('x' not in t for t in component_formulas) == True and all('n('
154             not in s for s in component_formulas) == True:
155             factors = []
156             for z in range(len(component_formulas)):
157                 if '(' and ')' in component_formulas[z]:
158                     pin = component_formulas[z].index('(')
159                     if component_formulas[z][:pin].isdigit() == True:
160                         factors.append(int(component_formulas[z][:pin]))
161                     else:
162                         factors.append(Decimal(component_formulas[z][:pin]))
163             else:
164                 factors.append(1)
165             normalized_factors = []
166             for y in range(len(factors)):
167                 normalized_factors.append(str(factors[y]/min(factors)))
168             if factors != normalized_factors:
169                 for x in range(len(normalized_factors)):
170                     if normalized_factors[x] == '1':
171                         if '(' and ')' in component_formulas[x]:
172                             start = component_formulas[x].index('(')
173                             end = component_formulas[x].index(')')
174                             component_formulas[x] =
175                             component_formulas[x][start+1:end]
176                     else:
177                         if '(' and ')' in component_formulas[x]:
178                             start = component_formulas[x].index('(')
179                             component_formulas[x] = normalized_factors[x] +
180                             component_formulas[x][start]
181                     else:
182                         component_formulas[x] = normalized_factors[x] + "("
183                         + component_formulas[x] + ")"
184
185             pre_formulas2 = []
186             for e in range(len(component_formulas)):
```

```

176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
if (' and ') in component_formulas[e]:
    start = component_formulas[e].index('(')
    end = component_formulas[e].index(')')
    letter_formula = component_formulas[b][start+1:end]
    elements = letter_formula.split(' ')
else:
    elements = component_formulas[e].split(' ')
deuteriums = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'D']
if deuteriums != []:
    if len(deuteriums) != 1:
        print entry_reader2[d].identifier
    hydrogens = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'H']
    if hydrogens == []:
        spot = elements.index(deuteriums[0])
        deuteriums[0] = deuteriums[0].replace('D', 'H')
        elements[spot] = deuteriums[0]
    else:
        if len(hydrogens) != 1:
            print entry_reader2[d].identifier
        spot = elements.index(deuteriums[0])
        spot2 = elements.index(hydrogens[0])
        stoichiometry1 = re.sub(r'[A-Z]+', '', deuteriums[0])
        stoichiometry2 = re.sub(r'[A-Z]+', '', hydrogens[0])
        total = int(stoichiometry1) + int(stoichiometry2)
        hydrogens[0] = hydrogens[0].replace(stoichiometry2,
                                             str(total))
        elements[spot2] = hydrogens[0]
        elements.remove(elements[spot])
    if '(' and ')' in component_formulas[b]:
        letter_formula = ''.join(elements)
        component_formulas[e] =
            component_formulas[b].replace(component_formulas[e][start+1:en
d], letter_formula)
    else:
        component_formulas[e] = ''.join(elements)
    pre_formulas2.append(hash(component_formulas[e]))
pre_formulas2 = list(set(pre_formulas2))
pre_formulas2 = sorted(pre_formulas2)
formulas.append(pre_formulas2)
positions.append(d+718)
if hash(entry_reader2[d].identifier) == hash('MTYHFB03'):
    numbers.append(29)
else:
    numbers.append(entry_reader2[d].crystal.spacegroup_number_and_setting
[0])
#A while loop is used to iterate through the list of formulas for hydrate
structures without an entry SMILES string
while None_formulas != []:
    print len(None_formulas)
    formulas_to_remove = []
    positions_to_write = []
    numbers_to_compare = []
    formulas_to_remove.append(formulas[0])
    positions_to_write.append(positions[0])
    numbers_to_compare.append(numbers[positions[0]])
    #The formula of the first hydrate structure without an entry SMILES
    #string is compared to all the other hydrate structure formulas
    #Matching formulas are appended to two lists: one has their position in
    #the corresponding input file, the other has their spacegroup number
    #The formula list includes all hydrate structures, the beginning
    #formulas are for hydrate structures without an entry SMILES string
    for c in range(1, len(formulas)):
        if formulas[0] == formulas[c]:

```

```

233     positions_to_write.append(positions[c])
234     numbers_to_compare.append(numbers[positions[c]])
235
236     #If there are no matches, the length of the numbers_to_compare list
237     #will only contain the hydrate structure being compared
238     #Otherwise, the stoichiometry of the water molecules in each hydrate
239     #structure with that formula is determined
240     if len(numbers_to_compare) != 1:
241         #Only cases where the spacegroup number is the same for at least
242         #two structures are checked for duplicates
243         if all(x == numbers_to_compare[0] for x in numbers_to_compare):
244             elif any(numbers_to_compare.count(x) > 1 for x in
245                     numbers_to_compare):
246                 positions_to_remove = []
247                 #Each potential duplicate is compared to all the other
248                 #structures in the list
249                 #If the potential duplicate does not match both the water
250                 #stoichiometry and spacegroup number of another structure in the
251                 #list it is placed in the lists to be removed
252                 while len(positions_to_write) != len(positions_to_remove):
253                     for j in range(len(positions_to_write)):
254                         if positions_to_write[j] not in positions_to_remove:
255                             positions_to_remove.append(positions_to_write[j])
256                             if positions_to_write[j] < 625:
257                                 match = 0
258                                 for k in range(len(positions_to_write)):
259                                     if positions_to_write[k] not in
260                                         positions_to_remove:
261                                         #If the spacegroup number is the same,
262                                         #the structures are put in lists to be
263                                         #analyzed in part2 of this script
264                                         if numbers_to_compare[j] ==
265                                         numbers_to_compare[k]:
266                                             if match == 0:
267                                                 new_lists_writer.write(entry_read
268                                                 er1[positions_to_write[j]])
269                                             if positions_to_write[k] < 625:
270                                                 new_lists_writer.write(entry_read
271                                                 er1[positions_to_write[k]])
272                                             else:
273                                                 new_lists_writer.write(entry_read
274                                                 er2[positions_to_write[k]-625])
275                                             match += 1
276                                             positions_to_remove.append(positions_
277                                                 to_write[k])
278                                             if match != 0:
279                                                 new_lists_writer.write(entry_reader1[position
280                                                 s_to_write[j]])
281
282                                         #The lists of formulas and positions is updated so structures that have
283                                         #been evaluated are removed
284                                         #The list of spacegroup numbers is generated for each analysis and
285                                         #therefore does not need to be updated
286                                         None_formulas = [i for i in None_formulas if i not in formulas_to_remove]
287                                         formulas = [i for i in formulas if i not in formulas_to_remove]
288                                         positions = [j for j in positions if j not in positions_to_write]
289
290                                         if __name__ == '__main__':
291                                             r = Runner()

```

277

r.run()

278

```

1 #7 Find potential duplicate none hydrate structures with same identifier prefix
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 from decimal import Decimal
9
10 #All the hydrate structures that share the same component stoichiometry and spacegroup
11 #number found in the previous script are analyzed here
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
13 Entries\hydrates_None_all_potential_duplicates.txt"
14
15 class Runner(argparse.ArgumentParser):
16
17     def __init__(self):
18         super(self.__class__, self).__init__(description=__doc__)
19         self.add_argument(
20             '-i', '--input', default=filepath1,
21             help='input database filepath')
22         self.add_argument(
23             '-o', '--output', default='hydrates_None_all_manual_check.gcd',
24             help='output file [hydrates_None_all_manual_check.gcd]')
25         self.add_argument(
26             '-m', '--maximum', default=0, type=int,
27             help='Maximum number of structures to find [all]')
28
29     args = self.parse_args()
30
31     self.args = args
32     self.settings = search.Search.Settings()
33     self.settings.max_hit_structures = self.args.maximum
34
35     def run(self):
36
37         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
38
39         with io.EntryWriter(self.args.output) as new_lists_writer:
40
41             #Lists of the hydrate identifiers and entry readers are generated
42             solvated_hydrates_hash = []
43             solvated_hydrates_entry = []
44             for a in range(len(entry_reader1)):
45                 solvated_hydrates_hash.append(hash(entry_reader1[a].identifier))
46                 solvated_hydrates_entry.append(entry_reader1[a])
47
48             #The previous output file printed the first hydrate structure identifier in
49             #the list of duplicates twice
50             #For example "ABCDEF, ABCDEF01, ABCDEF"
51             #That way the duplicates can be found by searching for the next occurrence
52             #of the first identifier
53             while len(solvated_hydrates_hash) != 0:
54                 print len(solvated_hydrates_hash)
55                 identifiers = []
56                 identifiers.append(solvated_hydrates_hash[0])
57                 for d in range(1,len(solvated_hydrates_hash)):
58                     if identifiers[0] != solvated_hydrates_hash[d]:
59                         identifiers.append(solvated_hydrates_hash[d])
60                     else:
61                         del solvated_hydrates_hash[d]
62                         del solvated_hydrates_entry[d]
63                         break

```

```

64     #The first six letters in each identifier are appended to a list
65     nomenclature = []
66     for e in range(len(identifiers)):
67         pin = solvated_hydrates_hash.index(identifiers[e])
68         nomenclature.append(solvated_hydrates_entry[pin].identifier[:6])
69
70     #Two structures with the same first six letters in their identifiers
71     #are considered potential duplicates
72     #If all the structures have the same first six letters in their
73     #identifiers they are written to the output file of potential duplicates
74     if all(x == x[0] for x in nomenclature):
75         pin1 = solvated_hydrates_hash.index(identifiers[0])
76         new_lists_writer.write(solvated_hydrates_entry[pin1])
77         for f in range(1, len(identifiers)):
78             pin2 = solvated_hydrates_hash.index(identifiers[f])
79             new_lists_writer.write(solvated_hydrates_entry[pin2])
80             new_lists_writer.write(solvated_hydrates_entry[pin1])
81     else:
82         #If not, they are iterated through to check for subsets of
83         #structures with the same first six letters in their identifier
84         remove_identifiers = []
85         while len(identifiers) != len(remove_identifiers):
86             for h in range(len(identifiers)):
87                 if identifiers[h] not in remove_identifiers:
88                     remove_identifiers.append(identifiers[h])
89                     if nomenclature.count(nomenclature[h]) > 1:
90                         pin3 = solvated_hydrates_hash.index(identifiers[h])
91                         new_lists_writer.write(solvated_hydrates_entry[pin3])
92                         #Any cases where at least two structures share the
93                         #same first six letters are written to the output
94                         #file of potential duplicates
95                         for i in range(len(identifiers)):
96                             if identifiers[i] not in remove_identifiers:
97                                 if nomenclature[h] == nomenclature[i]:
98                                     remove_identifiers.append(identifiers[i])
99                                     pin4 =
100                                     solvated_hydrates_hash.index(identifiers[i])
101                                     new_lists_writer.write(solvated_hydrates_
102                                         entry[pin4])
103                                     new_lists_writer.write(solvated_hydrates_entry[pin3])
104
105     #Each set of structures that has been analyzed is removed from the list
106     #of structures before the next set is tested
107     for j in range(len(identifiers)):
108         pin5 = solvated_hydrates_hash.index(identifiers[j])
109         solvated_hydrates_hash.remove(solvated_hydrates_hash[pin5])
100         solvated_hydrates_entry.remove(solvated_hydrates_entry[pin5])
101
102
103
104
105
106 if __name__ == '__main__':
107     r = Runner()
108     r.run()
109

```



```

hash('petrol'):'CCCCCC', hash('n-alkane'):'CCCCC',
hash('dichloromethanal'):'ClC(Cl)=O',
hash('1-fluoro-4-methylbenzene'):'Cc1ccc(F)cc1',
hash('1,2,3-trimethoxybenzene'):'COc1cccc(OC)c1OC',
hash('1,3,5-trimethoxybenzene'):'COc1cc(OC)cc(OC)c1',
hash('chloroform'):'C1C(Cl)Cl', hash('methanol'):'CO',
hash('hexane'):'CCCCCC', hash('ethanol'):'CCO',
hash('acetonitrile'):'CC#N', hash('toluene'):'Cc1ccccc1',
hash('tetrahydrofuran'):'C1CCOC1', hash('acetone'):'CC(C)=O',
hash('dichloromethane'):'ClCCl', hash('1,2-dichloroethane'):'ClCCl',
hash('dimethylsulfoxide'):'CS(C)=O', hash('1,4-dioxane'):'C1COCCO1',
hash('1,2-dimethoxyethane'):'COCCOC', hash('o-xylene'):'Cc1cccc1c',
hash('propanol'):'CCO', hash('dioxane'):'C1COCCO1',
hash('n-pentane'):'CCCCCC', hash('4-cyanopyridine'):'N#Cc1ccncc1',
hash("4,4'-bipyridine"):'c1cc(ccn1)clccncc1', hash('isopropanol'):'CC(C)O',
hash('1-butanol'):'CCCCO', hash('nitromethane'):'CN(=O)=O',
hash('n-hexane'):'CCCCCC', hash('dimethylformamide'):'CN(C)C=O',
hash('N,N-dimethylformamide'):'CN(C)C=O', hash('cyclohexane'):'C1CCCCC1',
hash('hydroquinone'):'Oc1ccc(O)cc1', hash('benzene'):'clcccc1',
hash('octane'):'CCCCCC', hash('cyclohexanone'):'O=C1CCCCC1',
hash('isobutanol'):'CC(C)CO', hash('p-xylene'):'Cc1ccc(C)cc1',
hash('pentane'):'CCCCCC', hash('pyridine'):'clccncc1',
hash('t-butanol'):'CC(C)(C)O', hash('m-xylene'):'Cc1cccc(C)c1',
hash('phenol'):'Oc1ccccc1', hash('o-cresol'):'Cc1ccccc1O',
hash('p-cresol'):'Cc1ccc(O)cc1',
hash('2,4,6-trimethylpyridine'):'Cc1cc(C)nc(C)c1',
hash('1-methylpyrrolidin-2-one'):'CN1CCCC1=O',
hash('1,1,2,2,3,3,4,4,5,5,6,6,7,7,8,8-hexadecafluoro-1,8-diiodooctane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)I',
hash('trifluoroethanol'):'OCC(F)(F)F', hash('dichloroethane'):'ClCCl',
hash('1,1,2,2-tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('2-propanol'):'CC(C)O', hash('chlorobenzene'):'Cc1cccc1',
hash('naphthalene'):'clccc2ccccc2c1', hash('propan-2-ol'):'CC(C)O',
hash('benzonitrile'):'N#Cc1ccccc1', hash('nitrobenzene'):'O=N(=O)c1ccccc1',
hash('dodecane'):'CCCCCCCCCCCC',
hash('1,2-dichlorobenzene'):'Clc1ccccc1Cl', hash('n-heptane'):'CCCCCC',
hash('methanesulfonate'):'CS(=O)(=O)[O-]',
hash('3-methylbutan-1-ol'):'CC(C)CCO',
hash('1,1,2-trichloroethane'):'ClCC(Cl)Cl', hash('n-butane'):'CCCC',
hash('dichlorine'):'ClCl', hash('o-dichlorobenzene'):'Clc1ccccc1Cl',
hash('N,N-dimethylacetamide'):'CN(C)C(C)=O',
hash('4-methylaniline'):'Cc1ccc(N)cc1', hash('morpholine'):'C1COCCN1',
hash('1,3-dimethylimidazolidin-2-one'):'CN1CCN(C)C1=O',
hash('1,3-dimethoxybenzene'):'COc1cccc(OC)c1', hash('n-propanol'):'CCCO',
hash('mesitylene'):'Cc1cc(C)cc(C)c1',
hash('acetylacetone'):'CC(=O)CC(C)=O', hash('diethylether'):'CCOCC',
hash('hydrazine'):'NN', hash('dimethylamine'):'CNC',
hash('1,2-difluorobenzene'):'Fc1ccccc1F', hash('cycloheptane'):'C1CCCCC1',
hash('formamide'):'NC=O', hash('2-butanol'):'CCC(C)O',
hash('ethylacetate'):'CCOC(C)=O', hash('butan-1-ol'):'CCCCO',
hash('bromomethane'):'CBr', hash('piperazine'):'C1CNCCN1',
hash('n-nonane'):'CCCCCCCCC',
hash('1,2,4-trimethoxybenzene'):'COc1ccc(OC)c(OC)c1',
hash('1,2-diaminoethane'):'NCCN', hash('n-butanol'):'CCCCO',
hash('N-methylformamide'):'CNC=O', hash('pyrrolidine'):'C1CCNCl',
hash('actone'):'CC(C)=O', hash('n-octane'):'CCCCCC',
hash('('S)-2-methyl-1-butanol'):'CCC(C)CO',
hash('N,N-diethylformamide'):'CCN(CC)C=O',
hash('tetracyanoethylene'):'N#CC(C#N)=C(C#N)C#N',
hash('fluorobenzene'):'Fc1ccccc1', hash('heptane'):'CCCCCC',
hash('tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('methylcyclohexane'):'CC1CCCCC1',
hash('tetrachloromethane'):'ClC(Cl)(Cl)Cl', hash('1-propanol'):'CCCO',
hash('acetamide'):'CC(N)=O', hash('propane'):'CCC',
hash('diethylamine'):'CCNCC',
hash('hexafluorophosphate'):'F[P-](F)(F)(F)(F)F',
hash('hexafluoropropan-2-ol'):'OC(C(F)(F)F)C(F)(F)F',

```

```

hash('1-chloro-4-methylbenzene'):'Cc1ccc(Cl)cc1',
hash('acetaldehyde'):'CC=O', hash('propiononitrile'):'CCC#N',
hash('bromobenzene'):'Brclcccc1'}
49
deuterated_dictionary = {hash('dideutero-dichloromethane'):'ClCCl',
hash('perdeutero-toluene'):'Cclcccc1', hash('octadeutero furan'):'C1CCOC1',
hash('deutero-ethanol'):'CCO', hash('deuterochloroform'):'C1C(Cl)Cl',
hash('deutero-chloroform'):'C1C(Cl)Cl', hash('deuterobenzene'):'c1cccc1',
hash('deuteromethanol'):'CO', hash('hexadeutero-benzene'):'c1cccc1',
hash('perdeuteroacetone'):'CC(C)=O',
hash('dideuterodichloromethane'):'ClCCl', hash('perdeuteromethanol'):'CO',
hash('hexadeutero benzene'):'c1cccc1'}
50
acids_dictionary = {hash('trifluoroacetic'):'OC(=O)C(F)(F)F',
hash('acetic'):'CC(O)=O', hash('formic'):'OC=O',
hash('sulfuric'):'OS(O)(=O)=O', hash('maleic'):'OC(=O)C=CC(O)=O',
hash('benzoic'):'OC(=O)c1cccc1', hash('succinic'):'OC(=O)CCC(O)=O'}
51
ethers_dictionary = {'t-butyl methyl ether':'COC(C)(C)C',
hash('di-isopropyl'):'CC(C)OC(C)C', hash('dipropyl'):'CCCOCCCC',
hash('diethyl'):'CCOCC', hash('diisopropyl'):'CC(C)OC(C)C',
hash('isopropyl'):'CC(C)OC(C)C', hash('t-butyl'):'COC(C)(C)C',
hash('acetic'):'CCOC(C)=O'}
52
acetates_dictionary = {hash('ethyl'):'CCOC(C)=O',
hash('amyl'):'CCCCOC(C)=O', hash('methyl'):'COC(C)=O',
hash('butyl'):'CCCCOC(C)=O', hash('isobutyl'):'CC(C)COC(C)=O'}
53
glycol_dictionary = {hash('ethylene'):'OCCO'}
54
oxalates_dictionary = {hash('dipropyl'):'CCCO(=O)C(=O)OCCC'}
55
acid_neutral_dictionary = {hash('oxalic'):'OC(=O)C(O)=O',
hash('nitric'):'ON(=O)=O', hash('trans-but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('fumaric'):'OC(=O)C=CC(O)=O', hash('phosphonic'):'OP(O)=O',
hash('perchloric'):'O[Cl](=O)(=O)=O', hash('phosphoric'):'OP(O)(O)=O',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(O)=O',
hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(O)=O)c(cc1C(O)=O)C(O)=O',
hash('trifluoroacetic'):'OC(=O)C(F)(F)F'}
56
acid_charged_dictionary = {hash('oxalic'):'OC(=O)C(=O)[O-]',
hash('nitric'):'O=N(=O)[O-]', hash('trans-but-2-enedioic'):'OC(=O)C=CC(=O)[O-]',
hash('but-2-enedioic'):'OC(=O)C=CC(=O)[O-]', hash('fumaric'):'OC(=O)C=CC(=O)[O-]', hash('phosphonic'):'OP(=O)[O-]',
hash('perchloric'):'O=[Cl](=O)(=O)[O-]', hash('phosphoric'):'OP(O)(=O)[O-]',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(=O)[O-]', hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(=O)[O-])c(cc1C(=O)[O-])C(O)=O',
hash('trifluoroacetatic'):'FC(F)(F)C(=O)[O-]'}
57
58
partial_name = [hash('ketone'), hash('propanoate'), hash('carbonate'),
hash('formate'), hash('iodine'), hash('tetrachloride'), hash('sulfoxide'),
hash('dichloride'), hash('chloride'), hash('benzoate'), hash('disulfide'),
hash('oxide')]
59
full_solvate_name_dictionary = {'methyl isobutyl ketone':'CC(C)CC(C)=O',
'deuterium oxide':'O', 'ethyl propanoate':'CCOC(=O)CC', 'dimethyl
carbonate':'COC(=O)OC', 'ethyl formate':'CCOC=O', '12]cycloparaphenylene
iodine':'II', 'carbon tetrachloride':'C1C(Cl)(Cl)Cl', 'dimethyl
sulfoxide':'CS(C)=O', 'methylene dichloride':'ClCCl', 'methylene
chloride':'ClCCl', 'ethyl benzoate':'CCOC(=O)c1cccc1', 'carbon
disulfide':'S=C=S'}
60
61
#There are ten instances of two distinct molecule SMILES strings giving the
same hash value
62
#They are accounted for here in this dictionary
63
#Any time SMILES strings are hashed, the dictionary is checked so the
following ten SMILES strings are represented by the numbers 1-10 rather
than repeat hash values
64
duplicate_smiles_dictionary =
{'CC(C)(C)c1c(O)c(O)c(C1)c2[NH+]3CCN(CC3)c12':1,
'CC1=C(C(CC2=NS(=O)(=O)c3cccc23)c2cccc2C)C(=O)N(N1)c1cccc1':2,
'OC(=O)COc1cccc2cccc12':3, 'Cc1c[nH+]cc(C)c1':4,
'Brclcccc(NC(=O)c2cccn2)cc1':5, 'CSC(N)=NNC(C)=CC(=O)c1cccc1':6,

```

```

'CC1(CO1)C(O)(CBr)CBr':7, 'OCC12CCCC(C1)(C(O)=O)C1(CCSCC1)O2':8,
'COC1CCC2CC3[NH2+]CCC3(C)c2c1':9,
'CC(=O)C(P(c1cccc1)c1cccc1)=C(C)Nc1c(F)cccc1F':10}

65
66     solvated_hydrates_entry = []
67     solvated_hydrates_hash = []
68     solvated_hydrates_SMILES = []

69
70     entry1 = entry_reader1
71     for z in range(len(entry1)):
72         solvated_hydrates_entry.append(entry1[z])
73         solvated_hydrates_hash.append(hash(entry1[z].identifier))

74
75     #Each hydrate structure is checked for missing solvent SMILES strings
76     #The chemical name of each hydrate structure is checked for in each solvent
77     #dictionary
78     #Any match to the solvent dictionary checks the entry SMILES string for the
79     #corresponding solvent SMILES string
80     entry1 = entry_reader1
81     tag = 12
82     for a in range(len(entry1)):
83         component_SMILES = []
84         entry_SMILES = entry1[a].molecule.smiles.split('.')
85         for b in range(len(entry_SMILES)):
86             if entry_SMILES[b] not in duplicate_smiles_dictionary:
87                 component_SMILES.append(hash(entry_SMILES[b]))
88             else:
89
90                 component_SMILES.append(duplicate_smiles_dictionary.get(entry_SMILES[b]))
91         component_SMILES = list(set(component_SMILES))
92         #Some hydrate structures require manual corrections to their SMILES
93         #string
94         #Example is xylene, orientation of the methyl group cannot be
95         #determined from chemical name or formula when SMILES string is missing
96         if entry1[a].identifier == 'BULVEL':
97             if hash('[O-][n+]1cccc1') not in component_SMILES:
98                 component_SMILES.append(hash('[O-][n+]1cccc1'))
99             if hash('[O]n1cccc1') in component_SMILES:
100                 component_SMILES.remove(hash('[O]n1cccc1'))
101         elif entry1[a].identifier == 'CEHKOS':
102             component_SMILES.append(hash('II'))
103         elif entry1[a].identifier == 'GOTJEE' or entry1[a].identifier ==
104             'HASVEF' or entry1[a].identifier == 'MTFBTZ10':
105             component_SMILES.append(tag)
106             tag += 1
107         elif entry1[a].identifier == 'JUHJIF':
108             if hash('BrBr') not in component_SMILES:
109                 component_SMILES.append(hash('BrBr'))
110             if hash('[Br]') in component_SMILES:
111                 component_SMILES.remove(hash('[Br]'))
112         elif entry1[a].identifier == 'LOKXUF':
113             if hash('OC(=O)C(F)(F)F') not in component_SMILES:
114                 component_SMILES.append(hash('OC(=O)C(F)(F)F'))
115             if hash('[O]C(=O)C(F)(F)F') in component_SMILES:
116                 component_SMILES.remove(hash('[O]C(=O)C(F)(F)F'))
117         elif entry1[a].identifier == 'OJUNEO':
118             if hash('ClCl') not in component_SMILES:
119                 component_SMILES.append(hash('ClCl'))
120             if hash('[Cl]') in component_SMILES:
121                 component_SMILES.remove(hash('[Cl]'))
122         elif entry1[a].identifier == 'ZILFILM':
123             component_SMILES.append(hash('CCO'))
124         else:
125             if entry1[a].chemical_name != None:
126                 pieces = entry1[a].chemical_name.split(' ')
127                 formula_pieces = entry1[a].formula.split(',')

```

```

122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
for c in range(len(pieces)):
    if hash(pieces[c]) == hash('xylene'):
        if entry1[a].identifier == 'QOWNEV':
            if hash('Cclcccc1C') not in component_SMILES:
                component_SMILES.append(hash('Cclcccc1C'))
        elif entry1[a].identifier == 'TEYCEF' or
entry1[a].identifier == 'TEYCEF01':
            if hash('Cclcccc(C)c1') not in component_SMILES:
                component_SMILES.append(hash('Cclcccc(C)c1'))
        elif entry1[a].identifier == 'MAMNAR':
            if hash('Cclccc(C)cc1') not in component_SMILES:
                component_SMILES.append(hash('Cclccc(C)cc1'))
    elif hash(pieces[c]) == hash('hydrochloride') or
hash(pieces[c]) == hash('bis(hydrochloride)'):
        for z in range(len(formula_pieces)):
            if formula_pieces[z] == 'Cl1 1-' or '(Cl1 1-)' in
formula_pieces[z]:
                if hash('[Cl-]') not in component_SMILES:
                    component_SMILES.append(hash('[Cl-]'))
            if formula_pieces[z] == 'H1 Cl1' or '(H1 Cl1)' in
formula_pieces[z]:
                if hash('Cl') not in component_SMILES:
                    component_SMILES.append(hash('Cl'))
    elif hash(pieces[c]) == hash('unidentified') or
hash(pieces[c]) == hash('unknown'):
        component_SMILES.append(tag)
        tag += 1
    elif hash(pieces[c]) == hash('glycol'):
        if hash(pieces[c-1]) in glycol_dictionary:
            if hash(glycol_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
                component_SMILES.append(hash(glycol_dictionary.get(hash(pieces[c-1]))))
    elif hash(pieces[c]) == hash('oxalate'):
        if hash(pieces[c-1]) in oxalates_dictionary:
            if hash(oxalates_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
                component_SMILES.append(hash(oxalates_dictionary.get(hash(pieces[c-1]))))
    elif hash(pieces[c]) == hash('acetate'):
        if hash(pieces[c-1]) in acetates_dictionary:
            if hash(acetates_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
                component_SMILES.append(hash(acetates_dictionary.get(hash(pieces[c-1]))))
    elif hash(pieces[c]) == hash('ether'):
        if hash(pieces[c-1]) == hash('petroleum'):
            component_SMILES.append(11)
        elif pieces[c-2] + ' ' + pieces[c-1] + ' ' + pieces[c] in ethers_dictionary:
            if hash(ethers_dictionary.get(pieces[c-2] + ' ' +
pieces[c-1] + ' ' + pieces[c])) not in component_SMILES:
                component_SMILES.append(hash(ethers_dictionary.get(pieces[c-2] + ' ' +
pieces[c-1] + ' ' + pieces[c])))
    elif hash(pieces[c-1]) in ethers_dictionary:
        if hash(ethers_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
            component_SMILES.append(hash(ethers_dictionary.get(hash(pieces[c-1]))))
    elif hash(pieces[c]) in partial_name:

```

```

166     if pieces[c] == 'ketone':
167         if pieces[c-2] + ' ' + pieces[c-1] + ' ' +
168             pieces[c] in full_solvate_name_dictionary:
169             if
170                 hash(full_solvate_name_dictionary.get(pieces[c-2]
171                     + ' ' + pieces[c-1] + ' ' + pieces[c])) not in
172                     component_SMILES:
173
174             component_SMILES.append(hash(full_solvate_name_
175                 dictionary.get(pieces[c-2] + ' ' +
176                     pieces[c-1] + ' ' + pieces[c])))
177
178     else:
179         if pieces[c-1] + ' ' + pieces[c] in
180             full_solvate_name_dictionary:
181             if
182                 hash(full_solvate_name_dictionary.get(pieces[c-1]
183                     + ' ' + pieces[c])) not in component_SMILES:
184
185             component_SMILES.append(hash(full_solvate_name_
186                 dictionary.get(pieces[c-1] + ' ' +
187                     pieces[c])))
188
189     elif hash(pieces[c]) == hash('acid'):
190         if hash(pieces[c-1]) in acids_dictionary:
191             if hash(acids_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
192
193             component_SMILES.append(hash(acids_dictionary.get(
194                 hash(pieces[c-1]))))
195
196     elif hash(pieces[c-1]) in acid_neutral_dictionary:
197         if hash(pieces[c-1]) == hash('oxalic'):
198             for z in range(len(formula_pieces)):
199                 if formula_pieces[z] == 'C2 H2 O4' or '(C2
200                     H2 O4)' in formula_pieces[z]:
201                     if
202                         hash(acid_neutral_dictionary.get(hash(pie
203                             ces[c-1]))) not in component_SMILES:
204
205             component_SMILES.append(hash(acid_neu
206                 tral_dictionary.get(hash(pieces[c-1]))))
207
208         if formula_pieces[z] == 'C2 H1 O4 1-' or
209             '(C2 H1 O4 1-)' in formula_pieces[z]:
210             if
211                 hash(acid_charged_dictionary.get(hash(pie
212                     ces[c-1]))) not in component_SMILES:
213
214             component_SMILES.append(hash(acid_cha
215                 rged_dictionary.get(hash(pieces[c-1]))))
216
217     elif hash(pieces[c-1]) == hash('nitric'):
218         for z in range(len(formula_pieces)):
219             if formula_pieces[z] == 'H1 N1 O3' or '(H1
220                 N1 O3)' in formula_pieces[z]:
221                 if
222                     hash(acid_neutral_dictionary.get(hash(pie
223                         ces[c-1]))) not in component_SMILES:
224
225             component_SMILES.append(hash(acid_neu
226                 tral_dictionary.get(hash(pieces[c-1]))))
227
228         if formula_pieces[z] == 'N1 O3 1-' or '(N1
229             O3 1-)' in formula_pieces[z]:
230             if
231                 hash(acid_charged_dictionary.get(hash(pie
232                     ces[c-1]))) not in component_SMILES:
233
234             component_SMILES.append(hash(acid_cha
235                 rged_dictionary.get(hash(pieces[c-1]))))

```

```

195                                         rged_dictionary.get(hash(pieces[c-1])
196                                         )))
197                                         elif hash(pieces[c-1]) ==
198                                         hash('trans-but-2-enedioic') or hash(pieces[c-1])
199                                         == hash('but-2-enedioic') or hash(pieces[c-1]) ==
200                                         hash('fumaric'):
201                                         for z in range(len(formula_pieces)):
202                                         if formula_pieces[z] == 'C4 H4 O4' or '(C4
203                                         H4 O4)' in formula_pieces[z]:
204                                         if
205                                         hash(acid_neutral_dictionary.get(hash(pie
206                                         ces[c-1]))) not in component_SMILES:
207                                         component_SMILES.append(hash(acid_neu
208                                         tral_dictionary.get(hash(pieces[c-1])))
209                                         )))
210                                         elif hash(pieces[c-1]) == hash('phosphonic'):
211                                         for z in range(len(formula_pieces)):
212                                         if formula_pieces[z] == 'H3 O3 P1' or '(H3
213                                         O3 P1)' in formula_pieces[z]:
214                                         if
215                                         hash(acid_neutral_dictionary.get(hash(pie
216                                         ces[c-1]))) not in component_SMILES:
217                                         component_SMILES.append(hash(acid_neu
218                                         tral_dictionary.get(hash(pieces[c-1])))
219                                         )))
220                                         elif hash(pieces[c-1]) == hash('perchloric'):

```

```

221             if formula_pieces[z] == 'H3 O4 P1' or '(H3
222                                         O4 P1)' in formula_pieces[z]:
223                 if
224                     hash(acid_neutral_dictionary.get(hash(pie
225                                         ces[c-1]))) not in component_SMILES:
226                         component_SMILES.append(hash(acid_neu
227                                         tral_dictionary.get(hash(pieces[c-1]))
228                                         )))
229             if formula_pieces[z] == 'H2 O4 P1 1-' or
230                 '(H2 O4 P1 1-)' in formula_pieces[z]:
231                 if
232                     hash(acid_charged_dictionary.get(hash(pie
233                                         ces[c-1]))) not in component_SMILES:
234                         component_SMILES.append(hash(acid_cha
235                                         rged_dictionary.get(hash(pieces[c-1]))
236                                         )))
237             elif hash(pieces[c-1]) ==
238                 hash('pyridine-2,6-dicarboxylic'):
239                 for z in range(len(formula_pieces)):
240                     if formula_pieces[z] == 'C7 H5 N1 O4' or
241                         '(C7 H5 N1 O4)' in formula_pieces[z]:
242                         if
243                             hash(acid_neutral_dictionary.get(hash(pie
244                                         ces[c-1]))) not in component_SMILES:
245                                 component_SMILES.append(hash(acid_neu
246                                         tral_dictionary.get(hash(pieces[c-1]))
247                                         )))
248             elif hash(pieces[c-1]) ==
249                 hash('benzene-1,2,4,5-tetracarboxylic'):
250                 for z in range(len(formula_pieces)):
251                     if formula_pieces[z] == 'C10 H6 O8' or
252                         '(C10 H6 O8)' in formula_pieces[z]:
253                         if
254                             hash(acid_neutral_dictionary.get(hash(pie
255                                         ces[c-1]))) not in component_SMILES:
256                                 component_SMILES.append(hash(acid_neu
257                                         tral_dictionary.get(hash(pieces[c-1]))
258                                         )))
259             elif hash(pieces[c-1]) ==
260                 hash('trifluoroacetatic'):
261                 for z in range(len(formula_pieces)):
262                     if formula_pieces[z] == 'C2 H1 F3 O2' or
263                         '(C2 H1 F3 O2)' in formula_pieces[z]:
264                         if
265                             hash(acid_neutral_dictionary.get(hash(pie
266                                         ces[c-1]))) not in component_SMILES:

```

```

248                     component_SMILES.append(hash(acid_neu
249                     tral_dictionary.get(hash(pieces[c-1]))
250                     )))
251             if formula_pieces[z] == 'C2 F3 O2 1-' or
252                 '(C2 F3 O2 1-)' in formula_pieces[z]:
253                 if
254                     hash(acid_charged_dictionary.get(hash(pie
255                     ces[c-1]))) not in component_SMILES:
256                         component_SMILES.append(hash(acid_cha
257                         rged_dictionary.get(hash(pieces[c-1]))
258                         )))
259             elif hash(pieces[c]) in deuterated_dictionary:
260                 if hash(deuterated_dictionary.get(hash(pieces[c]))) not
261                     in component_SMILES:
262                         component_SMILES.append(hash(deuterated_dictionar
263                         y.get(hash(pieces[c]))))
264             if hash(pieces[c]) == hash('dideutero-dichloromethane')
265             or hash(pieces[c]) == hash('dideuterodichloromethane'):
266                 if hash('Cl[C]Cl') in component_SMILES:
267                     component_SMILES.remove(hash('Cl[C]Cl'))
268             elif hash(pieces[c]) == hash('perdeutero-toluene'):
269                 if hash('[C]c1[c][c][c][c][c]1') in component_SMILES:
270                     component_SMILES.remove(hash('[C]c1[c][c][c][c]
271                     [c]1'))
272             elif hash(pieces[c]) == hash('deutero-ethanol'):
273                 if hash('[C][C][O]') in component_SMILES:
274                     component_SMILES.remove(hash('[C][C][O]'))
275             elif hash(pieces[c]) == hash('deuterochloroform') or
276                 hash(pieces[c]) == hash('deutero-chloroform'):
277                 if hash('Cl[C](Cl)Cl') in component_SMILES:
278                     component_SMILES.remove(hash('Cl[C](Cl)Cl'))
279             elif hash(pieces[c]) == hash('hexadeutero-benzene') or
280                 hash(pieces[c]) == hash('deuterobenzene') or
281                 hash(pieces[c]) == hash('hexadeuterobenzene'):
282                 if hash('[c]1[c][c][c][c][c]1') in component_SMILES:
283                     component_SMILES.remove(hash('[c]1[c][c][c][c]
284                     [c]1'))
285             elif hash(pieces[c]) == hash('deuteromethanol') or
286                 hash(pieces[c]) == hash('perdeuteromethanol'):
287                 if hash('[C][O]') in component_SMILES:
288                     component_SMILES.remove(hash('[C][O]'))
289             elif hash(pieces[c]) in solvents_dictionary:
290                 if hash(solvents_dictionary.get(hash(pieces[c]))) not
291                     in component_SMILES:
292                         component_SMILES.append(hash(solvents_dictionary.get(
293                             hash(pieces[c]))))
294
295             if hash('O') in component_SMILES:
296                 component_SMILES.remove(hash('O'))
297             if hash('[O]') in component_SMILES:
298                 component_SMILES.remove(hash('[O]'))
299
300             component_SMILES = sorted(component_SMILES)
301             solvated_hydrates_SMILES.append(component_SMILES)
302
303
304             print len(solvated_hydrates_SMILES)
305             print len(solvated_hydrates_hash)
306             print len(solvated_hydrates_entry)
307
308             while len(solvated_hydrates_SMILES) != 0:

```

```

290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
    print len(solvated_hydrates_SMILES)
    identifiers = []
    number = []
    identifiers.append(solvated_hydrates_hash[0])

    number.append(solvated_hydrates_entry[0].crystal.spacegroup_number_and_se-
tting[0])
    #The first hydrate structure entry SMILES string is compared against
    all the other hydrate structure entry SMILES strings
    for d in range(1,len(solvated_hydrates_SMILES)):
        #Hydrate structures that match the entry SMILES string are put into
        a list with the first hydrate structure
        #The spacegroup number of all matches are recorded in another list
        if solvated_hydrates_SMILES[0] == solvated_hydrates_SMILES[d]:
            identifiers.append(solvated_hydrates_hash[d])

            number.append(solvated_hydrates_entry[d].crystal.spacegroup_numbe-
r_and_setting[0])

    if len(identifiers) != 1:
        #The water stoichiometry of each match is determined the same way
        it was for hydrate structures without entry SMILES strings
        waters = []
        for g in range(len(identifiers)):
            pin = solvated_hydrates_hash.index(identifiers[g])
            title = solvated_hydrates_entry[pin].formula
            pieces = title.split(',')
            water_molecule = []
            for e in range(len(pieces)):
                if '(H2 O1)' in pieces[e] or pieces[e] == 'H2 O1' or '(D2
O1)' in pieces[e] or pieces[e] == 'D2 O1':
                    water_molecule.append(pieces[e])
            if water_molecule == [] or len(water_molecule) > 1:
                waters.append('undefined')
            elif water_molecule[0] == 'H2 O1' or water_molecule[0] == 'D2
O1':
                waters.append(1)
            else:
                water = str(water_molecule[0])
                end = water.index('(')
                water = water[:end]
                if 'x' in water:
                    waters.append('x')
                elif 'n' in water:
                    waters.append('n')
                elif water.isdigit() == False:
                    waters.append(Decimal(water))
                else:
                    waters.append(int(water))
        #Once again, any cases where the water stoichiometry or spacegroup
        number match undergo further analysis
        if any(waters.count(x) > 1 for x in waters) or all(x == number[0]
for x in number) or any(number.count(y) > 1 for y in number):
            identifiers_to_remove = []
            while len(identifiers) != len(identifiers_to_remove):
                for h in range(len(identifiers)):
                    if identifiers[h] not in identifiers_to_remove:
                        identifiers_to_remove.append(identifiers[h])
                    if waters.count(waters[h]) > 1:
                        match = 0
                        for i in range(len(identifiers)):
                            if identifiers[i] not in
                                identifiers_to_remove:
                                if waters[h] == waters[i]:
                                    if number[h] == number[i]:
                                        if match == 0:
                                            pin1 =

```

```

345                                         solvated_hydrates_hash.index(
346                                         identifiers[h])
347                                         duplicates_writer.write(solva-
348                                         ted_hydrates_entry[pin1])
349                                         pin2 =
350                                         solvated_hydrates_hash.index(iden-
351                                         tifiers[i])
352                                         duplicates_writer.write(solvated_
353                                         hydrates_entry[pin2])
354                                         match += 1
355                                         identifiers_to_remove.append(iden-
356                                         tifiers[i])
357                                         if match != 0:
358                                         duplicates_writer.write(solvated_hydrates_ent-
359                                         ry[pin1])
360                                         #The overall list of hydrates is updated to remove the set of
361                                         #structures that was analyzed
362                                         for j in range(len(identifiers)):
363                                         pin = solvated_hydrates_hash.index(identifiers[j])
364                                         solvated_hydrates_hash.remove(solvated_hydrates_hash[pin])
365                                         solvated_hydrates_entry.remove(solvated_hydrates_entry[pin])
366                                         solvated_hydrates_SMILES.remove(solvated_hydrates_SMILES[pin])
367                                         print 'solvated hydrates list is empty'
368                                         if __name__ == '__main__':
369                                         # This runs the script
370                                         r = Runner()
371                                         r.run()

```

```

1 #9 Find duplicate SMILES hydrate structures with packing similarity
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 import itertools
9 from decimal import Decimal
10 from ccdc.crystal import PackingSimilarity
11 import time
12
13 similarity_engine = PackingSimilarity()
14 similarity_engine.settings.ignore_hydrogen_positions = True
15
16 #This script looks at the potential duplicates identified in the previous script
17 #This script uses the packing similarity tool which gets stuck at certain entries as
18 #identified in a later comment
19 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_smiles_all_potential_duplicates.txt"
20
21 class Runner(argparse.ArgumentParser):
22
23     def __init__(self):
24         super(self.__class__, self).__init__(description=__doc__)
25         self.add_argument(
26             '-i', '--input', default=filepath1,
27             help='input database filepath')
28         self.add_argument(
29             '-o', '--output', default='hydrates_smiles_all_first_occurrence.gcd',
30             help='output file [hydrates_smiles_all_first_occurrence.gcd]')
31         self.add_argument(
32             '-m', '--maximum', default=0, type=int,
33             help='Maximum number of structures to find [all]')
34
35     args = self.parse_args()
36
37     self.args = args
38     self.settings = search.Search.Settings()
39     self.settings.max_hit_structures = self.args.maximum
40
41     def run(self):
42
43         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
44
45         with io.EntryWriter(self.args.output) as lone_wolf_writer:
46             with io.EntryWriter("hydrates_smiles_all_duplicates.gcd") as writer1:
47                 with io.EntryWriter("hydrates_smiles_all_manual_check.gcd") as writer2:
48
49                     solvated_hydrates_hash = []
50                     solvated_hydrates_entry = []
51
52                     #Split input file at entries: 1362, 1365, 2329, 2332
53                     #The following entries get stuck at analysis calculation: GUQBOJ,
54                     #NISNAF
55                     #Removing these potential duplicates of these entries from the
56                     #input file for a manual check before running the script should
57                     #prevent stalling
58                     for a in range(2332, len(entry_reader1)):
59                         solvated_hydrates_hash.append(hash(entry_reader1[a].identifier))
60                         solvated_hydrates_entry.append(entry_reader1[a])
61
62                     while len(solvated_hydrates_hash) != 0:
63                         print len(solvated_hydrates_hash)

```

```

63     identifiers = []
64     identifiers.append(solvated_hydrates_hash[0])
65     #Once again the set of potential duplicates is found by
66     #checking for the second occurrence of the starting identifier
67     for d in range(1,len(solvated_hydrates_hash)):
68         if identifiers[0] != solvated_hydrates_hash[d]:
69             identifiers.append(solvated_hydrates_hash[d])
70         else:
71             del solvated_hydrates_hash[d]
72             del solvated_hydrates_entry[d]
73             break
74     remove_identifiers = []
75     while len(identifiers) != len(remove_identifiers):
76         for s in range(len(identifiers)):
77             count = 0
78             if identifiers[s] not in remove_identifiers:
79                 pin1 = solvated_hydrates_hash.index(identifiers[s])
80                 reference = solvated_hydrates_entry[pin1].crystal
81                 remove_identifiers.append(identifiers[s])
82                 #Here each hydrate structure in the set of
83                 #potential duplicates is compared to subsequent
84                 #hydrate structures
85                 #The script iterates through the list so cases with
86                 #polymorphs are isolated
87                 #For example, if my list of matches has [A, B, C,
88                 #D] and A and C are duplicates and B and D are
89                 #duplicates
90                 #When A is checked for duplicates it will match C,
91                 #A and C will then be removed from the list so it is
92                 #now [B, D]
93                 #Then B will be checked for duplicates and D will
94                 #return as a match
95                 for f in range(len(identifiers)):
96                     if identifiers[f] not in remove_identifiers:
97                         pin2 =
98                         solvated_hydrates_hash.index(identifiers[f])
99                         sample =
100                         solvated_hydrates_entry[pin2].crystal
101                         #The packing similarity between two hydrate
102                         #structures is calculated
103                         analysis =
104                         similarity_engine.compare(reference, sample)
105                         #If the analysis returns None, the hydrate
106                         #structures are written to a list for a
107                         #manual check
108                         if analysis == None:
109                             if count == 0:
110
111                                 writer2.write(solvated_hydrates_entry
112 [pin1])
113                                 count += 1
114
115                                 writer2.write(solvated_hydrates_entry[pin
116 2])
117                                 remove_identifiers.append(identifiers[f])
118                             else:
119                                 #If the packing similarity returns 1
120                                 #(15 out of 15 packing points are the
121                                 #same) then the two hydrate structures
122                                 #are considered duplicates
123                                 if
124                                     analysis.nmatched_molecules/analysis.pack
125                                     ing_shell_size == 1:
126                                         #The first hydrate structure in the
127                                         #set is written to a file that
128                                         #contains the hydrate structures to
129                                         #keep

```

```

103                                         #The second hydrate structure in
104                                         the set is written to a file that
105                                         contains the hydrate structures to
106                                         remove
107                                         lone_wolf_writer.write(solvated_hydrates_entry[pin1])
108                                         writer1.write(solvated_hydrates_entry[pin2])
109                                         remove_identifiers.append(identifiers[f])
110                                         if count != 0:
111                                         writer2.write(solvated_hydrates_entry[pin1])
112                                         #The overall list of hydrate structures is updated to remove
113                                         those already analyzed
114                                         for j in range(len(identifiers)):
115                                         pin = solvated_hydrates_hash.index(identifiers[j])
116                                         solvated_hydrates_hash.remove(solvated_hydrates_hash[pin])
117                                         solvated_hydrates_entry.remove(solvated_hydrates_entry[pin])
118                                         print 'solvated hydrates list is empty'
119                                         if __name__ == '__main__':
120                                         # This runs the script
121                                         r = Runner()
122                                         r.run()

```

```

1 #10 Find potential duplicate none waterless forms with formula match
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 from decimal import Decimal
9
10 #This script is the same as the part1 script to find duplicates for hydrate structures
11 #without entry SMILES strings
12 #The only difference is that the water stoichiometry of each entry does not need to be
13 #determined
14 #In order for structures with the same formula to be considered potential duplicates
15 #there needs to be at least two structures with the same spacegroup number
16 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
17 Forms\complete_None_waterless_forms.txt"
18 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
19 Forms\complete_smiles_waterless_forms.txt"
20
21 class Runner(argparse.ArgumentParser):
22
23     def __init__(self):
24         super(self.__class__, self).__init__(description=__doc__)
25         self.add_argument(
26             '-i', '--input', default=filepath1,
27             help='input database filepath'
28         )
29         self.add_argument(
30             '-o', '--output',
31             default='waterless_forms_None_all_potential_duplicates.gcd',
32             help='output file [waterless_forms_None_all_potential_duplicates.gcd]'
33         )
34         self.add_argument(
35             '-m', '--maximum', default=0, type=int,
36             help='Maximum number of structures to find [all]'
37         )
38
39     def run(self):
40
41         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
42         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
43
44         with io.EntryWriter(self.args.output) as new_lists_writer:
45
46             None_formulas = []
47             formulas = []
48             positions = []
49             numbers = []
50             for a in range(len(entry_reader1)):
51                 component_formulas = entry_reader1[a].formula.split(',')
52                 pre_formulas = []
53                 for b in range(len(component_formulas)):
54                     if ' and ' in component_formulas[b]:
55                         start = component_formulas[b].index('(')
56                         end = component_formulas[b].index(')')
57                         component_formulas[b] = component_formulas[b][start+1:end]
58                         elements = component_formulas[b].split(' ')
59                         deuteriums = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'D']
60                         if deuteriums != []:
61                             if len(deuteriums) != 1:

```

```

62     print entry_reader1[a].identifier
63     hydrogens = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'H']
64     if hydrogens == []:
65         spot = elements.index(deuteriums[0])
66         deuteriums[0] = deuteriums[0].replace('D', 'H')
67         elements[spot] = deuteriums[0]
68     else:
69         if len(hydrogens) != 1:
70             print entry_reader1[a].identifier
71             spot = elements.index(deuteriums[0])
72             spot2 = elements.index(hydrogens[0])
73             stoichiometry1 = re.sub(r'[A-Z]+', '', deuteriums[0])
74             stoichiometry2 = re.sub(r'[A-Z]+', '', hydrogens[0])
75             total = int(stoichiometry1) + int(stoichiometry2)
76             hydrogens[0] = hydrogens[0].replace(stoichiometry2,
77             str(total))
78             elements[spot2] = hydrogens[0]
79             elements.remove(elements[spot])
80             component_formulas[b] = ' '.join(elements)
81             pre_formulas.append(hash(component_formulas[b]))
82             pre_formulas = list(set(pre_formulas))
83             pre_formulas = sorted(pre_formulas)
84             None_formulas.append(pre_formulas)
85             formulas.append(pre_formulas)
86             positions.append(a)
87             if hash(entry_reader1[a].identifier) == hash('KECYBU15'):
88                 numbers.append(87)
89             else:
90
91             numbers.append(entry_reader1[a].crystal.spacegroup_number_and_setting
92             [0])
93
94             for d in range(len(entry_reader2)):
95                 component_formulas = entry_reader2[d].formula.split(',')
96                 pre_formulas2 = []
97                 for e in range(len(component_formulas)):
98                     if ' and ' in component_formulas[e]:
99                         start = component_formulas[e].index('(')
100                        end = component_formulas[e].index(')')
101                        component_formulas[e] = component_formulas[e][start+1:end]
102                        elements = component_formulas[e].split(' ')
103                        deuteriums = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'D']
104                        if deuteriums != []:
105                            if len(deuteriums) != 1:
106                                print entry_reader2[d].identifier
107                                hydrogens = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'H']
108                                if hydrogens == []:
109                                    spot = elements.index(deuteriums[0])
110                                    deuteriums[0] = deuteriums[0].replace('D', 'H')
111                                    elements[spot] = deuteriums[0]
112                                else:
113                                    if len(hydrogens) != 1:
114                                        print entry_reader2[d].identifier
115                                        spot = elements.index(deuteriums[0])
116                                        spot2 = elements.index(hydrogens[0])
117                                        stoichiometry1 = re.sub(r'[A-Z]+', '', deuteriums[0])
118                                        stoichiometry2 = re.sub(r'[A-Z]+', '', hydrogens[0])
119                                        total = int(stoichiometry1) + int(stoichiometry2)
120                                        hydrogens[0] = hydrogens[0].replace(stoichiometry2,
121                                        str(total))
122                                        elements[spot2] = hydrogens[0]
123                                        elements.remove(elements[spot])
124                                        component_formulas[e] = ' '.join(elements)
125                                        pre_formulas2.append(hash(component_formulas[e]))
126                                        pre_formulas2 = list(set(pre_formulas2))

```

```

123     pre_formulas2 = sorted(pre_formulas2)
124     formulas.append(pre_formulas2)
125     positions.append(d+5450)
126     if hash(entry_reader2[d].identifier) == hash('MTYHFB03'):
127         numbers.append(29)
128     else:
129
130         numbers.append(entry_reader2[d].crystal.spacegroup_number_and_setting
131             [0])
132
133     while None_formulas != []:
134         print len(None_formulas)
135         formulas_to_remove = []
136         positions_to_write = []
137         numbers_to_compare = []
138         formulas_to_remove.append(formulas[0])
139         positions_to_write.append(positions[0])
140         numbers_to_compare.append(numbers[positions[0]])
141         for c in range(1, len(formulas)):
142             if formulas[0] == formulas[c]:
143                 positions_to_write.append(positions[c])
144                 numbers_to_compare.append(numbers[positions[c]])
145
146         if len(numbers_to_compare) != 1:
147             if all(x == numbers_to_compare[0] for x in numbers_to_compare) or
148                 any(numbers_to_compare.count(x) > 1 for x in numbers_to_compare):
149                 positions_to_remove = []
150                 while len(positions_to_write) != len(positions_to_remove):
151                     for j in range(len(positions_to_write)):
152                         if positions_to_write[j] not in positions_to_remove:
153                             positions_to_remove.append(positions_to_write[j])
154                             if positions_to_write[j] < 5450:
155                                 match = 0
156                                 for k in range(len(positions_to_write)):
157                                     if positions_to_write[k] not in
158                                         positions_to_remove:
159                                         if numbers_to_compare[j] ==
160                                         numbers_to_compare[k]:
161                                         if match == 0:
162                                             new_lists_writer.write(entry_read
163                                             er1[positions_to_write[j]])
164                                         if positions_to_write[k] < 5450:
165                                             new_lists_writer.write(entry_read
166                                             er1[positions_to_write[k]])
167                                         else:
168                                             new_lists_writer.write(entry_read
169                                             er2[positions_to_write[k]-5450])
170                                         match += 1
171                                         positions_to_remove.append(positions_
172                                         to_write[k])
173                                         if match != 0:
174                                             new_lists_writer.write(entry_reader1[position
175                                             s_to_write[j]])
176
177     None_formulas = [i for i in None_formulas if i not in formulas_to_remove]
178     formulas = [i for i in formulas if i not in formulas_to_remove]
179     positions = [j for j in positions if j not in positions_to_write]
180
181     if __name__ == '__main__':

```

```
175     r = Runner()  
176     r.run()  
177
```

```

1 #11 Find potential duplicate none waterless forms with same identifier prefix
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 from decimal import Decimal
9
10 #This script is the same as the part2 script to find duplicates for hydrate structures
11 #without entry SMILES strings
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_None_all_potential_duplicates.txt"
13
14 class Runner(argparse.ArgumentParser):
15
16     def __init__(self):
17         super(self.__class__, self).__init__(description=__doc__)
18         self.add_argument(
19             '-i', '--input', default=filepath1,
20             help='input database filepath')
21         self.add_argument(
22             '-o', '--output', default='waterless_forms_None_all_manual_check.gcd',
23             help='output file [waterless_forms_None_all_manual_check.gcd]')
24         self.add_argument(
25             '-m', '--maximum', default=0, type=int,
26             help='Maximum number of structures to find [all]')
27
28
29     args = self.parse_args()
30
31     self.args = args
32     self.settings = search.Search.Settings()
33     self.settings.max_hit_structures = self.args.maximum
34
35
36     def run(self):
37
38         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
39
40         with io.EntryWriter(self.args.output) as new_lists_writer:
41
42             duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
43             'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
44             'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
45
46             solvated_waterless_hash = []
47             solvated_waterless_entry = []
48             for a in range(len(entry_reader1)):
49                 if entry_reader1[a].identifier in duplicate_identifier_dictionary:
50
51                     solvated_waterless_hash.append(duplicate_identifier_dictionary.get(en
52                     try_reader1[a].identifier))
53                 else:
54                     solvated_waterless_hash.append(hash(entry_reader1[a].identifier))
55                     solvated_waterless_entry.append(entry_reader1[a])
56
57             while len(solvated_waterless_hash) != 0:
58                 print len(solvated_waterless_hash)
59                 identifiers = []
60                 identifiers.append(solvated_waterless_hash[0])
61                 for d in range(1,len(solvated_waterless_hash)):
62                     if identifiers[0] != solvated_waterless_hash[d]:
63                         identifiers.append(solvated_waterless_hash[d])
64                     else:
65                         del solvated_waterless_hash[d]

```

```
62     del solvated_waterless_entry[d]
63     break
64
65     nomenclature = []
66     for e in range(len(identifiers)):
67         pin = solvated_waterless_hash.index(identifiers[e])
68         nomenclature.append(solvated_waterless_entry[pin].identifier[:6])
69
70     if all(x == x[0] for x in nomenclature):
71         pin1 = solvated_waterless_hash.index(identifiers[0])
72         new_lists_writer.write(solvated_waterless_entry[pin1])
73         for f in range(1, len(identifiers)):
74             pin2 = solvated_waterless_hash.index(identifiers[f])
75             new_lists_writer.write(solvated_waterless_entry[pin2])
76             new_lists_writer.write(solvated_waterless_entry[pin1])
77     else:
78         remove_identifiers = []
79         while len(identifiers) != len(remove_identifiers):
80             for h in range(len(identifiers)):
81                 if identifiers[h] not in remove_identifiers:
82                     remove_identifiers.append(identifiers[h])
83                     if nomenclature.count(nomenclature[h]) > 1:
84                         pin3 = solvated_waterless_hash.index(identifiers[h])
85
86                         new_lists_writer.write(solvated_waterless_entry[pin3])
87                         )
88                         for i in range(len(identifiers)):
89                             if identifiers[i] not in remove_identifiers:
90                                 if nomenclature[h] == nomenclature[i]:
91                                     remove_identifiers.append(identifiers[i])
92                                     pin4 =
93                                     solvated_waterless_hash.index(identifiers
94                                     [i])
95
96                                     new_lists_writer.write(solvated_waterless
97                                     _entry[pin4])
98
99
100
101 if __name__ == '__main__':
102     r = Runner()
103     r.run()
```

```

1 #12 Find potential duplicate SMILES waterless forms with formula match
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script is different from the part1 script to find duplicates for hydrate
10 structures with entry SMILES strings
11 #The work that script does for hydrates had to be broken into two for anhydrous forms
12 since there are far more anhydrous forms with entry SMILES strings than hydrates with
13 SMILES strings
14 #This script carries out the first part which is checking the anhydrous forms for
15 formula matches
16 #This is done the same way as in the part1 script to find duplicates for anhydrous
17 forms without entry SMILES strings
18 #This was done due to the large amount of time required to check all the anhydrous
19 forms in the solvent dictionary
20 #This way time won't be wasted looking in the solvent dictionary for structures that do
21 not match the formula of any other anhydrous forms pulled from the CSD
22 #Structures that don't match the formula of any other structure cannot have a duplicate
23 in the list obtained for this work
24 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
25 Forms\complete_smiles_waterless_forms.txt"
26
27 class Runner(argparse.ArgumentParser):
28
29     def __init__(self):
30         super(self.__class__, self).__init__(description=__doc__)
31         self.add_argument(
32             '-i', '--input', default=filepath1,
33             help='input database filepath')
34
35         self.add_argument(
36             '-o', '--output', default='waterless_forms_smiles_all_repeated_formula.gcd',
37             help='output file [waterless_forms_smiles_all_repeated_formula.gcd]')
38
39         self.add_argument(
40             '-m', '--maximum', default=0, type=int,
41             help='Maximum number of structures to find [all]')
42
43     args = self.parse_args()
44
45     self.args = args
46     self.settings = search.Search.Settings()
47     self.settings.max_hit_structures = self.args.maximum
48
49     def run(self):
50
51         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
52
53         with io.EntryWriter(self.args.output) as duplicates_writer:
54
55             duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
56             'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
57             'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
58
59             formulas = []
60             positions = []
61             numbers = []
62             for a in range(len(entry_reader1)):
63                 print a
64                 component_formulas = entry_reader1[a].formula.split(',')
65                 pre_formulas = []
66                 for b in range(len(component_formulas)):
67
68

```

```

57     if '(' and ')' in component_formulas[b]:
58         start = component_formulas[b].index('(')
59         end = component_formulas[b].index(')')
60         component_formulas[b] = component_formulas[b][start+1:end]
61     elements = component_formulas[b].split(' ')
62     deuteriums = [i for i in elements if re.sub(r'[0-9]+', '', i) == 'D']
63     if deuteriums != []:
64         if len(deuteriums) != 1:
65             #These lines terminate the script and act as a safety for
66             #cases that are not accounted for
67             #With such a large set of structures, extraneous cases are
68             #hard to predict
69             #To save time, these lines are put in to catch extraneous
70             #cases, that way when there are not any time is not wasted
71             #accounting for them
72             print entry_reader1[a].identifier
73             print deuteriums[90]
74     hydrogens = [i for i in elements if re.sub(r'[0-9]+', '', i) ==
75     'H']
76     if hydrogens == []:
77         spot = elements.index(deuteriums[0])
78         deuteriums[0] = deuteriums[0].replace('D', 'H')
79         elements[spot] = deuteriums[0]
80     else:
81         if len(hydrogens) != 1:
82             print entry_reader1[a].identifier
83             print hydrogens[90]
84         spot = elements.index(deuteriums[0])
85         spot2 = elements.index(hydrogens[0])
86         stoichiometry1 = re.sub(r'[A-Z]+', '', deuteriums[0])
87         stoichiometry2 = re.sub(r'[A-Z]+', '', hydrogens[0])
88         total = int(stoichiometry1) + int(stoichiometry2)
89         hydrogens[0] = hydrogens[0].replace(stoichiometry2,
90         str(total))
91         elements[spot2] = hydrogens[0]
92         elements.remove(elements[spot])
93         component_formulas[b] = ' '.join(elements)
94         pre_formulas.append(hash(component_formulas[b]))
95     pre_formulas = list(set(pre_formulas))
96     pre_formulas = sorted(pre_formulas)
97     formulas.append(pre_formulas)
98     positions.append(a)
99     if hash(entry_reader1[a].identifier) == hash('KECYBU15'):
100        numbers.append(87)
101    elif hash(entry_reader1[a].identifier) == hash('MTYHFB03'):
102        numbers.append(29)
103    else:
104
105        numbers.append(entry_reader1[a].crystal.spacegroup_number_and_setting
106        [0])
107
108    while formulas != []:
109        print len(formulas)
110        formulas_to_remove = []
111        positions_to_write = []
112        numbers_to_compare = []
113        formulas_to_remove.append(formulas[0])
114        positions_to_write.append(positions[0])
115        numbers_to_compare.append(numbers[positions[0]])
116        for c in range(1, len(formulas)):
117            if formulas[0] == formulas[c]:
118                positions_to_write.append(positions[c])
119                numbers_to_compare.append(numbers[positions[c]])
120
121        if len(numbers_to_compare) != 1:
122            if all(x == numbers_to_compare[0] for x in numbers_to_compare) or
123            any(numbers_to_compare.count(x) > 1 for x in numbers_to_compare):

```

```

115     positions_to_remove = []
116     for h in range(len(positions_to_write)):
117         times1 = numbers_to_compare.count(numbers_to_compare[h])
118         if times1 == 1:
119             positions_to_remove.append(positions_to_write[h])
120     positions_to_check = [z for z in positions_to_write if z not in
121     positions_to_remove]
122     if positions_to_check != []:
123         for i in range(len(positions_to_check)):
124             duplicates_writer.write(entry_reader1[positions_to_check[i]])
125             duplicates_writer.write(entry_reader1[positions_to_check[0]])
126     formulas = [i for i in formulas if i not in formulas_to_remove]
127     positions = [j for j in positions if j not in positions_to_write]
128
129     print 'waterless formulas list is empty'
130
131
132 if __name__ == '__main__':
133     r = Runner()
134     r.run()
135

```



```

hash('methyl-cyclohexane'):'CC1CCCCC1',
hash('pyridine-3-carboxamide'):'NC(=O)c1cccnc1',
hash('3,7-dimethyl-3,7-dihydro-1H-purine-2,6-dione'):'CN1C=NC2=C1C(=O)NC(=O)N2C', hash('nitroethane'):'CCN(=O)=O',
hash('dimethylaminobenzaldehyde'):'CN(C)c1ccc(C=O)cc1',
hash('3-(1H-pyrrol-1-yl)propanenitrile'):'C1=CN(C=C1)CCC#N',
hash('trifluoromethanol'):'OC(F)(F)F', hash('methoxyethane'):'CCOC',
hash('methylanisole'):'COc1cccc1C', hash('xylene'):'Cc1ccc(C)cc1',
hash('petrol'):'CCCCCC', hash('n-alkane'):'CCCCC',
hash('dichloromethanal'):'ClC(Cl)=O',
hash('1-fluoro-4-methylbenzene'):'Cc1ccc(F)cc1',
hash('1,2,3-trimethoxybenzene'):'COc1cccc(OC)c1OC',
hash('1,3,5-trimethoxybenzene'):'COc1cc(OC)cc(OC)c1',
hash('chloroform'):'C1C(Cl)Cl', hash('methanol'):'CO',
hash('hexane'):'CCCCCC', hash('ethanol'):'CCO',
hash('acetonitrile'):'CC#N', hash('toluene'):'Cc1cccc1',
hash('tetrahydrofuran'):'C1CCOC1', hash('acetone'):'CC(C)=O',
hash('dichloromethane'):'ClCCl', hash('1,2-dichloroethane'):'ClCCCCl',
hash('dimethylsulfoxide'):'CS(C)=O', hash('1,4-dioxane'):'C1COCCO1',
hash('1,2-dimethoxyethane'):'COCCOC', hash('o-xylene'):'Cc1cccc1c',
hash('propanol'):'CCO', hash('dioxane'):'C1COCCO1',
hash('n-pentane'):'CCCCCC', hash('4-cyanopyridine'):'N#Cc1ccncc1',
hash("4,4'-bipyridine"):'c1cc(ccn1)c1ccncc1', hash('isopropanol'):'CC(C)O',
hash('1-butanol'):'CCCCO', hash('nitromethane'):'CN(=O)=O',
hash('n-hexane'):'CCCCCC', hash('dimethylformamide'):'CN(C)C=O',
hash('N,N-dimethylformamide'):'CN(C)C=O', hash('cyclohexane'):'C1CCCCC1',
hash('hydroquinone'):'Oc1ccc(O)cc1', hash('benzene'):'c1cccc1',
hash('octane'):'CCCCCC', hash('cyclohexanone'):'O=C1CCCCC1',
hash('isobutanol'):'CC(C)CO', hash('p-xylene'):'Cc1ccc(C)cc1',
hash('pentane'):'CCCCCC', hash('pyridine'):'c1ccncc1',
hash('t-butanol'):'CC(C)(C)O', hash('m-xylene'):'Cc1cccc(C)c1',
hash('phenol'):'Oc1cccc1', hash('o-cresol'):'Cc1cccc1O',
hash('p-cresol'):'Cc1ccc(O)cc1',
hash('2,4,6-trimethylpyridine'):'Cc1cc(C)nc(C)c1',
hash('1-methylpyrrolidin-2-one'):'CN1CCCC1=O',
hash('1,1,2,2,3,3,4,4,5,5,6,6,7,7,8,8-hexadecafluoro-1,8-diiodooctane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)I',
hash('trifluoroethanol'):'OCC(F)(F)F', hash('dichloroethane'):'ClCCCCl',
hash('1,1,2,2-tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('2-propanol'):'CC(C)O', hash('chlorobenzene'):'Clc1cccc1',
hash('naphthalene'):'c1ccc2cccc2c1', hash('propan-2-ol'):'CC(C)O',
hash('benzonitrile'):'N#Cc1cccc1', hash('nitrobenzene'):'O=N(=O)c1cccc1',
hash('dodecane'):'CCCCCCCCCCCC',
hash('1,2-dichlorobenzene'):'Clc1cccc1Cl', hash('n-heptane'):'CCCCCC',
hash('methanesulfonate'):'CS(=O)(=O)[O-]',
hash('3-methylbutan-1-ol'):'CC(C)CCO',
hash('1,1,2-trichloroethane'):'ClCC(Cl)Cl', hash('n-butane'):'CCCC',
hash('dichlorine'):'ClCl', hash('o-dichlorobenzene'):'Clc1cccc1Cl',
hash('N,N-dimethylacetamide'):'CN(C)C(C)=O',
hash('4-methylaniline'):'Cc1ccc(N)cc1', hash('morpholine'):'C1COCCN1',
hash('1,3-dimethylimidazolidin-2-one'):'CN1CCN(C)C1=O',
hash('1,3-dimethoxybenzene'):'COc1cccc(OC)c1', hash('n-propanol'):'CCCO',
hash('mesitylene'):'Cc1cc(C)cc(C)c1',
hash('acetylacetone'):'CC(=O)CC(C)=O', hash('diethylether'):'CCOCC',
hash('hydrazine'):'NN', hash('dimethylamine'):'CNC',
hash('1,2-difluorobenzene'):'Fc1cccc1F', hash('cycloheptane'):'C1CCCCC1',
hash('formamide'):'NC=O', hash('2-butanol'):'CCC(C)O',
hash('ethylacetate'):'CCOC(C)=O', hash('butan-1-ol'):'CCCCO',
hash('bromomethane'):'CBr', hash('piperazine'):'C1CNCCN1',
hash('n-nonane'):'CCCCCC',
hash('1,2,4-trimethoxybenzene'):'COc1ccc(OC)c(OC)c1',
hash('1,2-diaminoethane'):'NCCN', hash('n-butanol'):'CCCCO',
hash('N-methylformamide'):'CNC=O', hash('pyrrolidine'):'C1CCNCl1',
hash('actone'):'CC(C)=O', hash('n-octane'):'CCCCCC',
hash('S)-2-methyl-1-butanol'):'CCC(C)CO',
hash('N,N-diethylformamide'):'CCN(CC)C=O',
hash('tetracyanoethylene'):'N#CC(C#N)=C(C#N)C#N',

```

```

hash('2-methyl-1,3,5-trinitrobenzene'):'Cc1c(cc(cc1N(=O)=O)N(=O)=O)N(=O)=O',
hash('furan'):'O=C=CC=C1', hash('2-methylpentane'):'CCCC(C)C',
hash('isohexane'):'CCCC(C)C', hash('anisole'):'COc1ccccc1',
hash('butanol'):'CCCCO', hash('propan-1-ol'):'CCCO',
hash('dimethoxyethane'):'COCCOC', hash('cyclopentane'):'C1CCCC1',
hash('methoxybenzene'):'COc1ccccc1', hash('toulene'):'Cc1ccccc1',
hash('pyrene'):'c1cc2ccc3cccc4ccc(c1)c2c34', hash('butan-2-ol'):'CCC(C)O',
hash('3-methylpyridine'):'Cclcccncl', hash('triethylamine'):'CCN(CC)CC',
hash('dioxan'):'C1COCCO1', hash('2-methoxyethanol'):'COCCO',
hash("4,4'-bipyridyl"):'clcc(ccn1)c1ccncc1',
hash('dimethylacetamide'):'CN(C)C(C)=O', hash('thiophene'):'S1C=CC=C1',
hash('1,3-dioxolane'):'C1COCO1',
hash('1-chloronaphthalene'):'Clc1cccc2ccccc12',
hash('pentan-1-ol'):'CCCCCO', hash('hexan-1-ol'):'CCCCCC',
hash('2-phenylethanol'):'OCCc1ccccc1',
hash('1-chloro-4-methylbenzene'):'Cclccc(Cl)cc1',
hash('acetaldehyde'):'CC=O', hash('propiononitrile'):'CCC#N',
hash('bromobenzene'):'Brclccccc1'

53 deuterated_dictionary = {hash('dideutero-dichloromethane'):'ClcCl',
hash('perdeutero-toluene'):'Cclccccc1', hash('octadeutero-furan'):'C1CCOC1',
hash('deutero-ethanol'):'CCO', hash('deuterochloroform'):'C1C(Cl)Cl',
hash('deutero-chloroform'):'C1C(Cl)Cl', hash('deuterobenzene'):'c1ccccc1',
hash('deuteromethanol'):'CO', hash('hexadeutero-benzene'):'clccccc1',
hash('perdeuteroacetone'):'CC(C)=O',
hash('dideuterodichloromethane'):'ClCCl', hash('perdeuteromethanol'):'CO',
hash('hexadeutero-benzene'):'c1ccccc1'}

54 acids_dictionary = {hash('trifluoroacetic'):'OC(=O)C(F)(F)F',
hash('acetic'):'CC(O)=O', hash('formic'):'OC=O',
hash('sulfuric'):'OS(O)(=O)=O', hash('maleic'):'OC(=O)C=CC(O)=O',
hash('benzoic'):'OC(=O)c1ccccc1', hash('succinic'):'OC(=O)CCC(O)=O'}
55 ethers_dictionary = {'t-butyl methyl ether':'COC(C)(C)C',
hash('di-isopropyl'):'CC(C)OC(C)C', hash('dipropyl'):'CCCOCCCC',
hash('diethyl'):'CCOCC', hash('diisopropyl'):'CC(C)OC(C)C',
hash('isopropyl'):'CC(C)OC(C)C', hash('t-butyl'):'COC(C)(C)C',
hash('acetic'):'CCOC(C)=O'}

56 acetates_dictionary = {hash('ethyl'):'CCOC(C)=O',
hash('amyl'):'CCCCOC(C)=O', hash('methyl'):'COC(C)=O',
hash('butyl'):'CCCCOC(C)=O', hash('isobutyl'):'CC(C)OCOC(C)=O'}

57 glycol_dictionary = {hash('ethylene'):'OCOCO'}
58 oxalates_dictionary = {hash('dipropyl'):'CCCOOC(=O)C(=O)OCOC(=O)'}
59 acid_neutral_dictionary = {hash('oxalic'):'OC(=O)C(O)=O',
hash('nitric'):'ON(=O)=O', hash('trans-but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('fumaric'):'OC(=O)C=CC(O)=O', hash('phosphonic'):'OP(O)=O',
hash('perchloric'):'O[Cl](=O)(=O)=O', hash('phosphoric'):'OP(O)(O)=O',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(O)=O',
hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(O)=O)c(cc1C(O)=O)C(O)=O',
hash('trifluoroacetic'):'OC(=O)C(F)(F)F'}
60 acid_charged_dictionary = {hash('oxalic'):'OC(=O)C(=O)[O-]',
hash('nitric'):'O=N(=O)[O-]', hash('trans-but-2-enedioic'):'OC(=O)C=CC(=O)[O-]',
hash('but-2-enedioic'):'OC(=O)C=CC(=O)[O-]', hash('fumaric'):'OC(=O)C=CC(=O)[O-]', hash('phosphonic'):'OP(=O)[O-]', hash('perchloric'):'O=[Cl](=O)(=O)[O-]', hash('phosphoric'):'OP(O)(=O)[O-]', hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(=O)[O-]', hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(=O)[O-])c(cc1C(=O)[O-])C(O)=O', hash('trifluoroacetic'):'FC(F)(F)C(=O)[O-]'}

61 partial_name = [hash('ketone'), hash('propanoate'), hash('carbonate'),
hash('formate'), hash('iodine'), hash('tetrachloride'), hash('sulfoxide'),
hash('dichloride'), hash('chloride'), hash('benzoate'), hash('disulfide'),
hash('oxide')]

62 full_solvate_name_dictionary = {'methyl isobutyl ketone':'CC(C)CC(C)=O',
'deuterium oxide':'O', 'ethyl propanoate':'CCOC(=O)CC', 'dimethyl
carbonate':'COC(=O)OC', 'ethyl formate':'CCOC=O', '[12]cycloparaphenylene
iodine':'III', 'carbon tetrachloride':'ClC(Cl)(Cl)Cl', 'dimethyl

```

```

sulfoxide':'CS(C)=O', 'methylene dichloride':'ClCCl', 'methylene
chloride':'ClCCl', 'ethyl benzoate':'CCOC(=O)c1ccccc1', 'carbon
disulfide':'S=C=S' }

64
65     #The appearance of the same hash value for two distinct strings happens for
66     #15 of the anhydrous form identifiers
67     #These identifiers are also in a dictionary that assigns them a number
68     #between 1-15 rather than the same hash value as another structure identifier
69     duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
70     'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
71     'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
72     duplicate_smiles_dictionary =
73     {'CC(C)(C)c1c(O)c(O)c(Cl)c2[NH+]3CCN(CC3)c12':1,
74     'CC1=C(C(CC2=NS(=O)(=O)c3cccc23)c2cccc2C)C(=O)N(N1)c1cccc1':2,
75     'OC(=O)COc1cccc2cccc21':3, 'Cc1c[nH+]cc(C)c1':4,
76     'BrC1ccc(NC(=O)c2cccn2)cc1':5, 'CSC(N)=NC(C)=C(=O)c1cccc1':6,
77     'CC1(CO1)C(O)(CBr)CBr':7, 'OCC12CCCC(C1)(C(O)=O)C1(CCSCC1)O2':8,
78     'COc1ccc2CC3[NH2+]CCC3(C)c2c1':9,
79     'CC(=O)C(P(c1cccc1)c1cccc1)=C(C)Nc1c(F)cccc1F':10}

80     solvated_waterless_entry = []
81     solvated_waterless_hash = []
82     solvated_waterless_SMILES = []

83     #Split input file at entries: 25020, 50017, 75024, 100016, 125022, and 142002
84     entry1 = entry_reader1
85     for z in range(142002, len(entry1)):
86         solvated_waterless_entry.append(entry1[z])
87         if entry1[z].identifier in duplicate_identifier_dictionary:
88             solvated_waterless_hash.append(duplicate_identifier_dictionary.get(en-
89             try1[z].identifier))
90         else:
91             solvated_waterless_hash.append(hash(entry1[z].identifier))

92     while len(solvated_waterless_hash) != 0:
93         print len(solvated_waterless_hash)
94         entries1 = []
95         hashes1 = []
96         entries1.append(solvated_waterless_entry[0])
97         hashes1.append(solvated_waterless_hash[0])
98         for y in range(1, len(solvated_waterless_hash)):
99             if solvated_waterless_hash[0] != solvated_waterless_hash[y]:
100                 entries1.append(solvated_waterless_entry[y])
101                 hashes1.append(solvated_waterless_hash[y])
102             else:
103                 del solvated_waterless_hash[y]
104                 del solvated_waterless_entry[y]
105                 break
106
107             tag = 12
108             for a in range(len(entries1)):
109                 component_SMILES = []
110                 entry_SMILES = entries1[a].molecule.smiles.split('.')
111                 for b in range(len(entry_SMILES)):
112                     if entry_SMILES[b] not in duplicate_smiles_dictionary:
113                         component_SMILES.append(hash(entry_SMILES[b]))
114                     else:
115                         component_SMILES.append(duplicate_smiles_dictionary.get(entry_
116                         _SMILES[b]))
117                 component_SMILES = list(set(component_SMILES))
118                 if entries1[a].identifier == 'BULVEL':
119                     if hash('[O-][n+]'1cccc1') not in component_SMILES:
120                         component_SMILES.append(hash('[O-][n+]'1cccc1'))
121                     if hash('[O]n1cccc1') in component_SMILES:
122                         component_SMILES.remove(hash('[O]n1cccc1'))
```

```

113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167

    elif entries1[a].identifier == 'CEHKOS':
        component_SMILES.append(hash('III'))
    elif entries1[a].identifier == 'GOTJEE' or entries1[a].identifier == 'HASVEF' or entries1[a].identifier == 'MTFBTZ10':
        component_SMILES.append(tag)
        tag += 1
    elif entries1[a].identifier == 'JUHJIF':
        if hash('BrBr') not in component_SMILES:
            component_SMILES.append(hash('BrBr'))
        if hash('[Br]') in component_SMILES:
            component_SMILES.remove(hash('[Br]'))
    elif entries1[a].identifier == 'LOKXUF':
        if hash('OC(=O)C(F)(F)F') not in component_SMILES:
            component_SMILES.append(hash('OC(=O)C(F)(F)F'))
        if hash('[O]C(=O)C(F)(F)F') in component_SMILES:
            component_SMILES.remove(hash('[O]C(=O)C(F)(F)F'))
    elif entries1[a].identifier == 'OJUNEO':
        if hash('ClCl') not in component_SMILES:
            component_SMILES.append(hash('ClCl'))
        if hash('[Cl]') in component_SMILES:
            component_SMILES.remove(hash('[Cl]'))
    elif entries1[a].identifier == 'ZILFILM':
        component_SMILES.append(hash('CCO'))
    else:
        if entries1[a].chemical_name != None:
            pieces = entries1[a].chemical_name.split(' ')
            formula_pieces = entries1[a].formula.split(',')
            for c in range(len(pieces)):
                if hash(pieces[c]) == hash('xylene'):
                    if entries1[a].identifier == 'QOWNEV':
                        if hash('Cc1ccccclC') not in component_SMILES:
                            component_SMILES.append(hash('Cc1ccccclC'))
                    elif entries1[a].identifier == 'TEYCEF' or entries1[a].identifier == 'TEYCEF01':
                        if hash('Cc1cccc(C)c1') not in component_SMILES:
                            component_SMILES.append(hash('Cc1cccc(C)c1'))
                    elif entries1[a].identifier == 'MAMNAR':
                        if hash('Cc1ccc(C)cc1') not in component_SMILES:
                            component_SMILES.append(hash('Cc1ccc(C)cc1'))
                    elif hash(pieces[c]) == hash('hydrochloride') or hash(pieces[c]) == hash('bis(hydrochloride)'):
                        for z in range(len(formula_pieces)):
                            if formula_pieces[z] == 'Cl1 1-' or '(Cl1 1-)' in formula_pieces[z]:
                                if hash('[Cl-]') not in component_SMILES:
                                    component_SMILES.append(hash('[Cl-]'))
                            if formula_pieces[z] == 'H1 Cl1' or '(H1 Cl1)' in formula_pieces[z]:
                                if hash('Cl') not in component_SMILES:
                                    component_SMILES.append(hash('Cl'))
                    elif hash(pieces[c]) == hash('unidentified') or hash(pieces[c]) == hash('unknown'):
                        component_SMILES.append(tag)
                        tag += 1
                    elif hash(pieces[c]) == hash('glycol'):
                        if hash(pieces[c-1]) in glycol_dictionary:
                            if hash(glycol_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:
                                component_SMILES.append(hash(glycol_dictionary.get(hash(pieces[c-1]))))
                            elif hash(pieces[c]) == hash('oxalate'):
                                if hash(pieces[c-1]) in oxalates_dictionary:
                                    if hash(oxalates_dictionary.get(hash(pieces[c-1]))) not in component_SMILES:

```

168

```
component_SMILES.append(hash(oxalates_dictionary
    .get(hash(pieces[c-1]))))
```

169

```
elif hash(pieces[c]) == hash('acetate'):
```

170

```
if hash(pieces[c-1]) in acetates_dictionary:
```

171

```
if
```

172

```
hash(acetates_dictionary.get(hash(pieces[c-1])))  
not in component_SMILES:
```

173

```
component_SMILES.append(hash(acetates_dictionary
    .get(hash(pieces[c-1]))))
```

174

```
elif hash(pieces[c]) == hash('ether'):
```

175

```
if hash(pieces[c-1]) == hash('petroleum'):
```

176

```
component_SMILES.append(11)
```

177

```
elif pieces[c-2] + ' ' + pieces[c-1] + ' ' +
pieces[c] in ethers_dictionary:
```

178

```
if hash(ethers_dictionary.get(pieces[c-2] + ' '
+ pieces[c-1] + ' ' + pieces[c])) not in
component_SMILES:
```

179

```
component_SMILES.append(hash(ethers_dictionary
    .get(pieces[c-2] + ' ' + pieces[c-1] + ' '
    + pieces[c])))
```

180

```
elif hash(pieces[c-1]) in ethers_dictionary:
```

181

```
if
```

182

```
hash(ethers_dictionary.get(hash(pieces[c-1])))
```

183

```
not in component_SMILES:
```

184

```
component_SMILES.append(hash(ethers_dictionary
    .get(hash(pieces[c-1]))))
```

185

```
elif hash(pieces[c]) in partial_name:
```

186

```
if pieces[c] == 'ketone':
```

187

```
if pieces[c-2] + ' ' + pieces[c-1] + ' ' +
pieces[c] in full_solvate_name_dictionary:
```

188

```
if
```

189

```
hash(full_solvate_name_dictionary.get(pieces[c-
    2] + ' ' + pieces[c-1] + ' ' +
    pieces[c]))) not in component_SMILES:
```

190

```
component_SMILES.append(hash(full_solvate_
    _name_dictionary.get(pieces[c-2] + ' '
    + pieces[c-1] + ' ' + pieces[c])))
```

191

```
else:
```

192

```
if pieces[c-1] + ' ' + pieces[c] in
full_solvate_name_dictionary:
```

193

```
if
```

194

```
hash(full_solvate_name_dictionary.get(pieces[c-
    1] + ' ' + pieces[c]))) not in component_SMILES:
```

195

```
component_SMILES.append(hash(full_solvate_
    _name_dictionary.get(pieces[c-1] + ' '
    + pieces[c])))
```

196

```
elif hash(pieces[c]) == hash('acid'):
```

197

```
if hash(pieces[c-1]) in acids_dictionary:
```

198

```
if
```

199

```
hash(acids_dictionary.get(hash(pieces[c-1])))
```

200

```
not in component_SMILES:
```

```
component_SMILES.append(hash(acids_dictionary
    .get(hash(pieces[c-1]))))
```

```
elif hash(pieces[c-1]) in acid_neutral_dictionary:
```

```
if hash(pieces[c-1]) == hash('oxalic'):
```

```
for z in range(len(formula_pieces)):
```

```
if formula_pieces[z] == 'C2 H2 O4' or
```

```
'(C2 H2 O4)' in formula_pieces[z]:
```

```
if
```

```

hash(acid_neutral_dictionary.get(hash
(pieces[c-1]))) not in
component_SMILES:
200
    component_SMILES.append(hash(acid
    _neutral_dictionary.get(hash(piec
    es[c-1]))))
201    if formula_pieces[z] == 'C2 H1 O4 1-' or '(C2 H1 O4 1-)' in formula_pieces[z]:
202        if
203            hash(acid_charged_dictionary.get(hash
            (pieces[c-1]))) not in
            component_SMILES:
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222

```

```

223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
' (H3 O3 P1)' in formula_pieces[z]:
    if
        hash(acid_neutral_dictionary.get(hash(
            pieces[c-1]))) not in
            component_SMILES:
                component_SMILES.append(hash(acid_
                    _neutral_dictionary.get(hash(piec
                        es[c-1]))))
    if formula_pieces[z] == 'H2 O3 P1 1-' or
        '(H2 O3 P1 1-)' in formula_pieces[z]:
            if
                hash(acid_charged_dictionary.get(hash(
                    pieces[c-1]))) not in
                        component_SMILES:
                            component_SMILES.append(hash(acid_
                                _charged_dictionary.get(hash(piec
                                    es[c-1]))))
    elif hash(pieces[c-1]) == hash('perchloric'):
        for z in range(len(formula_pieces)):
            if formula_pieces[z] == 'H1 Cl1 O4' or
                '(H1 Cl1 O4)' in formula_pieces[z]:
                    if
                        hash(acid_neutral_dictionary.get(hash(
                            pieces[c-1]))) not in
                                component_SMILES:
                                    component_SMILES.append(hash(acid_
                                        _neutral_dictionary.get(hash(piec
                                            es[c-1]))))
            if formula_pieces[z] == 'Cl1 O4 1-' or
                '(Cl1 O4 1-)' in formula_pieces[z]:
                    if
                        hash(acid_charged_dictionary.get(hash(
                            pieces[c-1]))) not in
                                component_SMILES:
                                    component_SMILES.append(hash(acid_
                                        _charged_dictionary.get(hash(piec
                                            es[c-1]))))
            elif hash(pieces[c-1]) == hash('phosphoric'):
                for z in range(len(formula_pieces)):
                    if formula_pieces[z] == 'H3 O4 P1' or
                        '(H3 O4 P1)' in formula_pieces[z]:
                            if
                                hash(acid_neutral_dictionary.get(hash(
                                    pieces[c-1]))) not in
                                    component_SMILES:
                                        component_SMILES.append(hash(acid_
                                            _neutral_dictionary.get(hash(piec
                                                es[c-1]))))
            if formula_pieces[z] == 'H2 O4 P1 1-' or
                '(H2 O4 P1 1-)' in formula_pieces[z]:
                    if
                        hash(acid_charged_dictionary.get(hash(
                            pieces[c-1]))) not in
                            component_SMILES:
                                component_SMILES.append(hash(acid_
                                    _neutral_dictionary.get(hash(piec
                                        es[c-1]))))
            elif hash(pieces[c-1]) ==
                hash('pyridine-2,6-dicarboxylic'):
                    for z in range(len(formula_pieces)):
                        if formula_pieces[z] == 'C7 H5 N1 O4'

```

```
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999
999
```

```

269             if hash(deuterated_dictionary.get(hash(pieces[c]))) :
270                 not_in component_SMILES:
271
272                     component_SMILES.append(hash(deuterated_dictionary.get(hash(pieces[c]))))
273
274             if hash(pieces[c]) == hash('dideutero-dichloromethane') or
275                 hash(pieces[c]) == hash('dideuterodichloromethane'):
276                 if hash('Cl[Cl]C1') in component_SMILES:
277                     component_SMILES.remove(hash('Cl[Cl]C1'))
278
279             elif hash(pieces[c]) == hash('perdeutero-toluene'):
280                 if hash('[C]c1[c][c][c][c][c]1') in
281                     component_SMILES:
282
283                         component_SMILES.remove(hash('[C]c1[c][c][c][c]1'))
284
285             elif hash(pieces[c]) == hash('deutero-ethanol'):
286                 if hash('[C][C][O]') in component_SMILES:
287                     component_SMILES.remove(hash('[C][C][O]'))
288
289             elif hash(pieces[c]) == hash('deuterochloroform')
290                 or hash(pieces[c]) == hash('deutero-chloroform'):
291
292                 if hash('Cl[C](Cl)Cl') in component_SMILES:
293                     component_SMILES.remove(hash('Cl[C](Cl)Cl'))
294
295             elif hash(pieces[c]) == hash('hexadeutero-benzene')
296                 or hash(pieces[c]) == hash('deuterobenzene') or
297                 hash(pieces[c]) == hash('hexadeuterobenzene'):
298
299                 if hash('[c]1[c][c][c][c]1') in
300                     component_SMILES:
301
302                         component_SMILES.remove(hash('[c]1[c][c][c][c]1'))
303
304             elif hash(pieces[c]) == hash('deuteromethanol') or
305                 hash(pieces[c]) == hash('perdeuteromethanol'):
306                 if hash('[C][O]') in component_SMILES:
307                     component_SMILES.remove(hash('[C][O]'))
308
309             elif hash(pieces[c]) in solvents_dictionary:
310                 if hash(solvents_dictionary.get(hash(pieces[c]))) :
311                     not_in component_SMILES:
312
313                     component_SMILES.append(hash(solvents_dictionary.get(hash(pieces[c]))))
314
315             component_SMILES = sorted(component_SMILES)
316             solvated_waterless_SMILES.append(component_SMILES)
317
318         while len(solvated_waterless_SMILES) != 0:
319             number = []
320             identifiers = []
321             identifiers.append(hashes1[0])
322             if hashes1[0] == hash('KECYBU15'):
323                 number.append(87)
324             elif hashes1[0] == hash('MTYHFB03'):
325                 number.append(29)
326             else:
327
328                 number.append(entries1[0].crystal.spacegroup_number_and_setting[0])
329
330             for d in range(1,len(solvated_waterless_SMILES)):
331                 if solvated_waterless_SMILES[0] == solvated_waterless_SMILES[d]:
332                     identifiers.append(hashes1[d])
333                     if hashes1[d] == hash('KECYBU15'):
334                         number.append(87)
335                     elif hashes1[d] == hash('MTYHFB03'):
336                         number.append(29)
337                     else:
338
339                         number.append(entries1[d].crystal.spacegroup_number_and_setting[0])

```

```

315         etting[0])
316
317     if len(identifiers) != 1:
318         if all(x == number[0] for x in number) or any(number.count(y) >
319             1 for y in number):
320             identifiers_to_remove = []
321             for h in range(len(identifiers)):
322                 times1 = number.count(number[h])
323                 if times1 == 1:
324                     identifiers_to_remove.append(identifiers[h])
325             identifiers_to_check = [z for z in identifiers if z not in
326             identifiers_to_remove]
327             if identifiers_to_check != []:
328                 for i in range(len(identifiers_to_check)):
329                     pin =
330                     solvated_waterless_hash.index(identifiers_to_check[i])
331                     )
332                     duplicates_writer.write(solvated_waterless_entry[pin]
333                     )
334                     pin =
335                     solvated_waterless_hash.index(identifiers_to_check[0])
336                     duplicates_writer.write(solvated_waterless_entry[pin])
337                     for j in range(len(identifiers)):
338                         pin = solvated_waterless_hash.index(identifiers[j])
339                         solvated_waterless_hash.remove(solvated_waterless_hash[pin])
340                         pin2 = hashes1.index(identifiers[j])
341                         hashes1.remove(hashes1[pin2])
342                         entries1.remove(entries1[pin2])
343
344                         solvated_waterless_entry.remove(solvated_waterless_entry[pin])
345                         )
346
347                     solvated_waterless_SMILES.remove(solvated_waterless_SMILES[pi
348                     n])
349
350     else:
351         solvated_waterless_hash.remove(hashes1[0])
352         hashes1.remove(hashes1[0])
353         solvated_waterless_entry.remove(entries1[0])
354         entries1.remove(entries1[0])
355         solvated_waterless_SMILES.remove(solvated_waterless_SMILES[0])
356
357         print 'solvated waterless list is empty'
358
359
360     if __name__ == '__main__':
361         # This runs the script
362         r = Runner()
363         r.run()
364
365

```

```

1 #14 Find duplicate SMILES waterless forms with packing similarity
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8 import itertools
9 from decimal import Decimal
10 from ccdc.crystal import PackingSimilarity
11
12 similarity_engine = PackingSimilarity()
13 similarity_engine.settings.ignore_hydrogen_positions = True
14
15 #This script is written the same as the part2 script to find duplicates for hydrate
16 #structures with entry SMILES strings
16 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_smiles_all_potential_duplicates.txt"
17
18 class Runner(argparse.ArgumentParser):
19
20     def __init__(self):
21         super(self.__class__, self).__init__(description=__doc__)
22         self.add_argument(
23             '-i', '--input', default=filepath1,
24             help='input database filepath'
25         )
26         self.add_argument(
27             '-o', '--output', default='waterless_forms_smiles_all_first_occurrence.gcd',
28             help='output file [waterless_forms_smiles_all_first_occurrence.gcd]'
29         )
30         self.add_argument(
31             '-m', '--maximum', default=0, type=int,
32             help='Maximum number of structures to find [all]'
33         )
34
35     args = self.parse_args()
36
37     self.args = args
38     self.settings = search.Search.Settings()
39     self.settings.max_hit_structures = self.args.maximum
40
41     def run(self):
42
43         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
44
45         with io.EntryWriter(self.args.output) as lone_wolf_writer:
46             with io.EntryWriter("waterless_forms_smiles_all_duplicates.gcd") as writer1:
47                 with io.EntryWriter("waterless_forms_smiles_all_manual_check.gcd") as
writer2:
48
49                     duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2,
50                                         'OBASAN':3, 'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7,
51                                         'VOBXOZ':8, 'VUFJAI':9, 'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12,
52                                         'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
53
54                     solvated_waterless_hash = []
55                     solvated_waterless_entry = []
56
57                     #Split input file at entries: 3301, 3306, 8071, 8077, 13514, 13518,
58                     17337, 17340, 25752, 25756, 26203, 26207, 28814, 28820, 30802,
59                     30805, 33952, 33955, 35040, 35043, 37480, 37483, 38315, 38320,
60                     40159, 40162, 40293, 40296, 40626, 40629, 42708, 42711, 43261,
61                     43264, 43472, 43475, 44646, 44649, 45347, 45350, 45944, 45947,
62                     47758, 47769, 48118, 48129, 48550, 48561, 49064, 49070
63
64                     #The following entries get stuck at analysis calculation: HPHBNZ,
65                     AMAXOD, IFAWAP, QEHLUL, FALZEX, DAZPUS, EBIGUR, FEFTUI, HELXUR,

```

```

ICUVUZ, KOLWEN, LISZAQ, NOCYAH, NUGWIW, OLEGIV, REFLIX, SAXNUA,
SEYZII, TPCMOM, VEDFAO, WILNIO, CUXHOT, CUXLAJ, PUWYUD, PUXCOC,
ZEZJAS
56
57     for a in range(49070, len(entry_reader1)):
58         if entry_reader1[a].identifier in
59             duplicate_identifier_dictionary:
60
61                 solvated_waterless_hash.append(duplicate_identifier_dictionary.get(entry_reader1[a].identifier))
62
63             else:
64
65                 solvated_waterless_hash.append(hash(entry_reader1[a].identifier))
66
67             solvated_waterless_entry.append(entry_reader1[a])
68
69     while len(solvated_waterless_hash) != 0:
70         print len(solvated_waterless_hash)
71         identifiers = []
72         identifiers.append(solvated_waterless_hash[0])
73         for d in range(1, len(solvated_waterless_hash)):
74             if identifiers[0] != solvated_waterless_hash[d]:
75                 identifiers.append(solvated_waterless_hash[d])
76             else:
77                 del solvated_waterless_hash[d]
78                 del solvated_waterless_entry[d]
79                 break
80
81         remove_identifiers = []
82         while len(identifiers) != len(remove_identifiers):
83             for s in range(len(identifiers)):
84                 count = 0
85                 if identifiers[s] not in remove_identifiers:
86                     pin1 = solvated_waterless_hash.index(identifiers[s])
87                     reference = solvated_waterless_entry[pin1].crystal
88                     remove_identifiers.append(identifiers[s])
89                     for f in range(len(identifiers)):
90                         if identifiers[f] not in remove_identifiers:
91                             pin2 =
92                                 solvated_waterless_hash.index(identifiers[f])
93                                 sample =
94                                 solvated_waterless_entry[pin2].crystal
95                                 analysis =
96                                 similarity_engine.compare(reference, sample)
97                                 if analysis == None:
98                                     if count == 0:
99                                         writer2.write(solvated_waterless_entry[pin1])
99                                         count += 1
100
101                                         writer2.write(solvated_waterless_entry[pin2])
102                                         remove_identifiers.append(identifiers[f])
103                                         else:
104                                             if
105                                             analysis.nmatched_molecules/analysis.packing_shell_size == 1:
106
107                                                 lone_wolf_writer.write(solvated_waterless_entry[pin1])
108
109                                                 writer1.write(solvated_waterless_entry[pin2])
110
111                                                 remove_identifiers.append(identifiers[f])
112
113                                                 if count != 0:

```

```
100 writer2.write(solvated_waterless_entry[pin1])
101
102     for j in range(len(identifiers)):
103         pin = solvated_waterless_hash.index(identifiers[j])
104         solvated_waterless_hash.remove(solvated_waterless_hash[pin])
105         solvated_waterless_entry.remove(solvated_waterless_entry[pin])
106
107     print 'solvated waterless list is empty'
108
109
110 if __name__ == '__main__':
111     # This runs the script
112     r = Runner()
113     r.run()
114
```

```

1 #15 Find duplicates that fail packing similarity check
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #Potential duplicates that cannot be analyzed with packing similarity tool are
10 accounted for here
11 #Structures with unit cell parameters that are within 1 angstrom (for cell lengths) and
12 1 degree (for cell angles) are considered probable duplicates
13 #It is possible that two structures with similar unit cell parameters are polymorphs
14 and not duplicates
15 #Unfortunately one of the polymorphs in these cases will be lost here
16 #This method was the least manually intensive way to ensure that duplicate structures
17 are removed
18
19 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_Non_all_manual_check.txt"
20 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_smiles_all_manual_check.txt"
21 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_Non_all_manual_check.txt"
22 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_smiles_all_manual_check.txt"
23
24 class Runner(argparse.ArgumentParser):
25
26     def __init__(self):
27         super(self.__class__, self).__init__(description=__doc__)
28         self.add_argument(
29             '-i', '--input', default=filepath1,
30             help='input database filepath1'
31         )
32         self.add_argument(
33             '-o', '--output', default='hydrates_Non_entry_to_keep.gcd',
34             help='output file [hydrates_Non_entry_to_keep.gcd]'
35         )
36         self.add_argument(
37             '-m', '--maximum', default=0, type=int,
38             help='Maximum number of structures to find [all]'
39         )
40
41         args = self.parse_args()
42
43     def run(self):
44
45         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
46         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
47         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
48         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
49
50         with io.EntryWriter(self.args.output) as new_lists_writer:
51             with io.EntryWriter("hydrates_Non_duplicates_to_remove.gcd") as writer1:
52                 with io.EntryWriter("hydrates_smiles_entry_to_keep.gcd") as writer2:
53                     with io.EntryWriter("hydrates_smiles_duplicates_to_remove.gcd") as
writer3:
54                         with io.EntryWriter("waterless_forms_Non_entry_to_keep.gcd")
as writer4:
55                             with
io.EntryWriter("waterless_forms_Non_duplicates_to_remove.gcd")
" as writer5:

```

```

56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
with
io.EntryWriter("waterless_forms_smiles_entry_to_keep.gcd"
) as writer6:
    with
        io.EntryWriter("waterless_forms_smiles_duplicates_to_
remove.gcd") as writer7:
            loop = 0
            entry1 = entry_reader1
            while loop < 4:
                identifiers = []
                entries = []
                positions = []
                lengths = []
                angles = []
                #The unit cell lengths and angles for each
                structure are recorded in two separate lists
                for a in range(len(entry1)):
                    pre_lengths = []
                    pre_angles = []

                    identifiers.append(hash(entry1[a].identif
ier))
                    entries.append(entry1[a])
                    positions.append(a)
                    for b in range(3):

                        pre_lengths.append(entry1[a].crystal.
cell_lengths[b])

                        pre_angles.append(entry1[a].crystal.c
ell_angles[b])
                    lengths.append(pre_lengths)
                    angles.append(pre_angles)

                while len(identifiers) != 0:
                    print len(identifiers)
                    print positions[0]
                    compare_identifiers = []
                    compare_lengths = []
                    compare_angles = []

                    compare_identifiers.append(identifiers[0]
)
                    compare_lengths.append(lengths[0])
                    compare_angles.append(angles[0])
                    #The input files are organized so that
                    sets of potential duplicates are
                    written as A, B, C, ArgumentParser
                    #The second occurrence of A indicates
                    the end of one set and tells the script
                    to compare A, B, and C
                    for c in range(1, len(identifiers)):
                        if identifiers[0] != identifiers[c]:
                            compare_identifiers.append(identifi
ers[c])
                            compare_lengths.append(lengths[c]
)
                            compare_angles.append(angles[c])
                        else:
                            del identifiers[c]
                            del entries[c]
                            del positions[c]
                            del lengths[c]
                            del angles[c]

```



```

matched axes and angles will not be compared against new axes and angles
125                                         #For
example, if the a-axis and alpha angle of structure A match the c-axis and gamma angle
126                                         #Then the
of structure B
127                                         #Then the
b-axis and beta angle of structure A will only be compared to the a-axis and alpha
angle and b-axis and beta angle of structure B
128                                         if g not in
found:
129                                         if
abs(compare_lengths[d][1]-compare_lengths[e][g]) < 1 and
130                                         if
abs(compare_angles[d][1]-compare_angles[e][g]) < 1:
131                                         break
132                                         if count == 2:
133                                         for h in
134                                         if h
range(3):
135                                         if
not in found:
136                                         if
abs(compare_lengths[d][2]-compare_lengths[e][h]) < 1 and
137                                         if
abs(compare_angles[d][2]-compare_angles[e][h]) < 1:
138                                         break
139                                         #If all
three axis and angle pairs are matched between the two structures, they are considered
duplicates and written to two output files
140                                         #One output
file will have the first occurrence of the structure, the other will have all the
duplicate occurrences
141                                         if count ==
3:
142                                         if loop
== 0:
143                                         new_lists_writer.write(entries[d])
144                                         if loop
writer1.write(entries[e])
145                                         if loop
== 1:
146                                         writer2.write(entries[d])
147                                         if loop
writer3.write(entries[e])
148                                         if loop
== 2:
149                                         writer4.write(entries[d])
150                                         if loop
writer5.write(entries[e])
151                                         if loop
== 3:
152                                         writer6.write(entries[d])
153                                         if loop
writer7.write(entries[e])
154                                         #The list of entries to analyze is
identifiers_to_remove.append(compare_identifiers[e])
155                                         updated and the next set is evaluated
156                                         for y in range(len(compare_identifiers)):
                                         del identifiers[0]

```

```
157  
158  
159  
160  
161  
162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180 if __name__ == '__main__':  
181     r = Runner()  
182     r.run()  
183
```

```

1 #16 Remove overlapping duplicate structures
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\hydrates_smiles_all_first_occurrence.txt"
10 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\hydrates_smiles_all_duplicates.txt"
11 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\hydrates_None_entry_to_keep.txt"
12 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\hydrates_None_duplicates_to_remove.txt"
13
14 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\waterless_forms_smiles_all_first_occurrence.txt"
15 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\waterless_forms_smiles_all_duplicates.txt"
16 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\waterless_forms_None_entry_to_keep.txt"
17 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate Entries\waterless_forms_None_duplicates_to_remove.txt"
18
19 class Runner(argparse.ArgumentParser):
20
21     def __init__(self):
22         super(self.__class__, self).__init__(description=__doc__)
23         self.add_argument(
24             '-i', '--input', default=filepath1,
25             help='input database filepath1'
26         )
27         self.add_argument(
28             '-o', '--output',
29             default='complete_hydrates_merged_all_first_occurrence.gcd',
30             help='output file [complete_hydrates_merged_all_first_occurrence.gcd]'
31         )
32         self.add_argument(
33             '-m', '--maximum', default=0, type=int,
34             help='Maximum number of structures to find [all]'
35         )
36
37     args = self.parse_args()
38
39     self.args = args
40     self.settings = search.Search.Settings()
41     self.settings.max_hit_structures = self.args.maximum
42
43     def run(self):
44
45         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
46         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
47         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
48         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
49         entry_reader5 = io.EntryReader(filepath5, format='identifiers')
50         entry_reader6 = io.EntryReader(filepath6, format='identifiers')
51         entry_reader7 = io.EntryReader(filepath7, format='identifiers')
52         entry_reader8 = io.EntryReader(filepath8, format='identifiers')
53
54         with io.EntryWriter(self.args.output) as new_lists_writer:
55             with io.EntryWriter("complete_hydrates_merged_all_duplicates.gcd") as writer1:
56                 with
57                     io.EntryWriter("complete_waterless_forms_merged_all_first_occurrence.gcd")
58                 as writer2:

```

```

56
57     with
58         io.EntryWriter("complete_waterless_forms_merged_all_duplicates.gcd") as writer3:
59
60             duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2,
61             'OBASAN':3, 'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7,
62             'VOBXOZ':8, 'VUFJAI':9, 'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12,
63             'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
64
65             loop = 0
66             entry1 = entry_reader1
67             entry2 = entry_reader2
68             entry3 = entry_reader3
69             entry4 = entry_reader4
70             while loop < 2:
71                 #Since the first occurrence is based on the results from
72                 #the structures with entry SMILES strings
73                 #All the first occurrence structures are written to the new
74                 #output file of first occurrences
75                 first_occurrence = []
76                 for a in range(len(entry1)):
77                     if entry1[a].identifier in
78                         duplicate_identifier_dictionary:
79
80                         first_occurrence.append(duplicate_identifier_dictionary.get(entry1[a].identifier))
81
82                     else:
83                         first_occurrence.append(hash(entry1[a].identifier))
84                     if loop == 0:
85                         new_lists_writer.write(entry1[a])
86                     if loop == 1:
87                         writer2.write(entry1[a])
88
89                 #The same is done for all the duplicate structures from the
90                 #structures with entry SMILES strings
91                 duplicate = []
92                 for b in range(len(entry2)):
93                     if entry2[b].identifier in
94                         duplicate_identifier_dictionary:
95
96                         duplicate.append(duplicate_identifier_dictionary.get(entry2[b].identifier))
97
98                     else:
99                         duplicate.append(hash(entry2[b].identifier))
100                    if loop == 0:
101                        writer1.write(entry2[b])
102                    if loop == 1:
103                        writer3.write(entry2[b])
104
105                #First occurrence and duplicate results from structures
106                #without entry SMILES strings are screened
107                for c in range(len(entry3)):
108                    if entry3[c].identifier in
109                        duplicate_identifier_dictionary:
110
111                        query1 =
112                            duplicate_identifier_dictionary.get(entry3[c].identifier)
113
114                    else:
115                        query1 = hash(entry3[c].identifier)
116                    if entry4[c].identifier in
117                        duplicate_identifier_dictionary:
118
119                            query2 =
120                                duplicate_identifier_dictionary.get(entry4[c].identifier)
121
122                        else:
123                            query2 = hash(entry4[c].identifier)

```

```

102 #All structures from the entries to keep files are
103 #without an entry SMILES string
104 #Therefore they will not be present in the
105 #first_occurrence and duplicate lists because these only
106 #contain structures with entry SMILES strings
107 #Only the structures in the duplicates to remove files
108 #are compared to these lists
109 #If the structure from the duplicates to remove file is
110 #not in the first_occurrence and duplicate lists, the
111 #structure is part of a new duplicate pair
112 #The structures are written to the appropriate output
113 #files and the pair is added to the first_occurrence and
114 #duplicate lists
115 if query2 not in first_occurrence and query2 not in
116 duplicate:
117     if loop == 0:
118         new_lists_writer.write(entry3[c])
119         writer1.write(entry4[c])
120     if loop == 1:
121         writer2.write(entry3[c])
122         writer3.write(entry4[c])
123         first_occurrence.append(query1)
124         duplicate.append(query2)
125 else:
126     #If the structure from the duplicates to remove
127     #file is in the first_occurrence list, the pair
128     #needs to be re-written
129     #For example, A is in the entries to keep file and
130     #B is its partner in the duplicates to remove file
131     #The first occurrence list has B paired with C in
132     #the duplicates list
133     #Therefore the pair A and B needs to be swapped so
134     #B is the first occurrence and A is the duplicate
135     if query2 in first_occurrence:
136         #Here the duplicate to remove is written to the
137         #first occurrence output file and the entry to
138         #keep is written to the duplicate file
139         if loop == 0:
140             new_lists_writer.write(entry4[c])
141             writer1.write(entry3[c])
142         if loop == 1:
143             writer2.write(entry4[c])
144             writer3.write(entry3[c])
145             first_occurrence.append(query2)
146             duplicate.append(query1)
147     #If the structure from the duplicates to remove
148     #file is in the duplicates list, the entry to keep
149     #needs to be paired with the appropriate first occurrence
150     #For example, A is in the entries to keep file and
151     #B is its partner in the duplicates to remove file
152     #The duplicates list has B paired with C in the
153     #first occurrence list
154     #This means A is also a duplicate of C
155     #An additional instance of C needs to be added to
156     #first_occurrence and A needs to be added to
157     #duplicates
158     elif query2 in duplicate:
159         #Here the location of B in duplicates is used
160         #to find where C is in first_occurrence
161         #Here the entry to keep is written to the
162         #duplicate output file and the first_occurrence
163         #(i.e. C) is written to the first_occurrence file
164         place = duplicate.index(query2)
165         if loop == 0:
166             new_lists_writer.write(entry1[place])
167             writer1.write(entry3[c])
168         if loop == 1:

```

```
144 writer2.write(entry1[place])
145 writer3.write(entry3[c])
146 first_occurrence.append(first_occurrence[place])
147 duplicate.append(query1)
148
149 loop += 1
150 if loop == 1:
151     entry1 = entry_reader5
152     entry2 = entry_reader6
153     entry3 = entry_reader7
154     entry4 = entry_reader8
155
156
157
158
159
160
161
162 if __name__ == '__main__':
163     r = Runner()
164     r.run()
165
```

```

1 #17 Find stoichiometrically distinct duplicate structures
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds duplicate matches that have different stoichiometries for the
molecular components
10
11 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_first_occurrence.txt"
12 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_duplicates.txt"
13 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_first_occurrence.txt"
14 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_duplicates.txt"
15
16 class Runner(argparse.ArgumentParser):
17
18     def __init__(self):
19         super(self.__class__, self).__init__(description=__doc__)
20         self.add_argument(
21             '-i', '--input', default=filepath1,
22             help='input database filepath'
23         )
24         self.add_argument(
25             '-o', '--output', default='check_formula_hydrate_entry.gcd',
26             help='output file [check_formula_hydrate_entry.gcd]'
27         )
28         self.add_argument(
29             '-m', '--maximum', default=0, type=int,
30             help='Maximum number of structures to find [all]'
31         )
32
33     args = self.parse_args()
34
35     self.args = args
36     self.settings = search.Search.Settings()
37     self.settings.max_hit_structures = self.args.maximum
38
39     def run(self):
40
41         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
42         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
43         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
44         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
45
46         with io.EntryWriter(self.args.output) as new_lists_writer:
47             with io.EntryWriter("check_formula_hydrate_duplicate.gcd") as writer1:
48                 with io.EntryWriter("check_formula_waterless_form_entry.gcd") as writer2:
49                     with io.EntryWriter("check_formula_waterless_form_duplicate.gcd")
as writer3:
50
51             loop = 0
52             entry1 = entry_reader1
53             entry2 = entry_reader2
54             while loop < 2:
55                 for a in range(len(entry1)):
56                     #The formula of each structure in the match is used to compare
                     the stoichiometry of each molecule in the structure
57                     formulas1 = entry1[a].formula.split(',')
58                     #Once again, any formulas containing deuterium are modified to
                     contain the equivalent number of hydrogen atoms
59                     for b in range(len(formulas1)):
```

```

60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
if 'H' in formulas1[b] and 'D' in formulas1[b]:
    place = formulas1[b].index('D')
    place2 = formulas1[b].index('H')
    #Structures with more than 9 or less than 10 hydrogen
    or deuterium atoms are specified below so the script is
    able to pull out the correct number of atoms in each case
    #Here there are less than 10 deuterium atoms and less
    than 10 hydrogen atoms
    #The formula ends with the number of deuterium atoms in
    this case
    #In all other cases the number of hydrogen and
    deuterium atoms are not at the end of the formula and
    the location of a space can be used to determine the
    place value (one or ten)
    if len(str(formulas1[b])) == place+2:
        if formulas1[b][place2+2] == ' ':
            total = int(formulas1[b][place+1]) +
            int(formulas1[b][place2+1])
            formulas1[b] =
            formulas1[b].replace(formulas1[b][place2:place2+2],
            str('H') + str(total))
            formulas1[b] = formulas1[b][:place-1]
        else:
            print 'line 58'
            print entryl[a].identifier
            #Here there are less than 10 deuterium atoms and less
            than 10 hydrogen atoms
    elif formulas1[b][place2+2] == ' ' and
    formulas1[b][place+2] == ' ':
        total = int(formulas1[b][place+1]) +
        int(formulas1[b][place2+1])
        formulas1[b] =
        formulas1[b].replace(formulas1[b][place2:place2+2],
        str('H') + str(total))
        formulas1[b] = formulas1[b][:place] +
        formulas1[b][place+3:]
        #Here there are less than 10 deuterium atoms and more
        than 9 hydrogen atoms
    elif entryl[a].identifier == 'ROBYEM01' or
    entryl[a].identifier == 'ELOCIR' or
    entryl[a].identifier == 'DOMQOL' or
    entryl[a].identifier == 'VISZAZ' or
    entryl[a].identifier == 'UNATAE':
        total = int(formulas1[b][place+1]) +
        int(formulas1[b][place2+1:place2+3])
        formulas1[b] =
        formulas1[b].replace(formulas1[b][place2:place2+3],
        str('H') + str(total))
        formulas1[b] = formulas1[b][:place] +
        formulas1[b][place+3:]
        #Here there are more than 9 deuterium atoms and more
        than 9 hydrogen atoms
    elif entryl[a].identifier == 'SIGHOF01':
        total = int(formulas1[b][place+1:place+3]) +
        int(formulas1[b][place2+1:place2+3])
        formulas1[b] =
        formulas1[b].replace(formulas1[b][place2:place2+3],
        str('H') + str(total))
        formulas1[b] = formulas1[b][:place] +
        formulas1[b][place+3:]
        #Here there are more than 9 deuterium atoms and less
        than 10 hydrogen atoms
    elif entryl[a].identifier == 'QQQAEJ01':
        total = int(formulas1[b][place+1:place+3]) +
        int(formulas1[b][place2+1])
        formulas1[b] =
        formulas1[b].replace(formulas1[b][place2:place2+2],
        str('H') + str(total))

```

```

94
95     str('H') + str(total))
96     formulas1[b] = formulas1[b][:place] +
97     formulas1[b][place+3:]
98 else:
99     print 'line 61'
100    print entry1[a].identifier
101 elif 'D' in formulas1[b]:
102     formulas1[b] = re.sub(r'[D]', 'H', formulas1[b])
103 formulas2 = entry2[a].formula.split(',')
104 #The same procedure is carried out on the duplicate structures
105 for c in range(len(formulas2)):
106     if 'H' in formulas2[c] and 'D' in formulas2[c]:
107         place = formulas2[c].index('D')
108         place2 = formulas2[c].index('H')
109         if len(str(formulas2[c])) == place+2:
110             if formulas2[c][place2+2] == ' ':
111                 total = int(formulas2[c][place+1]) +
112                 int(formulas2[c][place2+1])
113                 formulas2[c] =
114                 formulas2[c].replace(formulas2[c][place2:place2+2],
115                                     str('H') + str(total))
116             elif entry2[a].identifier == 'ANTMEU04' or
117                 entry2[a].identifier == 'HATYUW01' or
118                 entry2[a].identifier == 'GASGEM01':
119                 total = int(formulas2[c][place+1]) +
120                 int(formulas2[c][place2+1:place2+3])
121                 formulas2[c] =
122                 formulas2[c].replace(formulas2[c][place2:place2+3],
123                                     str('H') + str(total))
124             else:
125                 print 'line 78'
126                 print entry2[a].identifier
127                 formulas2[c] = formulas2[c][:place-1]
128             elif formulas2[c][place2+2] == ' ' and
129                 formulas2[c][place+2] == ' ':
130                 total = int(formulas2[c][place+1]) +
131                 int(formulas2[c][place2+1])
132                 formulas2[c] =
133                 formulas2[c].replace(formulas2[c][place2:place2+2],
134                                     str('H') + str(total))
135                 formulas2[c] = formulas2[c][:place] +
136                 formulas2[c][place+3:]
137             elif entry2[a].identifier == 'DIBZAQ03' or
138                 entry2[a].identifier == 'DIBZAQ04' or
139                 entry2[a].identifier == 'ROBYEM02' or
140                 entry2[a].identifier == 'REXFIM01' or
141                 entry2[a].identifier == 'GUQXIB01' or
142                 entry2[a].identifier == 'LUSLOA01' or
143                 entry2[a].identifier == 'MAQWIM09' or
144                 entry2[a].identifier == 'MAQWIM10' or
145                 entry2[a].identifier == 'MAQWIM11' or
146                 entry2[a].identifier == 'MAQWIM12' or
147                 entry2[a].identifier == 'MAQWIM14' or
148                 entry2[a].identifier == 'MAQWIM06' or
149                 entry2[a].identifier == 'MAQWIM07' or
150                 entry2[a].identifier == 'MAQWIM08' or
151                 entry2[a].identifier == 'MAQWIM13' or
152                 entry2[a].identifier == 'OROQET' or
153                 entry2[a].identifier == 'ICUXUC01' or
154                 entry2[a].identifier == 'NALWUS03' or
155                 entry2[a].identifier == 'KOFCUE05' or
156                 entry2[a].identifier == 'NEZRUF01' or
157                 entry2[a].identifier == 'RUJCII01' or
158                 entry2[a].identifier == 'UPIVIA01' or
159                 entry2[a].identifier == 'VIWQUO' or
160                 entry2[a].identifier == 'GESXOT01' or
161                 entry2[a].identifier == 'MPRGOM01' or

```

```

122     entry2[a].identifier == 'TUTBUG01' or
123     entry2[a].identifier == 'KOFCUE04':
124         total = int(formulas2[c][place+1]) +
125             int(formulas2[c][place2+1:place2+3])
126         formulas2[c] =
127             formulas2[c].replace(formulas2[c][place2:place2+3],
128                 str('H') + str(total))
129         formulas2[c] = formulas2[c][:place] +
130             formulas2[c][place+3:]
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163 if __name__ == '__main__':
164     r = Runner()
165     r.run()
166

```

```

1 #18 Find distinct identifier prefix duplicate structures
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds duplicate matches that have a different set of the first six letters
10 #of the refcode
11 #Structures with the same first six letters in their refcode are multiple entries of
12 #the same structural components
13 #It is plausible that structures with distinct sets of the first six letters are
14 #composed of structural components that do not share the same stereochemistry
15 #Therefore, these matches were isolated to determine if there are any chiral centers
16 #that do not match up (R and S configurations)
17
18 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_first_occurrence.txt"
19 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_duplicates.txt"
20 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_first_occurrence.txt"
21 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_duplicates.txt"
22
23 class Runner(argparse.ArgumentParser):
24
25     def __init__(self):
26         super(self.__class__, self).__init__(description=__doc__)
27         self.add_argument(
28             '-i', '--input', default=filepath1,
29             help='input database filepath')
30
31         self.add_argument(
32             '-o', '--output', default='check_chirality_hydrates_entry.gcd',
33             help='output file [check_chirality_hydrates_entry.gcd]')
34
35         self.add_argument(
36             '-m', '--maximum', default=0, type=int,
37             help='Maximum number of structures to find [all]')
38
39     args = self.parse_args()
40
41     self.args = args
42     self.settings = search.Search.Settings()
43     self.settings.max_hit_structures = self.args.maximum
44
45     def run(self):
46
47         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
48         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
49         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
50         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
51
52         with io.EntryWriter(self.args.output) as new_lists_writer:
53             with io.EntryWriter("check_chirality_hydrates_duplicate.gcd") as writer1:
54                 with io.EntryWriter("check_chirality_waterless_forms_entry.gcd") as
writer2:
55                     with
56                         io.EntryWriter("check_chirality_waterless_forms_duplicate.gcd") as
writer3:
57
58                         loop = 0
59                         entry1 = entry_reader1
60                         entry2 = entry_reader2

```

```

57     while loop < 2:
58         for a in range(len(entry1)):
59             #The first six letters of each identifier are compared
60             #Any cases where the first six letters do not match are
61             #written to output file where they are flagged as having
62             #distinct six letter identifiers
63             if entry1[a].identifier[:6] != entry2[a].identifier[:6]:
64                 if loop == 0:
65                     new_lists_writer.write(entry1[a])
66                     writer1.write(entry2[a])
67                 if loop == 1:
68                     writer2.write(entry1[a])
69                     writer3.write(entry2[a])
70             loop += 1
71             if loop == 1:
72                 entry1 = entry_reader3
73                 entry2 = entry_reader4
74
75
76
77
78
79
80
81
82
83
84
85 if __name__ == '__main__':
86     r = Runner()
87     r.run()
88

```

```

1 #19 Check chirality of duplicates with distinct identifier prefixes
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script determines which matches with different sets of the first six letters in
10 their refcodes have distinct stereochemistry
11 #The hydrate structures and waterless form structures were evaluated with this script
12 #The input files and output file names were changed when switching from hydrate
13 structures to waterless forms
14 #i.e. Every instance of the word "hydrate" was changed to "waterless form" in the
15 filenames
16
17 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_chirality_hydrate_entry.txt"
18 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_chirality_hydrate_duplicate.txt"
19
20 #filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_chirality_waterless_forms_entry.txt"
21 #filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_chirality_waterless_forms_duplicate.txt"
22
23 class Runner(argparse.ArgumentParser):
24
25     def __init__(self):
26         super(self.__class__, self).__init__(description=__doc__)
27         self.add_argument(
28             '-i', '--input', default=filepath1,
29             help='input database filepath')
30
31         self.add_argument(
32             '-o', '--output', default='same_chirality_hydrate_entry.gcd',
33             help='output file [same_chirality_hydrate_entry.gcd]')
34
35         self.add_argument(
36             '-m', '--maximum', default=0, type=int,
37             help='Maximum number of structures to find [all]')
38
39     args = self.parse_args()
40
41     self.args = args
42     self.settings = search.Search.Settings()
43     self.settings.max_hit_structures = self.args.maximum
44
45     def run(self):
46
47         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
48         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
49
50         with io.EntryWriter(self.args.output) as chiral_writer:
51             with io.EntryWriter("same_chirality_hydrate_duplicate.gcd") as writer1:
52                 with io.EntryWriter("distinct_chirality_hydrate_entry.gcd") as writer2:
53                     with io.EntryWriter("distinct_chirality_hydrate_duplicate.gcd") as
writer3:
54                         with
55                             io.EntryWriter("achiral_names_one_missing_stereo_hydrate_entry.gc
d") as writer4:
56                             with
57                                 io.EntryWriter("achiral_names_one_missing_stereo_hydrate_duplicat
e.gcd") as writer5:
58                                     with
59                                         io.EntryWriter("one_chiral_name_one_missing_stereo_hydrat
e.gcd")

```



```
hash('FEFNOT10'), hash('FEFNOT10'),
hash('FEFNOT11'), hash('FEFNOT11'),
hash('FEFNOT11'), hash('FUSWOH04'),
hash('FUSWOH04'), hash('FUSWOH04'),
hash('FUSWOH04'), hash('FUSWOH04'),
hash('FUSWOH04'), hash('FUSWOH04'),
hash('FUSWOH04'), hash('IKALIR'),
hash('IKALIR'), hash('IKALIRO1'),
hash('IKALIRO1'), hash('IKALIRO2'),
hash('IKALIRO2'), hash('IKALIRO3'),
hash('IKALIRO3'), hash('IKAMEO'),
hash('IKAMEO'), hash('IKAMEO01'),
hash('IKAMEO01'), hash('IKAMOY'),
hash('IKAMOY'), hash('IKAMOY01'),
hash('IKAMOY01'), hash('IKAMOY02'),
hash('IKAMOY02'), hash('IKAMOY03'),
hash('IKAMOY03'), hash('IKAMOY04'),
hash('IKAMOY04'), hash('IKAMOY05'),
hash('IKAMOY05'), hash('IKAVOF'),
hash('IKAVOF'), hash('IKAVOF'),
hash('IKAVOF'), hash('IKAVOF'),
hash('IKAVOF'), hash('IKAVOF'),
hash('JARBUA'), hash('KUYGOA'),
hash('KUYGUG'), hash('KUYGUG'),
hash('KUYGUG'), hash('KUYGUG'),
hash('KUYGUG'), hash('KUYGUG'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHAN'), hash('KUYHAN'),
hash('KUYHER'), hash('KUYHER'),
hash('MUSDEL'), hash('MUSDOV'),
hash('NAHCIJ'), hash('NAHCIJ'),
hash('NAHCOP'), hash('NAHCOP'),
hash('NAHCOP'), hash('NAHCOP'),
hash('NAHCOP'), hash('NAHCOP'),
hash('NAHCOP'), hash('NAHCOP'),
hash('QQQCIY10'), hash('QQQCIY10'),
hash('QQQCIY10'), hash('QQQCIY10'),
hash('QQQCIY10'), hash('QQQCIY10'),
hash('QQQCIY10'), hash('RITMEN'),
hash('RITMEN'), hash('SOGUAN06'),
hash('XEBQAZ'), hash('ZZZPPI01'),
hash('ZZZWQA01'), hash('ZZZWQA01'),
hash('FEFNOT'), hash('FEFNOT'),
hash('FEFNOT'), hash('FEFNOT')
```

```
hash('CBMZPN02'), hash('CBMZPN03'),
hash('CBMZPN10'), hash('CBMZPN11'),
hash('CBMZPN12'), hash('CBMZPN13'),
hash('CBMZPN14'), hash('CBMZPN16'),
hash('CBMZPN17'), hash('CBMZPN18'),
hash('CBMZPN20'), hash('CBMZPN21'),
hash('CBMZPN22'), hash('CBMZPN23'),
hash('CBMZPN27'), hash('CBMZPN28'),
hash('CBMZPN29'), hash('CBMZPN30'),
hash('CBMZPN01'), hash('CBMZPN02'),
hash('CBMZPN03'), hash('CBMZPN10'),
hash('CBMZPN11'), hash('CBMZPN12'),
hash('CBMZPN13'), hash('CBMZPN14'),
hash('CBMZPN16'), hash('CBMZPN17'),
hash('CBMZPN18'), hash('CBMZPN20'),
hash('CBMZPN21'), hash('CBMZPN22'),
hash('CBMZPN23'), hash('CBMZPN27'),
hash('CBMZPN28'), hash('CBMZPN29'),
hash('CBMZPN30'), hash('CBMZPN01'),
hash('CBMZPN02'), hash('CBMZPN03'),
hash('CBMZPN10'), hash('CBMZPN11'),
hash('CBMZPN12'), hash('CBMZPN13'),
hash('CBMZPN14'), hash('CBMZPN16'),
hash('CBMZPN17'), hash('CBMZPN18'),
hash('CBMZPN20'), hash('CBMZPN21'),
hash('CBMZPN22'), hash('CBMZPN23'),
hash('CBMZPN27'), hash('CBMZPN28'),
hash('CBMZPN29'), hash('CBMZPN30'),
hash('EFUMAU'), hash('EFUMAU01'),
hash('EFUMAU02'), hash('EFUMAU03'),
hash('EFUMAU04'), hash('EFUMAU05'),
hash('EFUMAU06'), hash('EFUMAU07'),
hash('IKAKOW'), hash('IKAKOW01'),
hash('ETHANE01'), hash('ETHANE04'),
hash('ETHANE05'), hash('ETHANE06'),
hash('ETHANE07'), hash('ETHANE08'),
hash('ETHANE09'), hash('ETHANE10'),
hash('ETHANE11'), hash('JARBOUTI'),
hash('ZZZITY01'), hash('KAMTAW'),
hash('KAMTAW01'), hash('KAMTAW02'),
hash('KAMTAW03'), hash('KAMTAW04'),
hash('KAMTAW05'), hash('KAMTAW06'),
hash('DCLBEN'), hash('DCLBEN01'),
hash('DCLBEN02'), hash('DCLBEN03'),
hash('DCLBEN04'), hash('DCLBEN05'),
hash('DCLBEN06'), hash('DCLBEN07'),
hash('DCLBEN11'), hash('CISTON01'),
hash('TAQYUG'), hash('DUHJIB'),
hash('DUHJIB'), hash('ACETYL02'),
hash('ACETYL03'), hash('JAYDUI'),
hash('JAYDUI01'), hash('JAYDUI02'),
hash('JAYDUI03'), hash('JAYDUI04'),
hash('JAYDUI05'), hash('JAYDUI06'),
hash('JAYDUI07'), hash('CEKGUU'),
hash('CEKGUU01'), hash('CEKGUU02'),
hash('CEKGUU03'), hash('CEKGUU04'),
```

```

68     hash('DAYFUH'), hash('LUXYIM'),
69     hash('LUYMAT'), hash('LUYMAT01'),
70     hash('NIWFEE02'), hash('NIWFEE02'),
71     hash('NIWFEE04'), hash('NIWFEE04'),
72     hash('SABZOK'), hash('SOCTOT'),
73     hash('THIOUR06'), hash('THIOUR06'),
74     hash('THIOUR06'), hash('THIOUR06'),
75     hash('PUSBOV')]
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
#This script was adapted from the one
#used to find hydrate-anhydrate pairs
#with distinct stereochemistry
#The hydrates list contains the
#information for the entry to keep
#(first occurrence)
#The waterless list contains the
#informaiton for the duplicate
entry1 = entry_reader1
hydrates = []
waterless = []
count1 = 0
count2 = 0
count3 = 0
#The formula for the entry will always
#be complete unless a solvent molecule
#is squeezed out
#The formula of a molecule in an entry
#can be incomplete or overdone
#An example would be a disordered
#structure where a hydrogen atom 3D
#coordinate is not reported (incomplete)
#or a carbon atom position has three
#possibilities and is listed three times
#(overdone)
#This part of the script compares the
#entry formula pieces to each molecule
#formula to find cases where they do not
#match
#Cases where one formula is deuterated
#and the other is not are considered a
#match
for a in range(len(entry1)):
    print a
    formula_error = 0
    missing_components = 0
    loop2 = 0
    while loop2 < 2:
        letter_formulas = []
        formulas =
        entry1[a].formula.split(',')
        for b in range(len(formulas)):
            if '(' in formulas[b] and
               ')' in formulas[b]:
                if '(H2 O1)' not in
                   formulas[b] and '(D2
                   O1)' not in formulas[b]:
                    pin =
                    formulas[b].index('('
)
                    letter_formulas.append(
d(formulas[b][pin+1:-1]))
            else:
                if formulas[b] != 'H2
O1' and formulas[b] !=
'D2 O1':

```

```

98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
letter_formulas.append
d(formulas[b])
molecule1 = entry1[a].molecule
#Checked contains all the
molecule formulas that have
been matched to a formula in
the entry formula
#Tags contains the location of
the molecule formula that has
been found in the list of
molecule.components
checked = []
tags = []
#Error records the number of
cases where a molecule formula
does not match any entry formulas
#Count records the number of
molecule formulas that have
matched formulas in the entry
formula
error = 0
count = 0
#Molecule.components generates
the molecule formulas
for g in
range(len(molecule1.components)):
    if
        molecule1.components[g].formula
        not in checked and
        molecule1.components[g].formula
        != 'H2 O1':
            if
                molecule1.components[g].f
                ormula not in
                letter_formulas:
                    if 'H' in
                        molecule1.components[
                            g].formula:
                            if
                                molecule1.compone
                                nts[g].formula.re
                                place('H', 'D')
                                not in
                                letter_formulas:
partial_deuterated = 0
for t in
    if 'H'
        for
            s in range(len(pieces)):
                if 'H' in pieces[s]:
                    tag1 = pieces[s].index('H')
                    number1 = int(pieces[s][tag1+1:])
                if 'D' in pieces[s]:
                    tag2 = pieces[s].index('D')
                    number2 = int(pieces[s][tag2+1:])

```

```

125     partial_deuterated += 1
126
127     partial_deuterated == 1:
128         = molecule1.components[g].formula.split(' ')
129
130     in range(len(pieces2)):
131
132     'H' in pieces2[u]:
133
134     tag3 = pieces2[u].index('H')
135
136     number3 = int(pieces2[u][tag3+1:])
137
138     number1+number2 == number3:
139
140     checked.append(molecule1.components[g].formula)
141
142     tags.append(g)
143
144     count += 1
145
146     checked.append(molecule1.components[g].formula)
147
148     tags.append(g)
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163

```

if  
pieces2  
for u  
if  
if  
else:  
else:  
error  
else:  
count += 1  
else:  
checked.append(molecule1.components[g].formula)  
tags.append(g)  
error += 1  
else:  
checked.append(molecule1.components[g].formula)  
tags.append(g)  
count +=  
1  
if loop2 == 0:  
 hydrate\_tags = tags  
 hydrates\_checked = checked  
else:  
 waterless\_tags = tags  
 waterless\_checked = checked  
if error != 0:  
 formula\_error += 1

```

165 #If the number of molecule
166 formulas found does not match
167 the total number of formulas in
168 the entry formula, then one of
169 the molecule formulas is
170 missing from that entry
171 if count != len(letter_formulas):
172     missing_components += 1
173 loop2 += 1
174 entry1 = entry_reader2

175 entry2 = entry_reader1
176 entry3 = entry_reader2

177 #Any cases where there is an error
178 or components are missing go
179 through a secondary refinement
180 #If one structure (ex. entry) is
181 missing a molecule, then that
182 molecule is removed from the other
183 structure (ex. duplicate) since
184 they can now no longer be compared
185 if formula_error != 0 or
186 missing_components != 0:
187     remove_hydrate_tags = []
188     for c in
189         range(len(hydrate_tags)):
190             if
191                 entry2[a].molecule.components
192                     [hydrate_tags[c]].formula
193                         not in waterless_checked:
194
195                     remove_hydrate_tags.append(
196                         hydrate_tags[c])
197             for e in
198                 range(len(remove_hydrate_tags)):
199
200                 hydrate_tags.remove(remove_hydrate_tags[e])
201             remove_waterless_tags = []
202             for d in
203                 range(len(waterless_tags)):
204                     if
205                         entry3[a].molecule.components
206                             [waterless_tags[d]].formula
207                                 not in hydrates_checked:
208
209                         remove_waterless_tags.append(
210                             waterless_tags[d])
211             for f in
212                 range(len(remove_waterless_tags)):
213
214                 waterless_tags.remove(remove_waterless_tags[f])

215 #If the length of molecules for
216 each structure does not match, the
217 identifiers of the two structures
218 are printed for a manual inspection
219 if len(hydrate_tags) !=
220     len(waterless_tags):
221     print entry2[a].identifier
222     print entry3[a].identifier

223 #There are several possibilities

```

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

that can come from this script  
 #Structures that have chiral  
 centers may not have that  
 information recorded for the atoms  
 of the chiral molecule  
 #Structures that have chiral  
 centers may not have any indication  
 in the chemical name  
 #The chemical name is used as a  
 preliminary check before  
 investigating the atom chiral  
 configurations of each molecule  
 #If both structures have the same  
 name and the name indicates the  
 chirality, the match is true  
 #If the structures have distinct  
 names, both of which indicate  
 chirality, the match is false  
 #If either or both names do not  
 indicate the chirality, then the  
 chiral configuration of the atoms  
 in the molecules are checked  
 same\_names = 0  
 possibly\_the\_same\_names = 0  
 distinct\_names = 0  
 #Double R = both names indicate  
 chirality but are not the same  
 #Double blank = same names with no  
 chirality indicated  
 #RandBlank = one name has chirality  
 and the other does not  
 DoubleR = 0  
 DoubleBlank = 0  
 RandBlank = 0  
 #The chemical name of each  
 structure, along with any synonyms  
 are placed into a list  
 #The list for each structure are  
 compared against each other to see  
 if any of the names are the same  
 hydrate\_titles =  
 list(entry2[a].synonyms)  
 if entry2[a].identifier == 'WACHIR':  
 hydrate\_titles =  
 hydrate\_titles[0].split(',')  
 hydrate\_titles.append(entry2[a].chemical\_name)  
 for u in range(len(hydrate\_titles)):  
 caps = [i for i in  
 hydrate\_titles[u] if i.isupper()]  
 caps = list(set(caps))  
 for q in range(len(caps)):  
 if caps[q] != 'D' and  
 caps[q] != 'L' and caps[q]  
 != 'R' and caps[q] != 'S'  
 and caps[q] != 'E' and  
 caps[q] != 'Z':  
 hydrate\_titles[u] =  
 re.sub(caps[q],  
 caps[q].lower(),  
 hydrate\_titles[u])  
 else:  
 places = [j for j,  
 letter in  
 enumerate(hydrate\_titles[u]) if letter == caps[q]]  
 86

```

225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
for k in
range(len(places)):
    if places[k]+1 == len(hydrate_titles[u]):
        :
            hydrate_titles[u] =
            hydrate_titles[u].replace(hydrate_
titles[u][places[
k]], hydrate_titles[u][places[k]].lower(
()))
    elif
hydrate_titles[u][pla
ces[k]+1] != ')'
and
hydrate_titles[u][pla
ces[k]+1] != '-'
and
hydrate_titles[u][pla
ces[k]+1] != ',': hydrate_titles[u] =
hydrate_titles[u].replace(hydrate_
titles[u][places[
k]], hydrate_titles[u][places[k]].lower(
()))
    hydrate_pieces =
hydrate_titles[u].split(' ')
hydrate_deuterated = [s for s
in hydrate_pieces if 'deutero'
in s]
if len(hydrate_deuterated) == 1:
    clip =
hydrate_deuterated[0].index('
deutero')
if 'deutero-' in
hydrate_deuterated[0]:
    hydrate_replacement =
[hydrate_deuterated[0][cl
ip+8:]]
else:
    hydrate_replacement =
[hydrate_deuterated[0][cl
ip+7:]]
for t in
range(len(hydrate_pieces)):
    if hydrate_pieces[t] ==
hydrate_deuterated[0]:
        hydrate_pieces[t] =
hydrate_replacement[0
]
water = [r for r in
hydrate_pieces if 'hydrate' in
r and 'perhydrate' not in r]
D2O = 0
if len(water) == 0:
    if 'deuterium oxide
solvate' in
hydrate_titles[u].lower() 87

```



265

```

    elif
        waterless_titles[u][p
        laces2[k]+1] != ')'
        and
        waterless_titles[u][p
        laces2[k]+1] != '-'
        and
        waterless_titles[u][p
        laces2[k]+1] != ',', :

```

266

```

            waterless_titles[
            u] =
            waterless_titles[
            u].replace(waterl
            ess_titles[u][pla
            ces2[k]],
            waterless_titles[
            u][places2[k]].lo
            wer())

```

267

```

waterless_pieces =
waterless_titles[u].split(' ')
waterless_deuterated = [s for s
in waterless_pieces if
'deutero' in s.lower()]
if len(waterless_deuterated) == 1:

```

268

```

    clip =
    waterless_deuterated[0].lower
    ().index('deutero')
    if 'deutero-' in
    waterless_deuterated[0].lower
    () :
        waterless_replacement =
        [waterless_deuterated[0][
        clip+8:]]

```

269

```

else:
    waterless_replacement =
    [waterless_deuterated[0][
    clip+7:]]

```

270

```

for t in
range(len(waterless_pieces)):
    if waterless_pieces[t] ==
    == waterless_deuterated[0]:
        waterless_pieces[t] =
        waterless_replacement
        [0]

```

271

```

waterless_titles[u] =
'.join(waterless_pieces)

```

272

#Words that indicate chirality are searched for in each chemical name (ex. rac, D-, R-, etc.)

273

#Racemates are searched for first, followed by enantiomers, then cases where there is no chirality in either name, and cases where there is chirality in one name

274

```

if any('rac-' in x for x in
hydrate_titles) and any('rac-' in y
for y in waterless_titles) or
any('DL' in x for x in
hydrate_titles) and any('DL' in y
for y in waterless_titles) or
any('RS' in x for x in
hydrate_titles) and any('RS' in y
for y in waterless_titles) or

```

275

276

277

278

279

280

281

```

any('RS' in x for x in
hydrate_titles) and any('SR' in y
for y in waterless_titles) or
any('SR' in x for x in
hydrate_titles) and any('RS' in y
for y in waterless_titles) or
any('SR' in x for x in
hydrate_titles) and any('SR' in y
for y in waterless_titles) or
any('+-' in x for x in
hydrate_titles) and any('+-' in y
for y in waterless_titles):
    same_names += 1
    #If the two names are the same,
    #the two structures are written
    #to output files for matches
    #with the same chirality
    if bool(set(hydrate_titles) &
set(waterless_titles)) == True:
        chiral_writer.write(entry2[a]
)
        writer1.write(entry3[a])
    else:
        possibly_the_same_names += 1
        DoubleR += 1

elif any('D' in x for x in
hydrate_titles) and any('D' in y
for y in waterless_titles) or
any('L' in x for x in
hydrate_titles) and any('L' in y
for y in waterless_titles) or
any('R' in x for x in
hydrate_titles) and any('R' in y
for y in waterless_titles) or
any('S' in x for x in
hydrate_titles) and any('S' in y
for y in waterless_titles) or
any('(+)-' in x for x in
hydrate_titles) and any('(+)-' in y
for y in waterless_titles) or
any('(-)-' in x for x in
hydrate_titles) and any('(-)-' in y
for y in waterless_titles) or
any('cis-' in x for x in
hydrate_titles) and any('cis-' in y
for y in waterless_titles) or
any('trans-' in x for x in
hydrate_titles) and any('trans-' in y
for y in waterless_titles) or
any('meso-' in x for x in
hydrate_titles) and any('meso-' in y
for y in waterless_titles) or
any('E' in x for x in
hydrate_titles) and any('E' in y
for y in waterless_titles) or
any('Z' in x for x in
hydrate_titles) and any('Z' in y
for y in waterless_titles):
    same_names += 1
    if bool(set(hydrate_titles) &
set(waterless_titles)) == True:
        chiral_writer.write(entry2[a]
)
        writer1.write(entry3[a])

```

```

296
297
298
299
300
301
302
303
304
else:
    possibly_the_same_names += 1
    DoubleR += 1
elif all('rac-' not in x for x in
hydrate_titles) and all('rac-' not
in y for y in waterless_titles) and
all('+-' not in x for x in
hydrate_titles) and all('+-' not in
y for y in waterless_titles) and
all('D' not in x for x in
hydrate_titles) and all('D' not in
y for y in waterless_titles) and
all('L' not in x for x in
hydrate_titles) and all('L' not in
y for y in waterless_titles) and
all('R' not in x for x in
hydrate_titles) and all('R' not in
y for y in waterless_titles) and
all('S' not in x for x in
hydrate_titles) and all('S' not in
y for y in waterless_titles) and
all('(+)-' not in x for x in
hydrate_titles) and all('(+)-' not
in y for y in waterless_titles) and
all('(-)-' not in x for x in
hydrate_titles) and all('(-)-' not
in y for y in waterless_titles) and
all('cis-' not in x for x in
hydrate_titles) and all('cis-' not
in y for y in waterless_titles) and
all('trans-' not in x for x in
hydrate_titles) and all('trans-' not
in y for y in waterless_titles) and
all('meso-' not in x for x in
hydrate_titles) and all('meso-' not
in y for y in waterless_titles) and
all('E' not in x for x in
hydrate_titles) and all('E' not in
y for y in waterless_titles) and
all('Z' not in x for x in
hydrate_titles) and all('Z' not in
y for y in waterless_titles):
    possibly_the_same_names += 1
    DoubleBlank += 1
else:
    if any('D' in x for x in
hydrate_titles) or any('L' in x
for x in hydrate_titles) or
any('R' in x for x in
hydrate_titles) or any('S' in x
for x in hydrate_titles) or
any('(+)-' in x for x in
hydrate_titles) or any('(-)-'
in x for x in hydrate_titles)
or any('cis-' in x for x in
hydrate_titles) or any('trans-'
in x for x in hydrate_titles)
or any('E' in x for x in
hydrate_titles) or any('Z' in x
for x in hydrate_titles):
        if all('rac-' not in y for
y in waterless_titles) and
all('+-' not in y for y in
waterless_titles) and
all('D' not in y for y in
waterless_titles) and
all('L' not in y for y in
waterless_titles) and
all('R' not in y for y in
waterless_titles) and
all('S' not in y for y in
waterless_titles) and
all('E' not in y for y in
waterless_titles) and
all('meso-' not in y for y in
waterless_titles) and
all('cis-' not in y for y in
waterless_titles) and
all('trans-' not in y for y in
waterless_titles) and
all('(+)-' not in y for y in
waterless_titles) and
all('(-)-' not in y for y in
waterless_titles) and
all('Z' not in y for y in
waterless_titles):
            possibly_the_same_names += 1
            DoubleBlank += 1

```

```

waterless_titles) and
all('R' not in y for y in
waterless_titles) and
all('S' not in y for y in
waterless_titles) and
all('(+)-' not in y for y
in waterless_titles) and
all('(-)-' not in y for y
in waterless_titles) and
all('cis-' not in y for y
in waterless_titles) and
all('trans-' not in y for y
in waterless_titles) and
all('meso-' not in y for y
in waterless_titles) and
all('E' not in y for y in
waterless_titles) and
all('Z' not in y for y in
waterless_titles):
    possibly_the_same_names
    += 1
    RandBlank += 1
else:
    distinct_names += 1
    writer2.write(entry2[a])
    writer3.write(entry3[a])
elif any('D' in y for y in
waterless_titles) or any('L' in
y for y in waterless_titles) or
any('R' in y for y in
waterless_titles) or any('S' in
y for y in waterless_titles) or
any('(+)-' in y for y in
waterless_titles) or any('(-)-'
in y for y in waterless_titles)
or any('cis-' in y for y in
waterless_titles) or
any('trans-' in y for y in
waterless_titles) or any('E' in
y for y in waterless_titles) or
any('Z' in y for y in
waterless_titles):
    if all('rac-' not in x for
x in hydrate_titles) and
all('+-' not in x for x in
hydrate_titles) and all('D'
not in x for x in
hydrate_titles) and all('L'
not in x for x in
hydrate_titles) and all('R'
not in x for x in
hydrate_titles) and all('S'
not in x for x in
hydrate_titles) and
all('(+)-' not in x for x
in hydrate_titles) and
all('(-)-' not in x for x
in hydrate_titles) and
all('cis-' not in x for x
in hydrate_titles) and
all('trans-' not in x for x
in hydrate_titles) and
all('meso-' not in x for x
in hydrate_titles) and
all('E' not in x for x in
hydrate_titles) and all('Z'
not in x for x in
hydrate_titles):
        possibly_the_same_names
        += 1
        RandBlank += 1

```

```

313                                         hydrate_titles):
314                                         possibly_the_same_names
315                                         += 1
316                                         RandBlank += 1
317                                         else:
318                                         distinct_names += 1
319                                         writer2.write(entry2[a])
320                                         writer3.write(entry3[a])
321                                         #Structures with distinct
322                                         chirality in their names are
323                                         found here and are written to
324                                         output files for matches with
325                                         distinct chirality
326                                         else:
327                                         distinct_names += 1
328                                         writer2.write(entry2[a])
329                                         writer3.write(entry3[a])
330                                         #This script finds the cases where
331                                         none of the molecule components
332                                         have formulas that match the entry
333                                         formula
334                                         #The first part looks at structures
335                                         where there are no chiral
336                                         configurations of the molecules
337                                         either
338                                         if possibly_the_same_names != 0 and
hydrate_tags == []:
339                                         if hash(entry2[a].identifier)
340                                         in
341                                         hydrate_tags_empty_no_stereo_hydr
342                                         ates:
343                                         indices = [i for i, x in
344                                         enumerate(hydrate_tags_empty_
345                                         no_stereo_hydrates) if x ==
346                                         hash(entry2[a].identifier)]
347                                         if
348                                         any(hydrate_tags_empty_no_stereo_
349                                         _waterless_forms[j] ==
350                                         hash(entry3[a].identifier)
351                                         for j in indices) == True:
352                                         #If the two structures
353                                         have names where both
354                                         structures indicate the
355                                         chirality but are not
356                                         the same, the
357                                         structures are written
358                                         to two output files to
359                                         be manually checked
360                                         if DoubleR != 0:
361                                         writer8.write(entry2[
362                                         a])
363                                         writer9.write(entry3[
364                                         a])
365                                         #Otherwise they are
366                                         written to two output
367                                         files as matches with
368                                         the same chirality
369                                         else:
370                                         chiral_writer.write(e
371                                         ntry2[a])
372                                         writer1.write(entry3[
373                                         a])

```

```

339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
else:
    print
        entry2[a].identifier
    print
        entry3[a].identifier
#In this case there is only one
structure with chiral
configurations of the atoms
elif hash(entry2[a].identifier)
in
hydrate_tags_empty_one_stereo_hyd
rates and
hydrate_tags_empty_one_stereo_wat
erless_forms[hydrate_tags_empty_o
ne_stereo_hydrates.index(hash(ent
ry2[a].identifier))] ==
hash(entry3[a].identifier):
    #If the two structures have
    names where both structures
    indicate the chirality but
    are not the same, the
    structures are written to
    two output files to be
    manually checked
    if DoubleR != 0:
        writer8.write(entry2[a])
        writer9.write(entry3[a])
    else:
        #Otherwise they are
        written to two output
        files as matches with
        the same chirality
chiral_writer.write(entry
2[a])
writer1.write(entry3[a])
#If both names have no
chirality, the
structures are written
to additional output
files for structures
with achiral names and
only one structure with
a chiral configuration
if DoubleBlank != 0:
    writer4.write(entry2[
a])
    writer5.write(entry3[
a])
#If only one name has
chirality, the
structures are written
to additional output
files for structures
with only one chiral
name structure and only
one structure with a
chiral configuration
if RandBlank != 0:
    writer6.write(entry2[
a])
    writer7.write(entry3[
a])

```



388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

a]))

writer5.write(entry3[a])

if RandBlank != 0:

writer6.write(entry2[a])

writer7.write(entry3[a])

if hydrate\_chirality == [] and waterless\_chirality == []:

if DoubleR != 0:

writer8.write(entry2[a])

writer9.write(entry3[a])

else:

chiral\_writer.write(entry2[a])

writer1.write(entry3[a])

#If there are chiral configurations recorded for both structures, they are compared against each other #This is done by finding the same molecule for each structure #Once that is done, the atoms in that molecule are compared to find the same atom #The element is first compared, then the bonds to that element are compared in a branching out manner until only one atom remains that matches the target atom of the other structure #If these matching atoms have the same chirality, the next atom in the list is checked

if hydrate\_chirality != [] and waterless\_chirality != []:

#same molecules = have the same chiral configuration (R, S, Mixed, or none)

#possibly the same molecules = one has a chiral configuration for at least one atom that the other has no chiral configuration for

#distinct molecules = at least one atom has a different chiral configuration in the molecule (R for one, S for the other)

same\_molecules = 0

possibly\_the\_same\_molecules = 0

distinct\_molecules = 0

for c in

range(len(hydrate\_tags)):

for d in

range(len(waterless\_tags)):

):

```

415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
if
entry2[a].molecule.co
mponents[hydrate_tags
[c]].formula ==
entry3[a].molecule.co
mponents[waterless_ta
gs[d]].formula:
ordered_hydrate_c
hiral_centers =
[]
ordered_waterless
_chiral_centers
= []
atomic_symbols
= []
for e in
range(len(entry2[
a].molecule.compo
nents[hydrate_tag
s[c]].atoms)):
if
entry2[a].molecule.components[hydrate_tags[c]].atoms[e].atomic_symbol not in
atomic_symbols:
atomic_symbols.append(entry2[a].molecule.components[hydrate_tags[c]].atoms[e].atomic_symb
ol)
same_elements = [x for x in entry2[a].molecule.components[hydrate_tags[c]].atoms if
x.atomic_symbol == entry2[a].molecule.components[hydrate_tags[c]].atoms[e].atomic_symbol]
chiral_centers = []
for f
in range(len(same_elements)):
if
len(same_elements[f].bonds) == 4:
string1 = str(same_elements[f].bonds)
string1 = string1.replace(str(same_elements[f]), '')
if ') )' in string1:
string1 = string1.replace(') )', ')')
if ' Atom' in string1:
string1 = string1.replace(' Atom', ' Atom')
string1 = re.sub(r'[0-9]+', '', string1)
capitals = [i for i in range(len(string1)) if string1[i].isupper() == True]
extra = [j for j in capitals if j-1 in capitals]
if extra != []:
if len(extra) == 1:
string2 = string1[:extra[0]]
omit_to = string1[extra[0]:].index(')')
string2 = string2 + string1[extra[0]+omit_to:]
else:

```

```

441     for k in range(len(extra)):
442         if k == 0:
443             string2 = string1[:extra[k]]
444             omit_to = string1[extra[k]:].index(' ')
445         elif k == len(extra)-1:
446             string2 = string2 + string1[extra[k-1]+omit_to:extra[k]]
447             omit_to = string1[extra[k]:].index(' ')
448             string2 = string2 + string1[extra[k]+omit_to:]
449         else:
450             string2 = string2 + string1[extra[k-1]+omit_to:extra[k]]
451             omit_to = string1[extra[k]:].index(' ')
452         if string2.count('Atom(H)') <= 1 and string2.count('Atom(D)') == 0 or
453             string2.count('Atom(H)') == 0 and string2.count('Atom(D)') <= 1:
454             chiral_centers.append(same_elements[f])
455         else:
456             if string1.count('Atom(H)') <= 1 and string1.count('Atom(D)') == 0 or
457                 string1.count('Atom(H)') == 0 and string1.count('Atom(D)') <= 1:
458                 chiral_centers.append(same_elements[f])
459             ordered_hydrate_chiral_centers = ordered_hydrate_chiral_centers + chiral_centers
460             start = 0
461             waterless_matches = []
462             while len(ordered_hydrate_chiral_centers) != len(ordered_waterless_chiral_centers):
463                 if
464                     waterless_matches == []:
465                     symbol_matches = [x for x in entry3[a].molecule.components[waterless_tags[d]].atoms if
466                         x.atomic_symbol == ordered_hydrate_chiral_centers[start].atomic_symbol]
467                     for t
468                     in range(len(symbol_matches)):
469                         if len(symbol_matches[t].bonds) == 4:
470                             string3 = str(symbol_matches[t].bonds)
471                             string3 = string3.replace(str(symbol_matches[t]), '')
472                             if ') )' in string3:
473                                 string3 = string3.replace(') )', ')')
474                             if ' Atom' in string3:
475

```

```

        string3 = string3.replace(' Atom', ' Atom')
471
string3 = re.sub(r'[0-9]+', '', string3)
472
capitals = [i for i in range(len(string3)) if string3[i].isupper() == True]
473
extra = [j for j in capitals if j-1 in capitals]
474
if extra != []:
475
    if len(extra) == 1:
476
        string4 = string3[:extra[0]]
477
        omit_to = string3[extra[0]:].index(' ')
478
        string4 = string4 + string3[extra[0]+omit_to:]
479
    else:
480
        for k in range(len(extra)):
481
            if k == 0:
482
                string4 = string3[:extra[k]]
483
                omit_to = string3[extra[k]:].index(' ')
484
            elif k == len(extra)-1:
485
                string4 = string4 + string3[extra[-1]+omit_to:extra[k]]
486
                omit_to = string3[extra[k]:].index(' ')
487
                string4 = string4 + string3[extra[k]+omit_to:]
488
            else:
489
                string4 = string4 + string3[extra[-1]+omit_to:extra[k]]
490
                omit_to = string3[extra[k]:].index(' ')
491
        if string4.count('Atom(H)') <= 1 and string4.count('Atom(D)') == 0 or
string4.count('Atom(H)') == 0 and string4.count('Atom(D)') <= 1:
492
            waterless_matches.append(symbol_matches[t])
493
        else:
494
            if string3.count('Atom(H)') <= 1 and string3.count('Atom(D)') == 0 or
string3.count('Atom(H)') == 0 and string3.count('Atom(D)') <= 1:
495
            waterless_matches.append(symbol_matches[t])
496
waterless_bonds = []
497
updated_waterless_bonds = []
498
for f
in range(len(waterless_matches)):
499
pre_waterless_bonds = []
500
for
y in range(len(waterless_matches[f].bonds)):
501
string2 = str(waterless_matches[f].bonds[y])
502
string2 = string2.replace(str(waterless_matches[f]), '')

```

```

503
504     if ') )' in string2:
505
506         string2 = string2.replace(') )', ')')
507
508     if ' Atom' in string2:
509
510         string2 = string2.replace(' Atom', ' Atom')
511
512     string2 = re.sub(r'[0-9]+', '', string2)
513
514     pre_waterless_bonds.append(string2)
515
516     pre_waterless_bonds = sorted(pre_waterless_bonds)
517
518     waterless_bonds.append(pre_waterless_bonds)
519
520     updated_waterless_bonds.append(pre_waterless_bonds)
521
522     else:
523
524         waterless_bonds = list(updated_waterless_bonds)
525
526     len(waterless_matches) == 1:
527
528         ordered_waterless_chiral_centers.append(waterless_matches[0])
529
530         waterless_matches.remove(waterless_matches[0])
531
532         start
533
534         += 1
535
536     else:
537
538
539         hydrate_spider_atoms = []
540
541         waterless_spider_atoms = []
542
543         found_hydrate_atoms = []
544
545         hydrate_bonds = []
546
547         for z
548
549             in range(len(ordered_hydrate_chiral_centers[start].bonds)):
550
551                 string1 = str(ordered_hydrate_chiral_centers[start].bonds[z])
552
553                 string1 = string1.replace(str(ordered_hydrate_chiral_centers[start]), '')
554
555                 ') )' in string1:
556
557                     string1 = string1.replace(') )', ')')
558
559                     if
560
561                         ' Atom' in string1:
562
563                             string1 = string1.replace(' Atom', ' Atom')
564
565                             string1 = re.sub(r'[0-9]+', '', string1)
566
567                             hydrate_bonds.append(string1)
568
569                             hydrate_bonds = sorted(hydrate_bonds)
570
571                             hydrate_neighbours = ordered_hydrate_chiral_centers[start].neighbours
572
573                             found_hydrate_atoms.append(ordered_hydrate_chiral_centers[start])
574
575                             numberless_hydrate_neighbours = []
576
577                             numberless_waterless_neighbours = []

```

```

538     found_match = 0
539
540     found_match == 0:
541
542     bond_matches = []
543     atom_matches = []
544
545     numberless_hydrate_neighbours == []:
546
547     for w in range(len(waterless_bonds)):
548
549         if waterless_bonds[w] == hydrate_bonds:
550
551             bond_matches.append(waterless_bonds[w])
552
553             atom_matches.append(w)
554
555         else:
556
557             if hydrate_bonds == [] and all(x == [] for x in waterless_bonds):
558
559                 bond_matches.append(waterless_bonds[w])
560
561                 atom_matches.append(w)
562
563             else:
564
565                 for w in range(len(numberless_waterless_neighbours)):
566
567                     off_the_table_waterless = []
568
569                     off_the_table_hydrate = []
570
571                     match = 0
572
573                     for x in range(len(numberless_waterless_neighbours[w])):
574
575                         if len(waterless_bonds[w]) == len(hydrate_bonds):
576
577                             if waterless_neighbours[w][x] not in off_the_table_waterless:
578
579                                 go = 0
580
581                                 for u in range(len(numberless_hydrate_neighbours)):
582
583                                     if go == 0:
584
585                                         if hydrate_neighbours[u] not in off_the_table_hydrate:
586
587                                             if numberless_waterless_neighbours[w][x] ==
588                                                 numberless_hydrate_neighbours[u]:
589
590                                                 waterless_to_compare = []
591
592                                                 hydrate_to_compare = []
593
594                                                 if x == 0:
595
596                                                     for t in range(waterless_spider_atoms[w][x]):
597
598                                                         waterless_to_compare.append(waterless_bonds[w][t])
599
600                                                     else:
601
602                                                         for t in range(waterless_spider_atoms[w][x-1],
603
604                                                 waterless_spider_atoms[w][x]):
```

```

570
571     waterless_to_compare.append(waterless_bonds[w][t])
572         if u == 0:
573             for s in range(hydrate_spider_atoms[u]):
574                 hydrate_to_compare.append(hydrate_bonds[s])
575         else:
576             for s in range(hydrate_spider_atoms[u-1],
577                             hydrate_spider_atoms[u]):
578                 hydrate_to_compare.append(hydrate_bonds[s])
579
580         if waterless_to_compare == hydrate_to_compare:
581
582             off_the_table_waterless.append(waterless_neighbours[w][x])
583
584             off_the_table_hydrate.append(hydrate_neighbours[u])
585                 match += 1
586
587                 go += 1
588
589                 if match == len(numberless_hydrate_neighbours):
590                     bond_matches.append(waterless_bonds[w])
591
592                     atom_matches.append(w)
593
594
595 #When the length of bond matches is one that means the target atom only has one match,
596 making that matching atom its equivalent in the other structure
597
598         if len(bond_matches) == 1:
599
600             pin2 = waterless_bonds.index(bond_matches[0])
601
602             ordered_waterless_chiral_centers.append(waterless_matches[pin2])
603
604             waterless_matches.remove(waterless_matches[pin2])
605
606             updated_waterless_bonds.remove(updated_waterless_bonds[pin2])
607
608             start += 1
609
610             found_match += 1
611
612         else:
613
614             if numberless_hydrate_neighbours != []:
615
616                 hydrate_neighbours = list(new_hydrate_neighbours)
617
618                 numberless_hydrate_neighbours = []
619
620                 hydrate_spider_atoms = []
621
622                 for g in range(len(hydrate_neighbours)):
623
624                     string1 = str(hydrate_neighbours[g])

```

```

        string2 = re.sub(r'[0-9]+', '', string1)
602    numberless_hydrate_neighbours.append(string2)
603
604    new_hydrate_neighbours = []
605    hydrate_atom_branches = 0
606
607    for j in range(len(hydrate_neighbours)):
608
608        if hydrate_neighbours[j] not in found_hydrate_atoms:
609
609            for k in range(len(hydrate_neighbours[j].neighbours)):
610
610                if hydrate_neighbours[j].neighbours[k] not in found_hydrate_atoms:
611
611                    new_hydrate_neighbours.append(hydrate_neighbours[j].neighbours[k])
612
612                    hydrate_atom_branches += 1
613
613    hydrate_spider_atoms.append(hydrate_atom_branches)
614
614    hydrate_bonds = []
615
615    for l in range(len(hydrate_neighbours)):
616
616        pre_hydrate_bonds = []
617
617        for z in range(len(hydrate_neighbours[l].bonds)):
618
618            if all(str(x) not in str(hydrate_neighbours[l].bonds[z]) for x in
619            found_hydrate_atoms):
620
620                string1 = str(hydrate_neighbours[l].bonds[z])
621
621                string1 = string1.replace(str(hydrate_neighbours[l]), '')
622
622                if ') )' in string1:
623
623                    string1 = string1.replace(') )', ')')
624
624                if ' Atom' in string1:
625
625                    string1 = string1.replace(' Atom', ' Atom')
626
626                string1 = re.sub(r'[0-9]+', '', string1)
627
627                pre_hydrate_bonds.append(string1)
628
628    pre_hydrate_bonds = sorted(pre_hydrate_bonds)
629
629    hydrate_bonds = hydrate_bonds + pre_hydrate_bonds
630
630    if numberless_waterless_neighbours != []:
631
631        waterless_neighbours = list(new_waterless_neighbours)
632
632        waterless_neighbours_to_remove = []
633
633        for y in range(len(waterless_neighbours)):
634
634            if y not in atom_matches:
635
635                waterless_neighbours_to_remove.append(waterless_neighbours[y])
636
636        for z in range(len(waterless_neighbours_to_remove)):
637
637

```

```

        waterless_neighbours.remove(waterless_neighbours_to_remove[z])
635
636     numberless_waterless_neighbours = []
637     waterless_spider_atoms = []
638
639     else:
640
641         found_waterless_atoms = []
642
643         waterless_neighbours = []
644
645         for y in range(len(waterless_matches)):
646
647             waterless_neighbours.append(waterless_matches[y].neighbours)
648
649             if str(waterless_matches[y]) != 'Atom(H)':
650
651                 found_waterless_atoms.append([waterless_matches[y]])
652
653         for f in range(len(waterless_neighbours)):
654
655             pre_numberless_waterless_neighbours = []
656
657             for e in range(len(waterless_neighbours[f])):
658
659                 string1 = str(waterless_neighbours[f][e])
660
661                 string2 = re.sub(r'[0-9]+', '', string1)
662
663                 pre_numberless_waterless_neighbours.append(string2)
664
665             numberless_waterless_neighbours.append(pre_numberless_waterless_neighbours)
666
667             new_waterless_neighbours = []
668
669             for m in range(len(waterless_neighbours)):
670
671                 waterless_atom_branches = 0
672
673                 pre_new_waterless_neighbours = []
674
675                 pre_waterless_spider_atoms = []
676
677                 for n in range(len(waterless_neighbours[m])):
678
679                     if waterless_neighbours[m][n] not in found_waterless_atoms[m]:
680
681                         for o in range(len(waterless_neighbours[m][n].neighbours)):
682
683                             if waterless_neighbours[m][n].neighbours[o] not in
684                             found_waterless_atoms[m]:
685
686                                 pre_new_waterless_neighbours.append(waterless_neighbours[m][n].neighbours[o])
687
688                                 waterless_atom_branches += 1
689
690                                 pre_waterless_spider_atoms.append(waterless_atom_branches)
691
692             new_waterless_neighbours.append(pre_new_waterless_neighbours)
693
694             waterless_spider_atoms.append(pre_waterless_spider_atoms)
695
696             waterless_bonds = []
697
698             for p in range(len(waterless_neighbours)):

```



```

699
each one is
recorded

700
hydrate_chirality
= []

701
waterless_chirali
ty = []
for e in
range(len(ordered
_hydrate_chiral_c
enters)):
    if

702
ordered_hydrate_chiral_centers[e].is_chiral == True:
703
hydrate_chirality.append(ordered_hydrate_chiral_centers[e].chirality)
    else:

704
705
hydrate_chirality.append(' ')
706

707
ordered_waterless_chiral_centers[f].is_chiral == True:
708
waterless_chirality.append(ordered_waterless_chiral_centers[f].chirality)
    else:
709
710
waterless_chirality.append(' ')
711
    if
hydrate_chirality
 ==
waterless_chirali
ty:
    elif [s for s
in
hydrate_chirality
if s != ' ']
== [] or [t for
t in
waterless_chirali
ty if t != ' ']
== []:
    else:
        possibly_the_same_molecules += 1
    else:
        distinct_molecules += 1
    else:
        #No distinct molecules but
        possibly the same molecules
        relies on the name as a
        secondary check as has been
        seen before
        if distinct_molecules == 0:
            if
possibly_the_same_molecul
es == 0:
                if DoubleR != 0:
                    writer8.write(ent
ry2[a])
                    writer9.write(en

```

```

724     ry3[a])
725
726         chiral_writer.write(entry2[a])
727
728         writer1.write(entry3[a])
729
730     #No distinct molecules
731     #and no possibly the
732     #same molecules also
733     #relies on the name as a
734     #secondary check as has
735     #been seen before
736     else:
737         if DoubleR != 0:
738             writer8.write(entry2[a])
739
740             writer9.write(entry3[a])
741
742         else:
743             chiral_writer.write(entry2[a])
744
745             writer1.write(entry3[a])
746
747             writer4.write(entry2[a])
748
749             writer5.write(entry3[a])
750
751             if RandBlank != 0:
752                 writer6.write(entry2[a])
753
754                 writer7.write(entry3[a])
755
756             else:
757                 #Any distinct molecules
758                 #implies that the
759                 #chirality of each
760                 #structure is distinct,
761                 #the structures are
762                 #written to output files
763                 #for matches with
764                 #distinct chirality
765                 writer2.write(entry2[a])
766                 writer3.write(entry3[a])
767
768             print 'another one bites the dust'
769
770         if __name__ == '__main__':
771             # This runs the script
772             r = Runner()
773             r.run()

```

```

1 #20 Remove duplicates with distinct chirality from list of duplicates
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script was used interchangably to update the list of first occurrence and
10 #duplicate structures
11 #Filepath1 and filepath2 are always the list of first occurrence and duplicate
12 #structures in its current state
13 #Filepath3 and filepath4 are added to the list of first occurrence and duplicate
14 #structures (this was done after script #14)
15 #Filepath5 and filepath6 are removed from the list of first occurrence and duplicate
16 #structures (this was done after script #18)
17 #The same script was used for hydrate and waterless forms, substituting for each in the
18 #output filenames
19
20 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_first_occurrence.txt"
21 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_all_duplicates.txt"
22
23 #filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_None_entry_to_keep.txt"
24 #filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_None_duplicates_to_remove.txt"
25 #filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_smiles_entry_to_keep.txt"
26 #filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\hydrates_smiles_duplicates_to_remove.txt"
27
28 #filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_formula_hydrate_entry.txt"
29 #filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\check_formula_hydrate_duplicate.txt"
30
31 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step4 Find Pairs and Duplicates with
Distinct Chirality\distinct_chirality_hydrate_entry.txt"
32 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step4 Find Pairs and Duplicates with
Distinct Chirality\distinct_chirality_hydrate_duplicate.txt"
33
34 #filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_first_occurrence.txt"
35 #filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_all_duplicates.txt"
36 #filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_None_entry_to_keep.txt"
37 #filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_None_duplicates_to_remove.txt"
38 #filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_smiles_entry_to_keep.txt"
39 #filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
Entries\waterless_forms_smiles_duplicates_to_remove.txt"
40
41 class Runner(argparse.ArgumentParser):
42

```

```

43     def __init__(self):
44         super(self.__class__, self).__init__(description=__doc__)
45         self.add_argument(
46             '-i', '--input', default=filepath1,
47             help='input database filepath1'
48         )
49         self.add_argument(
50             '-o', '--output', default='hydrates_all_first_occurrence.gcd',
51             help='output file [hydrates_all_first_occurrence.gcd]'
52         )
53         self.add_argument(
54             '-m', '--maximum', default=0, type=int,
55             help='Maximum number of structures to find [all]'
56     )
57
58     args = self.parse_args()
59
60     self.args = args
61     self.settings = search.Search.Settings()
62     self.settings.max_hit_structures = self.args.maximum
63
64     def run(self):
65
66         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
67         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
68         #entry_reader3 = io.EntryReader(filepath3, format='identifiers')
69         #entry_reader4 = io.EntryReader(filepath4, format='identifiers')
70         entry_reader5 = io.EntryReader(filepath5, format='identifiers')
71         entry_reader6 = io.EntryReader(filepath6, format='identifiers')
72
73         with io.EntryWriter(self.args.output) as new_lists_writer:
74             with io.EntryWriter("hydrates_all_duplicates.gcd") as writer1:
75
76                 duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
77                 'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
78                 'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14,
79                 'ZOVFEV':15}
80
81                 #Create lists of first occurrences and duplicates that need to be removed
82                 first_occurrence = []
83                 for a in range(len(entry_reader5)):
84                     if entry_reader5[a].identifier in duplicate_identifier_dictionary:
85
86                         first_occurrence.append(duplicate_identifier_dictionary.get(entry_reader5[a].identifier))
87                     else:
88                         first_occurrence.append(hash(entry_reader5[a].identifier))
89
90                 duplicate = []
91                 for b in range(len(entry_reader6)):
92                     if entry_reader6[b].identifier in duplicate_identifier_dictionary:
93
94                         duplicate.append(duplicate_identifier_dictionary.get(entry_reader6[b].identifier))
95                     else:
96                         duplicate.append(hash(entry_reader6[b].identifier))
97
98                 #Check whether each match in the current list of first occurrences and
99                 #duplicates is on the list to be removed
100                for a in range(len(entry_reader1)):
101                    skip = 0
102                    if entry_reader1[a].identifier in duplicate_identifier_dictionary:
103                        query =
104                            duplicate_identifier_dictionary.get(entry_reader1[a].identifier)
105                    else:
106                        query = hash(entry_reader1[a].identifier)
107                    if entry_reader2[a].identifier in duplicate_identifier_dictionary:
108
109
```

```

101         query2 =
102         duplicate_identifier_dictionary.get(entry_reader2[a].identifier)
103     else:
104         query2 = hash(entry_reader2[a].identifier)
105     #If both structures are in the corresponding lists for removal
106     #Whether or not these two structures appear on the same line
107     #dictates if they are removed
108     if query in first_occurrence and query2 in duplicate:
109         indice = [i for i, x in enumerate(first_occurrence) if x ==
110             query and duplicate[i] == query2]
111         if len(indice) == 1:
112             skip += 1
113
114         #The new output file updates the first occurrence and duplicate
115         #list so it does not contain any matches that needed to be removed
116         if skip == 0:
117             new_lists_writer.write(entry_reader1[a])
118             writer1.write(entry_reader2[a])
119
120         #Matches that need to be added are simply appended to the output files
121         #for the new list of first occurrences and duplicates
122         #for c in range(len(entry_reader3)):
123             #new_lists_writer.write(entry_reader3[c])
124             #writer1.write(entry_reader4[c])
125
126
127
128
129
130
131
132
133 if __name__ == '__main__':
134     r = Runner()
135     r.run()
136

```

```

1 #21 Find hydrates with waterless forms using SMILES string method
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds hydrate-anhydrate pairs from structures with an entry SMILES string
10
11 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\smiles_solvated_hydrates.txt"
12 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\smiles_solvated_waterless_forms.txt"
13 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\smiles_hydrates.txt"
14 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
Forms\smiles_waterless_forms.txt"
15
16 class Runner(argparse.ArgumentParser):
17
18     def __init__(self):
19         super(self.__class__, self).__init__(description=__doc__)
20         self.add_argument(
21             '-i', '--input', default=filepath1,
22             help='input database filepath'
23         )
24         self.add_argument(
25             '-o', '--output', default='potential_hydrates_with_waterless_form.gcd',
26             help='output file [potential_hydrates_with_waterless_form.gcd]'
27         )
28         self.add_argument(
29             '-m', '--maximum', default=0, type=int,
30             help='Maximum number of structures to find [all]'
31         )
32
33     args = self.parse_args()
34
35     self.args = args
36     self.settings = search.Search.Settings()
37     self.settings.max_hit_structures = self.args.maximum
38
39     def run(self):
40
41         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
42         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
43         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
44         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
45
46         total = 0
47         count = 0
48
49         with io.EntryWriter(self.args.output) as pairing_writer:
50             with io.EntryWriter("potential_waterless_forms_with_hydrate.gcd") as writer1:
51                 with io.EntryWriter("potential_hydrates_without_waterless_form.gcd") as
writer2:
52                     with
53                         io.EntryWriter("potential_waterless_forms_without_hydrate.gcd") as
writer3:
54
solvents_dictionary = {hash('but-2-ene'):'CC=CC',
hash('1,1,2-dichloroethane'):'ClCC(Cl)Cl',
hash('4-aminopyridine'):'Nc1ccncc1',
hash('N-methylurea'):'CNC(N)=O',
hash('2-methylphenol'):'Cc1cccc1O',
hash('tetraglyme'):'COCCOCCOCOCOCOC',
hash('isophorone'):'CC1=CC(=O)CC(C)(C)C1',

```

```

hash('L-carvone'):'CC(=C)C1CC=C(C)C(=O)C1',
hash('2-cyanopyridine'):'N#Cc1ccccc1',
hash('biuret'):'NC(=O)NC(N)=O', hash('styrene'):'C=Cc1cccc1',
hash('carbondisulfide'):'S=S=S',
hash('tricyanomethanide'):'N#C[C-](C#N)C#N',
hash('methyl-cyclohexane'):'CC1CCCCC1',
hash('pyridine-3-carboxamide'):'NC(=O)c1cccncl',
hash('3,7-dimethyl-3,7-dihydro-1H-purine-2,6-dione'):'CN1C=NC2=C1
C(=O)NC(=O)N2C', hash('nitroethane'):'CCN(=O)=O',
hash('dimethylaminobenzaldehyde'):'CN(C)c1ccc(C=O)cc1',
hash('3-(1H-pyrrol-1-yl)propanenitrile'):'Cl=CN(C=C1)CCC#N',
hash('trifluoromethanol'):'OC(F)(F)F',
hash('methoxyethane'):'CCOC',
hash('methylanisole'):'COc1ccccc1C',
hash('xylene'):'Cc1ccc(C)cc1', hash('petrol'):'CCCCCC',
hash('n-alkane'):'CCCCC', hash('dichloromethanal'):'ClC(Cl)=O',
hash('1-fluoro-4-methylbenzene'):'Cc1ccc(F)cc1',
hash('1,2,3-trimethoxybenzene'):'COc1ccccc(OC)c1OC',
hash('1,3,5-trimethoxybenzene'):'COc1cc(OC)cc(OC)c1',
hash('chloroform'):'C1C(Cl)Cl', hash('methanol'):'CO',
hash('hexane'):'CCCCCC', hash('ethanol'):'CCO',
hash('acetonitrile'):'CC#N', hash('toluene'):'Cc1ccccc1',
hash('tetrahydrofuran'):'C1CCOC1', hash('acetone'):'CC(C)=O',
hash('dichloromethane'):'ClCCl',
hash('1,2-dichloroethane'):'ClCCCCl',
hash('dimethylsulfoxide'):'CS(C)=O',
hash('1,4-dioxane'):'C1COCCO1',
hash('1,2-dimethoxyethane'):'COCCOC',
hash('o-xylene'):'Cc1ccccc1C', hash('propanol'):'CCO',
hash('dioxane'):'C1COCCO1', hash('n-pentane'):'CCCCCC',
hash('4-cyanopyridine'):'N#Cc1ccncc1',
hash("4,4'-bipyridine"):'c1cc(ccn1)c1ccncc1',
hash('isopropanol'):'CC(C)O', hash('1-butanol'):'CCCCO',
hash('nitromethane'):'CN(=O)=O', hash('n-hexane'):'CCCCCC',
hash('dimethylformamide'):'CN(C)C=O',
hash('N,N-dimethylformamide'):'CN(C)C=O',
hash('cyclohexane'):'C1CCCCC1',
hash('hydroquinone'):'Oc1ccc(O)cc1',
hash('benzene'):'c1ccccc1', hash('octane'):'CCCCCC',
hash('cyclohexanone'):'O=C1CCCCC1',
hash('isobutanol'):'CC(C)CO', hash('p-xylene'):'Cc1ccc(C)cc1',
hash('pentane'):'CCCCCC', hash('pyridine'):'c1ccncc1',
hash('t-butanol'):'CC(C)(C)O', hash('m-xylene'):'Cc1ccccc(C)c1',
hash('phenol'):'Oc1ccccc1', hash('o-cresol'):'Cc1ccccc1O',
hash('p-cresol'):'Cc1ccc(O)cc1',
hash('2,4,6-trimethylpyridine'):'Cc1cc(C)nc(C)c1',
hash('1-methylpyrrolidin-2-one'):'CN1CCCC1=O',
hash('1,1,2,2,3,3,4,4,5,5,6,6,7,7,8,8-hexadecafluoro-1,8-diiodooctane'):'FC(F)(I)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)C(F)(F)I', hash('trifluoroethanol'):'OCC(F)(F)F',
hash('dichloroethane'):'ClCCl',
hash('1,1,2,2-tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('2-propanol'):'CC(C)O',
hash('chlorobenzene'):'Clc1ccccc1',
hash('naphthalene'):'c1ccc2ccccc2c1',
hash('propan-2-ol'):'CC(C)O',
hash('benzonitrile'):'N#Cc1ccccc1',
hash('nitrobenzene'):'O=N(=O)c1ccccc1',
hash('dodecane'):'CCCCCCCCCCCC',
hash('1,2-dichlorobenzene'):'Clc1ccccc1Cl',
hash('n-heptane'):'CCCCCC',
hash('methanesulfonate'):'CS(=O)(=O)[O-]',
hash('3-methylbutan-1-ol'):'CC(C)CCO',
hash('1,1,2-trichloroethane'):'ClCC(Cl)Cl',
hash('n-butane'):'CCCC', hash('dichlorine'):'ClCl',
hash('o-dichlorobenzene'):'Clc1ccccc1Cl',
hash('N,N-dimethylacetamide'):'CN(C)C(C)=O',

```

```

hash('diethylether'):'CCOCC', hash('hydrazine'):'NN',
hash('dimethylamine'):'CNC',
hash('1,2-difluorobenzene'):'Fc1cccc1F',
hash('cycloheptane'):'C1CCCCC1', hash('formamide'):'NC=O',
hash('2-butanol'):'CCC(C)O', hash('ethylacetate'):'CCOC(C)=O',
hash('butan-1-ol'):'CCCCO', hash('bromomethane'):'CBr',
hash('piperazine'):'C1CNCCN1', hash('n-nonane'):'CCCCCC',
hash('1,2,4-trimethoxybenzene'):'COc1ccc(OC)c(OC)c1',
hash('1,2-diaminoethane'):'NCCN', hash('n-butanol'):'CCCCO',
hash('N-methylformamide'):'CNC=O',
hash('pyrrolidine'):'C1CCNC1', hash('actone'):'CC(C)=O',
hash('n-octane'):'CCCCCC',
hash('S)-2-methyl-1-butanol'):'CCC(C)CO',
hash('N,N-diethylformamide'):'CCN(CC)C=O',
hash('tetracyanoethylene'):'N#CC(C#N)=C(C#N)C#N',
hash('fluorobenzene'):'Fc1cccc1', hash('heptane'):'CCCCCC',
hash('tetrachloroethane'):'ClC(Cl)C(Cl)Cl',
hash('methylcyclohexane'):'CC1CCCC1',
hash('tetrachloromethane'):'ClC(Cl)(Cl)Cl',
hash('1-propanol'):'CCCO', hash('acetamide'):'CC(N)=O',
hash('propane'):'CCC', hash('diethylamine'):'CCNCC',
hash('hexafluorophosphate'):'F[P-](F)(F)(F)F',
hash('hexafluoropropan-2-ol'):'OC(C(F)(F)F)C(F)(F)F',
hash('bis(hexafluorophosphate)'):'F[P-](F)(F)(F)F',
hash('2-methyl-1,3,5-trinitrobenzene'):'Cc1c(cc(cc1N(=O)=O)=O)=O',
hash('furan'):'O1C=CC=C1',
hash('2-methylpentane'):'CCCC(C)C',
hash('isohexane'):'CCCC(C)C', hash('anisole'):'COc1cccc1',
hash('butanol'):'CCCO', hash('propan-1-ol'):'CCCO',
hash('dimethoxyethane'):'COCCOC',
hash('cyclopentane'):'C1CCCC1',
hash('methoxybenzene'):'COc1cccc1',
hash('toluene'):'Cc1cccc1',
hash('pyrene'):'c1cc2cc3cccc4ccc(c1)c2c34',
hash('butan-2-ol'):'CCC(C)O',
hash('3-methylpyridine'):'Cc1cccnc1',
hash('triethylamine'):'CCN(CC)CC', hash('dioxan'):'C1COCCO1',
hash('2-methoxyethanol'):'COCCO',
hash('4,4'-bipyridyl'):'c1cc(ccn1)c1ccncc1',
hash('dimethylacetamide'):'CN(C)C(C)=O',
hash('thiophene'):'S1C=CC=C1', hash('1,3-dioxolane'):'C1COCO1',
hash('1-chloronaphthalene'):'Clc1cccc2cccc12',
hash('pentan-1-ol'):'CCCCO', hash('hexan-1-ol'):'CCCCCO',
hash('2-phenylethanol'):'OCCc1cccc1',
hash('1-chloro-4-methylbenzene'):'Cc1ccc(Cl)cc1',
hash('acetaldehyde'):'CC=O', hash('propiononitrile'):'CCC#N',
hash('bromobenzene'):'Brclcccc1'}
deuterated_dictionary =
{hash('dideutero-dichloromethane'):'ClCCl',
hash('perdeutero-toluene'):'Cc1cccc1',
hash('octadeuterofuran'):'C1CCOC1',
hash('deutero-ethanol'):'CCO',
hash('deuterochloroform'):'C1C(Cl)Cl',
hash('deutero-chloroform'):'C1C(Cl)Cl',
hash('deuterobenzene'):'c1cccc1',
hash('deuteromethanol'):'CO',
hash('hexadeutero-benzene'):'c1cccc1',
hash('perdeuteroacetone'):'CC(C)=O',
hash('dideuterodichloromethane'):'ClCCl',
hash('perdeuteromethanol'):'CO',
hash('hexadeuterobenzene'):'c1cccc1'}

acids_dictionary = {hash('trifluoroacetic'):'OC(=O)C(F)(F)F',
hash('acetic'):'CC(O)=O', hash('formic'):'OC=O',
hash('sulfuric'):'OS(O)(=O)=O',
hash('maleic'):'OC(=O)C=CC(O)=O',
hash('benzoic'):'OC(=O)c1cccc1',

```

```

57
hash('succinic'):'OC(=O)CCC(O)=O',
ethers_dictionary = {'t-butyl methyl ether':'COC(C)(C)C',
hash('di-isopropyl'):'CC(C)OC(C)C', hash('dipropyl'):'CCCOCCCC',
hash('diethyl'):'CCOCC', hash('diisopropyl'):'CC(C)OC(C)C',
hash('isopropyl'):'CC(C)OC(C)C', hash('t-butyl'):'COC(C)(C)C',
hash('acetic'):'CCOC(C)=O'},
58
acetates_dictionary = {hash('ethyl'):'CCOC(C)=O',
hash('amyl'):'CCCCCCOC(C)=O', hash('methyl'):'COC(C)=O',
hash('butyl'):'CCCCOC(C)=O', hash('isobutyl'):'CC(C)COCC(C)=O'},
glycol_dictionary = {hash('ethylene'):'OCCO'}},
59
oxalates_dictionary = {hash('dipropyl'):'CCCOOC(=O)C(=O)OCCC'}},
60
acid_neutral_dictionary = {hash('oxalic'):'OC(=O)C(O)=O',
hash('nitric'):'ON(=O)=O',
hash('trans-but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('but-2-enedioic'):'OC(=O)C=CC(O)=O',
hash('fumaric'):'OC(=O)C=CC(O)=O',
hash('phosphonic'):'OP(O)=O',
hash('perchloric'):'O[Cl](=O)(=O)=O',
hash('phosphoric'):'OP(O)(O)=O',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(O)=O',
hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(O)=O)c(cc1C(O)=O)C(O)=O',
hash('trifluoroacetic'):'OC(=O)C(F)(F)F'},
61
acid_charged_dictionary = {hash('oxalic'):'OC(=O)C(=O)[O-]',
hash('nitric'):'O=N(=O)[O-]',
hash('trans-but-2-enedioic'):'OC(=O)C=CC(=O)[O-]',
hash('but-2-enedioic'):'OC(=O)C=CC(=O)[O-]',
hash('fumaric'):'OC(=O)C=CC(=O)[O-]',
hash('phosphonic'):'OP(=O)[O-]',
hash('perchloric'):'O=[Cl](=O)(=O)[O-]',
hash('phosphoric'):'OP(O)(=O)[O-]',
hash('pyridine-2,6-dicarboxylic'):'OC(=O)c1cccc(n1)C(=O)[O-]',
hash('benzene-1,2,4,5-tetracarboxylic'):'OC(=O)c1cc(C(=O)[O-])c(cc1C(=O)[O-])C(O)=O',
hash('trifluoroacetic'):'FC(F)(F)C(=O)[O-]'}
62
63
64
partial_name = [hash('ketone'), hash('propanoate'),
hash('carbonate'), hash('formate'), hash('iodine'),
hash('tetrachloride'), hash('sulfoxide'), hash('dichloride'),
hash('chloride'), hash('benzoate'), hash('disulfide'),
hash('oxide')]}
65
full_solvate_name_dictionary = {'methyl isobutyl
ketone':'CC(C)CC(C)=O', 'deuterium oxide':'O', 'ethyl
propanoate':'CCOC(=O)CC', 'dimethyl carbonate':'COC(=O)OC',
'ethyl formate':'CCOC=O', '12]cycloparaphenylene iodine':'II',
'carbon tetrachloride':'ClC(Cl)(Cl)Cl', 'dimethyl
sulfoxide':'CS(C)=O', 'methylene dichloride':'ClCCl',
'methylene chloride':'ClCCl', 'ethyl
benzoate':'CCOC(=O)c1ccccc1', 'carbon disulfide':'S=C=S'}
66
67
hydrates_SMILES = []
68
waterless_forms_SMILES = []
69
hydrates_formulas = []
70
waterless_forms_formulas = []
71
72
duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2,
'OBASAN':3, 'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7,
'VOBOXZ':8, 'VUFJAI':9, 'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12,
'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
73
duplicate_name_dictionary = {"N,N'-bis((Diethoxyphosphinoyl)phenylmethyl)-1,4-diaminobenzene":1,
'bis(n-Hexyl)':2, 'Benzoylmethyl':3,
'(((4-methylphenyl)sulfonyl)amino)(phenyl)acetic':4,
'1,1-diphenylethyl':5, 'hemikis(1-naphthylamine)':6,
'(E)-1-(2-Nitrobenzylidene)-4-phenylthiosemicarbazide':7}
74
duplicate_smiles_dictionary = {'CC(C)(C)c1c(O)c(O)c(Cl)c2[NH+]3CCN(CC3)c12':1,
'CC1=C(C(CC2=NS(=O)(=O)c3cccc23)c2cccc2C)C(=O)N(N1)c1ccccc1':2}

```

```

'OC(=O)COc1cccc2ccccc12':3, 'Cc1c[nH+]cc(C)c1':4,
'Brclccc(NC(=O)c2cccn2)cc1':5,
'CSC(N)=NNC(C)=CC(=O)c1cccc1':6, 'CC1(CO1)C(O)(CBr)CBr':7,
'OCC12CCCC(C1)(C(O)=O)C1(CCSC1)O2':8,
'COc1ccc2CC3[NH2+]CCC3(C)c2c1':9,
'CC(=O)C(P(c1cccc1)c1cccc1)=C(C)Nc1c(F)cccc1F':10}

75
76    loop = 0
77    entryl = entry_reader1
78    tag = 12
79    #The solvent dictionaries are used to correct incomplete SMILES
80    #strings for solvent molecules that are missing
81    while loop < 4:
82        for a in range(len(entryl)):
83            component_SMILES = []
84            entry_SMILES = entryl[a].molecule.smiles.split('.')
85            for b in range(len(entry_SMILES)):
86                if entry_SMILES[b] not in
87                    duplicate_smiles_dictionary:
88                        component_SMILES.append(hash(entry_SMILES[b]))
89                else:
90
91                    component_SMILES.append(duplicate_smiles_dictionary.get(entry_SMILES[b]))
92            component_SMILES =
93            list(set(component_SMILES))
94            if entryl[a].identifier == 'BULVEL':
95                if hash('[O-][n+]1cccc1') not in component_SMILES:
96                    component_SMILES.append(hash('[O-][n+]1cccc1'))
97                if hash('[O]n1cccc1') in component_SMILES:
98                    component_SMILES.remove(hash('[O]n1cccc1'))
99
100           elif entryl[a].identifier == 'CEHKOS':
101               component_SMILES.append(hash('II'))
102           elif entryl[a].identifier == 'GOTJEE' or
103               entryl[a].identifier == 'HASVEF' or
104               entryl[a].identifier == 'MTFBTZ10':
105                   component_SMILES.append(tag)
106                   tag += 1
107           elif entryl[a].identifier == 'JUHJIF':
108               if hash('BrBr') not in component_SMILES:
109                   component_SMILES.append(hash('BrBr'))
110               if hash('[Br]') in component_SMILES:
111                   component_SMILES.remove(hash('[Br]'))
112           elif entryl[a].identifier == 'LOKXUF':
113               if hash('OC(=O)C(F)(F)F') not in component_SMILES:
114                   component_SMILES.append(hash('OC(=O)C(F)(F)F'))
115               if hash('[O]C(=O)C(F)(F)F') in component_SMILES:
116                   component_SMILES.remove(hash('[O]C(=O)C(F)(F)F'))
117           elif entryl[a].identifier == 'OJUNEO':
118               if hash('ClCl') not in component_SMILES:
119                   component_SMILES.append(hash('ClCl'))
120               if hash('[Cl]') in component_SMILES:
121                   component_SMILES.remove(hash('[Cl]'))
122           elif entryl[a].identifier == 'ZILFILM':
123               component_SMILES.append(hash('CCO'))
124           else:
125               if entryl[a].chemical_name != None:
126                   pieces = entryl[a].chemical_name.split(' ')
127                   formula_pieces = entryl[a].formula.split(',')
128                   for c in range(len(pieces)):
129                       if hash(pieces[c]) == hash('xylene'):
130                           if entryl[a].identifier == 'QOWNEV':
131                               if hash('Cc1cccc1C') not in
132                                   component_SMILES:
133
134                                       component_SMILES.append(hash('Cc1cccc1C'))
```

```

126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154

        elif entry1[a].identifier == 'TEYCEF'
        or entry1[a].identifier ==
        'TEYCEF01':
            if hash('Cclcccc(C)c1') not in
            component_SMILES:
                component_SMILES.append(hash('Ccl
                cccc(C)c1'))
        elif entry1[a].identifier == 'MAMNAR':
            if hash('Cclccc(C)ccl') not in
            component_SMILES:
                component_SMILES.append(hash('Ccl
                ccc(C)ccl'))
        elif hash(pieces[c]) ==
        hash('hydrochloride') or hash(pieces[c]) ==
        hash('bis(hydrochloride)'):
            for z in range(len(formula_pieces)):
                if formula_pieces[z] == 'Cl1 1-' or
                '(Cl1 1-)' in formula_pieces[z]:
                    if hash('[Cl-]') not in
                    component_SMILES:
                        component_SMILES.append(hash(
                        '[Cl-]'))
                if formula_pieces[z] == 'H1 Cl1' or
                '(H1 Cl1)' in formula_pieces[z]:
                    if hash('Cl') not in
                    component_SMILES:
                        component_SMILES.append(hash(
                        'Cl'))
        elif hash(pieces[c]) ==
        hash('unidentified') or hash(pieces[c]) ==
        hash('unknown'):
            component_SMILES.append(tag)
            tag += 1
        elif hash(pieces[c]) == hash('glycol'):
            if hash(pieces[c-1]) in
            glycol_dictionary:
                if
                hash(glycol_dictionary.get(hash(pie
                ces[c-1]))) not in component_SMILES:
                    component_SMILES.append(hash(glyc
                    ol_dictionary.get(hash(pieces[c-1
                    ]))))
        elif hash(pieces[c]) == hash('oxalate'):
            if hash(pieces[c-1]) in
            oxalates_dictionary:
                if
                hash(oxalates_dictionary.get(hash(pie
                ces[c-1]))) not in component_SMILES:
                    component_SMILES.append(hash(oxal
                    ates_dictionary.get(hash(pieces[c
                    -1]))))
        elif hash(pieces[c]) == hash('acetate'):
            if hash(pieces[c-1]) in
            acetates_dictionary:
                if
                hash(acetates_dictionary.get(hash(pie
                ces[c-1]))) not in component_SMILES:
                    component_SMILES.append(hash(acet
                    ates_dictionary.get(hash(pieces[c
                    -1]))))

```

```

155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
    elif hash(pieces[c]) == hash('ether'):
        if hash(pieces[c-1]) ==
            hash('petroleum'):
                component_SMILES.append(11)
            elif pieces[c-2] + ' ' + pieces[c-1] +
                ' ' + pieces[c] in ethers_dictionary:
                if
                    hash(ethers_dictionary.get(pieces[c-2]
                        ] + ' ' + pieces[c-1] + ' ' +
                        pieces[c])) not in component_SMILES:
                        component_SMILES.append(hash(ethe
                            rs_dictionary.get(pieces[c-2] +
                                ' ' + pieces[c-1] + ' ' +
                                pieces[c])))
            elif hash(pieces[c-1]) in
                ethers_dictionary:
                if
                    hash(ethers_dictionary.get(hash(piece
                        s[c-1]))) not in component_SMILES:
                        component_SMILES.append(hash(ethe
                            rs_dictionary.get(hash(pieces[c-1]
                                ))))
            elif hash(pieces[c]) in partial_name:
                if pieces[c] == 'ketone':
                    if pieces[c-2] + ' ' + pieces[c-1]
                        + ' ' + pieces[c] in
                            full_solvate_name_dictionary:
                            if
                                hash(full_solvate_name_dictionary
                                    .get(pieces[c-2] + ' ' +
                                        pieces[c-1] + ' ' + pieces[c])) not in component_SMILES:
                                    component_SMILES.append(hash(
                                        full_solvate_name_dictionary.
                                            get(pieces[c-2] + ' ' +
                                                pieces[c-1] + ' ' +
                                                pieces[c])))
                else:
                    if pieces[c-1] + ' ' + pieces[c] in
                        full_solvate_name_dictionary:
                        if
                            hash(full_solvate_name_dictionary
                                .get(pieces[c-1] + ' ' +
                                    pieces[c])) not in
                                        component_SMILES:
                                            component_SMILES.append(hash(
                                                full_solvate_name_dictionary.
                                                    get(pieces[c-1] + ' ' +
                                                        pieces[c])))
            elif hash(pieces[c]) == hash('acid'):
                if hash(pieces[c-1]) in acids_dictionary:
                    if
                        hash(acids_dictionary.get(hash(pieces
                            [c-1]))) not in component_SMILES:
                            component_SMILES.append(hash(acid
                                s_dictionary.get(hash(pieces[c-1]
                                    ))))
            elif hash(pieces[c-1]) in
                acid_neutral_dictionary:
                if hash(pieces[c-1]) ==
                    hash('oxalic'):
                        for z in

```

```

180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
range(len(formula_pieces)):
    if formula_pieces[z] == 'C2
        H2 O4' or '(C2 H2 O4)' in
            formula_pieces[z]:
                if
                    hash(acid_neutral_diction
                        ary.get(hash(pieces[c-1]))
                    )) not in
                        component_SMILES:
                            component_SMILES.append(hash(acid_neutral_
                                dictionary.get(hash(p
                                    ieces[c-1]))))
    if formula_pieces[z] == 'C2
        H1 O4 1-' or '(C2 H1 O4
        1-)' in formula_pieces[z]:
        if
            hash(acid_charged_diction
                ary.get(hash(pieces[c-1]))
            )) not in
                component_SMILES:
                    component_SMILES.append(hash(acid_charged_
                        dictionary.get(hash(p
                            ieces[c-1]))))
    elif hash(pieces[c-1]) ==
        hash('nitric'):
            for z in
                range(len(formula_pieces)):
                    if formula_pieces[z] == 'H1
                        N1 O3' or '(H1 N1 O3)' in
                            formula_pieces[z]:
                                if
                                    hash(acid_neutral_diction
                                        ary.get(hash(pieces[c-1]))
                                    )) not in
                                        component_SMILES:
                                            component_SMILES.append(hash(acid_neutral_
                                                dictionary.get(hash(p
                                                    ieces[c-1]))))
    if formula_pieces[z] == 'N1
        O3 1-' or '(N1 O3 1-)' in
            formula_pieces[z]:
                if
                    hash(acid_charged_diction
                        ary.get(hash(pieces[c-1]))
                    )) not in
                        component_SMILES:
                            component_SMILES.append(hash(acid_charged_
                                dictionary.get(hash(p
                                    ieces[c-1]))))
    elif hash(pieces[c-1]) ==
        hash('trans-but-2-enedioic') or
        hash(pieces[c-1]) ==
        hash('but-2-enedioic') or
        hash(pieces[c-1]) == hash('fumaric'):
            for z in
                range(len(formula_pieces)):
                    if formula_pieces[z] == 'C4
                        H4 O4' or '(C4 H4 O4)' in
                            formula_pieces[z]:

```

197

```

    if
        hash(acid_neutral_dictionary.get(hash(pieces[c-1])))
    )) not in
        component_SMILES:
            component_SMILES.append(hash(acid_neutral_dictionary.get(hash(pieces[c-1])))))
    if formula_pieces[z] == 'C4
        H3 O4 1-' or '(C4 H3 O4
    1-) ' in formula_pieces[z]:
        if
            hash(acid_charged_dictionary.get(hash(pieces[c-1])))
        )) not in
            component_SMILES:
                component_SMILES.append(hash(acid_charged_dictionary.get(hash(pieces[c-1])))))
    elif hash(pieces[c-1]) ==
        hash('phosphonic'):
            for z in
                range(len(formula_pieces)):
                    if formula_pieces[z] == 'H3
                        O3 P1' or '(H3 O3 P1) ' in
                            formula_pieces[z]:
                                if
                                    hash(acid_neutral_dictionary.get(hash(pieces[c-1])))
                                )) not in
                                    component_SMILES:
                                        component_SMILES.append(hash(acid_neutral_dictionary.get(hash(pieces[c-1])))))
    if formula_pieces[z] == 'H2
        O3 P1 1-' or '(H2 O3 P1
    1-) ' in formula_pieces[z]:
        if
            hash(acid_charged_dictionary.get(hash(pieces[c-1])))
        )) not in
            component_SMILES:
                component_SMILES.append(hash(acid_charged_dictionary.get(hash(pieces[c-1])))))
    elif hash(pieces[c-1]) ==
        hash('perchloric'):
            for z in
                range(len(formula_pieces)):
                    if formula_pieces[z] == 'H1
                        Cl1 O4' or '(H1 Cl1 O4) ' in
                            formula_pieces[z]:
                                if
                                    hash(acid_neutral_dictionary.get(hash(pieces[c-1])))
                                )) not in
                                    component_SMILES:
                                        component_SMILES.append(hash(acid_neutral_dictionary.get(hash(pieces[c-1])))))

```

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

```

215                                         nd(hash(acid_neutral_
216                                         dictionary.get(hash(p
217                                         ieces[c-1]))))
218                                         if formula_pieces[z] ==
219                                         'C11 O4 1-' or '(C11 O4
220                                         1-) ' in formula_pieces[z]:
221                                         if
222                                         hash(acid_charged_diction
223                                         ary.get(hash(pieces[c-1]))
224                                         )) not in
225                                         component_SMILES:
226                                         component_SMILES.append(
227                                         hash(acid_charged_
228                                         dictionary.get(hash(p
229                                         ieces[c-1]))))
230                                         elif hash(pieces[c-1]) ==
231                                         hash('pyridine-2,6-dicarboxylic'):
232                                         for z in

```



```

250                                         nd(hash(acid_charged_
251                                         dictionary.get(hash(p
252                                         ieces[c-1]))))
253
254     elif hash(pieces[c]) in
255         deuterated_dictionary:
256             if
257                 hash(deuterated_dictionary.get(hash(piece
258                                         s[c]))) not in component_SMILES:
259
260                 component_SMILES.append(hash(deuterat
261                                         ed_dictionary.get(hash(pieces[c]))))
262
263             if hash(pieces[c]) ==
264                 hash('dideutero-dichloromethane') or
265                 hash(pieces[c]) ==
266                 hash('dideuterodichloromethane'):
267                     if hash('Cl[C]Cl') in
268                         component_SMILES:
269
270                         component_SMILES.remove(hash('Cl[
271                                         C]Cl'))
272
273             elif hash(pieces[c]) ==
274                 hash('perdeutero-toluene'):
275                 if hash('[C]c1ccccc1') in
276                     component_SMILES:
277
278                         component_SMILES.remove(hash(' [C]
279                                         c1ccccc1'))
280
281             elif hash(pieces[c]) ==
282                 hash('deutero-ethanol'):
283                 if hash('[C][C][O]') in
284                     component_SMILES:
285
286                         component_SMILES.remove(hash(' [C]
287                                         [C][O]'))
288
289             elif hash(pieces[c]) ==
290                 hash('deuterochloroform') or
291                 hash(pieces[c]) ==
292                 hash('deutero-chloroform'):
293                     if hash('Cl[C](Cl)Cl') in
294                         component_SMILES:
295
296                         component_SMILES.remove(hash('Cl[
297                                         C](Cl)Cl'))
298
299             elif hash(pieces[c]) ==
300                 hash('hexadeutero-benzene') or
301                 hash(pieces[c]) ==
302                 hash('deuterobenzene') or
303                 hash(pieces[c]) ==
304                 hash('hexadeuterobenzene'):
305                     if hash('[c]1ccccc1') in
306                         component_SMILES:
307
308                         component_SMILES.remove(hash(' [c]
309                                         1ccccc1'))
310
311             elif hash(pieces[c]) ==
312                 hash('deuteromethanol') or
313                 hash(pieces[c]) ==
314                 hash('perdeuteromethanol'):
315                     if hash('[C][O]') in
316                         component_SMILES:
317
318                         component_SMILES.remove(hash(' [C]
319                                         [O]'))
320
321             elif hash(pieces[c]) in solvents_dictionary:
322                 if
323                     hash(solvents_dictionary.get(hash(pie
324                                         ces[c]))) not in component_SMILES:

```

```

component_SMILES.append(hash(solvents
                             _dictionary.get(hash(pieces[c]))))

#The formula of each structure is also recorded
formulas = entry1[a].formula.split(',')
letter_formulas = []
for b in range(len(formulas)):
    if 'D' in formulas[b] and 'H' in formulas[b]:
        pieces = formulas[b].split(' ')
        for c in range(len(pieces)):
            if 'H' in pieces[c]:
                pin1 = pieces[c].index('H')
                piece1 = pieces[c]
                stoich1 = int(pieces[c][pin1+1:])
            elif 'D' in pieces[c]:
                pin2 = pieces[c].index('D')
                piece2 = pieces[c]
                stoich2 = int(pieces[c][pin2+1:])
        stoich = stoich1 + stoich2
        pieces.remove(piece1)
        pieces.remove(piece2)
        pieces.append('H' + str(stoich))
        pieces = sorted(pieces)
        formulas[b] = ' '.join(pieces)
    elif 'D' in formulas[b]:
        formulas[b] = str(formulas[b]).replace('D', 'H')
if loop == 0 or loop == 2:
    if '(' in formulas[b] and ')' in formulas[b]:
        if '(H2 O1)' not in formulas[b]:
            pin = formulas[b].index('(')
            letter_formulas.append(formulas[b][pin+1:-1])
    else:
        if formulas[b] != 'H2 O1':
            letter_formulas.append(formulas[b])
else:
    if '(' in formulas[b] and ')' in formulas[b]:
        pin = formulas[b].index('(')
        letter_formulas.append(formulas[b][pin+1:-1])
    else:
        letter_formulas.append(formulas[b])

if loop == 0 or loop == 2:
    hydrates_formulas.append(letter_formulas)
else:
    waterless_forms_formulas.append(letter_formulas)

if loop == 0 or loop == 2:
    if hash('O') in component_SMILES:
        component_SMILES.remove(hash('O'))
    if hash('[O]') in component_SMILES:
        component_SMILES.remove(hash('[O]'))
    component_SMILES = sorted(component_SMILES)
    hydrates_formulas.append(component_SMILES)
if loop == 1 or loop == 3:
    component_SMILES = sorted(component_SMILES)
    waterless_forms_formulas.append(component_SMILES)
loop += 1
if loop == 1:
    entry1 = entry_reader2
if loop == 2:
    entry1 = entry_reader3
if loop == 3:
    entry1 = entry_reader4

```

```

336     waterless_forms_recorded = []
337
338     print 'made it to comparison'
339
340     for d in range(len(hydrates_SMILES)):
341         count = 0
342         for e in range(len(waterless_forms_SMILES)):
343             #If the SMILES string and formula of the hydrate
344             #structure match those of the waterless form, the two
345             #structures are written to output files as
346             #hydrate-anhydride pairs
347             if hydrates_SMILES[d] == waterless_forms_SMILES[e] and
348             hydrates_formulas[d] == waterless_forms_formulas[e]:
349                 count += 1
350                 #There are 28,388 structures in the solvated
351                 #waterless forms input file
352                 #These numbers serve as place holders so the
353                 #correct structures are compared and written to
354                 #output files
355                 if e < 28388:
356                     if entry_reader2[e].identifier not in
357                     duplicate_identifier_dictionary:
358                         #Each waterless form that matches a hydrate
359                         #is recorded in a list
360                         waterless_forms_recorded.append(hash(entry_re
361                         ader2[e].identifier))
362
363             else:
364                 waterless_forms_recorded.append(duplicate_ide
365                 ntifier_dictionary.get(entry_reader2[e].ident
366                 ifier))
367
368             else:
369                 if entry_reader4[e-28388].identifier not in
370                 duplicate_identifier_dictionary:
371
372                 waterless_forms_recorded.append(hash(entry_re
373                 ader4[e-28388].identifier))
374
375             else:
376                 waterless_forms_recorded.append(duplicate_ide
377                 ntifier_dictionary.get(entry_reader4[e-28388]
378                 .identifier))
379
380             #There are 3,499 structures in the solvated
381             #hydrates input file
382             if d < 3439 and e < 28388:
383                 pairing_writer.write(entry_reader1[d])
384                 writer1.write(entry_reader2[e])
385             elif d < 3439 and e >= 28388:
386                 pairing_writer.write(entry_reader1[d])
387                 writer1.write(entry_reader4[e-28388])
388             elif d >= 3439 and e < 28388:
389                 pairing_writer.write(entry_reader3[d-3439])
390                 writer1.write(entry_reader2[e])
391             elif d >= 3439 and e >= 28388:
392                 pairing_writer.write(entry_reader3[d-3439])
393                 writer1.write(entry_reader4[e-28388])
394
395             #Hydrates that do not match any waterless forms are written
396             #as hydrates without known anhydride forms
397             if count == 0:
398                 if d < 3439:
399                     writer2.write(entry_reader1[d])
400                 else:
401                     writer2.write(entry_reader3[d-3439])
402
403             print d

```

```
381 print 'made it to generating files with unpaired structures'
382
383 #Any waterless forms that did not match a hydrate structure are
384 #written as anhydrous forms without a known hydrate form
385 for f in range(len(entry_reader2)):
386     if entry_reader2[f].identifier not in
387         duplicate_identifier_dictionary:
388             if hash(entry_reader2[f].identifier) not in
389                 waterless_forms_recorded:
390                     writer3.write(entry_reader2[f])
391             else:
392                 if
393                     duplicate_identifier_dictionary.get(entry_reader2[f].iden
394                         tifier) not in waterless_forms_recorded:
395                             writer3.write(entry_reader2[f])
396             for g in range(len(entry_reader4)):
397                 if entry_reader4[g].identifier not in
398                     duplicate_identifier_dictionary:
399                         if hash(entry_reader4[g].identifier) not in
400                             waterless_forms_recorded:
401                                 writer3.write(entry_reader4[g])
402             else:
403                 if
404                     duplicate_identifier_dictionary.get(entry_reader4[g].iden
405                         tifier) not in waterless_forms_recorded:
406                         writer3.write(entry_reader4[g])
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999
999
999
```

```

1 #22 Find hydrates with waterless forms using None method
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds hydrates-anhydrate pairs where either one or both structure is
10 without an entry SMILES string
11 #The only automated piece that is done by the script is finding all the hydrate
12 structures without an entry SMILES string that match the formula of a waterless format
13 #It also finds all the waterless forms without an entry SMILES string that match the
14 formula of a hydrate
15 #The output files from this script were manually screened to find true
16 hydrate-anhydrate pairs
17
18 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
19 Forms\None_solvated_hydrates.txt"
20 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
21 Forms\None_solvated_waterless_forms.txt"
22 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
23 Forms\None_hydrates.txt"
24 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
25 Forms\None_waterless_forms.txt"
26 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
27 Forms\smiles_solvated_hydrates.txt"
28 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
29 Forms\smiles_solvated_waterless_forms.txt"
30 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
31 Forms\smiles_hydrates.txt"
32 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
33 Forms\smiles_waterless_forms.txt"
34
35 class Runner(argparse.ArgumentParser):
36
37     def __init__(self):
38         super(self.__class__, self).__init__(description=__doc__)
39         self.add_argument(
40             '-i', '--input', default=filepath1,
41             help='input database filepath1'
42         )
43         self.add_argument(
44             '-o', '--output', default='None_hydrates_with_formula_match.gcd',
45             help='output file [None_hydrates_with_formula_match.gcd]'
46         )
47         self.add_argument(
48             '-m', '--maximum', default=0, type=int,
49             help='Maximum number of structures to find [all]'
50         )
51
52         args = self.parse_args()
53
54         self.args = args
55         self.settings = search.Search.Settings()
56         self.settings.max_hit_structures = self.args.maximum
57
58     def run(self):
59
60         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
61         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
62         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
63         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
64         entry_reader5 = io.EntryReader(filepath5, format='identifiers')
65         entry_reader6 = io.EntryReader(filepath6, format='identifiers')
66         entry_reader7 = io.EntryReader(filepath7, format='identifiers')
67         entry_reader8 = io.EntryReader(filepath8, format='identifiers')

```

```

56
57     total = 0
58     count = 0
59
60     with io.EntryWriter(self.args.output) as pairing_writer:
61         with io.EntryWriter("waterless_forms_that_match_None_hydrate_formula.gcd") as writer1:
62             with io.EntryWriter("None_waterless_forms_with_formula_match.gcd") as writer2:
63                 with io.EntryWriter("hydrates_that_match_None_waterless_forms.gcd") as writer3:
64
65                     None_hydrates_formulas = []
66                     None_waterless_forms_formulas = []
67                     hydrates_formulas = []
68                     waterless_forms_formulas = []
69
70                     loop = 0
71                     entry1 = entry_reader1
72                     #Generate a list of formulas for hydrates and waterless forms
73                     #with and without an entry SMILES string
74                     while loop < 8:
75                         for a in range(len(entry1)):
76                             component_formulas = entry1[a].formula.split(',')
77                             for b in range(len(component_formulas)):
78                                 if ' ( and )' in component_formulas[b]:
79                                     start = component_formulas[b].index('(')
80                                     end = component_formulas[b].index(')')
81                                     component_formulas[b] =
82                                         component_formulas[b][start+1:end]
83                                     component_formulas[b] =
84                                     hash(component_formulas[b])
85
86                                     if loop == 0 or loop == 2:
87                                         if hash('H2 O1') in component_formulas:
88                                             component_formulas.remove(hash('H2 O1'))
89                                         if hash('D2 O1') in component_formulas:
90                                             component_formulas.remove(hash('D2 O1'))
91                                         component_formulas = sorted(component_formulas)
92                                         None_hydrates_formulas.append(component_formulas)
93
94                                     if loop == 1 or loop == 3:
95                                         component_formulas = sorted(component_formulas)
96
97                                         None_waterless_forms_formulas.append(component_formulas)
98
99                                     if loop == 4 or loop == 6:
100                                         if hash('H2 O1') in component_formulas:
101                                             component_formulas.remove(hash('H2 O1'))
102                                         if hash('D2 O1') in component_formulas:
103                                             component_formulas.remove(hash('D2 O1'))
104                                         component_formulas = sorted(component_formulas)
105                                         hydrates_formulas.append(component_formulas)
106
107                                         if loop == 5 or loop == 7:
108                                             component_formulas = sorted(component_formulas)
109
110                                         waterless_forms_formulas.append(component_formulas)
111
112                                         loop += 1
113                                         if loop == 1:
114                                             entry1 = entry_reader2
115                                         if loop == 2:
116                                             entry1 = entry_reader3
117                                         if loop == 3:
118                                             entry1 = entry_reader4
119                                         if loop == 4:
120                                             entry1 = entry_reader5
121                                         if loop == 5:
122                                             entry1 = entry_reader6

```

```

113
114
115
116
117
118         if loop == 6:
119             entry1 = entry_reader7
120         if loop == 7:
121             entry1 = entry_reader8
122
123             #If a hydrate formula matches a waterless form, the two
124             structures are written to output files as potential
125             hydrate-anhydrite pairs
126             #Hydrates without an entry SMILES string are first compared to
127             waterless forms without an entry SMILES string
128             for d in range(len(None_hydrates_formulas)):
129                 for e in range(len(None_waterless_forms_formulas)):
130                     if None_hydrates_formulas[d] ==
131                         None_waterless_forms_formulas[e]:
132                         #There are 192 solvated hydrates without an entry
133                         SMILES string
134                         #There are 1,555 solvated waterless forms without
135                         an entry SMILES string
136                         if d < 192 and e < 1555:
137                             pairing_writer.write(entry_reader1[d])
138                             writer1.write(entry_reader2[e])
139                         elif d < 192 and e >= 1555:
140                             pairing_writer.write(entry_reader1[d])
141                             writer1.write(entry_reader4[e-1555])
142                         elif d >= 192 and e < 1555:
143                             pairing_writer.write(entry_reader3[d-192])
144                             writer1.write(entry_reader2[e])
145                         elif d >= 192 and e >= 1555:
146                             pairing_writer.write(entry_reader3[d-192])
147                             writer1.write(entry_reader4[e-1555])
148
149             #Hydrates without an entry SMILES string are then compared
150             to waterless forms with an entry SMILES string
151             for f in range(len(waterless_forms_formulas)):
152                 if None_hydrates_formulas[d] ==
153                     waterless_forms_formulas[f]:
154                         if d < 192 and f < 28388:
155                             pairing_writer.write(entry_reader1[d])
156                             writer1.write(entry_reader6[f])
157                         elif d < 192 and f >= 28388:
158                             pairing_writer.write(entry_reader1[d])
159                             writer1.write(entry_reader8[f-28388])
160                         elif d >= 192 and f < 28388:
161                             pairing_writer.write(entry_reader3[d-192])
162                             writer1.write(entry_reader6[f])
163                         elif d >= 192 and f >= 28388:
164                             pairing_writer.write(entry_reader3[d-192])
165                             writer1.write(entry_reader8[f-28388])
166
167             #Waterless forms without an entry SMILES string are compared to
168             hydrates with an entry SMILES string
169             for g in range(len(None_waterless_forms_formulas)):
170                 for h in range(len(hydrates_formulas)):
171                     if None_waterless_forms_formulas[g] ==
172                         hydrates_formulas[h]:
173                         if g < 1555 and h < 3439:
174                             writer2.write(entry_reader2[g])
175                             writer3.write(entry_reader5[h])
176                         elif g < 1555 and h >= 3439:
177                             writer2.write(entry_reader2[g])
178                             writer3.write(entry_reader7[h-3439])
179                         elif g >= 1555 and h < 3439:
180                             writer2.write(entry_reader4[g-1555])
181                             writer3.write(entry_reader5[h])
182                         elif g >= 1555 and h >= 3439:
183                             writer2.write(entry_reader4[g-1555])
184                             writer3.write(entry_reader7[h-3439])

```

```
170  
171  
172 if __name__ == '__main__':  
173     # This runs the script  
174     r = Runner()  
175     r.run()  
176
```

```

1 #23 Find hydrate-anhydrate pairs with distinct chirality
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds hydrate-anhydrate pairs with distinct chirality.
10 #This script is the same as the one used to find duplicate structures with distinct
11 #chirality
12 #The only difference is that water is not included in the molecule formula check
13 #against the entry formula for hydrate structures
14
15 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
16 Waterless Forms\set1_potential_hydrates_with_waterless_form.txt"
17 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
18 Waterless Forms\set1_potential_waterless_forms_with_hydrate.txt"
19 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
20 Waterless Forms\set2_potential_hydrates_with_waterless_form.txt"
21 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
22 Waterless Forms\set2_potential_waterless_forms_with_hydrate.txt"
23 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
24 Waterless Forms\set3_potential_hydrates_with_waterless_form.txt"
25 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
26 Waterless Forms\set3_potential_waterless_forms_with_hydrate.txt"
27 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
28 Waterless Forms\potential_None_hydrates_with_waterless_form.txt"
29 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
30 Waterless Forms\potential_None_waterless_forms_with_hydrate.txt"
31
32 class Runner(argparse.ArgumentParser):
33
34     def __init__(self):
35         super(self.__class__, self).__init__(description=__doc__)
36         self.add_argument(
37             '-i', '--input', default=filepath1,
38             help='input database filepath')
39
40         self.add_argument(
41             '-o', '--output', default='same_chirality_hydrates_with_waterless_form.gcd',
42             help='output file [same_chirality_hydrates_with_waterless_form.gcd]')
43
44         self.add_argument(
45             '-m', '--maximum', default=0, type=int,
46             help='Maximum number of structures to find [all]')
47
48     args = self.parse_args()
49
50     self.args = args
51     self.settings = search.Search.Settings()
52     self.settings.max_hit_structures = self.args.maximum
53
54     def run(self):
55
56         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
57         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
58         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
59         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
60         entry_reader5 = io.EntryReader(filepath5, format='identifiers')
61         entry_reader6 = io.EntryReader(filepath6, format='identifiers')
62         entry_reader7 = io.EntryReader(filepath7, format='identifiers')
63         entry_reader8 = io.EntryReader(filepath8, format='identifiers')
64
65         with io.EntryWriter(self.args.output) as chiral_writer:
66             with io.EntryWriter("same_chirality_waterless_forms_with_hydrate.gcd") as

```

```

writer1:
    with
        io.EntryWriter("distinct_chirality_hydrates_with_enantiomer_waterless_for
m.gcd") as writer2:
            with
                io.EntryWriter("distinct_chirality_waterless_forms_with_enantiomer_hy
drate.gcd") as writer3:
                    with
                        io.EntryWriter("achiral_names_one_missing_stereo_hydrates_with_wa
terless_form.gcd") as writer4:
                            with
                                io.EntryWriter("achiral_names_one_missing_stereo_waterless_fo
rms_with_hydrate.gcd") as writer5:
                                    with
                                        io.EntryWriter("one_chiral_name_one_missing_stereo_hydrat
es_with_waterless_form.gcd") as writer6:
                                            with
                                                io.EntryWriter("one_chiral_name_one_missing_stereo_wa
terless_forms_with_hydrate.gcd") as writer7:
                                                    with
                                                        io.EntryWriter("chiral_names_check_hydrates.gcd")
as writer8:
                                                        with
                                                            io.EntryWriter("chiral_names_check_waterless_
forms.gcd") as writer9:
                                                                duplicate_identifier_dictionary =
                                                                {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
                                                                'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6,
                                                                'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
                                                                'VUFJAJ':10, 'WEPOQA':11, 'XADFUD':12,
                                                                'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
                                                                hydrate_tags_empty_one_stereo_hydrates =
                                                                [hash('BEXPOK'), hash('QQQCYA01'),
                                                                hash('VOLZON'), hash('VOMPEU'),
                                                                hash('WEGPED'), hash('AHLPRO'),
                                                                hash('YUPBER')]
                                                                hydrate_tags_empty_one_stereo_waterless_f
orms = [hash('BEXPIE'), hash('NAMNSB'),
hash('VOMNUI'), hash('VOMNUI'),
hash('WEGPIH'), hash('HOPROL'),
hash('VALINO01')]
                                                                hydrate_tags_empty_no_stereo_hydrates =
                                                                [hash('ADPRTY'), hash('BEWWIK'),
hash('BEWWIK'), hash('BUTAHC11'),
hash('DUGTOQ'), hash('ETHYLO'),
hash('ETHYLO'), hash('ETXDHY01'),
hash('ETXDHY01'), hash('FEFNOT01'),
hash('FEFNOT01'), hash('FEFNOT01'),
hash('FEFNOT05'), hash('FEFNOT05'),
hash('FEFNOT05'), hash('FEFNOT05')]

```

```
hash('FEFNOT05'), hash('FEFNOT10'),  
hash('FEFNOT10'), hash('FEFNOT10'),  
hash('FEFNOT11'), hash('FEFNOT11'),  
hash('FEFNOT11'), hash('FUSWOH04'),  
hash('FUSWOH04'), hash('FUSWOH04'),  
hash('FUSWOH04'), hash('FUSWOH04'),  
hash('FUSWOH04'), hash('FUSWOH04'),  
hash('FUSWOH04'), hash('FUSWOH04'),  
hash('FUSWOH04'), hash('IKALIR'),  
hash('IKALIR'), hash('IKALIRO1'),  
hash('IKALIRO1'), hash('IKALIRO2'),  
hash('IKALIRO2'), hash('IKALIRO3'),  
hash('IKALIRO3'), hash('IKAMEO'),  
hash('IKAMEO'), hash('IKAMEO01'),  
hash('IKAMEO01'), hash('IKAMOY'),  
hash('IKAMOY'), hash('IKAMOY01'),  
hash('IKAMOY01'), hash('IKAMOY02'),  
hash('IKAMOY02'), hash('IKAMOY03'),  
hash('IKAMOY03'), hash('IKAMOY04'),  
hash('IKAMOY04'), hash('IKAMOY05'),  
hash('IKAMOY05'), hash('IKAVOF'),  
hash('IKAVOF'), hash('IKAVOF'),  
hash('IKAVOF'), hash('IKAVOF'),  
hash('IKAVOF'), hash('IKAVOF'),  
hash('JARBUA'), hash('KUYGOA'),  
hash('KUYGUG'), hash('KUYGUG'),  
hash('KUYGUG'), hash('KUYGUG'),  
hash('KUYHAN'), hash('KUYHAN'),  
hash('KUYHAN'), hash('KUYHAN'),  
hash('KUYHAN'), hash('KUYHAN'),  
hash('KUYHAN'), hash('KUYHAN'),  
hash('KUYHER'), hash('KUYHER'),  
hash('MUSDEL'), hash('MUSDOV'),  
hash('NAHCIJ'), hash('NAHCIJ'),  
hash('NAHCOP'), hash('NAHCOP'),  
hash('NAHCOP'), hash('NAHCOP'),  
hash('NAHCOP'), hash('NAHCOP'),  
hash('QQQCIY10'), hash('QQQCIY10'),  
hash('QQQCIY10'), hash('QQQCIY10'),  
hash('QQQCIY10'), hash('QQQCIY10'),  
hash('QQQCIY10'), hash('RITMEN'),  
hash('RITMEN'), hash('SOGUAN06'),  
hash('XEBQAZ'), hash('ZZZPPI01'),  
hash('ZZZWQA01'), hash('ZZZWQA01'),  
hash('FEFNOT'), hash('FEFNOT'),  
hash('FEFNOT'), hash('FEFNOT'),  
hash('FEFNOT'), hash('FEFNOT'),  
hash('FEFNOT'), hash('FEFNOT')
```

hash('CBMZPN16'), hash('CBMZPN17'),  
hash('CBMZPN18'), hash('CBMZPN20'),  
hash('CBMZPN21'), hash('CBMZPN22'),  
hash('CBMZPN23'), hash('CBMZPN27'),  
hash('CBMZPN28'), hash('CBMZPN29'),  
hash('CBMZPN30'), hash('CBMZPN01'),  
hash('CBMZPN02'), hash('CBMZPN03'),  
hash('CBMZPN10'), hash('CBMZPN11'),  
hash('CBMZPN12'), hash('CBMZPN13'),  
hash('CBMZPN14'), hash('CBMZPN16'),  
hash('CBMZPN17'), hash('CBMZPN18'),  
hash('CBMZPN20'), hash('CBMZPN21'),  
hash('CBMZPN22'), hash('CBMZPN23'),  
hash('CBMZPN27'), hash('CBMZPN28'),  
hash('CBMZPN29'), hash('CBMZPN30'),  
hash('CBMZPN01'), hash('CBMZPN02'),  
hash('CBMZPN03'), hash('CBMZPN10'),  
hash('CBMZPN11'), hash('CBMZPN12'),  
hash('CBMZPN13'), hash('CBMZPN14'),  
hash('CBMZPN16'), hash('CBMZPN17'),  
hash('CBMZPN18'), hash('CBMZPN20'),  
hash('CBMZPN21'), hash('CBMZPN22'),  
hash('CBMZPN23'), hash('CBMZPN27'),  
hash('CBMZPN28'), hash('CBMZPN29'),  
hash('CBMZPN30'), hash('CBMZPN01'),  
hash('CBMZPN02'), hash('CBMZPN03'),  
hash('CBMZPN10'), hash('CBMZPN11'),  
hash('CBMZPN12'), hash('CBMZPN13'),  
hash('CBMZPN14'), hash('CBMZPN16'),  
hash('CBMZPN17'), hash('CBMZPN18'),  
hash('CBMZPN20'), hash('CBMZPN21'),  
hash('CBMZPN22'), hash('CBMZPN23'),  
hash('CBMZPN27'), hash('CBMZPN28'),  
hash('CBMZPN29'), hash('CBMZPN30'),  
hash('EFUMAU'), hash('EFUMAU01'),  
hash('EFUMAU02'), hash('EFUMAU03'),  
hash('EFUMAU04'), hash('EFUMAU05'),  
hash('EFUMAU06'), hash('EFUMAU07'),  
hash('IKAKOW'), hash('IKAKOW01'),  
hash('ETHANE01'), hash('ETHANE04'),  
hash('ETHANE05'), hash('ETHANE06'),  
hash('ETHANE07'), hash('ETHANE08'),  
hash('ETHANE09'), hash('ETHANE10'),  
hash('ETHANE11'), hash('JARBOU'),  
hash('ZZZITY01'), hash('KAMTAW'),  
hash('KAMTAW01'), hash('KAMTAW02'),  
hash('KAMTAW03'), hash('KAMTAW04'),  
hash('KAMTAW05'), hash('KAMTAW06'),  
hash('DCLBEN'), hash('DCLBEN01'),  
hash('DCLBEN02'), hash('DCLBEN03'),  
hash('DCLBEN04'), hash('DCLBEN05'),  
hash('DCLBEN06'), hash('DCLBEN07'),  
hash('DCLBEN11'), hash('CISTON01'),  
hash('TAQYUG'), hash('DUHJIB'),  
hash('DUHJIB'), hash('ACETYL02'),  
hash('ACETYL03'), hash('JAYDUI')

```

hash('BUNJAV01'), hash('CBMZPN01'),
hash('CBMZPN02'), hash('CBMZPN03'),
hash('CBMZPN10'), hash('CBMZPN11'),
hash('CBMZPN12'), hash('CBMZPN13'),
hash('CBMZPN14'), hash('CBMZPN16'),
hash('CBMZPN17'), hash('CBMZPN18'),
hash('CBMZPN20'), hash('CBMZPN21'),
hash('CBMZPN22'), hash('CBMZPN23'),
hash('CBMZPN27'), hash('CBMZPN28'),
hash('CBMZPN29'), hash('CBMZPN30'),
hash('PINCOL'), hash('PINCOL01'),
hash('PINCOL02'), hash('PINCOL03'),
hash('DAYFUH'), hash('LUXYIM'),
hash('LUYMAT'), hash('LUYMAT01'),
hash('NIWFEE02'), hash('NIWFEE02'),
hash('NIWFEE04'), hash('NIWFEE04'),
hash('SABZOK'), hash('SOCTOT'),
hash('THIOUR06'), hash('THIOUR06'),
hash('THIOUR06'), hash('THIOUR06'),
hash('PUSBOV')]

72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109

loop = 2
while loop < 4:
    if loop == 0:
        entry1 = entry_reader1
    if loop == 1:
        entry1 = entry_reader3
    if loop == 2:
        entry1 = entry_reader5
    if loop == 3:
        entry1 = entry_reader7
hydrates = []
waterless = []
count1 = 0
count2 = 0
count3 = 0
for a in range(len(entry1)):
    print a
    loop2 = 0
    if loop2 == 0:
        if loop == 0:
            entry1 = entry_reader1
        if loop == 1:
            entry1 = entry_reader3
        if loop == 2:
            entry1 = entry_reader5
        if loop == 3:
            entry1 = entry_reader7
formula_error = 0
missing_components = 0
while loop2 < 2:
    letter_formulas = []
    formulas =
    entry1[a].formula.split(',')
    for b in
        range(len(formulas)):
            if '(' in formulas[b]
            and ')' in formulas[b]:
                if '(H2 O1)' not in
                    formulas[b] and
                    '(D2 O1)' not in
                    formulas[b]:
                        pin =
                            formulas[b].index
                            ('(')
                        letter_formulas.append(pin)
    count1 += len(letter_formulas)
    count2 += len(letter_formulas)
    count3 += len(letter_formulas)
    loop2 += 1
    if count1 == 0 and count2 == 0 and count3 == 0:
        break

```

```

110                                         ppPEND(formulas[b]
111                                         [pin+1:-1])
112                                         else:
113                                         if formulas[b] != 'H2 O1' and
114                                         formulas[b] != 'D2
115                                         O1':
116                                         letter_formulas.a
117                                         ppPEND(formulas[b]
118                                         )
119                                         molecule1 =
120                                         entryl[a].molecule
121                                         checked = []
122                                         tags = []
123                                         error = 0
124                                         count = 0
125                                         for g in
126                                         range(len(molecule1.component
127                                         s)):
128                                         if
129                                         molecule1.components[g].f
130                                         ormula not in checked
131                                         and
132                                         molecule1.components[g].f
133                                         ormula != 'H2 O1':
134                                         if
135                                         molecule1.components[g].formula not in
136                                         letter_formulas:
137                                         if 'H' in
138                                         molecule1.compone
139                                         nts[g].formula:
140                                         if
141                                         molecule1.components[g].formula.replace('H', 'D') not in letter_formulas:
142                                         partial_deuterated = 0
143                                         for t
144                                         in range(len(formulas)):
145                                         if 'H' in formulas[t] and 'D' in formulas[t]:
146                                         pieces = formulas[t].split(' ')
147                                         for s in range(len(pieces)):
148                                         if 'H' in pieces[s]:
149                                         tag1 = pieces[s].index('H')
150                                         number1 = int(pieces[s][tag1+1:])
151                                         if 'D' in pieces[s]:
152                                         tag2 = pieces[s].index('D')
153                                         number2 = int(pieces[s][tag2+1:])
154                                         partial_deuterated += 1
155                                         if
156                                         partial_deuterated == 1:
157                                         pieces2 = molecule1.components[g].formula.split(' ')
158                                         for u
159                                         in range(len(pieces2)):
160                                         if 'H' in pieces2[u]:

```

```

139     tag3 = pieces2[u].index('H')
140
141     number3 = int(pieces2[u][tag3+1:])
142
143     if number1+number2 == number3:
144
145         checked.append(molecule1.components[g].formula)
146
147         tags.append(g)
148
149         count += 1
150
151     else:
152
153         checked.append(molecule1.components[g].formula)
154
155         tags.append(g)
156
157         error += 1
158
159     else:
160
161         checked.append(molecule1.components[g].formula)
162
163         tags.append(g)
164
165         count += 1
166
167         error += 1
168
169     else:
170
171         checked.append(molecule1.components[g].formula)
172
173         tags.append(g)
174
175         count += 1
176
177         error += 1
178
179     if loop2 == 0:
180
181         hydrate_tags = tags
182         hydrates_checked =
183         checked
184
185     else:
186
187         waterless_tags = tags
188         waterless_checked =
189         checked
190
191     if error != 0:
192
193         formula_error += 1
194
195     if count != len(letter_formulas):
196
197         missing_components += 1
198
199     loop2 += 1
200
201     if loop == 0:
202
203         entry1 = entry_reader2
204
205     if loop == 1:

```

```

179                                         entry1 = entry_reader4
180                                         if loop == 2:
181                                             entry1 = entry_reader6
182                                         if loop == 3:
183                                             entry1 = entry_reader8
184
185                                         if loop == 0:
186                                             entry2 = entry_reader1
187                                             entry3 = entry_reader2
188                                         if loop == 1:
189                                             entry2 = entry_reader3
190                                             entry3 = entry_reader4
191                                         if loop == 2:
192                                             entry2 = entry_reader5
193                                             entry3 = entry_reader6
194                                         if loop == 3:
195                                             entry2 = entry_reader7
196                                             entry3 = entry_reader8
197
198                                         if formula_error != 0 or
199                                         missing_components != 0:
200                                             remove_hydrate_tags = []
201                                             for c in
202                                                 range(len(hydrate_tags)):
203                                                 if
204                                                     entry2[a].molecule.compon
205                                                     ents[hydrate_tags[c]].for
206                                                     mula not in
207                                                     waterless_checked:
208
209                                         remove_hydrate_tags.append(hydrate_tags[c])
210
211                                         for e in
212                                             range(len(remove_hydrate_tags)):
213
214                                         hydrate_tags.remove(remove
215                                         _hydrate_tags[e])
216                                         remove_waterless_tags = []
217                                         for d in
218                                             range(len(waterless_tags)):
219                                                 if
220                                                     entry3[a].molecule.compon
221                                                     ents[waterless_tags[d]].for
222                                                     mula not in
223                                                     hydrates_checked:
224
225                                         remove_waterless_tags.append(waterless_tags[d])
226                                         for f in
227                                             range(len(remove_waterless_ta
228                                             gs)):
229
230                                         waterless_tags.remove(remove
231                                         _waterless_tags[f])
232
233                                         if len(hydrate_tags) !=
234                                         len(waterless_tags):
235                                             print(entry2[a].identifier)
236                                             print(entry3[a].identifier)
237                                             print(waterless_tags[90])
238
239                                         same_names = 0
240                                         possibly_the_same_names = 0
241                                         distinct_names = 0

```

```

220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
709
710
711
712
713
714
715
716
717
718
719
719
720
721
722
723
724
725
726
727
728
729
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
749
750
751
752
753
754
755
756
757
758
759
759
760
761
762
763
764
765
766
767
768
769
769
770
771
772
773
774
775
776
777
778
779
779
780
781
782
783
784
785
786
787
788
789
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
809
810
811
812
813
814
815
816
817
818
819
819
820
821
822
823
824
825
826
827
828
829
829
830
831
832
833
834
835
836
837
838
839
839
840
841
842
843
844
845
846
847
848
849
849
850
851
852
853
854
855
856
857
858
859
859
860
861
862
863
864
865
866
867
868
869
869
870
871
872
873
874
875
876
877
878
879
879
880
881
882
883
884
885
886
887
888
889
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
909
910
911
912
913
914
915
916
917
918
919
919
920
921
922
923
924
925
926
927
928
929
929
930
931
932
933
934
935
936
937
938
939
939
940
941
942
943
944
945
946
947
948
949
949
950
951
952
953
954
955
956
957
958
959
959
960
961
962
963
964
965
966
967
968
969
969
970
971
972
973
974
975
976
977
978
979
979
980
981
982
983
984
985
986
987
988
989
989
990
991
992
993
994
995
996
997
998
999

```

```

244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
    if 'deutero-' in
hydrate_deuterated[0]:
    hydrate_replacement
    =
    [hydrate_deuterated[0
    ][clip+8:]]
else:
    hydrate_replacement
    =
    [hydrate_deuterated[0
    ][clip+7:]]
for t in
range(len(hydrate_pieces)
):
    if
hydrate_pieces[t]
==

hydrate_deuterated[0
]:
    hydrate_pieces[t]
    =
    hydrate_replacement[0]
water = [r for r in
hydrate_pieces if 'hydrate'
in r and 'perhydrate' not
in r]
D2O = 0
if len(water) == 0:
    if 'deuterium oxide
solvate' in
hydrate_titles[u].lower()
    or '(deuterium oxide)
solvate' in
hydrate_titles[u].lower():
        water = [r for r in
hydrate_pieces if
'deuterium' in r or
'oxide' in r or
'solvate' in r]
elif '(dideuterium
oxide)' in
hydrate_titles[u].lower():
        water = [r for r in
hydrate_pieces if
'dideuterium' in r
or 'oxide' in r]
elif 'deuterium oxide'
in
hydrate_titles[u].lower():
        water = [r for r in
hydrate_pieces if
'deuterium' in r or
'oxide' in r]
hydrate_pieces = [v for v
in hydrate_pieces if v not
in water]
hydrate_titles[u] =
'.join(hydrate_pieces)
waterless_titles =
list(entry3[a].synonyms)
waterless_titles.append(entry3[a]
```

```

264                                         .chemical_name)
265                                         for u in
266                                         range(len(waterless_titles)):
267                                         caps2 = [i for i in
268                                         waterless_titles[u] if
269                                         i.isupper()]
270                                         caps2 = list(set(caps2))
271                                         for q in range(len(caps2)):
272                                         if caps2[q] != 'D' and
273                                         caps2[q] != 'L' and
274                                         caps2[q] != 'R' and
275                                         caps2[q] != 'S' and
276                                         caps2[q] != 'E' and
277                                         caps2[q] != 'Z':
278                                         waterless_titles[u]
279                                         = re.sub(caps2[q],
280                                         caps2[q].lower(),
281                                         waterless_titles[u])
282                                         else:
283                                         places2 = [j for j,
284                                         letter in
285                                         enumerate(waterless_tit-
286                                         les[u]) if letter
287                                         == caps2[q]]
288                                         for k in
289                                         range(len(places2)):
290                                         if places2[k]+1
291                                         ==
292                                         len(waterless_tit-
293                                         les[u]):
294                                         waterless_titles[u] = waterless_titles[u].replace(waterless_titles[u][places2[k]],
295                                         waterless_titles[u][places2[k]].lower())
296                                         elif
297                                         waterless_titles[u][places2[k]+1]
298                                         != ')' and
299                                         waterless_titles[u][places2[k]+1]
300                                         != '-' and
301                                         waterless_titles[u][places2[k]+1]
302                                         != ',', :
303                                         waterless_titles[u] = waterless_titles[u].replace(waterless_titles[u][places2[k]],
304                                         waterless_titles[u][places2[k]].lower())
305                                         waterless_pieces =
306                                         waterless_titles[u].split(
307                                         ')
308                                         waterless_deuterated = [s
309                                         for s in waterless_pieces
310                                         if 'deutero' in s.lower()]
311                                         if
312                                         len(waterless_deuterated)
313                                         == 1:
314                                         clip =
315                                         waterless_deuterated[0].l-
316                                         ower().index('deutero')
317                                         if 'deutero-' in
318                                         waterless_deuterated[0].l-
319                                         ower():
320                                         waterless_replacement
321                                         =
322                                         [waterless_deuterated
323                                         [0][clip+8:]]
324                                         else:

```

```

        waterless_replacement
        =
        [waterless_deuterated
         [0][clip+7:]]
    for t in
        range(len(waterless_pieces)):
        if
            waterless_pieces[t]
            ==
            waterless_deuterated[
            0]:
                waterless_pieces[
                t] =
                waterless_replace
                ment[0]
                waterless_titles[u] =
                '.join(waterless_pieces)
        if any('rac-' in x for x in
            hydrate_titles) and any('rac-' in
            y for y in waterless_titles)
        or any('DL' in x for x in
            hydrate_titles) and any('DL' in
            y for y in waterless_titles) or
        any('RS' in x for x in
            hydrate_titles) and any('RS' in
            y for y in waterless_titles) or
        any('RS' in x for x in
            hydrate_titles) and any('SR' in
            y for y in waterless_titles) or
        any('SR' in x for x in
            hydrate_titles) and any('RS' in
            y for y in waterless_titles) or
        any('+-' in x for x in
            hydrate_titles) and any('+-' in
            y for y in waterless_titles):
            same_names += 1
        if bool(set(hydrate_titles)
            & set(waterless_titles)) ==
            True:
                chiral_writer.write(entry
                2[a])
                writer1.write(entry3[a])
        else:
            possibly_the_same_names
            += 1
            DoubleR += 1
    elif any('D' in x for x in
        hydrate_titles) and any('D' in
        y for y in waterless_titles) or
    any('L' in x for x in
        hydrate_titles) and any('L' in
        y for y in waterless_titles) or
    any('R' in x for x in
        hydrate_titles) and any('R' in
        y for y in waterless_titles) or
    any('S' in x for x in
        hydrate_titles) and any('S' in
        y for y in waterless_titles) or
    any('(+)-' in x for x in
        hydrate_titles) and any('(+)-' in
        y for y in waterless_titles):

```

```

hydrate_titles) and any('(+)-'
in y for y in waterless_titles)
or any('(-)-' in x for x in
hydrate_titles) and any('(-)-'
in y for y in waterless_titles)
or any('cis-' in x for x in
hydrate_titles) and any('cis-' in y for y in waterless_titles)
or any('trans-' in x for x in
hydrate_titles) and
any('trans-' in y for y in
waterless_titles) or
any('meso-' in x for x in
hydrate_titles) and any('meso-' in y for y in waterless_titles)
or any('E' in x for x in
hydrate_titles) and any('E' in y for y in waterless_titles) or
any('Z' in x for x in
hydrate_titles) and any('Z' in y for y in waterless_titles):
    same_names += 1
    if bool(set(hydrate_titles) & set(waterless_titles)) == True:
        299
        300
301
        chiral_writer.write(entry2[a])
        writer1.write(entry3[a])
    else:
        302
        303
        304
        possibly_the_same_names
        += 1
        DoubleR += 1
        305
        306
        elif all('rac-' not in x for x in hydrate_titles) and
            all('rac-' not in y for y in waterless_titles) and all('+-'
            not in x for x in hydrate_titles) and all('+-'
            not in y for y in waterless_titles) and all('D'
            not in x for x in hydrate_titles) and all('D' not
            in y for y in waterless_titles) and all('L'
            not in x for x in hydrate_titles) and all('L' not
            in y for y in waterless_titles) and all('R'
            not in x for x in hydrate_titles) and all('R' not
            in y for y in waterless_titles) and all('S'
            not in x for x in hydrate_titles) and all('S' not
            in y for y in waterless_titles) and all('(+)-'
            not in x for x in hydrate_titles) and
            all('(+)-' not in y for y in waterless_titles) and
            all('(-)-' not in x for x in hydrate_titles) and all('(-)-'
            not in y for y in waterless_titles) and
            all('cis-' not in x for x in hydrate_titles) and all('cis-' not
            in y for y in waterless_titles) and
            all('trans-' not in x for x in hydrate_titles) and
            all('trans-' not in y for y in waterless_titles):

```

307  
308  
309  
310

311

```
all('trans-' not in y for y in
waterless_titles) and
all('meso-' not in x for x in
hydrate_titles) and all('meso-' not in y for y in
waterless_titles) and all('E' not in x for x in
hydrate_titles) and all('E' not in y for y in waterless_titles)
possibly_the_same_names += 1
DoubleBlank += 1
else:
    if any('D' in x for x in
hydrate_titles) or any('L' in x for x in
hydrate_titles) or any('R' in x for x in
hydrate_titles) or any('S' in x for x in
hydrate_titles) or
any('(+)-' in x for x in
hydrate_titles) or
any('(-)-' in x for x in
hydrate_titles) or
any('cis-' in x for x in
hydrate_titles) or
any('trans-' in x for x in
hydrate_titles) or any('E' in x for x in
hydrate_titles) or any('Z' in x for x in
hydrate_titles):
        if all('rac-' not in y
for y in
waterless_titles) and
all('+-' not in y for y in waterless_titles)
and all('D' not in y
for y in
waterless_titles) and
all('L' not in y for y in waterless_titles)
and all('R' not in y
for y in
waterless_titles) and
all('S' not in y for y in waterless_titles)
and all('(+)-' not in y
for y in
waterless_titles) and
all('(-)-' not in y for y in waterless_titles)
and all('cis-' not in y
for y in
waterless_titles) and
all('trans-' not in y
for y in
waterless_titles) and
all('meso-' not in y
for y in
waterless_titles) and
all('E' not in y for y in waterless_titles)
and all('Z' not in y
```

```

312
313
314
315
316
317
318
319
320
for y in
waterless_titles):
    possibly_the_same_names += 1
    RandBlank += 1
else:
    distinct_names += 1
writer2.write(entry2[
a])
writer3.write(entry3[
a])
elif any('D' in y for y in
waterless_titles) or
any('L' in y for y in
waterless_titles) or
any('R' in y for y in
waterless_titles) or
any('S' in y for y in
waterless_titles) or
any('(+)-' in y for y in
waterless_titles) or
any('(-)-' in y for y in
waterless_titles) or
any('cis-' in y for y in
waterless_titles) or
any('trans-' in y for y in
waterless_titles) or
any('E' in y for y in
waterless_titles) or
any('Z' in y for y in
waterless_titles):
    if all('rac-' not in x
for x in
hydrate_titles) and
all('+-' not in x for x
in hydrate_titles) and
all('D' not in x for x
in hydrate_titles) and
all('L' not in x for x
in hydrate_titles) and
all('R' not in x for x
in hydrate_titles) and
all('S' not in x for x
in hydrate_titles) and
all('(+)-' not in x for
x in hydrate_titles)
and all('(-)-' not in x
for x in
hydrate_titles) and
all('cis-' not in x for
x in hydrate_titles)
and all('trans-' not in
x for x in
hydrate_titles) and
all('meso-' not in x
for x in
hydrate_titles) and
all('E' not in x for x
in hydrate_titles) and
all('Z' not in x for x
in hydrate_titles):
    possibly_the_same_names += 1

```

```

321                                         RandBlank += 1
322
323
324
325                                         writer2.write(entry2[a])
326                                         writer3.write(entry3[a])
327
328
329
330
331                                         if possibly_the_same_names != 0
332                                         and hydrate_tags == []:
333                                         if
334                                         hash(entry2[a].identifier)
335                                         in
336                                         hydrate_tags_empty_no_stereo_
337                                         hydrates:
338                                         indices = [i for i, x
339                                         in
340                                         enumerate(hydrate_tags_em-
341                                         pty_no_stereo_hydrates)
342                                         if x ==
343                                         hash(entry2[a].identifier)
344                                         ]
345                                         if
346                                         any(hydrate_tags_empty_no_
347                                         _stereo_waterless_forms[j]
348                                         ] ==
349                                         hash(entry3[a].identifier)
350                                         ) for j in indices) ==
351                                         True:
352                                         if DoubleR != 0:
353                                         writer8.write(ent-
354                                         ry2[a])
355                                         writer9.write(ent-
356                                         ry3[a])
357                                         else:
358                                         chiral_writer.wri-
359                                         te(entry2[a])
360                                         writer1.write(ent-
361                                         ry3[a])
362                                         else:
363                                         print
364                                         entry2[a].identifier
365                                         print
366                                         entry3[a].identifier
367                                         print
368                                         hydrate_tags[90]
369                                         elif
370                                         hash(entry2[a].identifier)
371                                         in
372                                         hydrate_tags_empty_one_stereo_
373                                         _hydrates and
374                                         hydrate_tags_empty_one_stereo_
375                                         _waterless_forms[hydrate_tags_
376                                         _empty_one_stereo_hydrates.in-
377                                         dex(hash(entry2[a].identifier)
378                                         ))] ==
379                                         hash(entry3[a].identifier)

```

```

346                                     if DoubleR != 0:
347                                         writer8.write(entry2[
348                                             a])
349                                         writer9.write(entry3[
350                                             a])
351                                     else:
352                                         chiral_writer.write(e
353                                             ntry2[a])
354                                         writer1.write(entry3[
355                                             a])
356                                         if DoubleBlank != 0:
357                                         writer4.write(ent
358                                             ry2[a])
359                                         writer5.write(ent
360                                             ry3[a])
361                                         if RandBlank != 0:
362                                         writer6.write(ent
363                                             ry2[a])
364                                         writer7.write(ent
365                                             ry3[a])
366                                     else:
367                                         print
368                                             entry2[a].identifier
369                                         print
370                                             entry3[a].identifier
371                                         print hydrate_tags[90]
372                                         if possibly_the_same_names != 0
373                                         and hydrate_tags != []:
374                                         hydrate_chirality = []
375                                         for c in
376                                         range(len(hydrat
377                                             e_tags)):
378                                         for e in
379                                         range(len(entry2[a].molec
380                                             ule.components[hydrate_ta
381                                             gs[c]].atoms)):
382                                         if
383                                             entry2[a].molecule.co
384                                             mponents[hydrate_tags[
385                                                 c]].atoms[e].is_chir
386                                             al == True:
387                                         hydrate_chirality
388                                             .append(entry2[a]
389                                             .molecule.compone
390                                             nts[hydrate_tags[
391                                                 c]].atoms[e].chir
392                                             ality)
393                                         waterless_chirality = []
394                                         for d in
395                                         range(len(waterless_tags)):
396                                         for f in
397                                         range(len(entry3[a].molec
398                                             ule.components[waterless_
399                                             tags[d]].atoms)):
400                                         if
401                                             entry3[a].molecule.co
402                                             mponents[waterless_ta
403                                             gs[d]].atoms[f].is_chir
404                                             ality)

```

373

iral == True:

374

375

```
waterless_chirality.append(entry3[a].molecule.components[waterless_tags[d]].atoms[f].chirality)
```

376

377

```
if hydrate_chirality == []
and waterless_chirality != []
or hydrate_chirality != []
and waterless_chirality ==
[]:
    if DoubleR != 0:
```

378

```
writer8.write(entry2[a])
```

379

380

```
writer9.write(entry3[a])
```

else:

381

```
chiral_writer.write(entry2[a])
```

382

383

```
writer1.write(entry3[a])
```

```
if DoubleBlank != 0:
```

384

```
writer4.write(entry2[a])
```

385

386

```
writer5.write(entry3[a])
```

```
if RandBlank != 0:
```

387

```
writer6.write(entry2[a])
```

```
writer7.write(entry3[a])
```

388

389

```
if hydrate_chirality == []
and waterless_chirality == []
[]:
```

390

391

```
    if DoubleR != 0:
```

392

```
writer8.write(entry2[a])
```

393

394

```
writer9.write(entry3[a])
```

395

else:

396

397

```
chiral_writer.write(entry2[a])
```

398

399

```
writer1.write(entry3[a])
```

```
if hydrate_chirality != []
and waterless_chirality != []
[]:
```

```
    same_molecules = 0
```

```
possibly_the_same_molecule
```

```

400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
        es = 0
        distinct_molecules = 0
        for c in
            range(len(hydrate_tags)):
                for d in
                    range(len(waterless_t
ags)):
                        if
                            entry2[a].molecul
e.components[hydr
ate_tags[c]].form
                            ula ==
                            entry3[a].molecul
e.components[wate
rless_tags[d]].fo
rmula:
                                for e in
                                    if
                                        entry2[a].molecu
le.components[hydr
ate_tags[c]].atoms[e].atomic_symbol not in
                                        atomic_symbols:
                                            atomic_symbols.append(entry2[a].molecule.components[hydrate_tags[c]].atoms[e].atomic_symbol)
                                same_elements = [x for x in entry2[a].molecule.components[hydrate_tags[c]].atoms if
                                x.atomic_symbol == entry2[a].molecule.components[hydrate_tags[c]].atoms[e].atomic_symbol]
                                chiral_centers = []
                                for
                                    f in range(len(same_elements)):
                                        if len(same_elements[f].bonds) == 4:
                                            string1 = str(same_elements[f].bonds)
                                            string1 = string1.replace(str(same_elements[f]), '')
                                            if ')' )' in string1:
                                                string1 = string1.replace(')' ), ')')
                                            if ' Atom' in string1:
                                                string1 = string1.replace(' Atom', ' Atom')
                                            string1 = re.sub(r'[0-9]+', '', string1)
                                            capitals = [i for i in range(len(string1)) if string1[i].isupper() == True]
                                            extra = [j for j in capitals if j-1 in capitals]
                                            if extra != []:
                                                if len(extra) == 1:
                                                    string2 = string1[:extra[0]]
                                                    omit_to = string1[extra[0]:].index(')')
                                                    string2 = string2 + string1[extra[0]+omit_to:]

```

```

428
429     else:
430
431         for k in range(len(extra)):
432
433             if k == 0:
434
435                 string2 = string1[:extra[k]]
436
437             elif k == len(extra)-1:
438
439                 string2 = string2 + string1[extra[k-1]+omit_to:extra[k]]
440
441                 omit_to = string1[extra[k]:].index(')')
442
443                 string2 = string2 + string1[extra[k]+omit_to:]
444
445             else:
446
447                 string2 = string2 + string1[extra[k-1]+omit_to:extra[k]]
448
449                 omit_to = string1[extra[k]:].index(')')
450
451         if string2.count('Atom(H)') <= 1 and string2.count('Atom(D)') == 0 or
452         string2.count('Atom(H)') == 0 and string2.count('Atom(D)') <= 1:
453
454             chiral_centers.append(same_elements[f])
455
456     else:
457
458         if string1.count('Atom(H)') <= 1 and string1.count('Atom(D)') == 0 or
459         string1.count('Atom(H)') == 0 and string1.count('Atom(D)') <= 1:
460
461             chiral_centers.append(same_elements[f])
462
463 ordered_hydrate_chiral_centers = ordered_hydrate_chiral_centers + chiral_centers
464                                         start = 0
465
466 waterless_matches = []
467                                         while
468 len(ordered_hydrate_chiral_centers) != len(ordered_waterless_chiral_centers):
469                                         if
470 waterless_matches == []:
471
472 symbol_matches = [x for x in entry3[a].molecule.components[waterless_tags[d]].atoms if
473 x.atomic_symbol == ordered_hydrate_chiral_centers[start].atomic_symbol]
474                                         for
475 t in range(len(symbol_matches)):
476
477 if len(symbol_matches[t].bonds) == 4:
478
479     string3 = str(symbol_matches[t].bonds)
480
481     string3 = string3.replace(str(symbol_matches[t]), ' ')
482
483     if ')' )' in string3:
484
485         string3 = string3.replace(')' ), ')')
486
487     if ' Atom' in string3:
488
489         string3 = string3.replace(' Atom', ' Atom')
490
491     string3 = re.sub(r'[0-9]+', ' ', string3)
492
493

```

```

    capitals = [i for i in range(len(string3)) if string3[i].isupper() == True]
461
462     extra = [j for j in capitals if j-1 in capitals]
463
464     if extra != []:
465
466         string4 = string3[:extra[0]]
467
468         omit_to = string3[extra[0]:].index(' ')
469
470         string4 = string4 + string3[extra[0]+omit_to:]
471
472     else:
473
474         for k in range(len(extra)):
475
476             if k == 0:
477
478                 string4 = string3[:extra[k]]
479
480                 omit_to = string3[extra[k]:].index(' ')
481
482             elif k == len(extra)-1:
483
484                 string4 = string4 + string3[extra[k-1]+omit_to:extra[k]]
485
486                 omit_to = string3[extra[k]:].index(' ')
487
488                 string4 = string4 + string3[extra[k]+omit_to:]
489
490             else:
491
492                 string4 = string4 + string3[extra[k-1]+omit_to:extra[k]]
493
494                 omit_to = string3[extra[k]:].index(' ')
495
496             if string4.count('Atom(H)') <= 1 and string4.count('Atom(D)') == 0 or
497             string4.count('Atom(H)') == 0 and string4.count('Atom(D)') <= 1:
498
499                 waterless_matches.append(symbol_matches[t])
500
501             else:
502
503                 if string3.count('Atom(H)') <= 1 and string3.count('Atom(D)') == 0 or
504                 string3.count('Atom(H)') == 0 and string3.count('Atom(D)') <= 1:
505
506                     waterless_matches.append(symbol_matches[t])
507
508     waterless_bonds = []
509
510     updated_waterless_bonds = []
511
512     f in range(len(waterless_matches)):
513
514         pre_waterless_bonds = []
515
516         for y in range(len(waterless_matches[f].bonds)):
517
518             string2 = str(waterless_matches[f].bonds[y])
519
520             string2 = string2.replace(str(waterless_matches[f]), ' ')
521
522             if ') )' in string2:
523
524                 string2 = string2.replace(') )', '))')

```

```

493     if ' Atom' in string2:
494
495         string2 = string2.replace(' Atom', ' Atom')
496
497         pre_waterless_bonds.append(string2)
498
499         pre_waterless_bonds = sorted(pre_waterless_bonds)
500
501         waterless_bonds.append(pre_waterless_bonds)
502
503         updated_waterless_bonds.append(pre_waterless_bonds)
504
505     else:
506
507         waterless_bonds = list(updated_waterless_bonds)
508
509         if len(waterless_matches) == 1:
510
511             ordered_waterless_chiral_centers.append(waterless_matches[0])
512
513             waterless_matches.remove(waterless_matches[0])
514
515             start += 1
516
517         else:
518
519             hydrate_spider_atoms = []
520
521             waterless_spider_atoms = []
522
523             found_hydrate_atoms = []
524
525             hydrate_bonds = []
526
527             for z in range(len(ordered_hydrate_chiral_centers[start].bonds)):
528
529                 string1 = str(ordered_hydrate_chiral_centers[start].bonds[z])
530
531                 string1 = string1.replace(str(ordered_hydrate_chiral_centers[start]), '')
532
533                 if ') )' in string1:
534
535                     string1 = string1.replace(') )', ')')
536
537                 if ' Atom' in string1:
538
539                     string1 = string1.replace(' Atom', ' Atom')
540
541                     string1 = re.sub(r'[0-9]+', '', string1)
542
543                     hydrate_bonds.append(string1)
544
545                     hydrate_bonds = sorted(hydrate_bonds)
546
547                     hydrate_neighbours = ordered_hydrate_chiral_centers[start].neighbours
548
549                     found_hydrate_atoms.append(ordered_hydrate_chiral_centers[start])
550
551                     numberless_hydrate_neighbours = []
552
553                     numberless_waterless_neighbours = []
554
555                     found_match = 0
556
557                     while found_match == 0:
558
559                         if

```

```

bond_matches = []
528
atom_matches = []
529
if numberless_hydrate_neighbours == []:
530
    for w in range(len(waterless_bonds)):
531
        if waterless_bonds[w] == hydrate_bonds:
532
            bond_matches.append(waterless_bonds[w])
533
            atom_matches.append(w)
534
else:
535
    if hydrate_bonds == [] and all(x == [] for x in waterless_bonds):
536
        bond_matches.append(waterless_bonds[w])
537
        atom_matches.append(w)
538
else:
539
    for w in range(len(numberless_waterless_neighbours)):
540
        off_the_table_waterless = []
541
        off_the_table_hydrate = []
542
        match = 0
543
        for x in range(len(numberless_waterless_neighbours[w])):
544
            if len(waterless_bonds[w]) == len(hydrate_bonds):
545
                if waterless_neighbours[w][x] not in off_the_table_waterless:
546
                    go = 0
547
                    for u in range(len(numberless_hydrate_neighbours)):
548
                        if go == 0:
549
                            if hydrate_neighbours[u] not in off_the_table_hydrate:
550
                                if numberless_waterless_neighbours[w][x] ==
numberless_hydrate_neighbours[u]:
551
                                    waterless_to_compare = []
552
                                    hydrate_to_compare = []
553
                                    if x == 0:
554
                                        for t in range(waterless_spider_atoms[w][x]):
555
                                            waterless_to_compare.append(waterless_bonds[w][t])
556
                                        else:
557
                                            for t in
range(waterless_spider_atoms[w][x-1], waterless_spider_atoms[w][x]):
558
                                                waterless_to_compare.append(waterless_bonds[w][t])

```

```

559                     if u == 0:
560                         for s in range(hydrate_spider_atoms[u]):
561
562                             hydrate_to_compare.append(hydrate_bonds[s])
563
564                     else:
565                         for s in range(hydrate_spider_atoms[u-1],
566 hydrate_spider_atoms[u]):
567
568                             hydrate_to_compare.append(hydrate_bonds[s])
569
570                         if waterless_to_compare == hydrate_to_compare:
571
572                             off_the_table_waterless.append(waterless_neighbours[w][x])
573
574                             off_the_table_hydrate.append(hydrate_neighbours[u])
575
576                             match += 1
577
578                             go += 1
579
580                             if match == len(numberless_hydrate_neighbours):
581
582                                 bond_matches.append(waterless_bonds[w])
583
584                                 atom_matches.append(w)
585
586
587                             if len(bond_matches) == 1:
588
589                                 pin2 = waterless_bonds.index(bond_matches[0])
590
591                                 ordered_waterless_chiral_centers.append(waterless_matches[pin2])
592
593                                 waterless_matches.remove(waterless_matches[pin2])
594
595                                 updated_waterless_bonds.remove(updated_waterless_bonds[pin2])
596
597                                 start += 1
598
599                                 found_match += 1
600
601                             else:
602
603                                 if numberless_hydrate_neighbours != []:
604
605                                     hydrate_neighbours = list(new_hydrate_neighbours)
606
607                                     numberless_hydrate_neighbours = []
608
609                                     hydrate_spider_atoms = []
610
611                                     for g in range(len(hydrate_neighbours)):
612
613                                         string1 = str(hydrate_neighbours[g])
614
615                                         string2 = re.sub(r'[0-9]+', ' ', string1)
616
617                                         numberless_hydrate_neighbours.append(string2)

```

```

590
591     new_hydrate_neighbours = []
592
593     hydrate_atom_branches = 0
594
595     for j in range(len(hydrate_neighbours)):
596
597         if hydrate_neighbours[j] not in found_hydrate_atoms:
598
599             for k in range(len(hydrate_neighbours[j].neighbours)):
600
601                 if hydrate_neighbours[j].neighbours[k] not in found_hydrate_atoms:
602
603                     new_hydrate_neighbours.append(hydrate_neighbours[j].neighbours[k])
604
605                     hydrate_atom_branches += 1
606
607             hydrate_spider_atoms.append(hydrate_atom_branches)
608
609     hydrate_bonds = []
610
611     for l in range(len(hydrate_neighbours)):
612
613         pre_hydrate_bonds = []
614
615         for z in range(len(hydrate_neighbours[l].bonds)):
616
617             if all(str(x) not in str(hydrate_neighbours[l].bonds[z]) for x in
618                  found_hydrate_atoms):
619
620                 string1 = str(hydrate_neighbours[l].bonds[z])
621
622                 string1 = string1.replace(str(hydrate_neighbours[l]), ' ')
623
624                 if ')' )' in string1:
625
626                     string1 = string1.replace(')' ), ' )')
627
628                 if ' Atom' in string1:
629
630                     string1 = string1.replace(' Atom', ' Atom')
631
632                     string1 = re.sub(r'[0-9]+', ' ', string1)
633
634                     pre_hydrate_bonds.append(string1)
635
636             pre_hydrate_bonds = sorted(pre_hydrate_bonds)
637
638             hydrate_bonds = hydrate_bonds + pre_hydrate_bonds
639
640     if numberless_waterless_neighbours != []:
641
642         waterless_neighbours = list(new_waterless_neighbours)
643
644         waterless_neighbours_to_remove = []
645
646         for y in range(len(waterless_neighbours)):
647
648             if y not in atom_matches:
649
650                 waterless_neighbours_to_remove.append(waterless_neighbours[y])
651
652         for z in range(len(waterless_neighbours_to_remove)):
653
654             waterless_neighbours.remove(waterless_neighbours_to_remove[z])
655
656     numberless_waterless_neighbours = []

```

```

623         waterless_spider_atoms = []
624
625     else:
626
627         found_waterless_atoms = []
628
629         waterless_neighbours = []
630
631         for y in range(len(waterless_matches)):
632
633             waterless_neighbours.append(waterless_matches[y].neighbours)
634
635             if str(waterless_matches[y]) != 'Atom(H)':
636
637                 found_waterless_atoms.append([waterless_matches[y]])
638
639         for f in range(len(waterless_neighbours)):
640
641             pre_numberless_waterless_neighbours = []
642
643             for e in range(len(waterless_neighbours[f])):
644
645                 string1 = str(waterless_neighbours[f][e])
646
647                 string2 = re.sub(r'[0-9]+', '', string1)
648
649                 pre_numberless_waterless_neighbours.append(string2)
650
651             numberless_waterless_neighbours.append(pre_numberless_waterless_neighbours)
652
653             new_waterless_neighbours = []
654
655             for m in range(len(waterless_neighbours)):
656
657                 waterless_atom_branches = 0
658
659                 pre_new_waterless_neighbours = []
660
661                 pre_waterless_spider_atoms = []
662
663                 for n in range(len(waterless_neighbours[m])):
664
665                     if waterless_neighbours[m][n] not in found_waterless_atoms[m]:
666
667                         for o in range(len(waterless_neighbours[m][n].neighbours)):
668
669                             if waterless_neighbours[m][n].neighbours[o] not in
670                             found_waterless_atoms[m]:
671
672
673                             pre_new_waterless_neighbours.append(waterless_neighbours[m][n].neighbours[o])
674
675                             waterless_atom_branches += 1
676
677                             pre_waterless_spider_atoms.append(waterless_atom_branches)
678
679             new_waterless_neighbours.append(pre_new_waterless_neighbours)
680
681             waterless_spider_atoms.append(pre_waterless_spider_atoms)
682
683             waterless_bonds = []
684
685             for p in range(len(waterless_neighbours)):
686
687                 pre_waterless_bonds = []

```

```

656         for q in range(len(waterless_neighbours[p])):
657             extra_pre_waterless_bonds = []
658             for z in range(len(waterless_neighbours[p][q].bonds)):
659                 if all(str(x) not in str(waterless_neighbours[p][q].bonds[z]) for x in
660 found_waterless_atoms[p]):
661                     string2 = str(waterless_neighbours[p][q].bonds[z])
662                     string2 = string2.replace(str(waterless_neighbours[p][q]), '')
663                     if ')' )' in string2:
664                         string2 = string2.replace(')' ), ')')
665                     if ' Atom' in string2:
666                         string2 = string2.replace(' Atom', ' Atom')
667                     string2 = re.sub(r'[0-9]+', '', string2)
668                     extra_pre_waterless_bonds.append(string2)
669                     extra_pre_waterless_bonds = sorted(extra_pre_waterless_bonds)
670                     pre_waterless_bonds = pre_waterless_bonds + extra_pre_waterless_bonds
671                     waterless_bonds.append(pre_waterless_bonds)
672             for j in range(len(hydrate_neighbours)):
673                 if str(hydrate_neighbours[j]) != 'Atom(H)':
674                     found_hydrate_atoms.append(hydrate_neighbours[j])
675             for m in range(len(waterless_neighbours)):
676                 for n in range(len(waterless_neighbours[m])):
677                     if str(waterless_neighbours[m][n]) != 'Atom(H)':
678                         found_waterless_atoms[m].append(waterless_neighbours[m][n])
679
680         if len(waterless_neighbours) == 0:
681             ordered_waterless_chiral_centers.append(waterless_matches[0])
682             waterless_matches.remove(waterless_matches[0])
683             updated_waterless_bonds.remove(updated_waterless_bonds[0])
684             start += 1
685             found_match += 1
686
687         hydrate_chirality = []
688         waterless_chirality = []
689         for e in
690         if

```

```

ordered_hydrate_chiral_centers[e].is_chiral == True:
689
690     hydrate_chirality.append(ordered_hydrate_chiral_centers[e].chirality)
691
692     hydrate_chirality.append(' ')
693
694     for f in range(len(ordered_waterless_chiral_centers)):
695         if ordered_waterless_chiral_centers[f].is_chiral == True:
696             waterless_chirality.append(ordered_waterless_chiral_centers[f].chirality)
697
698         else:
699             waterless_chirality.append(' ')
700
701         if hydrate_chirality == waterless_chirality:
702             same_molecules += 1
703
704             elif [s for s in hydrate_chirality if s != ' '] == [] or [t for t in waterless_chirality if t != ' '] == []:
705
706             possibly_the_same_molecules += 1
707
708             if distinct_molecules == 0:
709                 if possibly_the_same_molecules == 0:
710                     if DoubleR != 0:
711                         writer8.write(entry2[a])
712
713                         writer9.write(entry3[a])
714
715                         chiral_writer.write(entry2[a])
716
717                         writer1.write(entry3[a])
718
719                         writer8.write(entry2[a])
720
721                         writer9.write(entry3[a])
722
723                         chiral_writer.write(entry2[a])
724
725                         writer1.write(entry3[a])
726
727                         if DoubleBlank != 0:
728                             writer4.write(entry2[a])
729
730                             writer5.write(entry3[a])
731
732                             RandBlank != 0:
733
734                             writer6.write(entry2[a])
735
736                             writer7.write(entry3[a])

```

```
725  
726  
727  
728  
729  
730  
731  
732  
733  
734 if __name__ == '__main__':  
    # This runs the script  
    r = Runner()  
    r.run()  
735  
736  
737  
738  
  
    else:  
        writer2.write(entry2[a])  
        writer3.write(entry3[a])  
  
    loop += 1  
    print 'another loop bites the dust'
```

```

1 #24 Create the three main classes
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script creates the three main classes of structures: (1) hydrate-anhydride pairs,
10 #hydrates without a known anhydrous form, and (3) anhydrous forms without a known
11 #hydrate-anhydride
12 #It also places hydrates and waterless forms from hydrate-anhydride pairs that have
13 #distinct chirality in classes 2 and 3 if there are no true hydrate-anhydride pairs
14 #remaining for those structures
15
16 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
17 Forms\None_solvated_hydrates.txt"
18 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
19 Forms\None_solvated_waterless_forms.txt"
20 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
21 Forms\None_hydrates.txt"
22 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
23 Forms\None_waterless_forms.txt"
24 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
25 Waterless Forms\potential_None_hydrates_with_waterless_form.txt"
26 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
27 Waterless Forms\potential_None_waterless_forms_with_hydrate.txt"
28 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
29 Waterless Forms\potential_hydrates_with_waterless_form.txt"
30 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
31 Waterless Forms\potential_waterless_forms_with_hydrate.txt"
32 filepath9 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
33 Waterless Forms\potential_hydrates_without_waterless_form.txt"
34 filepath10 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
35 Waterless Forms\potential_waterless_forms_without_hydrate.txt"
36 filepath11 = "C:\Users\jenwe\Hydrates Manuscript\Step4 Find Pairs and Duplicates with
37 Distinct Chirality\complete_distinct_chirality_hydrates_in_pairs.txt"
38 filepath12 = "C:\Users\jenwe\Hydrates Manuscript\Step4 Find Pairs and Duplicates with
39 Distinct Chirality\complete_distinct_chirality_waterless_forms_in_pairs.txt"
40
41 class Runner(argparse.ArgumentParser):
42
43     def __init__(self):
44         super(self.__class__, self).__init__(description=__doc__)
45         self.add_argument(
46             '-i', '--input', default=filepath1,
47             help='input database file')
48         self.add_argument(
49             '-o', '--output', default='class1_hydrates_with_waterless_forms.gcd',
50             help='output file [class1_hydrates_with_waterless_forms.gcd]')
51         self.add_argument(
52             '-m', '--maximum', default=0, type=int,
53             help='Maximum number of structures to find [all]')
54
55     args = self.parse_args()
56
57     self.args = args
58     self.settings = search.Search.Settings()
59     self.settings.max_hit_structures = self.args.maximum
60
61     def run(self):
62
63         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
64         entry_reader2 = io.EntryReader(filepath2, format='identifiers')

```

```

52 entry_reader3 = io.EntryReader(filepath3, format='identifiers')
53 entry_reader4 = io.EntryReader(filepath4, format='identifiers')
54 entry_reader5 = io.EntryReader(filepath5, format='identifiers')
55 entry_reader6 = io.EntryReader(filepath6, format='identifiers')
56 entry_reader7 = io.EntryReader(filepath7, format='identifiers')
57 entry_reader8 = io.EntryReader(filepath8, format='identifiers')
58 entry_reader9 = io.EntryReader(filepath9, format='identifiers')
59 entry_reader10 = io.EntryReader(filepath10, format='identifiers')
60 entry_reader11 = io.EntryReader(filepath11, format='identifiers')
61 entry_reader12 = io.EntryReader(filepath12, format='identifiers')
62
63 with io.EntryWriter(self.args.output) as check_writer:
64     with io.EntryWriter("class1_waterless_forms_with_hydrate.gcd") as writer1:
65         with
66             io.EntryWriter("class2_hydrates_with_no_reported_waterless_form.gcd")
67             as writer2:
68                 with
69                     io.EntryWriter("class3_waterless_forms_with_no_reported_hydrate.gcd")
70                     as writer3:
71
72             duplicate_hash_dictionary = {'GADKOL':1, 'MECYEY':2,
73                                         'OBASAN':3, 'OMOMAE':4, 'OTONUI':5, 'QOJWAQ':6, 'TZBFZO':7,
74                                         'VOBKOZ':8, 'VUFJAI':9, 'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12,
75                                         'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
76
77             #True hydrate-anhydride pairs are placed in checked lists
78             #False hydrate-anhydride pairs are placed in false lists
79             checked_HW = []
80             checked_WH = []
81             false_HW = []
82             false_WH = []
83
84             loop = 0
85             entry1 = entry_reader1
86             list1 = []
87             list = list1
88             #A list of structures from each input file is created
89             while loop < 12:
90                 for a in range(len(entry1)):
91                     if entry1[a].identifier not in duplicate_hash_dictionary:
92                         list.append(hash(entry1[a].identifier))
93                     else:
94                         list.append(duplicate_hash_dictionary.get(entry1[a].identifier))
95
96                     loop += 1
97                     if loop == 1:
98                         entry1 = entry_reader2
99                         list2 = []
100                        list = list2
101
102                     if loop == 2:
103                         entry1 = entry_reader3
104                         list3 = []
105                         list = list3
106
107                     if loop == 3:
108                         entry1 = entry_reader4
109                         list4 = []
110                         list = list4
111
112                     if loop == 4:
113                         entry1 = entry_reader5
114                         list5 = []
115                         list = list5
116
117                     if loop == 5:
118                         entry1 = entry_reader6
119                         list6 = []
120                         list = list6
121
122                     if loop == 6:
123                         entry1 = entry_reader7
124                         list7 = []
125                         list = list7
126
127                     if loop == 7:
128                         entry1 = entry_reader8
129                         list8 = []
130                         list = list8
131
132                     if loop == 8:
133                         entry1 = entry_reader9
134                         list9 = []
135                         list = list9
136
137                     if loop == 9:
138                         entry1 = entry_reader10
139                         list10 = []
140                         list = list10
141
142                     if loop == 10:
143                         entry1 = entry_reader11
144                         list11 = []
145                         list = list11
146
147                     if loop == 11:
148                         entry1 = entry_reader12
149                         list12 = []
150                         list = list12
151
152                     if loop == 12:
153                         entry1 = entry_reader1
154                         list13 = []
155                         list = list13
156
157
158
159
160

```

```

110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
    entry1 = entry_reader7
    list7 = []
    list = list7
    if loop == 7:
        entry1 = entry_reader8
        list8 = []
        list = list8
    if loop == 8:
        entry1 = entry_reader9
        list9 = []
        list = list9
    if loop == 9:
        entry1 = entry_reader10
        list10 = []
        list = list10
    if loop == 10:
        entry1 = entry_reader11
        list11 = []
        list = list11
    if loop == 11:
        entry1 = entry_reader12
        list12 = []
        list = list12

#Hydrates from the potential hydrate-anhydrate pairs list are
checked for in the distinct chirality list
for b in range(len(list7)):
    count = 0
    for c in range(len(list11)):
        #If the hydrate and corresponding anhydrate are one of
        #the pairs in the distinct chirality list they are
        #appended to the false lists
        if list11[c] == list7[b] and list12[c] == list8[b]:
            if list7[b] not in false_HW:
                false_HW.append(list7[b])
            if list8[b] not in false_WH:
                false_WH.append(list8[b])
            count += 1
    #Otherwise, the two structures are written to class 1
    #output files
    if count == 0:
        check_writer.write(entry_reader7[b])
        if list7[b] not in checked_HW:
            checked_HW.append(list7[b])
        writer1.write(entry_reader8[b])
        if list8[b] not in checked_WH:
            checked_WH.append(list8[b])

#Hydrates-anhydrate pairs found using the None method are
#appended to the lists for class 1
for d in range(len(list5)):
    check_writer.write(entry_reader5[d])
    if list5[d] not in checked_HW:
        checked_HW.append(list5[d])
    writer1.write(entry_reader6[d])
    if list6[d] not in checked_WH:
        checked_WH.append(list6[d])

#Hydrates in the false list that are not in class 1 are written
#to the class 2 output file
for z in range(len(false_HW)):
    if false_HW[z] not in checked_HW:
        pin = list7.index(false_HW[z])
        writer2.write(entry_reader7[pin])

#Waterless forms in the false list that are not in class 1 are
#written to the class 3 output file

```

```

170     for y in range(len(false_WH)):
171         if false_WH[y] not in checked_WH:
172             pin2 = list8.index(false_WH[y])
173             writer3.write(entry_reader8[pin2])
174
175             #Hydrates that are not in class 1 and not in the false list are
176             #written to the class 2 output file
177             for f in range(len(list1)):
178                 if list1[f] not in checked_HW and list1[f] not in false_HW:
179                     writer2.write(entry_reader1[f])
180                     false_HW.append(list1[f])
181
182             #Waterless forms that are not in class 1 and not in the false
183             #list are written to the class 3 output file
184             for g in range(len(list2)):
185                 if list2[g] not in checked_WH and list2[g] not in false_WH:
186                     writer3.write(entry_reader2[g])
187                     false_WH.append(list2[g])
188
189             #Hydrates that are not in class 1 and not in the false list are
190             #written to the class 2 output file
191             for h in range(len(list3)):
192                 if list3[h] not in checked_HW and list3[h] not in false_HW:
193                     writer2.write(entry_reader3[h])
194                     false_HW.append(list3[h])
195
196             #Waterless forms that are not in class 1 and not in the false
197             #list are written to the class 3 output file
198             for i in range(len(list4)):
199                 if list4[i] not in checked_WH and list4[i] not in false_WH:
200                     writer3.write(entry_reader4[i])
201                     false_WH.append(list4[i])
202
203             #Hydrates that are not in class 1 and not in the false list are
204             #written to the class 2 output file
205             for j in range(len(list9)):
206                 if list9[j] not in checked_HW and list9[j] not in false_HW:
207                     writer2.write(entry_reader9[j])
208                     false_HW.append(list9[j])
209
210             #Waterless forms that are not in class 1 and not in the false
211             #list are written to the class 3 output file
212             for k in range(len(list10)):
213                 if list10[k] not in checked_WH and list10[k] not in false_WH:
214                     writer3.write(entry_reader10[k])
215                     false_WH.append(list10[k])
216
217
218 if __name__ == '__main__':
219     # This runs the script
220     r = Runner()
221     r.run()
222

```

```

1 #25 Find stoichiometrically distinct hydrate-anhydrate pairs
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script finds hydrate-anhydrate pairs that have distinct stoichiometry
10 #This script is the same as the one used to find duplicate matches with distinct
11 #stoichiometry
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
13 Classes\class1_hydrates_with_waterless_forms.txt"
14 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
15 Classes\class1_waterless_forms_with_hydrate.txt"
16
17 class Runner(argparse.ArgumentParser):
18
19     def __init__(self):
20         super(self.__class__, self).__init__(description=__doc__)
21         self.add_argument(
22             '-i', '--input', default=filepath1,
23             help='input database filepath')
24         self.add_argument(
25             '-o', '--output', default='unique_stoichiometry_hydrates.gcd',
26             help='output file [unique_stoichiometry_hydrates.gcd]')
27         self.add_argument(
28             '-m', '--maximum', default=0, type=int,
29             help='Maximum number of structures to find [all]')
30
31     args = self.parse_args()
32
33     self.args = args
34     self.settings = search.Search.Settings()
35     self.settings.max_hit_structures = self.args.maximum
36
37     def run(self):
38
39         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
40         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
41
42         total = 0
43         count = 0
44
45         with io.EntryWriter(self.args.output) as sample_writer:
46             with io.EntryWriter("unique_stoichiometry_waterless_forms.gcd") as writer1:
47
48                 entry1 = entry_reader1
49                 entry2 = entry_reader2
50                 for a in range(len(entry_reader1)):
51                     formulas1 = entry1[a].formula.split(',')
52                     water = []
53                     for b in range(len(formulas1)):
54                         if '(H2 O1)' in formulas1[b] or '(D2 O1)' in formulas1[b] or
55                         formulas1[b] == 'H2 O1' or formulas1[b] == 'D2 O1':
56                             water.append(formulas1[b])
57                         elif 'H' in formulas1[b] and 'D' in formulas1[b]:
58                             place = formulas1[b].index('D')
59                             place2 = formulas1[b].index('H')
60                             if len(str(formulas1[b])) == place+2:
61                                 if formulas1[b][place2+2] == ' ':
62                                     total = int(formulas1[b][place+1]) +
63                                     int(formulas1[b][place2+1])

```

```

63         formulas1[b] =
64         formulas1[b].replace(formulas1[b][place2:place2+2],
65             str('H') + str(total))
66         formulas1[b] = formulas1[b][:place-1]
67     else:
68         print 'line 58'
69         print entry_reader1[a].identifier
70     elif formulas1[b][place2+2] == ' ' and
71         formulas1[b][place+2] == ' ':
72         total = int(formulas1[b][place+1]) +
73             int(formulas1[b][place2+1])
74         formulas1[b] =
75             formulas1[b].replace(formulas1[b][place2:place2+2],
76                 str('H') + str(total))
77         formulas1[b] = formulas1[b][:place] +
78             formulas1[b][place+3:]
79     elif entry1[a].identifier == 'ROBYEM01' or
80         entry1[a].identifier == 'ELOCIR' or entry1[a].identifier == 'DOMQOL' or
81         entry1[a].identifier == 'VISZAZ' or
82         entry1[a].identifier == 'UNATAE':
83         total = int(formulas1[b][place+1]) +
84             int(formulas1[b][place2+1:place2+3])
85         formulas1[b] =
86             formulas1[b].replace(formulas1[b][place2:place2+3],
87                 str('H') + str(total))
88         formulas1[b] = formulas1[b][:place] +
89             formulas1[b][place+3:]
90     elif entry1[a].identifier == 'SIGHOF01':
91         total = int(formulas1[b][place+1:place+3]) +
92             int(formulas1[b][place2+1:place2+3])
93         formulas1[b] =
94             formulas1[b].replace(formulas1[b][place2:place2+3],
95                 str('H') + str(total))
96         formulas1[b] = formulas1[b][:place] +
97             formulas1[b][place+3:]
98     else:
99         print 'line 61'
100        print entry1[a].identifier
101    elif 'D' in formulas1[b]:
102        formulas1[b] = re.sub(r'[D]', 'H', formulas1[b])
103    formulas1 = [x for x in formulas1 if x not in water]
104    formulas2 = entry2[a].formula.split(',')
105    for c in range(len(formulas2)):
106        if 'H' in formulas2[c] and 'D' in formulas2[c]:
107            place = formulas2[c].index('D')
108            place2 = formulas2[c].index('H')
109            if len(str(formulas2[c])) == place+2:
110                if formulas2[c][place2+2] == ' ':
111                    total = int(formulas2[c][place+1]) +
112                        int(formulas2[c][place2+1])
113                    formulas2[c] =
114                        formulas2[c].replace(formulas2[c][place2:place2+2],
115                            str('H') + str(total))
116            elif entry2[a].identifier == 'ANTMEU04' or
117                entry2[a].identifier == 'HATTYUW01' or
118                entry2[a].identifier == 'GASGEM01':
119                    total = int(formulas2[c][place+1]) +
120                        int(formulas2[c][place2+1:place2+3])

```

```

101
102     formulas2[c] =
103     formulas2[c].replace(formulas2[c][place2:place2+3],
104                           str('H') + str(total))
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
    else:
        print 'line 78'
        print entry_reader2[a].identifier
        formulas2[c] = formulas2[c][:place-1]
    elif formulas2[c][place2+2] == ' ' and
formulas2[c][place+2] == ' ':
        total = int(formulas2[c][place+1]) +
int(formulas2[c][place2+1])
        formulas2[c] =
        formulas2[c].replace(formulas2[c][place2:place2+2],
str('H') + str(total))
        formulas2[c] = formulas2[c][:place] +
formulas2[c][place+3:]
    elif entry2[a].identifier == 'DIBZAQ03' or
entry2[a].identifier == 'DIBZAQ04' or entry2[a].identifier ==
'ROBYEM02' or entry2[a].identifier == 'REXFIM01' or
entry2[a].identifier == 'GUQXIB01' or entry2[a].identifier ==
'LUSLOA01' or entry2[a].identifier == 'MAQWIM09' or
entry2[a].identifier == 'MAQWIM10' or entry2[a].identifier ==
'MAQWIM11' or entry2[a].identifier == 'MAQWIM12' or
entry2[a].identifier == 'MAQWIM14' or entry2[a].identifier ==
'MAQWIM06' or entry2[a].identifier == 'MAQWIM07' or
entry2[a].identifier == 'MAQWIM08' or entry2[a].identifier ==
'MAQWIM13' or entry2[a].identifier == 'OROQET' or
entry2[a].identifier == 'ICUXUC01' or entry2[a].identifier ==
'NALWUS03' or entry2[a].identifier == 'KOFCUE05' or
entry2[a].identifier == 'NEZRUF01' or entry2[a].identifier ==
'RUJCI01' or entry2[a].identifier == 'UPIVIA01' or
entry2[a].identifier == 'VIWQUO' or entry2[a].identifier ==
'GESXOT01' or entry2[a].identifier == 'MPRGOM01' or
entry2[a].identifier == 'TUTBUG01' or entry2[a].identifier ==
'KOFCUE04':
        total = int(formulas2[c][place+1]) +
int(formulas2[c][place2+1:place2+3])
        formulas2[c] =
        formulas2[c].replace(formulas2[c][place2:place2+3],
str('H') + str(total))
        formulas2[c] = formulas2[c][:place] +
formulas2[c][place+3:]
    elif entry2[a].identifier == 'SIGHOF02' or
entry2[a].identifier == 'LEPXOT01':
        total = int(formulas2[c][place+1:place+3]) +
int(formulas2[c][place2+1:place2+3])
        formulas2[c] =
        formulas2[c].replace(formulas2[c][place2:place2+3],
str('H') + str(total))
        formulas2[c] = formulas2[c][:place] +
formulas2[c][place+3:]
    elif entry2[a].identifier == 'YUHHUE02' or
entry2[a].identifier == 'NALWUS04' or entry2[a].identifier ==
'NALWUS05' or entry2[a].identifier == 'NALWUS06' or
entry2[a].identifier == 'NALWUS07':
        total = int(formulas2[c][place+1:place+3]) +
int(formulas2[c][place2+1])
        formulas2[c] =
        formulas2[c].replace(formulas2[c][place2:place2+2],
str('H') + str(total))
        formulas2[c] = formulas2[c][:place] +
formulas2[c][place+3:]
    else:
        print 'line 81'
        print entry2[a].identifier
    elif 'D' in formulas2[c]:
        formulas2[c] = re.sub(r'[D]', 'H', formulas2[c])

```

```
127     formulas1 = sorted(formulas1)
128     formulas2 = sorted(formulas2)
129     if len(formulas1) != 1 and len(formulas2) != 1:
130         if formulas1 != formulas2:
131             sample_writer.write(entry1[a])
132             writer1.write(entry2[a])
133
134
135
136
137 if __name__ == '__main__':
138     # This runs the script
139     r = Runner()
140     r.run()
141
```

```

1 #26 Determine numbers for flowchart
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import re
8
9 #This script refines the three classes that were generated previously to contain
10 structures without metals, duplicates, or mismatched stoichiometry
11 #Duplicates and metals can be removed before finding hydrate-anhydrite pairs as shown
12 in the flowchart of the paper
13 #Or they can be removed from everything simultaneously at this point which is what was
14 done in this work
15
16 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
17 Entries\metals_hydrates.txt"
18 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
19 Entries\metals_waterless_forms.txt"
20 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
21 Entries\hydrates_all_duplicates.txt"
22 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step3 Find Metal and Duplicate
23 Entries\waterless_forms_all_duplicates.txt"
24 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
25 Classes\class1_hydrates_with_waterless_forms.txt"
26 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
27 Classes\class1_waterless_forms_with_hydrate.txt"
28 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
29 Classes\class2_hydrates_with_no_reported_waterless_form.txt"
30 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
31 Classes\class3_waterless_forms_with_no_reported_hydrate.txt"
32 filepath9 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
33 Forms\complete_None_hydrates.txt"
34 filepath10 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
35 Forms\complete_None_waterless_forms.txt"
36 filepath11 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
37 Forms\complete_smiles_hydrates.txt"
38 filepath12 = "C:\Users\jenwe\Hydrates Manuscript\Step2 Find Hydrates and Waterless
39 Forms\complete_smiles_waterless_forms.txt"
40 filepath13 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
41 Waterless Forms\potential_hydrates_with_waterless_form.txt"
42 filepath14 = "C:\Users\jenwe\Hydrates Manuscript\Step5 Find Paired Hydrates and
43 Waterless Forms\potential_waterless_forms_with_hydrate.txt"
44 filepath15 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
45 Classes\unique_stoichiometry_hydrates.txt"
46 filepath16 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
47 Classes\unique_stoichiometry_waterless_forms.txt"
48
49 class Runner(argparse.ArgumentParser):
50
51     def __init__(self):
52         super(self.__class__, self).__init__(description=__doc__)
53         self.add_argument(
54             '-i', '--input', default=filepath1,
55             help='input database filepath')
56         self.add_argument(
57             '-o', '--output', default='chart_class1_hydrates_with_waterless_forms.gcd',
58             help='output file [chart_class1_hydrates_with_waterless_forms.gcd]')
59         self.add_argument(
60             '-m', '--maximum', default=0, type=int,
61             help='Maximum number of structures to find [all]')
62
63     args = self.parse_args()

```

```

49     self.args = args
50     self.settings = search.Search.Settings()
51     self.settings.max_hit_structures = self.args.maximum
52
53     def run(self):
54
55         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
56         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
57         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
58         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
59         entry_reader5 = io.EntryReader(filepath5, format='identifiers')
60         entry_reader6 = io.EntryReader(filepath6, format='identifiers')
61         entry_reader7 = io.EntryReader(filepath7, format='identifiers')
62         entry_reader8 = io.EntryReader(filepath8, format='identifiers')
63         entry_reader9 = io.EntryReader(filepath9, format='identifiers')
64         entry_reader10 = io.EntryReader(filepath10, format='identifiers')
65         entry_reader11 = io.EntryReader(filepath11, format='identifiers')
66         entry_reader12 = io.EntryReader(filepath12, format='identifiers')
67         entry_reader13 = io.EntryReader(filepath13, format='identifiers')
68         entry_reader14 = io.EntryReader(filepath14, format='identifiers')
69         entry_reader15 = io.EntryReader(filepath15, format='identifiers')
70         entry_reader16 = io.EntryReader(filepath16, format='identifiers')
71
72         total = 0
73         count = 0
74
75         with io.EntryWriter(self.args.output) as sample_writer:
76             with io.EntryWriter("chart_class1_waterless_forms_with_hydrate.gcd") as
77                 writer1:
78                     with
79                         io.EntryWriter("chart_class2_hydrates_without_waterless_forms.gcd") as
80                             writer2:
81                                 with
82                                     io.EntryWriter("chart_class3_waterless_forms_without_hydrate.gcd") as
83                                         writer3:
84                                             with
85                                                 io.EntryWriter("chart_complete_None_hydrates.gcd") as
86                                                 writer4:
87                                                     with
88                                                         io.EntryWriter("chart_complete_None_waterless_forms.gcd") as
89                                                         writer5:
90                                                             with
91                                                               io.EntryWriter("chart_complete_smiles_hydrates.gcd") as
92                                                               writer6:
93                                                                   with
94                                                                       io.EntryWriter("chart_complete_smiles_waterless_forms.gcd") as
95                                                                       writer7:
96                                                                           with
97                                                                             io.EntryWriter("chart_potential_pairs_hydrates.gcd") as
98                                                                             writer8:
99                                                                                 with
100                                                                 io.EntryWriter("chart_potential_pairs_waterless_forms.gcd") as
101                                                                 writer9:
102
103             duplicate_identifier_dictionary =
104             {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
105              'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6,
106              'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
107              'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12,
108              'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
109
110             loop = 0
111             entry1 = entry_reader1
112             entry2 = entry_reader2
113             dont_want_hydrates = []
114             dont_want_waterless_forms = []
115             #All the hydrate and waterless form
116             #structures with metals are placed into

```

95  
96  
97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

```

lists with a "dont_want" label
#All the hydrate and waterless form
structures that are duplicates are
placed in the lists with a "dont_want"
label
while loop < 2:
    for a in range(len(entry1)):
        dont_want_hydrates.append(hash(entry1[a].identifier))
    for b in range(len(entry2)):
        if entry2[b].identifier in
            duplicate_identifier_dictionary:
                dont_want_waterless_forms.append(hash(entry2[b].identifier))
        else:
            dont_want_waterless_forms.append(hash(entry2[b].identifier))
loop += 1
if loop == 1:
    entry1 = entry_reader3
    entry2 = entry_reader4

print 'lists of metals and duplicates
complete'
print len(dont_want_hydrates)
print len(dont_want_waterless_forms)

#Hydrate and waterless forms that do
not have the same stoichiometry are
placed into two lists
hydrates_not_pair = []
for c in range(len(entry_reader15)):

    hydrates_not_pair.append(hash(entry_reader15[c].identifier))

waterless_forms_not_pair = []
for d in range(len(entry_reader16)):
    if entry_reader16[d].identifier not
        in duplicate_identifier_dictionary:
        waterless_forms_not_pair.append(hash(entry_reader16[d].identifier))
    else:
        waterless_forms_not_pair.append(duplicate_identifier_dictionary.get(entry_reader16[d].identifier))

print 'list of stoichiometrically
distinct complete'
print len(hydrates_not_pair)
print len(waterless_forms_not_pair)

loop1 = 0
entry3 = entry_reader5
entry4 = entry_reader6
unpaired_hydrates = []
unpaired_waterless_forms = []

```



164  
165166  
167

168

169  
170  
171172  
173

174

175

176

177  
178  
179  
180181     partner\_hydrates.append(entry3[g])  
182     partner\_waterless\_forms.append(entry4[g])

183

184  
185  
186  
187

```

    ry4[f])

unpaired_waterless_forms.append(duplicate_identifier_dictionary.get(
    entry4[f].identifier))

if loop1 == 2:
    writer5.write(entry4[f])
if loop1 == 3:
    writer7.write(entry4[f])

#For lists of hydrate-anhydrate pairs, pairs that have distinct stoichiometry are removed from the pairs lists
if loop1 == 0 or loop1 == 4:
    for g in range(len(entry3)):
        if entry4[g].identifier not in duplicate_identifier_dictionary:
            count = 0
            if hash(entry3[g].identifier) in hydrates_not_pair:
                indices = [i for i,
                           x in enumerate(hydrates_no_t_pair) if x ==
                           hash(entry3[g].identifier)]
                for h in range(len(indices)):
                    if hash(entry4[g].id
                           entifier) == waterless_forms_no_t_pair[indices[h]]:
                        count += 1
            if count != 0:
                if loop1 == 0:
                    if count == 0:
                        if hash(entry3[g].identi
                               fier) not in dont_want_hydrates and
                           hash(entry4[g].identi
                               fier) not in dont_want_waterless_f
                               orms:
                            if loop1 == 0:

```

```

188 paired_hydrates.append(hash(entry3[g].identifier))
189 paired_waterless_forms.append(hash(entry4[g].identifier))
190
191 writer8.write(entry3[g])
192 writer9.write(entry4[g])
193
194 partner_hydrates.append(entry3[g])
195
196 partner_waterless_forms.append(entry4[g])
197
198
199
200
201
202
203
204
205
206
207
208
209

```

```

210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
    duplicate_identifier_
    dictionary.get(entry4
    [g].identifier) not
    in
    dont_want_waterless_f
    orms:
        if loop1 == 0:
            sample_writer.write(entry3[g])
            writer1.write(entry4[g])
            paired_hydrates.append(hash(entry3[g].identifier))
            paired_waterless_forms.append(duplicate_identifier_dictionary.get(entry4[g].identifier))
                if loop1 == 4:
                    writer8.write(entry3[g])
                    writer9.write(entry4[g])
                    if loop1 == 0:
                        if
                            hash(entry3[g].id
                            entifier) not
                            in
                            dont_want_hydrate
                            s and
                            duplicate_identif
                            ier_dictionary.ge
                            t(entry4[g].ident
                            ifier) in
                            dont_want_waterle
                            ss_forms:
                    partner_hydrates.append(entry3[g])
                    if
                        hash(entry3[g].id
                        entifier) in
                        dont_want_hydrate
                        s and
                        duplicate_identif
                        ier_dictionary.ge
                        t(entry4[g].ident
                        ifier) not in
                        dont_want_waterle
                        ss_forms:
            partner_waterless_forms.append(entry4[g])
            loop1 += 1
            if loop1 == 1:
                entry3 = entry_reader7
                entry4 = entry_reader8
            if loop1 == 2:
                entry3 = entry_reader9
                entry4 = entry_reader10
            if loop1 == 3:
                entry3 = entry_reader11
                entry4 = entry_reader12
            if loop1 == 4:
                entry3 = entry_reader13
                entry4 = entry_reader14
            print 'first phase of new list
            generation complete'
#Here hydrates and waterless forms that
#were part of a pair that did not share

```

```

240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261 if __name__ == '__main__':
262     # This runs the script
263     r = Runner()
264     r.run()
265
the same stoichiometry are checked for
a second appearance in the list of true
hydrate-anhydrite pairs
#If these structures do not appear in
the true hydrate-anhydrite pairs list,
then they are written to the true class
2 (hydrate) and class 3 (anhydrite)
lists for structures without a known
partner
partner_hydrates2 =
list(set(partner_hydrates))
partner_waterless_forms2 =
list(set(partner_waterless_forms))

for i in range(len(partner_hydrates2)):
    if
        hash(partner_hydrates2[i].identifier)
        not in paired_hydrates and
        hash(partner_hydrates2[i].identifier)
        not in unpaired_hydrates:
            writer2.write(partner_hydrates2[i]
                          )
for j in
range(len(partner_waterless_forms2)):
    if
        partner_waterless_forms2[j].identifie
        r not in
        duplicate_identifier_dictionary:
            if
                hash(partner_waterless_forms2[j].i
                    dentifier) not in
                paired_waterless_forms and
                hash(partner_waterless_forms2[j].i
                    dentifier) not in
                unpaired_waterless_forms:
                    writer3.write(partner_waterle
                                  ss_forms2[j])
            else:
                if
                    duplicate_identifier_dictionary.g
                    et(partner_waterless_forms2[j].id
                        entifier) not in
                    paired_waterless_forms and
                    duplicate_identifier_dictionary.g
                    et(partner_waterless_forms2[j].id
                        entifier) not in
                    unpaired_waterless_forms:
                        writer3.write(partner_waterle
                                      ss_forms2[j])
print 'second phase of new list
generation complete'

```

#27 List of hydrate-anhydrite pairs

Hydrates	Anhydrous Forms		
1. ACEQOQ	ACEQUW	49. AFEVOX	HISHEX
2. AVEPUN	AVEPEX	50. AFICAV	AFICEZ
3. BULMEA02	TFACET	51. AFUTUS	MODNUP
4. CYTOSM04	CYTSIN	52. AGEKEE	AGEKII
5. CYTOSM04	CYTSIN02	53. AGEQIO	TALVOU
6. DOGGUD	DOGHAK	54. AGOKUD	AGOKOX
7. ELEVOG	UNOGOT	55. AHEREK	JAYPUU
8. KAPZUY	JAWXUB	56. AHEREK	JAYPUU01
9. KAQBEL	WEMHOO	57. AHEREK	JAYPUU03
10. KEMBIO	XOHVEX	58. AHIMUA	AHIMOU
11. LUHHUS	LUHHOM	59. AHOXLH	MOYHAJ
12. NESKED	NESJUS	60. AHOXLH02	MOYHAJ
13. NUMVEW01	JEGMIR	61. AJISIW	AJISES
14. OXACDH27	OXALAC02	62. AJUPAW	VILRIS
15. OXACDH27	OXALAC04	63. AJUQEBC	RABTIX
16. OXACDH28	OXALAC02	64. AJUWIL	TESZIA
17. OXACDH28	OXALAC04	65. AKATEN	FIDSAN
18. PECYII	PECYEE	66. AKIHIN	AKIHAF
19. RAVBUL	NAFMEM	67. AKIJEK	NAFZEC
20. SEFGIW	NAFMEM	68. ALABEU	DIKSUM
21. TANCES	TANDET	69. ALUPAA	GAKQUF
22. TAWJEH	MUBHOG	70. AMBZCL	PEYNOZ
23. TUDDAZ	TUDBOL	71. AMCHCA	AMMCHC10
24. TUDDED	TUDBUR	72. AMEQAK	AMEPOX
25. VUZYIA	CELSOD	73. AMEVET	AMEVAP
26. WUVLOP	GEVYOY	74. AMPCIH01	AMCILL
27. XEHPIJ	ESTILE10	75. AMTETZ	EJIQEUE
28. BIUHYD01	ZZZLQC01	76. AMTETZ	EJIQEUE02
29. ABAZUB	ABEDUJ	77. ANDOON05	ANDOON
30. ABAZUB01	ABEDUJ	78. ANIPET	ANIPIX
31. ABEBIV	ABEDUJ	79. ANIPUJ	ANIQAQ
32. ABEBOB	ABEDUJ	80. ANORIG	TAFJUH
33. ABEBUH	ABEDUJ	81. ANSFON	DAPSUO03
34. ABIMUW	ABINAD	82. ANSFON	DAPSUO15
35. ABIPEJ	ABINUX	83. APICAE	TATBEX01
36. ABOQOZ	ABOQIT	84. AQAROZ	HEGLOV
37. ACCAAH	ACEDAC01	85. AQASAM	HEGMAI
38. ACCAAH	ACEDAC10	86. AQASEQ	HEGMEM
39. ACCAAH	ACEDAC11	87. AQASIU	HEGLUB
40. ACCAAH	ACEDAC15	88. AQOMAU	UNOGIN
41. ACCAAH	ACEDAC16	89. AQOMAU	UNOGIN02
42. ACCAAH	ACEDAC17	90. AQOMAU01	UNOGIN
43. ACCAAH	ACEDAC18	91. AQOMAU01	UNOGIN02
44. ACCAAH	ACEDAC19	92. AQOMEY	UNOGIN
45. ACCAAH	ACEDAC20	93. AQOMEY	UNOGIN02
46. ACCAAH	HEWCOE	94. AQOMEY01	UNOGIN
47. ACXMAC	CEMHEK	95. AQOMEY01	UNOGIN02
48. ADAMAU	NORJUA	96. ARGHCL10	LARGIN02
		97. ARGIND	TAQBIY
		98. ARIGAK	HUPWEV

99. ASAXUO	ASAYAV	151.	BARBAD	BARBAC
100. ASIMOG	ASIMIA	152.	BARBAD	BARBAC02
101. ASIMOG	ASIMIA01	153.	BARBAD02	BARBAC
102. ASPARM	VIKKEG	154.	BARBAD02	BARBAC02
103. ATOYIR	LOLYUG	155.	BCDEXD03	WEWTOJ
104. ATOYIR	LOLYUG01	156.	BCDEXD04	WEWTOJ
105. ATZDZM10	ATZTHD10	157.	BCDEXD05	WEWTOJ
106. AVEMUK	KUXPUP	158.	BCLMSN	WOYPAB
107. AVEMUK	KUXPUP01	159.	BECXEO	TEGHAQ
108. AVEQAU	AVENUL	160.	BEGVAL	CEXQEB
109. AVINEZ	AVINID	161.	BEKBOJ01	MEYGII01
110. AVOKIG	AVOJUR	162.	BENTAC	TANRIK
111. AVOKIG	AVOJUR02	163.	BEQWID	LEWBOE
112. AWUSEQ	AWUSAM	164.	BESLUG	BEGDIB01
113. AWUWAR	WAYBAB	165.	BEVNUM	DUYSUL
114. AXEGUH	AXEGOB	166.	BEVXUX	HOQLAB
115. AXEXIM	AXEXAE	167.	BEXMEX	WANJEB
116. AXEXIM	AXEXAE01	168.	BICMEF	ZASQOZ
117. AXEXOS	AXEXAE	169.	BICQUB	FADHEZ
118. AXEXOS	AXEXAE01	170.	BICQUB01	FADHEZ
119. AXOBEW	AXOBAS	171.	BIHNAH	BOGSEW
120. AXUGEG	AXUGIK	172.	BIHWAR	UFOLEJ
121. AXUJIM	QIMZOB	173.	BIJCOO	GAQQEV
122. AXUJIM	QIMZOB01	174.	BIJDON03	WUYNUA
123. AXUJUY	QIMZOB	175.	BIJDON04	WUYNUA
124. AXUJUY	QIMZOB01	176.	BIJDON05	WUYNUA
125. AYAXOP	AYAXUV	177.	BIJVIZ	BIJSAO
126. AYILAW	TEZVAW	178.	BIKPIV	GENQIB
127. AYIRUX	ZEDDOB	179.	BIKPIV	GENQIB01
128. AYOGIE	IHOQUS	180.	BIMYEB	BIMYAX
129. AYOGIE	IHOQUS01	181.	BIMYEB	BIMYAX01
130. AYUBEC	HICYEZ	182.	BIMYEB	BIMYAX02
131. AZTHPN	CIPWUT	183.	BINWEC	PARBAC
132. AZULAJ	SEKNOM	184.	BIRMEU	MEBQEQQ
133. BADGAR	QOQFAF	185.	BIRMEU	MEBQEQQ01
134. BAFDUI	BAFFAQ	186.	BIRMEU03	MEBQEQQ
135. BAFFIY	BAFFOE	187.	BIRMEU03	MEBQEQQ01
136. BAFHUL10	BAFJAT10	188.	BIUHYD	ZZZLQC01
137. BAFPAZ	CALDEY	189.	BIYSOR	BIYSIL
138. BAFXOV	BAFXIP	190.	BIZWOY	BIZWIS
139. BAGJEY10	AZTHYM10	191.	BNITRA10	VECGIV
140. BAHMAB	DOLBIR	192.	BOBMEM	BOCBAY
141. BAHMAB	DOLBIR07	193.	BOCNIQ	TICYOT
142. BAHMAB	DOLBIR08	194.	BOLDOX	LDOPAS03
143. BAHMAB	DOLBIR14	195.	BOLGAK	ACANLC10
144. BAHMAB	DOLBIR43	196.	BOLLIZ	RUYFEX01
145. BAHMAB	GLYCIN35	197.	BOPQAY	COFDUW10
146. BAHMAB	GLYCIN57	198.	BOPQAY01	COFDUW10
147. BAHMAB	GLYCIN67	199.	BOQNOM	BOQNIG
148. BAKJUS	POAZHN	200.	BOXGAY	ISINOP
149. BANAPQ10	BANAQP10	201.	BOYFOK03	ZULQAY
150. BAPJIL	BAPJEH	202.	BOYFOK04	ZULQAY

203.	BOZYUM	TIMZAQ	255.	CIMGUA	CIMETD04
204.	BOZYUM	TIMZAQ01	256.	CIRDIS	CIRDOY
205.	BRADOM	BIPADO	257.	CIRFEQ	CIPZIM
206.	BUDNEU	TOBHEZ	258.	CITARC	CITRAC10
207.	BUDNEU01	TOBHEZ	259.	CIWPIJ	CIWTUZ
208.	BUDPUM	FUHJOI	260.	CIWQAC	CIWPEF
209.	BUDTUP	BUDTOJ	261.	CIYPUX	KOMTUC
210.	BUDXII	LABZAP	262.	CIZGUN01	RABTIX
211.	BUFTAZ	GUFREG	263.	CIZQAD	MENMIB
212.	BUGJES01	JEYDEW	264.	CIZQAD	MENMIB01
213.	BUHSOM	DUPBOG	265.	CIZQAD	MENMIB02
214.	BULMAW10	TFACET	266.	CLANDH	CLANAC
215.	BUVSEQ01	WEWTOJ	267.	CLQUON01	HOJLOI
216.	CADVIM	EHIWEZ	268.	COGQEW	CUCVUS
217.	CADVUY	YUYSMOU	269.	COKCOV	COKCIP
218.	CAFINE	NIWFEE03	270.	COMPAT	COMPEA
219.	CAFINE	NIWFEE05	271.	COMPAT	COMPEA01
220.	CAKFAV	BUBTIB	272.	CONYIO	VETVOG01
221.	CAPWUL	HEYJOK	273.	COTVOX	COTVIR
222.	CARNID	YUYPAJ	274.	COTZIV	VETVOG01
223.	CARNID	YUYPAJ01	275.	COVMEF	BARWUM
224.	CASSEU	CUNRIO	276.	COVMEF	BARWUM01
225.	CAVREX	FEDLEG	277.	COVPIN	UHITOV
226.	CAYPIC	HUYWAB	278.	COWYAP	COWXUI
227.	CAZKUK	PATTOU	279.	COXPAH	YEPJOT
228.	CECQEJ	CECYIV	280.	COYFUS	COYFOM
229.	CECQIN	CECPIM	281.	COYGAZ	COYFOM
230.	CEDLIJ	YUTCEV	282.	COYGED	COYFOM
231.	CEDLIJ	YUTCEV01	283.	COYGIH	COYFOM
232.	CEDLIJ	YUTCEV02	284.	COZVIY	UWIJUH
233.	CEDLIJ	YUTCEV03	285.	CREATH	JOHJIB
234.	CEF MAD	ALOMUK	286.	CREATH	JOHJIB01
235.	CEHTAK10	MECWIC	287.	CREATH	JOHJIB02
236.	CEHTAK10	MECWIC03	288.	CUJKID	IJUVAL
237.	CEJGIK	CEJGAC	289.	CUJREG	CUHVUY
238.	CEJGIK	CEJGAC01	290.	CUWKAG	CUWJUZ
239.	CELXEY	ACACAK	291.	CUYWAU	CUYQUI
240.	CELYEZ	CELYAV	292.	CUYWAU	CUYQUI01
241.	CEMQUG	CEMQOA	293.	CYSTAC	CYSTE A
242.	CEMQUG	CEMQOA01	294.	CYSTIN10	CYSTBR
243.	CESVII	CESWIJ	295.	DAFNOO	RIZGEM
244.	CETZUY	PAKMAR	296.	DAFRAG	DAFXAM
245.	CEYQOO	CEYQII	297.	DAHXOZ	TPCYPO
246.	CIGDON	CIGDIH	298.	DAJWOA	DAJWUG
247.	CIJMEQ	BOXQUA	299.	DAMMIP	DAMNAI
248.	CIJYUQ	FEPJOB	300.	DAPRAQ	GAHBAV
249.	CIJYUQ	FEPJOB01	301.	DAZLEW	POPDOO
250.	CIKDOQ	MAJRIZ	302.	DEDBUJ	OCHTET
251.	CIKDOQ	MAJRIZ01	303.	DEDBUJ	OCHTET01
252.	CIKJAH	ZASQOZ	304.	DEDBUJ	OCHTET03
253.	CIMGUA	CIMETD	305.	DEDNIM	DEDMUX
254.	CIMGUA	CIMETD01	306.	DEDNIM	DEDMUX01

307.	DEDNIM	DEDMUX02	359.	DUSQOY	POFNOO
308.	DEHBUP	NDNHCL10	360.	DUTLEJ10	JAXGUK
309.	DEPDEH	DEPDAD	361.	DUZHIP	YERVAS
310.	DETHHC20	LIHJAQ	362.	EACLTH10	HIVROT
311.	DETHHC20	LIHJAQ01	363.	EACLTH10	HIVROT01
312.	DEVXEK	DALGON	364.	ECEFOJ	YEZLUN
313.	DEVXEK	DALGON01	365.	ECORAS	ECOQUL
314.	DEVXEK	DALGON02	366.	EDAJEA	MEQQEE
315.	DEVXEK	DALGON03	367.	EDIDUQ	QIMZOB
316.	DEYREF	DOPNED	368.	EDIDUQ	QIMZOB01
317.	DGLSLM10	QQQFGD02	369.	EDOKOX	RIZFEL
318.	DIFPEO	TAJPAW	370.	EDUHIU	EDUGOZ
319.	DIGXOJ	CASNEQ	371.	EFASOV	AFEDIY
320.	DILFAF	FIFGOQ	372.	EFIBUR	EFICAY
321.	DILFAF	FIFGOQ01	373.	EFIDUS01	UDAYUT
322.	DILNAP	DILNET	374.	EFIDUS01	UDAYUT01
323.	DIPFAJ01	JEYDEW	375.	EFIDUS01	UDAYUT02
324.	DIPFAJ10	JEYDEW	376.	EFIDUS02	UDAYUT
325.	DIPICA10	AFEBUI	377.	EFIDUS02	UDAYUT01
326.	DIPMUK	NILYAI	378.	EFIDUS02	UDAYUT02
327.	DITFUJ	DITFOD	379.	EFIDUS03	UDAYUT
328.	DIWNON	PIDJES	380.	EFIDUS03	UDAYUT01
329.	DLPROM01	QANRUT	381.	EFIDUS03	UDAYUT02
330.	DMAPYC	JUBKAS	382.	EFIFO01	KETXIR
331.	DOBLOV	CUVHAC01	383.	EFOBIL	EFOBEH
332.	DOBLOV	CUVHAC02	384.	EFOCOS	EFOCIM
333.	DOFTAT	DOFSUM	385.	EGOBAF	EGUHIZ
334.	DOFTAT	DOFSUM01	386.	EGUTUW	EDUGUF
335.	DOHXII	DOHXAA	387.	EJEVIA	RAKSIG
336.	DOLRUR	JEDXAR	388.	EJUPUU	DORBES
337.	DOPLIF	DOPLOL	389.	ELASUG	ELATAN
338.	DOPLIF	DOPLOL02	390.	ELASUG	ELATAN01
339.	DOPLIF	DOPLOL03	391.	ELOJUM	ELOJOG
340.	DOPLIF	DOPLOL05	392.	ELOJUM	ELOJOG01
341.	DOXVUH	BOWWIT	393.	ELOKIB	ELOKAT
342.	DOZYUO	DUBMAQ	394.	ELOKIB	ELOKEX
343.	DRPRDL	GIXXOB	395.	ENEBUW	UGIVAI
344.	DRPRDL	GIXXOB01	396.	ENEBUW	UGIVAI01
345.	DUBMEU	DUBMAQ	397.	ENECAD	UGIVAI
346.	DUCXEG	TICFUG	398.	ENECAD	UGIVAI01
347.	DUDGAL	BEWKUM	399.	ENODUH	UHITOV
348.	DUHKEW	BURLOP	400.	ENODUH01	UHITOV
349.	DUJPAB	BENPRL	401.	ENUJED	HOFQID
350.	DUJPAB	BENPRL01	402.	EPHEDH04	EPHEDR01
351.	DUJPAB	BENPRL02	403.	EPOCET	FEGCOM
352.	DUKTOS	CAPTAZ	404.	EQOTAG	FEPGUE
353.	DUKWIQ	MEBQEQQ	405.	EQOTUA	RUHTUL02
354.	DUKWIQ	MEBQEQQ01	406.	EQOVEM	RUHTOF02
355.	DULKAX	OMEJUN	407.	EQOVOW	RUHTIZ02
356.	DUNDAT	UGOVIX	408.	EQOWAJ	WOQRERB
357.	DUNDAT	UGOVIX01	409.	EQOWAJ	WOQRIF
358.	DUSHIJ	VUCXEX	410.	ERUHUV	VIFSAF

411.	ETAMCM	HIVROT	463.	FEFNOT02	CBMZPN01
412.	ETAMCM	HIVROT01	464.	FEFNOT02	CBMZPN03
413.	ETOHOM	AWUNAI	465.	FEFNOT02	CBMZPN11
414.	EVAGIT	EPACAB	466.	FEFNOT02	CBMZPN12
415.	EVAVIH	BBENAN	467.	FEFNOT02	CBMZPN16
416.	EVIXIQ	RIVKAJ	468.	FEFNOT02	CBMZPN20
417.	EVODOI	EVODUO	469.	FEFRIS	PUBMUU
418.	EVUQAO	BOBVIY	470.	FEFRIS	PUBMUU01
419.	EVUQAO	BOBVIY04	471.	FEFRIS	PUBMUU02
420.	EVUQAO	BOBVIY05	472.	FEFRIS	PUBMUU23
421.	EVUQAO	BOBVIY06	473.	FEFYEW	FEFYAS
422.	EVUQAO	BOBVIY07	474.	FEHMAJ	IZIVUK
423.	EVUQAO	BOBVIY08	475.	FEJBAA	JOFWIM
424.	EVUQAO	BOBVIY09	476.	FEJBAA	JOFWIM01
425.	EVUQAO	BOBVIY13	477.	FEJBAA	JOFWIM03
426.	EVUQAO	BOBVIY14	478.	FELCIJ	FELCOP
427.	EVUQAO	BOBVIY15	479.	FELVIE	TEDYOT
428.	EVUQAO	BOBVIY16	480.	FELWEB	TEDYOT
429.	EVUQAO01	BOBVIY	481.	FEPBAF	FENZUV
430.	EVUQAO01	BOBVIY04	482.	FEFOX	WOQBAF
431.	EVUQAO01	BOBVIY05	483.	FEFOX	WOQBAF02
432.	EVUQAO01	BOBVIY06	484.	FETXEI	IKASUK
433.	EVUQAO01	BOBVIY07	485.	FEWMID	OJOQOV01
434.	EVUQAO01	BOBVIY08	486.	FEZHAV	WICSUX
435.	EVUQAO01	BOBVIY09	487.	FIBTAM	FIBTEQ
436.	EVUQAO01	BOBVIY13	488.	FICZAU	BEJNEM
437.	EVUQAO01	BOBVIY14	489.	FIFFUV	ACPRET03
438.	EVUQAO01	BOBVIY15	490.	FIFWAV	FIFVOI
439.	EVUQAO01	BOBVIY16	491.	FIFWEZ	FIFVOI
440.	EWAMIZ	EWAMAR	492.	FILBEH	ROMMEM
441.	EWAMIZ	EWAMAR01	493.	FIMVOM	FESRIE
442.	EWAQEZ	EWAQAV	494.	FITYAI	SAGXOP
443.	EWAYAE	EWAYEI	495.	FITYAI	SAGXOP01
444.	EWIVOX	VOKWAT	496.	FITYAI	SAGXOP03
445.	EYELAX	ANISAC	497.	FIZVUF	DUMRAE
446.	EYELAX	ANISAC03	498.	FIZVUF	DUMRAF
447.	EYELAX	ANISAC04	499.	FOMPUT	FOMPON
448.	EYIKUU	FIFGOQ	500.	FONBUG	TUWYOB
449.	EYIKUU	FIFGOQ01	501.	FONHEW	BTCOAC
450.	EYOGIJ	YULJUM	502.	FOSLUW	ZOVQEI
451.	EYUQAS	CARPEQ	503.	FOSPUZ	UZIROM
452.	EZOBUQ	EZOCAX	504.	FOWMUA	HIMCAJ
453.	FADYIV	FADYER	505.	FOWTOB	HUCVEH
454.	FAFPEJ	YIXRON	506.	FOYTOC	YEHROS
455.	FAJSOY	FAJSIS	507.	FURGAA01	JICLAI
456.	FAMMOY	FAMPOB	508.	FURGAA04	JICLAI
457.	FAYVUW	FAYVOQ	509.	FUSWIB	EFUMAU
458.	FAYWUX	HEPOPH01	510.	FUSWIB	EFUMAU03
459.	FAYWUX	HEPOPH02	511.	FUSWIB	EFUMAU06
460.	FEDBUL	FEDBOF	512.	FUSWOH	EFUMAU
461.	FEDNEI	FEDNIM	513.	FUSWOH	EFUMAU03
462.	FEFBEX	NAFZEC	514.	FUSWOH	EFUMAU06

515.	GADYUI	GADXUH	567.	GOBMAN	YUPHUO
516.	GAFVIS	VUXBAR	568.	GOFWOP	QQQAUJ03
517.	GAFVOY	VUXBAR	569.	GOFWOP	QQQAUJ04
518.	GAHJIL	GAHJEH	570.	GOFWOP	QQQAUJ07
519.	GAJMIN	GAJMOT	571.	GOFWUV	QQQAUJ03
520.	GANCUV	GANDAC	572.	GOFWUV	QQQAUJ04
521.	GAQYON	ALETUG	573.	GOFWUV	QQQAUJ07
522.	GAQYON	ALETUG01	574.	GORFIC	ROBXEN
523.	GASNIY	GASNEU	575.	GOWCAW	GOWBOJ
524.	GAVWUW	GAWSED	576.	GUBLAS	RIFBAL
525.	GAXTII	USIWEZ	577.	GUBLAS	RIFBAL01
526.	GEDYUM	IJAQIV	578.	GUCNOJ	GUCNUP
527.	GEDYUM	IJAQIV01	579.	GUHLEC	FOWFOM
528.	GEDYUM	IJAQIV02	580.	GULTEM	DECGPU10
529.	GEHRUH	DAMGEF	581.	GULTOW	LAKVOI
530.	GEHXAV	GEHWUO	582.	GUMJIJ	NAGVAS
531.	GEJQAO	URICAC	583.	GUSHAD01	SLFNMB01
532.	GEMBEF	POSVEA	584.	GUSHAD01	SLFNMB02
533.	GEMLAN	GEMKUG	585.	GUSHAD01	SLFNMB05
534.	GEQXUX	GEQXOR	586.	GUSHAD01	SLFNMB06
535.	GEQXUX	GEQXOR01	587.	GUXZEG	RABTIX
536.	GERWUX	GERXEI	588.	GUYGUD01	GUYGOX03
537.	GERXAE	GERXEI	589.	HABJUP	PAQUCL
538.	GESCIS	GESBEN	590.	HABJUP01	PAQUCL
539.	GESCOY	VIDMAX	591.	HAFNUX	HAFNOR
540.	GESCOY	VIDMAX02	592.	HAJWAR	HAJWIZ
541.	GESHET	GESGAO	593.	HAKDUT	HAKDON
542.	GESKUK	GESKOE	594.	HATKET	ROMMEM
543.	GETLUN	TIZVEE	595.	HATQAX	MEQQOO
544.	GETLUN	TIZVEE01	596.	HAXBUD	AZAXEG
545.	GEVBIT	CASQUJ	597.	HAXBUD	LABJON01
546.	GEVNED	GEVNON	598.	HAXBUD01	AZAXEG
547.	GEVNIH	GEVNON	599.	HAXBUD01	LABJON01
548.	GEXXAI	HOJLOI	600.	HAXQAA	BADVIL10
549.	GEYVIS	GEYTIQ	601.	HBARBT	ALXANM01
550.	GICXOF	EQAHEL	602.	HEBZEW	HEBZAS
551.	GIHJOX	DAZXOS	603.	HEBZEW	HEBZAS01
552.	GIMSIF	GIMSOL	604.	HECMUX	VUNUY10
553.	GITVAI	PUJZEA	605.	HECMUX	VUNUY11
554.	GIVCEU	GIVCAQ	606.	HECNAE	VUNZEJ10
555.	GIVYAL	GIVXUE	607.	HEPDUE	NOTEST01
556.	GIXDAS	ETDIAM01	608.	HEPNAR	ZUTMAC
557.	GIXDAS	ETDIAM15	609.	HESWIM	HEQMUM
558.	GIXDAS	ETDIAM18	610.	HESYAF	HESXUY
559.	GIXDEW	HEXMDA	611.	HETHAQ	HETGUJ
560.	GIXDIA	QATVUC	612.	HEZWAK	WAGLUM
561.	GIXDOG	ROKZOG	613.	HIFVEY	TUFGEG
562.	GIXDUM	QATWAJ	614.	HIKGOZ	BILCUW
563.	GIZMAE	GIXXOB	615.	HIKGOZ	BILCUW01
564.	GIZMAE	GIXXOB01	616.	HIKNIY	BEURID10
565.	GLUCMH	GLUCSA	617.	HIKNIY	LYFURA
566.	GLUCMH	GLUCSA20	618.	HILXEG	BORSEG

619.	HILXEG	BORSEG01	671.	IFILIT	IFILAL
620.	HIMDOY	HIMDAK	672.	IFILIT	IFILAL01
621.	HINFAN	VAGYEI	673.	IFILIT	IFILAL02
622.	HIPKAS	ALITOL	674.	IFIZIG	VENLUW
623.	HIPKAS	GLUCIT	675.	IFIZIG04	VENLUW
624.	HIPKAS	GLUCIT02	676.	IFURIL	YOMDUA
625.	HIPKAS	GLUCIT03	677.	IGALIN	VACHAK
626.	HIRLUQ	PEXXEW	678.	IGATUH	HACVAL
627.	HIRZUF	UNOGIN	679.	IHAQE0	IHAPOX
628.	HIRZUF	UNOGIN02	680.	IHAQE0	IHAPOX01
629.	HIYJII	HIYJEE	681.	IHAQE0	IHAPOX02
630.	HODHUE	MIYDII	682.	IHATAN	XIPYIF
631.	HODRUO	JELSEY	683.	IHATAN	XIPYIF01
632.	HOFNAR	RUZWUE	684.	IHEGOT	IHEGIN
633.	HOGCOX	HOGCIR	685.	IHOTUW	COKBIN
634.	HOMPRO10	QAMVAB	686.	IHUQEJ	IHURUA
635.	HONXUD	REVZAT	687.	IJEQET	LASPRT
636.	HOQHAW	XERBEB	688.	IJESIZ	IJESAR
637.	HOSHEC	NANRAV	689.	IJESIZ01	IJESAR
638.	HOSHIG	NANRAV	690.	IJILET	XOQGIV
639.	HOTYIZ	GIDRAN	691.	IJISEB	IJISAX
640.	HOTYUL	GIDRAN	692.	IJUKAA	SUZTAJ
641.	HOWPEP	WAMFAS	693.	IKALAJ	IKAKOW
642.	HOWPEP	WAMFAS01	694.	IKEDIM	SAGYAC
643.	HOWPEP	WAMFAS02	695.	IKEXOM	VOLZEC
644.	HOWPEP	WAMFAS04	696.	IKEXOM	VOLZEC01
645.	HOXGEG	HOXGIK	697.	IKUROV	IKURUB
646.	HUCLIA	SAGMEU	698.	ILEJUG	PUBMUU
647.	HUCSUV	CURFUR	699.	ILEJUG	PUBMUU01
648.	HUFLUQ	XAHHOG	700.	ILEJUG	PUBMUU02
649.	HUHWEM	TIJLUT	701.	ILEJUG	PUBMUU23
650.	HULKUU	AVURAK	702.	ILINEW	GOLMIF
651.	HUMJEE	COTZAN06	703.	IMEGIR	NELPUP
652.	HUMJEE	HXACAN	704.	IMOPIJ	IMOPEF
653.	HUMJEE	HXACAN29	705.	INEFEM	ROMMEM
654.	HUMJEE	HXACAN39	706.	INOSND01	INOSIN10
655.	HURMUC	HURMOW	707.	INOSND01	INOSIN11
656.	HUSVEX	HUSVIB	708.	IPEHUH	ITESIK
657.	HUSVEX	HUSVIB01	709.	IQUICM	IQUINC
658.	HUVGUC	HUVHAJ	710.	IRELEX	YEGGIA
659.	HXBIIUR10	HBIURT10	711.	IRERUU	IRETAC
660.	HXMTHH	HXMTHAM	712.	IRIZAL	FARRUM
661.	IBALOP	UBACOS	713.	IRIZAL	FARRUM03
662.	IBOVED	IBOHIT	714.	ISEKEW	RUYFEX01
663.	IBUJOH	EWAYOS	715.	ISUSUK	UQUDOB
664.	IBUXIN	UJEMEB	716.	IVEGIA	UZOVAH03
665.	ICEGAZ	SITCEE	717.	IVOSAP	IVOQUH
666.	ICOSUR	AJISES01	718.	IVUQIZ	IVUQOF
667.	IDIWOI	SAZBIF	719.	IVUQIZ	IVUQOF01
668.	IDUCUF	IDUDAM	720.	IVUQIZ	IVUQOF02
669.	IDUYEM	MUYSUW	721.	IVUQIZ	IVUQOF03
670.	IFAXOE	FAKDUS	722.	IVUQIZ	IVUQOF04

723.	IVUSAU	IVURUN	775.	KADBAU	BIPTID01
724.	IWUDUZ	IFOXIM	776.	KAKHEL	KAKHIP
725.	IXISAL02	IXIRUE02	777.	KALZIG	SUNKAO01
726.	IYUXOQ	LUYQAX	778.	KAMSEZ	KAMROI
727.	IZAPUU	YASTUI	779.	KAMSEZ	KAMROI01
728.	JAGXOE	HIGTEZ	780.	KAMSEZ01	KAMROI
729.	JARGUH	SENJEC	781.	KAMSEZ01	KAMROI01
730.	JASSAY	VEVJOW	782.	KAPPUP	KAPPOJ
731.	JATDUD	JATFAL	783.	KARGUI	KARFIV
732.	JATJAP	JATHUH	784.	KASHAP	KASGIW
733.	JAYBOD	MOKBAR	785.	KATSEG	KATSAC
734.	JAYBOD	MOKBAR01	786.	KAVFUL	KAVGAS
735.	JEBJIM	JEBJOS	787.	KEFHUA	RAKSAZ
736.	JEBREO	CAQZOL	788.	KEKHOY	KEKHIS
737.	JEDKIN	GASNEU	789.	KEQYOX	KEQYIR
738.	JEDKIN01	GASNEU	790.	KEXZAQ	KEXZEU
739.	JEDKIN02	GASNEU	791.	KEXZIY	KEXZEU
740.	JEDTOB	TPEPHO	792.	KEXZUK	KEXZEU
741.	JEDTOB	TPEPHO01	793.	KEYBAT	KEXZEU
742.	JEDTOB	TPEPHO07	794.	KEYBOH	KEXZEU
743.	JEDTOB	TPEPHO12	795.	KEYTOX	KEYSOW
744.	JEDTOB01	TPEPHO	796.	KICCOO	WEMWEQ
745.	JEDTOB01	TPEPHO01	797.	KIFQEW	CEKHAC
746.	JEDTOB01	TPEPHO07	798.	KIGMET	KIGMAP
747.	JEDTOB01	TPEPHO12	799.	KIJFOB	ALKINA
748.	JEDTOB02	TPEPHO	800.	KIJFOB	ALKINA01
749.	JEDTOB02	TPEPHO01	801.	KIKCIR	OTOLUC
750.	JEDTOB02	TPEPHO07	802.	KIMDER	KIMDAN
751.	JEDTOB02	TPEPHO12	803.	KIPPUU	TEKBUI
752.	JEHRAP	KAHLEK	804.	KITMIL	AMUQOQ
753.	JEHRAP	KAHLEK01	805.	KITMIL	AMUQOQ01
754.	JEJXAX	YIRZOP	806.	KOFCAL	GAPMAN
755.	JEKGIP	JEKGEL	807.	KOGPIG	UXONAY
756.	JEMROL	KULBIE	808.	KOJZUE	LICKUG
757.	JEPJOD	BOBZAC	809.	KOMNIL	XINVIA
758.	JEPJOD	BOBZAC02	810.	KONTIQ	IJUMEG
759.	JERRUU	JERROO	811.	KONTIQ	IJUMEG04
760.	JETJEY	VEYSEY	812.	KONTIQ	IJUMEG06
761.	JEVQUX	TOYSUX	813.	KONTIQ01	IJUMEG
762.	JEXZER	YADSOL	814.	KONTIQ01	IJUMEG04
763.	JEYZIX	BOMDUC	815.	KONTIQ01	IJUMEG06
764.	JIDPUH	WETFAD	816.	KONTIQ03	IJUMEG
765.	JIDPUH	WETFAD01	817.	KONTIQ03	IJUMEG04
766.	JIWPIP	JIWPTEL	818.	KONTIQ03	IJUMEG06
767.	JIXCOI	VIDMAX	819.	KONTIQ04	IJUMEG
768.	JIXCOI	VIDMAX02	820.	KONTIQ04	IJUMEG04
769.	JIYWET	JOWWIB	821.	KONTIQ04	IJUMEG06
770.	JOGMUP	LOFFET	822.	KONTIQ06	IJUMEG
771.	JOGMUP01	LOFFET	823.	KONTIQ06	IJUMEG04
772.	JONNIL	JONNEH	824.	KONTIQ06	IJUMEG06
773.	JULBOH	HEZZER	825.	KOPDOK	KOPDAW
774.	JUXYIM	HEQMUM	826.	KOPDOK	KOPDAW01

827.	KOPDOK	KOPDAW02	879.	LEXPUA01	PUHCEB
828.	KOPDOK	KOPDAW03	880.	LEYKAE	PEWYUO
829.	KOPZEW	KOPZAS	881.	LEYKAE	PEWYUO01
830.	KUCSIM	KUCSAE	882.	LEYKAE	PEWYUO02
831.	KUCSIM	KUCSAE01	883.	LEYZOG	IFUPAB
832.	KUGKAZ	EVESIJ04	884.	LEZJEF	LEZJAB
833.	KUQTAS	PENDAM	885.	LEZJEF	LEZJAB01
834.	KUSYON	SEBBEG	886.	LEZKAC	FIQFER
835.	KUTSAU	DODAMB	887.	LIBPUK	LIBPOE
836.	KUVBEI	ELLAGC	888.	LIBWUQ	LUYWOR
837.	KUVYUW	AJEYAQ	889.	LIDNOE	LIDNUK
838.	KUVYUW	AJEYAQ01	890.	LIFNOE	BISMEV
839.	KUVYUW	AJEYAQ02	891.	LIFNOE	BISMEV01
840.	KUVYUW	AJEYAQ03	892.	LIFNOE	BISMEV03
841.	KUVYUW	AJEYAQ04	893.	LIFNOE	BISMEV04
842.	KUZKOG	NODTIJ	894.	LIFNOE	BISMEV07
843.	LACFON	LACGII	895.	LIFNOE	BISMEV14
844.	LACFON01	LACGII	896.	LIGXAC	RUKHAG
845.	LACTOS01	EYOCUQ	897.	LIGXAC	RUKHAG01
846.	LACTOS01	EYOCUQ01	898.	LIKUNK	LIKLAU
847.	LACTOS03	EYOCUQ	899.	LIPZOA	NINSOS
848.	LACTOS03	EYOCUQ01	900.	LIZWOI	CIRNEY
849.	LAFDIH	FIFVOI	901.	LOBSOL	COTXUG
850.	LAHQAO	SIYZIK	902.	LOBSOL01	COTXUG
851.	LAJGUA	PIWXKEY	903.	LOCCAI	YIVGIV
852.	LAKTUM	GULTIQ	904.	LOCCAI	YIVGIV01
853.	LAPVEE	LAPVAA	905.	LOCCAI	YIVGIV02
854.	LAPWOQ	MUNCOO	906.	LOCNEX	FILGEM
855.	LAQSON	NAVSUY	907.	LOCNEX	FILGEM01
856.	LAQSON	NAVSUY02	908.	LOCNEX	FILGEM03
857.	LAQSON01	NAVSUY	909.	LOFKOG	REWRUH
858.	LAQSON01	NAVSUY02	910.	LOFTEH	LOFTIL
859.	LARWIN	ZORTAB	911.	LOJJIE	ZEXLUJ
860.	LASYOV	LASYUB	912.	LOJJIE	ZEXLUJ02
861.	LASZEM	LASZAI	913.	LOLDEW	LOLDAS
862.	LATBOZ	DAVPEW	914.	LONGUR	NAVGAT
863.	LATPIH	UQUDOB	915.	LONGUR	NAVGAT01
864.	LAXDIA	WACGEO	916.	LONKIK	PUBHOK
865.	LAYPEI	LAYPAE	917.	LOPTOB	KUJDEA
866.	LCYSSCC	CYSTCL	918.	LORTUI	ZIFQEM
867.	LCYSSCC	CYSTCL01	919.	LOXSUC10	YIKNUD
868.	LEBJUX	VATSAK	920.	LOZREZ	PUBTIQ
869.	LEBJUX	VATSAK01	921.	LOZROJ	LOZRID
870.	LEBKEI	VATSAK	922.	LSERMH10	LSERIN01
871.	LEBKEI	VATSAK01	923.	LSERMH10	LSERIN16
872.	LEBLEK	LEBLAG	924.	LSERMH10	LSERIN21
873.	LEFSUM	BODVAS	925.	LSERMH10	LSERIN50
874.	LEGZII	QAMSOO	926.	LSERMH16	LSERIN01
875.	LEJLIX	TEGFIW	927.	LSERMH16	LSERIN16
876.	LEJYOP	GOKNOK	928.	LSERMH16	LSERIN21
877.	LETGAV	LETFUO	929.	LSERMH16	LSERIN50
878.	LEXPUA	PUHCEB	930.	LUDXOX	LUDYEO

931.	LUJGEE	LUJFUT	983.	MICFEM	MICFAI
932.	LUKTER	LUKTAN	984.	MIFVEC	RIBWAB
933.	LUKXIA	BZPPOS	985.	MIQXUH	JACYAQ
934.	LUPAND	LPNECL	986.	MISMOR	MISMIL
935.	LYSCLH	DLLYSC10	987.	MISRIQ	MISREM
936.	MAFSUI	ALITOL	988.	MITCDH01	MITOMC
937.	MAFSUI	DMANTL	989.	MIVNIO	TAQGEY
938.	MAFSUI	DMANTL01	990.	MIWJAE	YIZCOZ
939.	MAFSUI	DMANTL10	991.	MIXTOD	HOFGAK
940.	MAFSUI	GALACT	992.	MIYJUZ	CUBBIL
941.	MAHDOP	LACJEE	993.	MIZHAF	RABDAV
942.	MAHKOX	GEPNOG	994.	MODMAV	MODLUO
943.	MAHKOX	OHIVAE	995.	MOHCUI	MOPYEV
944.	MAJQIY	AHIHAC	996.	MOJZUI	TOZKEA
945.	MAQYOS	XIJMOS	997.	MOKRUA	XEBYUA
946.	MARLUN	MARLOH	998.	MOKXIU	CIRDOY
947.	MAVPUX	MAWDUM	999.	MOPSER	LERKUQ
948.	MAWMEF	MAWMIJ	1000.	MOPZOH	GAVVEF
949.	MAWMEF	MAWMIJ01	1001.	MORPHC	EFASAH
950.	MAWMEF	MAWMIJ03	1002.	MORPHM	MORPIN01
951.	MCYTIM10	RADKOX01	1003.	MOSPAM	VAXHEI
952.	MEBQUG	MEBQEQQ	1004.	MOSWIB	FOSRAI
953.	MEBQUG	MEBQEQQ01	1005.	MOVDTIB	EHOWIH
954.	MECLOV	LIFTIF	1006.	MOVDTIB	EHOWIH02
955.	MECXOJ	MECXID	1007.	MOVDTIB	EHOWIH03
956.	MECYUQ	FOYSUI	1008.	MOVDTIB	EHOWIH04
957.	MEDMAL	MEDLUE	1009.	MOVDTIB	EHOWIH05
958.	MEDMAL	MEDLUE01	1010.	MOVDTIB02	EHOWIH
959.	MEKPAW	MEKNOI	1011.	MOVDTIB02	EHOWIH02
960.	MELFUF	MELFIT	1012.	MOVDTIB02	EHOWIH03
961.	MELFUF	MELFIT01	1013.	MOVDTIB02	EHOWIH04
962.	MELFUF	MELFIT02	1014.	MOVDTIB02	EHOWIH05
963.	MELFUF	MELFIT04	1015.	MOZWOO	MOSXOI
964.	MELFUF	MELFIT05	1016.	MPIPCA	MPIPAN
965.	MELFUF	MELFIT06	1017.	MSULIM	POMDAW
966.	MELFUF	MELFIT07	1018.	MSULIM	POMDAW01
967.	MELFUF	MELFIT09	1019.	MUBPOP	NEQVOV
968.	MELFUF	MELFIT19	1020.	MUBPOP	NEQVOV01
969.	MELFUF01	MELFIT	1021.	MUDDUK	DLSERN
970.	MELFUF01	MELFIT01	1022.	MUDDUK	DLSERN19
971.	MELFUF01	MELFIT02	1023.	MUFNEG	IDIMAJ
972.	MELFUF01	MELFIT04	1024.	MUKYOI	BOBHOP
973.	MELFUF01	MELFIT05	1025.	MUSCOU	DUHJIB
974.	MELFUF01	MELFIT06	1026.	MUSCUA	DUHJIB
975.	MELFUF01	MELFIT07	1027.	MUWLIA	VULGUF
976.	MELFUF01	MELFIT09	1028.	MUYTOQ	HTRYPT10
977.	MELFUF01	MELFIT19	1029.	MYTOLD	EFURIH
978.	MENHIX	BELPAL	1030.	MYTOLD	EFURIH02
979.	MEPYRZ	MPYRAZ	1031.	MYTOLD	EPINOS
980.	MEXJOO	MEXJII	1032.	MYTOLD	FOPKOK
981.	MEYRAJ	MEYQUC	1033.	MYTOLD	IFAKAC
982.	MHEXDO10	DMHXDL	1034.	MYTOLD	MUINOS

1035.	MYTOLD	MYINOL	1087.	NURAMH01	NIMFOE01
1036.	MYTOLD	MYINOL01	1088.	NURAMH01	NIMFOE02
1037.	MYTOLD	MYINOL03	1089.	NURJEP	NURHOX
1038.	MYTOLD	QIKZAN	1090.	NURJEP01	NURHOX
1039.	MYTOLD	QIKZAN01	1091.	NUTWII	REZRAP
1040.	MYTOLD	QIKZAN02	1092.	NUYCUG	NUYCOA
1041.	MYTOLD	QIKZUH	1093.	OBEQAN	HEBFUR
1042.	MYTOLD	YEPNOW	1094.	OBEQAN	HEBFUR01
1043.	MYTOLD	ZIVVIL	1095.	OBEQAN01	HEBFUR
1044.	NABWET	NABWAP	1096.	OBEQAN01	HEBFUR01
1045.	NAJYIH	NEFFUA	1097.	OBILOA	ABEWAG
1046.	NATPEE	NEDXAX	1098.	OCADEE	YAWZAZ
1047.	NATPEE01	NEDXAX	1099.	OCANIQ	VURWAH
1048.	NAZDAV	NAZCUO	1100.	OCANIQ	VURWAH03
1049.	NBARBT	NBARBA	1101.	OCAZEX	MOKYER
1050.	NEBFIM	NEBFOS	1102.	OCAZEX	MOKYER01
1051.	NEFCUY	NEFCEI	1103.	OCEGIO	IGARAK
1052.	NEFTEY	MUJKEH	1104.	OCESEW	CAXMOD
1053.	NESVOZ	AFAZEM	1105.	OCESEW	CAXMOD01
1054.	NESVOZ	AFAZEM01	1106.	ODOBAK	KETXIR
1055.	NEXSIU	THIOUR	1107.	ODUSAJ	NOWFEM
1056.	NEXSIU	THIOUR19	1108.	OFEJOA	OFEJUG
1057.	NEXSUG	THIOUR	1109.	OFEJOA	OFEJUG01
1058.	NEXSUG	THIOUR19	1110.	OFIREC	OFIQUR
1059.	NEXXAP	SUYBAR	1111.	OFIREC	OFIQUR01
1060.	NIDPEY	NIDJIW	1112.	OHAVUQ	DUXSUK
1061.	NIDPEY01	NIDJIW	1113.	OHAVUQ	DUXSUK04
1062.	NIFHAL	NIFGUE	1114.	OJOGEA	YINFEI
1063.	NIGZEJ	POVFIP	1115.	OKATIF	LAKJUC
1064.	NIHPIF	KEMHAL	1116.	OKEJIZ	OKEJEV
1065.	NIHPIF	KEMHAL02	1117.	OKEMAT	WUYPOW
1066.	NIMRUX	FLUBIP01	1118.	OKEMAT	WUYPOW01
1067.	NIMRUX	FLUBIP02	1119.	OKESAY	LEBKUZ
1068.	NIMRUX	YACZIO	1120.	OREHOK	ITIZOA
1069.	NINSIM	NINSEI	1121.	OROXAV	BUCRIA
1070.	NISHUS	TOAZOC	1122.	OSEXIU	OSESEQ
1071.	NOCNOK	NOCPAY	1123.	OSIGEE	GADBUL
1072.	NOGKID	NOGKEZ	1124.	OSIVAO	OSITUG
1073.	NOJWAK	NOJVUD	1125.	OVAYUH	OVAYOB
1074.	NOMDUP	NITPOL	1126.	OVOMOC	OVOMIW
1075.	NOMDUP	NITPOL01	1127.	OXACDH09	OXALAC02
1076.	NOQZID	EVAZAD	1128.	OXACDH09	OXALAC04
1077.	NOTKEO	NUVFEQ	1129.	OXAYUJ	QUKHIP
1078.	NOVTID	NOVTEZ	1130.	OXAYUJ	QUKHIP01
1079.	NUBFAR	NUBDUJ	1131.	OXAYUJ	QUKHIP02
1080.	NUJLOU	NUJLIO	1132.	OXUHEW02	OXUHAS02
1081.	NUMFUY	CIZQUX	1133.	OXUHIA02	PEXXEW
1082.	NUMFUY	CIZQUX01	1134.	OYETET	TIJPAF
1083.	NURAMH	NIMFOE	1135.	OZOWOR	URIRET
1084.	NURAMH	NIMFOE01	1136.	PABPAL	PABNAJ
1085.	NURAMH	NIMFOE02	1137.	PACHEH	LAKFEJ
1086.	NURAMH01	NIMFOE	1138.	PAFBEG	ABIPIN

1139.	PAFBEG	ETEXIK	1191.	PEYSOD	CBMZPN01
1140.	PAFBEG	ETEXIK02	1192.	PEYSOD	CBMZPN03
1141.	PAFBEG	ETEXIK06	1193.	PEYSOD	CBMZPN11
1142.	PAFBEG	KAHVUL	1194.	PEYSOD	CBMZPN12
1143.	PAGYUS	ALXANM01	1195.	PEYSOD	CBMZPN16
1144.	PAHNAP	WENROX	1196.	PEYSOD	CBMZPN20
1145.	PAJDQU	HICRUI	1197.	PEZBEB	WUVKON
1146.	PAKOJM	PAYKOJ	1198.	PEZVOF	SODVUC
1147.	PANKEV	VAGVAA	1199.	PHBARM	PHBARB
1148.	PANYLB	ANPYAB	1200.	PHBARM	PHBARB05
1149.	PAPNIG	PAPNUS	1201.	PHBARM	PHBARB10
1150.	PAPNIG01	PAPNUS	1202.	PHBARM	PHBARB11
1151.	PASXAK	PASWUD	1203.	PHBARM	PHBARB12
1152.	PASXIS	PASWUD	1204.	PHBZAC	JOZZIH
1153.	PASXIS02	PASWUD	1205.	PHBZAC	JOZZIH01
1154.	PASXIS03	PASWUD	1206.	PHGLOH	PHGLOL
1155.	PAXQAH	GAVKOG	1207.	PHOLCL	CUZDIK
1156.	PAYWOD	VANXIS	1208.	PIGFET	PIGFAP
1157.	PAZGII	PAZGUU	1209.	PIKKAV	VAWDEC
1158.	PECXIE	PECXEA	1210.	PILBUJ	PILBOD
1159.	PEDNAQ	WASPEL	1211.	PILSIO	KIXWEV
1160.	PEFGOX	PEFGIR	1212.	PINOLH01	PINCOL
1161.	PEFGOX	PEFGIR01	1213.	PINOLH01	PINCOL01
1162.	PEFGOX	PEFGIR03	1214.	PIPACA	PIPACB
1163.	PEFGOX01	PEFGIR	1215.	PIPERH	ITIZOA
1164.	PEFGOX01	PEFGIR01	1216.	PIPTUF	PIPTOZ
1165.	PEFGOX01	PEFGIR03	1217.	PIRRAL	PIRREP
1166.	PEGWUU	PEGWII	1218.	POBRON	WEWTOJ
1167.	PEGWUU	PEGWII01	1219.	PODPEE	RACCEE
1168.	PEJREC	PEJRIG	1220.	POSTAS	POSTEW
1169.	PEKJUM	PEKJEW	1221.	POSTAS	POSTEW01
1170.	PEKJUM	PEKJEW01	1222.	POTPET	GLYGLY
1171.	PEKJUM01	PEKJEW	1223.	POTPET	GLYGLY01
1172.	PEKJUM01	PEKJEW01	1224.	POVSEZ	POVRUO
1173.	PEKQUT	PEKQON	1225.	PROAMH10	QOBHEW
1174.	PEKRIG	PEKREC	1226.	PROAMH11	QOBHEW
1175.	PELYOW	PELYUC	1227.	PUBMII	PUBMUU
1176.	PEMCOA	UDOMUW	1228.	PUBMII	PUBMUU01
1177.	PEMPIH	VEQHAB	1229.	PUBMII	PUBMUU02
1178.	PEMTEI	PEMTAE	1230.	PUBMII	PUBMUU23
1179.	PEPCIW	XODYOE	1231.	PUBMII01	PUBMUU12
1180.	PERCAR	POFYEP	1232.	PUDYUI	PUDYOC
1181.	PERKEF	PERKIJ	1233.	PUDYUI	PUDYOC01
1182.	PERKOP	PERKUV	1234.	PUFGUU	PUFJUX
1183.	PESHAX	QOQNEP	1235.	PUFJIL	PUFJUX
1184.	PETRIS	AFEBUI	1236.	PUHPOX	BOBHOP
1185.	PEWYOI	PEWYUO	1237.	PUKFAE	PUKBEE
1186.	PEWYOI	PEWYUO01	1238.	PUVMAU	SAQJEZ
1187.	PEWYOI	PEWYUO02	1239.	PUYSEI	EDEXOA
1188.	PEWYOI01	PEWYUO	1240.	PUZGAT	VETVOG01
1189.	PEWYOI01	PEWYUO01	1241.	PUZGAT01	VETVOG01
1190.	PEWYOI01	PEWYUO02	1242.	PYMDSD	PYMSUL10

1243.	PYMDSD	PYMSUL11	1295.	QOBJOH	MAWMIJ01
1244.	PYMELL10	KEFGUA	1296.	QOBJOH	MAWMIJ03
1245.	PYMELL12	KEFGUA	1297.	QOHYAO	TETDAM01
1246.	PYRTHA10	PYRDNA01	1298.	QOHYAO	TETDAM02
1247.	PYRTHA10	PYRDNA02	1299.	QOHYAO01	TETDAM01
1248.	PYZDCX	IYAWAG	1300.	QOHYAO01	TETDAM02
1249.	QABCII	QABCEE	1301.	QOHYES	TETDAM01
1250.	QAMBIR	IVEVIO	1302.	QOHYES	TETDAM02
1251.	QAMMUP	QAMNIE	1303.	QOQXAV	FAFWIS
1252.	QANBOW	UJAWUZ	1304.	QOQXAV	FAFWIS01
1253.	QANVAF	IJUTAJ	1305.	QOQXAV	FAFWIS02
1254.	QANVAF	IJUTAJ01	1306.	QOSQIY	GEMJIU
1255.	QANVAF	IJUTAJ02	1307.	QOZKOG	KIXFIH
1256.	QARVUB	QARVOV	1308.	QOZSAZ	LAGJOU
1257.	QATJUR	RUGPUF	1309.	QUBREL	PAKNOG
1258.	QATNEF	QUWYIR	1310.	QUFCUQ	QUFCOK
1259.	QEHSER	IKOBAN	1311.	QUINCX	QUINCB10
1260.	QEHTUT	RUDYIZ	1312.	QUKHUB	QUKHIP
1261.	QEHWOO	XAQXAQ	1313.	QUKHUB	QUKHIP01
1262.	QEKDUI	SAJFIT	1314.	QUKHUB	QUKHIP02
1263.	QEMGEV	XEBTEF	1315.	QUNHAJ	QUNGUC
1264.	QEMLOK	QEKWEJ	1316.	QUNHAJ	QUNGUC01
1265.	QERHOL	KINKOH	1317.	QURTIH	HACTPH10
1266.	QERHOL	QIFPEC	1318.	QURTIH	HACTPH12
1267.	QEVPPE	VATSAK	1319.	QURYEI	QURYAE
1268.	QEVPPE	VATSAK01	1320.	QUWCIV	QUWBBA
1269.	QEVPPE01	VATSAK	1321.	QUWXOW	RUKHAG
1270.	QEVPPE01	VATSAK01	1322.	RABBUN	NAGVUM
1271.	QEXRIV	WASNAD	1323.	RABCAU	NAGWAT
1272.	QEYNED	QEYNAM	1324.	RABCEY	NAGWAT
1273.	QIDCIP	WIQMUE	1325.	RABJOP	DUHJIB
1274.	QIFKAT	QIFJAS	1326.	RAFHUE	LABQUD
1275.	QIFKAT	QIFJAS01	1327.	RALDEN	XEHKOK
1276.	QIFKAT	QIFJAS02	1328.	RAPJIE	UMIQEO
1277.	QIHSIJ	CAWKEQ	1329.	RAPRUW	AZSTBB02
1278.	QIJZUE	RORQUE	1330.	RAQCOC	TEXNEP
1279.	QILHOK	APANIQ	1331.	RASBEV	DUHJIB
1280.	QIMKOM	QIMKIG	1332.	RASBIZ	DUHJIB
1281.	QIMKOM	QIMKIG02	1333.	RATBOD	RATBIX
1282.	QIMKOM	QIMKIG03	1334.	RAVFEA	HQOXDO
1283.	QIMKOM02	QIMKIG	1335.	RAVFOK	OCUSOU
1284.	QIMKOM02	QIMKIG02	1336.	RAVFUQ	OCUSOU
1285.	QIMKOM02	QIMKIG03	1337.	RAVGEB	PECYON
1286.	QIMXOB	QIMXER	1338.	RAVQAI	RAVPUB
1287.	QIMXUH	QIMXER	1339.	RAWBIA	RAWBEW
1288.	QIQLUX	QIQLIL	1340.	RAZWEX	RAZWIB
1289.	QIQLUX	QIQLIL01	1341.	REGRAY	HOXOCD
1290.	QIVTUK	ZZZEEU01	1342.	REGREA	HOXOCD
1291.	QIVTUK	ZZZEEU08	1343.	REGRIE	HOXOCD
1292.	QIYKEO	ZULQAY	1344.	RELBAM	REKZUD
1293.	QOBJEW	QOBJAS	1345.	RENQIN	RENQEJ
1294.	QOBJOH	MAWMIJ	1346.	RENQIN01	RENQEJ

1347.	REQWEP	REQWAL	1399.	SAQYAM	PUDXES02
1348.	REVCOL	SEVPEP01	1400.	SAQYAM	PUDXES04
1349.	REVGUV	REVGOP	1401.	SARJIF	FEDPUA
1350.	REVMIP	GAYJAT	1402.	SASYOD	SASXOC
1351.	REVQUF	REVRAM	1403.	SATSEL	DURDAV
1352.	RIGLAW	RIGLEA	1404.	SATSEL	DURDAV01
1353.	RIGMIE	SEQVAN	1405.	SAXDUR	PINCOL
1354.	RIKBIW	RIKBOC	1406.	SAXDUR	PINCOL01
1355.	RINXUH	RINWUG	1407.	SAXQAL	SAXPOY
1356.	RIPFON	RIPFIH	1408.	SAXQAL	SAXPOY01
1357.	RIPVAP	RIPTUH	1409.	SAYPIT	FADHUP
1358.	RIRBID	CTVHVH	1410.	SAYPOZ	FADHOJ
1359.	RIRBID	CTVHVH01	1411.	SAYQEP	SILTOW
1360.	RIRBID	VOKBIG	1412.	SAYQEP	SILTOW01
1361.	RIVJEM	RIVJAI	1413.	SAYQEP	SILTOW11
1362.	RIYYEE	WOVXIO	1414.	SAZQOA	PINCOL
1363.	RIYYEE03	WOVXIO	1415.	SAZQOA	PINCOL01
1364.	ROFFEY	RIZFAI	1416.	SAZXAS	QQQACY01
1365.	ROFNUW	ZOTMOM	1417.	SAZXAS01	QQQACY01
1366.	ROFXOA	KUTPOF	1418.	SECGIQ	SECGOW
1367.	ROHBEV	QUCKOP	1419.	SEDCAG	IVASON
1368.	ROHQOW	OXUXIQ	1420.	SEHKEW	TARTAC
1369.	ROJVIX	ROJVET	1421.	SEHKEW	TARTAC02
1370.	ROKSOZ	QOPYOK	1422.	SEHKEW	TARTAC24
1371.	ROLDAY	ROLCUR	1423.	SENPAF	SENNUX
1372.	ROLVOF	DIXFAR	1424.	SENRUZ	BUDMOD
1373.	ROMDUT	VETVOG01	1425.	SEPDAS	WEDSOO
1374.	ROPHUB	ROPJEN	1426.	SEQBOF	HARJOC
1375.	ROPJAJ	ROPJEN	1427.	SEQBOF	HARJOC01
1376.	RUCFOL	RUCFUR	1428.	SEWBIIH	GUWWUS
1377.	RUKGUZ	RUKHAG	1429.	SEYDIJ	OQUQEX
1378.	RUVPUT	RUVPON	1430.	SEYDIJ01	OQUQEX
1379.	RUWGEV	PROLIN	1431.	SEYDIJ02	OQUQEX
1380.	RUWGEV	PROLIN04	1432.	SEYDIM	SEYDOS
1381.	RUWKEZ	RUWKAV	1433.	SIDCOY	RIZGEM
1382.	RUYGUO	MADTAP	1434.	SIGBUG	TELCEU01
1383.	SABZUQ	SABZEA	1435.	SILFAX	SILDUP
1384.	SAJRUQ	SAJQUP	1436.	SILPEL	SILHON
1385.	SAKYOS10	SAKYUY10	1437.	SIMPOT	OJOCIB
1386.	SAMKIB	IWODUT	1438.	SITRIY	SITREU
1387.	SAMSEH	SAMVOU	1439.	SLBZAC01	RASYOZ
1388.	SAMSEH	SAMVOU02	1440.	SOBBUG	EYOJUZ
1389.	SAMVUA	SAMVOU	1441.	SOCMOO	ABIPIN
1390.	SAMVUA	SAMVOU02	1442.	SOCMOO	ETEXIK
1391.	SAMVUA01	SAMVOU	1443.	SOCMOO	ETEXIK02
1392.	SAMVUA01	SAMVOU02	1444.	SOCMOO	ETEXIK06
1393.	SANACM	AFAZEM	1445.	SOCMOO	KAHVUL
1394.	SANACM	AFAZEM01	1446.	SOGUAN03	ZZZAYP03
1395.	SANACM01	AFAZEM	1447.	SOGUAN20	ZZZAYP03
1396.	SANACM01	AFAZEM01	1448.	SOKPAL	SOKPEP
1397.	SAQQUZ	SAQQOT	1449.	SOMLAI	NEQLOM
1398.	SAQYAM	PUDXES	1450.	SOQNIW	VIWRRAV

1451.	SOQNOC	VITBEG	1503.	TEJMIE	LETBER
1452.	SOWSEC	QATWAJ	1504.	TEJPSEE	BOKTIF
1453.	SOWSIG	QATWUD	1505.	TEJTUX	XAZNES
1454.	SOWSUS	ETDIAM01	1506.	TESTOM	TESTON10
1455.	SOWSUS	ETDIAM15	1507.	TESTOM01	TESTON10
1456.	SOWSUS	ETDIAM18	1508.	TEVFAB	CEXQEB
1457.	SOWTIH	HEXMDA	1509.	TEVFEG	TEVFIK
1458.	SOYDAM	SOYCUF	1510.	TEVFEG	TEVFIK01
1459.	SUGBIF	YADSOL	1511.	TEVVOG	ETANIY
1460.	SUHHEI	OCUSOU	1512.	TEVVOG	ETANIY01
1461.	SULAMH10	SULAMD	1513.	TEXQOC	LOFFET
1462.	SULAMH10	SULAMD01	1514.	TEZMER	JACXUJ
1463.	SULAMH10	SULAMD02	1515.	TFMSAD	TFMSUL
1464.	SULAMH10	SULAMD11	1516.	TFMSPH	TFMSUL
1465.	SULSUX	GOVVEA	1517.	TFMSTH	TFMSUL
1466.	SULSUX	GOVVEA01	1518.	THEOPH	BAPLOT01
1467.	SUMMOC	SUMMIW	1519.	THEOPH	BAPLOT02
1468.	SUNGUD	SUNGOX	1520.	THEOPH	BAPLOT04
1469.	SUPKET	SUPKIX	1521.	THEOPH	DUWXEA
1470.	SURRIG	SURREC	1522.	THEOPH01	BAPLOT01
1471.	SUYVAJ	LITJED	1523.	THEOPH01	BAPLOT02
1472.	TABQIY	TABPUJ	1524.	THEOPH01	BAPLOT04
1473.	TAJRAX	QQQACY01	1525.	THEOPH01	DUWXEA
1474.	TAJRAY	TAJQUR	1526.	THEOPH02	BAPLOT01
1475.	TAKJIB	TAJSUV	1527.	THEOPH02	BAPLOT02
1476.	TAKJIB	TAKJUN	1528.	THEOPH02	BAPLOT04
1477.	TAKJIB	TAKKAU	1529.	THEOPH02	DUWXEA
1478.	TAMQUW	OSILAF	1530.	THEOPH03	BAPLOT01
1479.	TAMQUW	OSILAF01	1531.	THEOPH03	BAPLOT02
1480.	TAPYEP	WIRNIU	1532.	THEOPH03	BAPLOT04
1481.	TAQCOE	AQIYEG	1533.	THEOPH03	DUWXEA
1482.	TARTDL	TARTAC02	1534.	THEOPH04	BAPLOT01
1483.	TARTMM	TARTAM	1535.	THEOPH04	BAPLOT02
1484.	TARTMM01	TARTAM	1536.	THEOPH04	BAPLOT04
1485.	TATBIB	TATBAT	1537.	THEOPH04	DUWXEA
1486.	TAZMOW	EFURIH	1538.	THIAMC	UNEXOA
1487.	TAZMOW	EFURIH02	1539.	THIAMC13	UNEXOA
1488.	TAZMOW	EPINOS	1540.	THIMCH10	GEYXOX
1489.	TAZMOW	IFAKAC	1541.	THIRDN10	BELZEX
1490.	TAZMOW	MUINOS	1542.	THYMMH	THYMIN
1491.	TAZMOW	MYINOL	1543.	THYMMH	THYMIN02
1492.	TAZMOW	MYINOL01	1544.	TIDPOL	JAYZEO
1493.	TAZMOW	QIKZAN	1545.	TIKBIY	TIKBEU
1494.	TAZMOW	QIKZAN01	1546.	TIKBIY	TIKBEU01
1495.	TAZMOW	QIKZAN02	1547.	TIRLOX	UYEHAI
1496.	TAZMOW	YEPNOW	1548.	TIRMIS	QAHQOI
1497.	TAZMOW	ZIVVIL	1549.	TISJAH	PABHIJ
1498.	TECLAP	IKAKOW	1550.	TITGOT	LEJVED
1499.	TEDFAL	CEDKUT	1551.	TIZNEW	TIZNIA
1500.	TEGXEI	YERTOE	1552.	TLALAN10	ZZZIFQ01
1501.	TEHREE	TEHRII	1553.	TMAMSU02	NUNSIY
1502.	TEJGEX	TEJGAT	1554.	TMAMSU02	NUNSIY01

1555.	TOHGIG	OMIXUD	1607.	UJIXAO	SURREC
1556.	TOKYIC	BEQWUS	1608.	UJOQUF	KAMPIY
1557.	TOKYIC	NEDMIS	1609.	UKIZUK	UKOBEC
1558.	TOKYIC	NEDMIS03	1610.	UKIZUK01	UKOBEC
1559.	TOLCAY	IBUMAT	1611.	ULEKUS	DEYSAE
1560.	TOLCAY01	IBUMAT	1612.	ULEQOR	ULEQIL
1561.	TOPBUX	TOPBOR	1613.	ULICOJ	NUZSIM
1562.	TOQFAI	TOQFEM	1614.	ULICOJ	NUZSIM02
1563.	TORSEZ	OBAZOJ	1615.	ULIQAJ	ZEFTAI
1564.	TOSBOU	TOSBIO	1616.	ULOBAA	ULOBOO
1565.	TOSVAA	HUFXAJ	1617.	UMIPIR	XOKVEA
1566.	TOSVAA	HUFXAJ01	1618.	UNEBIY	HEWTAF
1567.	TREHAL01	DEKYEX01	1619.	UNEWUG	UNEWIU
1568.	TREHAL02	DEKYEX01	1620.	UNEWUG	UNEWIU01
1569.	TREHAL03	DEKYEX01	1621.	UNEXAN	UNEWIU
1570.	TREHAL10	DEKYEX01	1622.	UNEXAN	UNEWIU01
1571.	TREHAL11	DEKYEX	1623.	UNUKAQ	SIHLOL
1572.	TREHAL11	DEKYEX01	1624.	UPIZAW	KIZCUT
1573.	TREHAL12	DEKYEX	1625.	UPUKOH	CUFFUG
1574.	TREHAL12	DEKYEX01	1626.	UPUKUN	CUFFUG
1575.	TREHAL13	DEKYEX	1627.	UPUNEA	QERCEW
1576.	TREHAL14	DEKYEX	1628.	UPUNEA	QERCEW01
1577.	TRIMES10	BTCOAC	1629.	UQIFOR	HIBXOF
1578.	TRMHXD	DMHXDM	1630.	UQIFOR	HIBXOF01
1579.	TUBWAQ	ROZHEU	1631.	URUCER	URUCAN
1580.	TUFSOD	TUFSIX	1632.	UTOHOB	UQESAL
1581.	TUQYUZ	TUQYOT	1633.	UVOHET	GIVJUQ
1582.	TUQZAG	TUQYOT	1634.	UVOHET	GIVJUQ02
1583.	TUXLOP	VALQAD	1635.	UVOJUL	UVOKAS
1584.	TUZFIF	TUZFOL	1636.	UZOVIP	UZOVAH03
1585.	TUZFIF01	TUZFOL	1637.	VABKOA	TICFUG
1586.	TYRAMH	SENJEC	1638.	VACBOT	VACBUZ
1587.	UBAFEK	SALWEK	1639.	VACBOT	VACBUZ01
1588.	UBEGEPE	KUVLAP	1640.	VADVEC	DURDAV
1589.	UCEHUG	WAKBAM	1641.	VADVEC	DURDAV01
1590.	UCOJIG	MAJRIZ	1642.	VAJWUZ	ZEGBOE
1591.	UCOJIG	MAJRIZ01	1643.	VALHCL10	VALEHC10
1592.	UCUVET	CUHNEY10	1644.	VALHCL11	VALEHC10
1593.	UCUVET	CUHNEZ	1645.	VALMOL	BERYAZ
1594.	UDAJEP	GUFMAV	1646.	VALMOL	BERYAZ01
1595.	UDENOH	USIMUF	1647.	VAMCAO	VAMBUH
1596.	UFETUX	UFEVAF	1648.	VATKAD	TIRMIR
1597.	UHAHUG	ZOKYON	1649.	VAWWIB	VAWWEX
1598.	UHELAV	UHEKOI	1650.	VAWWOH	VAWWEX
1599.	UHIZES	BENPRL	1651.	VAWWUN	VAWWEX
1600.	UHIZES	BENPRL01	1652.	VAXDUT	DEMYEZ
1601.	UHIZES	BENPRL02	1653.	VAXDUT	DEMYEZ01
1602.	UHURUL	KAPNAQ	1654.	VAXDUT	DEMYEZ03
1603.	UIJIJEC	UIJIJAY	1655.	VAXFAB	DEMYEZ
1604.	UIJIPIN	GEXKID	1656.	VAXFAB	DEMYEZ01
1605.	UIJPOT	GEXKID	1657.	VAXFAB	DEMYEZ03
1606.	UIJIPUZ	GEXKID	1658.	VAYBOO	VABVAX

1659.	VEBREZ	GEYRAD	1711.	WECGUI	WECGIW02
1660.	VEFNAY	VEFNEC	1712.	WECHAP	WECGIW
1661.	VEFPUR	ZOYMOP06	1713.	WECHAP	WECGIW01
1662.	VEFPUR	ZOYMOP07	1714.	WECHAP	WECGIW02
1663.	VEFPUR01	ZOYMOP06	1715.	WEGPAA	YOGZAV
1664.	VEFPUR01	ZOYMOP07	1716.	WEPMOU	AQEMIT
1665.	VEHZAK	ETIWEJ	1717.	WEPMOU	AQEMIT02
1666.	VEJDAQ	VEJCUJ	1718.	WIBQOP	XAYCAB
1667.	VELQUZ	AVULEI	1719.	WIBWOW	KIGKUH
1668.	VELREM	QUKHIP	1720.	WIBWOW	KIGKUH01
1669.	VELREM	QUKHIP01	1721.	WIBWOW	KIGKUH02
1670.	VELREM	QUKHIP02	1722.	WICKUR	JACDAV
1671.	VELYUH	TAKRON	1723.	WIFCIY	WIFCEU
1672.	VERQOB	VERQER	1724.	WIGXUG	ICIBOL
1673.	VERUCH	VERUCE	1725.	WIJNEI	GESTOQ
1674.	VEXRUL	VEXREV	1726.	WIKDAV	MUMBEL
1675.	VICSIM	NAVSUY	1727.	WIKDAV	MUMBEL04
1676.	VICSIM	NAVSUY02	1728.	WINKUA	NETRUZ
1677.	VIMRIS	KUZDIS	1729.	WINKUA	NOBVEF
1678.	VIPCUS	ZOKSUN	1730.	WIRTIA	NUTSUQ
1679.	VIPKUA	YADSOL	1731.	WITHIQ	RUYFEX01
1680.	VIVQIB	OQUPAS	1732.	WIVGEO	RUHWUM
1681.	VIVQOH	OQUPAS	1733.	WIVGEO02	RUHWUM
1682.	VIVQOH01	OQUPAS	1734.	WIWBIP	WIVZUY
1683.	VOBVUE	VOBWAL	1735.	WIZMUO	SOLMOW
1684.	VOBWEP	VOBWoz	1736.	WOBNEG	RUJQOD
1685.	VODJII	HUJBUK	1737.	WOBQAG	DIPGIS10
1686.	VOHQAM	THEOPI01	1738.	WOBQAG	DIPGIS12
1687.	VOHWIZ	RITDUU	1739.	WOBQAG	DIPGIS13
1688.	VOJRIV	PAZADO10	1740.	WOCREL	SEPNUX
1689.	VOKNOA	MORPLN01	1741.	WODKOQ	WODKIK
1690.	VOLWEA	VOLWUQ	1742.	WODVAN	PHTHAC
1691.	VOYSIL	CEXQUR	1743.	WODZAS	WODZEW
1692.	VOYSOR	CEXQUR	1744.	WOGGIK	XOTWUY
1693.	VULGOZ	VULGUF	1745.	WOMBAB	WOLZUS
1694.	VUMDEN	VUMDAJ	1746.	WOVYEL	HIQWEJ
1695.	VUNFEQ	VUNFIU	1747.	WOVYEL01	HIQWEJ
1696.	VURTIM	VURSUX	1748.	WOZPAE	MECWIC
1697.	VUXREM	VUXRIQ	1749.	WOZPAE	MECWIC03
1698.	VUXREM	VUXRIQ01	1750.	WUDHOU	IVEVIO
1699.	WADROK	WACNOF	1751.	WUTVIQ	RENCUI
1700.	WAFNAT	COTZAN06	1752.	WUWJAA	UNEXOA
1701.	WAFNAT	HXACAN	1753.	WUWREM	WUWRIQ
1702.	WAFNAT	HXACAN29	1754.	WUWROW	ATPRCL01
1703.	WAFNAT	HXACAN39	1755.	WUXPUC	WUXQAJ
1704.	WAFAV	WAFNUN	1756.	WUZDIG	WUZDEC
1705.	WAGLOG	WAGLUM	1757.	XADXAE	XADWUX
1706.	WAKCAO	WAKCES	1758.	XAHGAR	KOVruk
1707.	WATDUR	WATFAZ	1759.	XAHWOS	GORKUV
1708.	WAYDEH	WAYDAD	1760.	XAHWOS	GORKUV01
1709.	WECGUI	WECGIW	1761.	XAHXAF	LEFHEL
1710.	WECGUI	WECGIW01	1762.	XAJTIM	XAJSUX

1763.	XANLOO	WIGMUV	1815.	XOZTAI	XOZTEM
1764.	XAQREN	XAQRAJ	1816.	XOZTAI	XOZTIQ
1765.	XASLOS	NIXZAX	1817.	XURHOH	IDIMAJ
1766.	XAVQEQQ	IYUQIC	1818.	XUSGUO	XUSGOI
1767.	XAYGEJ	VETVOG01	1819.	XUVLOP	EFEMUX01
1768.	XEBBAI	PAQLUP	1820.	XUVPAG	PAHBUW
1769.	XEGPUX	EXAXEG	1821.	XUVROW	KIDHEN
1770.	XEGQEIQ	EWUZED	1822.	XYFUAD10	BRADOS
1771.	XEJLON	XEJLUT	1823.	YABHOA	BITCEM11
1772.	XENTES	SEZREU	1824.	YABHOA	BITCEM13
1773.	XENTES	SEZREU01	1825.	YADZAH	CBMZPN01
1774.	XEPNUD	XEPNIR	1826.	YADZAH	CBMZPN03
1775.	XEPREQ	ADEZER	1827.	YADZAH	CBMZPN11
1776.	XESPUH	XESQAO	1828.	YADZAH	CBMZPN12
1777.	XEVRUO	XEVROI	1829.	YADZAH	CBMZPN16
1778.	XEVRUO	XEVSAY	1830.	YADZAH	CBMZPN20
1779.	XEYYOQ	PIHLAU	1831.	YAFCOY	IXEWIR
1780.	XEYZUX	CEXXIO	1832.	YAFCOY	IXEWIR01
1781.	XEZHER	XEZHOB	1833.	YAFCOY	IXEWIR02
1782.	XIBGIZ	XIBGEV	1834.	YAFCOY	IXEWIR03
1783.	XICPIK	XICPOQ	1835.	YAFCOY	IXEWIR04
1784.	XICRAE	OFUYIZ	1836.	YAFHUL	ECAQAD01
1785.	XIKPUF	XIKQAM	1837.	YAGFEU	SAQJAW
1786.	XIKPUF	XIKQAM01	1838.	YAKWAJ	BISMEV
1787.	XIPCUU	LONDOJ	1839.	YAKWAJ	BISMEV01
1788.	XIRKUG	ZOBCOK	1840.	YAKWAJ	BISMEV03
1789.	XIVHOA	ZIVKAQ	1841.	YAKWAJ	BISMEV04
1790.	XIYSAA	ALAHC	1842.	YAKWAJ	BISMEV07
1791.	XIZPEB	XIZNUP	1843.	YAKWAJ	BISMEV14
1792.	XIZPEB01	XIZNUP	1844.	YANCUM	YAKVAI
1793.	XODYIY	XODYOE	1845.	YANHIH	KOBBA
1794.	XOFSOA	BOFTAT	1846.	YAPMEJ	UHITOV
1795.	XOHDAB	XOHCUU	1847.	YASHIL	YASGOQ
1796.	XOHFEH	XOHFAD	1848.	YASHIL	YASGOQ01
1797.	XOJJAJ	XOJHUB	1849.	YASHIL	YASGOQ02
1798.	XOJMOA	XOJMIU	1850.	YASHIL	YASGOQ03
1799.	XOJMUG	XOJMIU	1851.	YECVUX	QOBGUL
1800.	XOJPUJ	XOJQUE	1852.	YECVUX	QOBHAS
1801.	XOKKIT	GIXXOB	1853.	YECWAE	LIHJAQ
1802.	XOKKIT	GIXXOB01	1854.	YECWAE	LIHJAQ01
1803.	XOKTEY	XOKTAU	1855.	YECWEI	QOBHIA
1804.	XOKVIE	XOKVEA	1856.	YEFBER	YASRIU
1805.	XOMBEH	XOMCAE	1857.	YEFLOL	BOJDOD
1806.	XOMWOL	COTZAN06	1858.	YEJNEG	QAMXUY
1807.	XOMWOL	HXACAN	1859.	YEJNEG	QAMXUY01
1808.	XOMWOL	HXACAN29	1860.	YEKHH	YEKHON
1809.	XOMWOL	HXACAN39	1861.	YESTOG	LEVSAJ
1810.	XOSHUI	JAYFOF	1862.	YETPIW	YETPES
1811.	XOXRAE	XOXQUX	1863.	YIJBID	YIJBEZ
1812.	XOYXUF	HUZKUI	1864.	YILGAB	XIPJOY
1813.	XOZSUB	XOZTEM	1865.	YILKUZ	SEBROG
1814.	XOZSUB	XOZTIQ	1866.	YIMZAV	HOXOCD

1867.	YINFAC	ZZZIVG02	1919.	ZIPWOM	GAMBUT01
1868.	YIQFEK	OGEKOB	1920.	ZIPWOM	GAMBUT04
1869.	YIXGUJ	YIXHAQ	1921.	ZIVVOR	EFURIH
1870.	YIXSUU	YIXSOO	1922.	ZIVVOR	EFURIH02
1871.	YIZQOO	PHTHAC	1923.	ZIVVOR	EPINOS
1872.	YODRUF	BOLBOT01	1924.	ZIVVOR	FOPKOK
1873.	YODYUL	ROMMEM	1925.	ZIVVOR	IFAKAC
1874.	YOTKEY	GAPSIB	1926.	ZIVVOR	MUINOS
1875.	YOWRAF	QONSET	1927.	ZIVVOR	MYINOL
1876.	YOYPIL	EGORIB	1928.	ZIVVOR	MYINOL01
1877.	YOYZIX	MAJRIZ	1929.	ZIVVOR	MYINOL03
1878.	YOYZIX	MAJRIZ01	1930.	ZIVVOR	QIKZAN
1879.	YUDFAF	COWSIR	1931.	ZIVVOR	QIKZAN01
1880.	YUDFAF	COWSIR01	1932.	ZIVVOR	QIKZAN02
1881.	YUFGIR	YUFGAJ	1933.	ZIVVOR	QIKZUH
1882.	YUFGIR	YUFGAJ01	1934.	ZIVVOR	YEPNOW
1883.	YUJNUM01	CUYCEF	1935.	ZIVVOR	ZIVVIL
1884.	YUJNUM01	LICWOM	1936.	ZIZYAI	TITLIS
1885.	YUJNUM01	LICWOM01	1937.	ZOQHEU	KIWTEP
1886.	YUJPAU01	CUYCEF	1938.	ZOYMUV	ZOYMOP
1887.	YUJPAU01	LICWOM	1939.	ZOYMUV	ZOYMOP01
1888.	YUJPAU01	LICWOM01	1940.	ZOYMUV	ZOYMOP02
1889.	YUKVAD	WOLXOK03	1941.	ZUDBUV	PAPCOA
1890.	YULFUH	YULGAO	1942.	ZUNHIB	RAMGOB
1891.	YUPBER	VALINO	1943.	ZUNHIB	RAMGOB02
1892.	YUVVAO	YUWYEW	1944.	ZUNHOH	RAMGOB
1893.	YUVVAO	YUWYEW01	1945.	ZUNHOH	RAMGOB02
1894.	YUXGUV	ZZZEEU01	1946.	ZUNHOH01	RAMGOB
1895.	YUXGUV	ZZZEEU08	1947.	ZUNHOH01	RAMGOB02
1896.	YUYROZ	YUYRIT	1948.	ZZZAMS03	OPENAN
1897.	YUYTIV	PICACC	1949.	ZZZAMS06	OPENAN
1898.	YUZTET	QIMKIG	1950.	ZZZPPI02	URICAC
1899.	YUZTET	QIMKIG02	1951.	ZZZPRW01	MAJRIZ
1900.	YUZTET	QIMKIG03	1952.	ZZZPRW01	MAJRIZ01
1901.	ZACSOO	JALYAZ	1953.	ZZZRLO01	ACRDIN
1902.	ZAGTUZ	ZAGZOZ	1954.	ZZZRLO01	ACRDIN01
1903.	ZAHJIB	CYURAC03	1955.	ZZZRLO01	ACRDIN05
1904.	ZAHQAA	PAQUCL	1956.	ZZZRLO01	ACRDIN06
1905.	ZARLEM	XAPTIU	1957.	ZZZRLO01	ACRDIN08
1906.	ZATTEV	ZATTIZ	1958.	ZZZSBA01	HUYBUY
1907.	ZEBWUA	ZEBWOU	1959.	ZZZSBA03	HUYBUY
1908.	ZEFMII	ZEFMUU	1960.	ZZZSQK02	BEVBAF
1909.	ZEFSUB	ZEFTAI	1961.	AEDTAH	EDTAXX01
1910.	ZEKPEM	ZEKPOW	1962.	AEDTAH	EDTAXX02
1911.	ZEKPEM	ZEKPOW01	1963.	AMPALX	AHPHAL
1912.	ZEKPIQ	ZEKPOW	1964.	DAWGOX	KETXIR
1913.	ZEKPIQ	ZEKPOW01	1965.	ETANOE	ETANIY
1914.	ZELTUI	ZELVAQ	1966.	ETANOE	ETANIY01
1915.	ZESQUL	ZESRAS	1967.	ETAPIA	ETANIY
1916.	ZEXYOQ	MEWSUE	1968.	ETAPIA	ETANIY01
1917.	ZICPAF	FIBYEW	1969.	FINFEN	SUNFIS
1918.	ZIFQOW	ZIFQEM	1970.	NEZMAJ	NEZTAQ

1971.	PASXIS01	PASWUD	2023.	QESVOA	QESVUG
1972.	PAXQAH01	GAVKOG	2024.	QIMXOB01	QIMXER
1973.	QESFID	LOBDUB	2025.	QQQCIY10	CEKGUU
1974.	SAKDUF	FUDGOA	2026.	QQQCIY10	CEKGUU02
1975.	SAYQEP01	SILTOW	2027.	QQQCIY10	CEKGUU06
1976.	SAYQEP01	SILTOW01	2028.	RITMEN	CIVXIO01
1977.	SAYQEP01	SILTOW11	2029.	ROPRAR	ROPQUK
1978.	TIZVUT	TIZWAA	2030.	VOMPEU	VOMNUI
1979.	TIZVUT	TIZWAA01	2031.	WEGPED	WEGPIH
1980.	UGUFOS	JERPIG	2032.	XEBQAZ	CARPEQ
1981.	VAHKIY	VAHKIY10	2033.	ZZZPPI01	URICAC
1982.	YECWIM	QOBHIA	2034.	ZZZWQA01	BUNJAV
1983.	SAGSUR	BODWAS	2035.	CUXYOK	CUXZAX
1984.	YUFCAE	EFEZAR	2036.	FABMOM	ESEZIM
1985.	ADPRTY	ADPTRR	2037.	HCPZHO	HCPZBO
1986.	BABYOV	BABYUB	2038.	KEMDIQ	KEMDEM
1987.	BEXPOK	BEXPIE	2039.	KEMDIQ	KEMDEM01
1988.	BUTAHC11	DUHJIB	2040.	NIJLEX	NEFPET
1989.	DUGTOQ	DUHJIB	2041.	PINOLH	PINCOL
1990.	ETHYLO	DUFBOV10	2042.	PINOLH	PINCOL01
1991.	ETXDHY01	DUFBOV10	2043.	QAFKAK	DAYFUH
1992.	FEFNOT01	CBMZPN01	2044.	QQQBKD01	PYRGAL02
1993.	FEFNOT01	CBMZPN03	2045.	RESMIL	DADPDS
1994.	FEFNOT01	CBMZPN11	2046.	SOWVAB	QOBHUM
1995.	FEFNOT01	CBMZPN12	2047.	SOWVOP	QOBHEW
1996.	FEFNOT01	CBMZPN16	2048.	TOLQUG	NOVGUA
1997.	FEFNOT01	CBMZPN20	2049.	XOZSIP	RUCDOJ
1998.	FUSWOH04	EFUMAU	2050.	LUXYAE	LUXYIM
1999.	FUSWOH04	EFUMAU03	2051.	YIHJOP	ETGUAN
2000.	FUSWOH04	EFUMAU06	2052.	AHEREK	JAYPUU05
2001.	GISWOW	UFADUC	2053.	NENMEZ	LUYMAT
2002.	IKALIR	IKAKOW	2054.	CAFINE	NIWFEE02
2003.	IKALIR01	IKAKOW	2055.	CAFINE	NIWFEE04
2004.	IKAMEO	IKAKOW	2056.	YINDUU	SABZOK
2005.	IKAMOY	IKAKOW	2057.	JOSKUZ	SALBUT01
2006.	IKAMOY01	IKAKOW	2058.	BIMGIP	SOCTOT
2007.	IKAVOF	ETHANE01	2059.	NEXSIU	THIOUR06
2008.	IKAVOF	ETHANE04	2060.	NEXSUG	THIOUR06
2009.	JARBUA	JARBOU	2061.	OCAZEZ	UBOTIQ
2010.	KUYGOA	ZZZITY01	2062.	PUSBIP	PUSBOV
2011.	KUYGUG	KAMTAW	2063.	QEVLAE	VUKVAY
2012.	KUYGUG	KAMTAW03	2064.	COSNUT	COSNON
2013.	KUYHAN	DCLBEN			
2014.	KUYHAN	DCLBEN02			
2015.	KUYHAN	DCLBEN03			
2016.	KUYHER	CISTON01			
2017.	LAFLOV	KAYSOT			
2018.	MUSDEL	DUHJIB			
2019.	MUSD OV	DUHJIB			
2020.	NAHCIJ	ACETYL02			
2021.	NAHCOP	JAYDUI			
2022.	QANGES	RABTIX			

```

1 #28 Determine water molecule stoichiometry
2
3 from __future__ import division
4 from ccdc import io, search
5 import argparse
6 import os
7 import glob
8 from decimal import*
9
10 #This script counts the number of water molecules in each hydrate entry
11
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
13 Classes\chart_class1_hydrates_with_waterless_forms.txt"
14 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
15 Classes\chart_class2_hydrates_without_waterless_forms.txt"
16
17 class Runner(argparse.ArgumentParser):
18     def __init__(self):
19         super(self.__class__, self).__init__(description=__doc__)
20         self.add_argument(
21             '-i', '--input', default=filepath1,
22             help='input database filepath')
23         self.add_argument(
24             '-o', '--output', default='waters.gcd',
25             help='output file [waters.gcd]')
26         self.add_argument(
27             '-m', '--maximum', default=0, type=int,
28             help='Maximum number of structures to find [all]')
29
30     args = self.parse_args()
31
32     self.args = args
33     self.settings = search.Search.Settings()
34     self.settings.max_hit_structures = self.args.maximum
35
36     def run(self):
37
38         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
39         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
40
41         with io.EntryWriter(self.args.output) as waters_writer:
42
43             duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
44             'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
45             'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
46
47             loop = 0
48             total_monohydrate = 0
49             total_dihydrate = 0
50             total_trihydrate = 0
51             total_tetrahydrate = 0
52             total_pentahydrate = 0
53             total_hexahydrate = 0
54             total_heptahydrate = 0
55             total_octahydrate = 0
56             total_nonahydrate = 0
57             total_decahydrate = 0
58             total_more_than_ten = 0
59             total_hemihydrate = 0
60             total_less_than_one = 0
61             total_more_than_one = 0
62             total_undefined = 0
63             grand_total = 0
64             entryl = entry_reader1

```

```

64 while loop < 2:
65     monohydrate = 0
66     dihydrate = 0
67     trihydrate = 0
68     tetrahydrate = 0
69     pentahydrate = 0
70     hexahydrate = 0
71     heptahydrate = 0
72     octahydrate = 0
73     nonahydrate = 0
74     decahydrate = 0
75     more_than_ten = 0
76     hemihydrate = 0
77     less_than_one = 0
78     more_than_one = 0
79     undefined = 0
80     sub_total = 0
81
82     #An already found list is used so hydrates in hydrate-anhydride pairs
83     #that have multiple anhydrous forms will not be counted multiple times
84     already_found = []
85     for a in range(len(entry1)):
86         sub_total += 1
87         grand_total += 1
88         skip = 0
89         if loop == 0:
90             if entry1[a].identifier in duplicate_identifier_dictionary:
91                 if
92                     duplicate_identifier_dictionary.get(entry1[a].identifier)
93                     not in already_found:
94
95                         already_found.append(duplicate_identifier_dictionary.get(
96                             entry1[a].identifier))
97                     else:
98                         skip += 1
99             else:
100                 if hash(entry1[a].identifier) not in already_found:
101                     already_found.append(hash(entry1[a].identifier))
102                 else:
103                     skip += 1
104
105             #Some structures require a manual specification of the number of
106             #water molecules
107             if entry1[a].identifier == 'BULMEA03' or entry1[a].identifier ==
108                 'LOYXAA':
109                 monohydrate += 1
110                 total_monomhydrate += 1
111                 if entry1[a].identifier == 'HEQQOJ':
112                     less_than_one += 1
113                     total_less_than_one += 1
114                 if entry1[a].identifier == 'KIJFAN':
115                     undefined += 1
116                     total_undefined += 1
117                 if skip == 0:
118                     #The formula of each entry is searched for water
119                     #The integer preceding the formula for water will be either an
120                     #integral or non-integral number
121                     #The number of water molecules were sorted up to 10, then any
122                     #structure with more than 10 waters was one category
123                     #Non-integral waters was restricted to 0.5 (hemihydrates), less
124                     #than 1, more than 1, and undefined (ex. x(H2 O1) or n(H2 O1))
125                     title = entry1[a].formula
126                     pieces = title.split(',')
127                     water = []
128                     for b in range(len(pieces)):
129                         if '(H2 O1)' in pieces[b] or pieces[b] == 'H2 O1' or '(D2
130                         O1)' in pieces[b] or pieces[b] == 'D2 O1':
131                             water.append(pieces[b])
132                     if len(water) == 1:

```

```

120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
    if '(' not in water[0]:
        monohydrate += 1
        total_monomhydrate += 1
    else:
        pin = water[0].find('(')
        water[0] = water[0][:pin]
        if 'x' in water[0] or 'n' in water[0]:
            undefined += 1
            total_undefined += 1
        elif '.' in water[0]:
            water[0] = Decimal(water[0])
            if water[0] == 0.5:
                hemihydrate += 1
                total_hemihydrate += 1
            elif water[0] < 1:
                less_than_one += 1
                total_less_than_one += 1
            else:
                more_than_one += 1
                total_more_than_one += 1
        else:
            water[0] = int(water[0])
            if water[0] == 2:
                dihydrate += 1
                total_dihydrate += 1
            elif water[0] == 3:
                trihydrate += 1
                total_trihydrate += 1
            elif water[0] == 4:
                tetrahydrate += 1
                total_tetrahydrate += 1
            elif water[0] == 5:
                pentahydrate += 1
                total_pentahydrate += 1
            elif water[0] == 6:
                hexahydrate += 1
                total_hexahydrate += 1
            elif water[0] == 7:
                heptahydrate += 1
                total_heptahydrate += 1
            elif water[0] == 8:
                octahydrate += 1
                total_octahydrate += 1
            elif water[0] == 9:
                nonahydrate += 1
                total_nonahydrate += 1
            elif water[0] == 10:
                decahydrate += 1
                total_decahydrate += 1
            else:
                more_than_ten += 1
                total_more_than_ten += 1

#The number of structures in each water category are printed followed
#by the percentage from the total number of structures in the input file
print 'start of one round'
print monohydrate
print monohydrate/sub_total
print dihydrate
print dihydrate/sub_total
print trihydrate
print trihydrate/sub_total
print tetrahydrate
print tetrahydrate/sub_total
print pentahydrate
print pentahydrate/sub_total
print hexahydrate

```

```

186     print hexahydrate/sub_total
187     print heptahydrate
188     print heptahydrate/sub_total
189     print octahydrate
190     print octahydrate/sub_total
191     print nonahydrate
192     print nonahydrate/sub_total
193     print decahydrate
194     print decahydrate/sub_total
195     print more_than_ten
196     print more_than_ten/sub_total
197     print hemihydrate
198     print hemihydrate/sub_total
199     print less_than_one
200     print less_than_one/sub_total
201     print more_than_one
202     print more_than_one/sub_total
203     print undefined
204     print undefined/sub_total
205     print 'end of one round'
206
207     loop += 1
208     if loop == 1:
209         entry1 = entry_reader2
210
211     #Numbers for the grand total are also recorded, these are determined by
212     #adding the numbers for class 1 and class 2 hydrates together
213     print 'start of one round'
214     print total_monohydrate
215     print total_monohydrate/grand_total
216     print total_dihydrate
217     print total_dihydrate/grand_total
218     print total_trihydrate
219     print total_trihydrate/grand_total
220     print total_tetrahydrate
221     print total_tetrahydrate/grand_total
222     print total_pentahydrate
223     print total_pentahydrate/grand_total
224     print total_hexahydrate
225     print total_hexahydrate/grand_total
226     print total_heptahydrate
227     print total_heptahydrate/grand_total
228     print total_octahydrate
229     print total_octahydrate/grand_total
230     print total_nonahydrate
231     print total_nonahydrate/grand_total
232     print total_decahydrate
233     print total_decahydrate/grand_total
234     print total_more_than_ten
235     print total_more_than_ten/grand_total
236     print total_hemihydrate
237     print total_hemihydrate/grand_total
238     print total_less_than_one
239     print total_less_than_one/grand_total
240     print total_more_than_one
241     print total_more_than_one/grand_total
242     print total_undefined
243     print total_undefined/grand_total
244     print 'end of one round'
245
246 if __name__ == '__main__':
247     r = Runner()
248     r.run()

```

```

1 #29 Find pairs with 1 and 2+ components and 2000 unpaired subsets
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 from decimal import*
8
9 #This script separates structures from each class into those with 1 organic component
10 and 2+ organic components
11
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
13 Classes\chart_class1_hydrates_with_waterless_forms.txt"
14 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
15 Classes\chart_class1_waterless_forms_with_hydrate.txt"
16 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
17 Classes\chart_class2_hydrates_without_waterless_forms.txt"
18 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
19 Classes\chart_class3_waterless_forms_without_hydrate.txt"
20
21
22 class Runner(argparse.ArgumentParser):
23     def __init__(self):
24         super(self.__class__, self).__init__(description=__doc__)
25         self.add_argument(
26             '-i', '--input', default=filepath1,
27             help='input database filepath')
28         self.add_argument(
29             '-o', '--output', default='one_organic_hydrates_with_waterless_forms.gcd',
30             help='output file [one_organic_hydrates_with_waterless_forms.gcd]')
31         self.add_argument(
32             '-m', '--maximum', default=0, type=int,
33             help='Maximum number of structures to find [all]')
34
35     args = self.parse_args()
36
37     self.args = args
38     self.settings = search.Search.Settings()
39     self.settings.max_hit_structures = self.args.maximum
40
41     def run(self):
42
43         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
44         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
45         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
46         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
47
48         with io.EntryWriter(self.args.output) as waters_writer:
49             with
50                 io.EntryWriter("more_than_one_organic_hydrates_with_waterless_forms.gcd")
51                 as writer1:
52                     with io.EntryWriter("one_organic_waterless_forms_with_hydrate.gcd") as
53                         writer2:
54                             with
55                                 io.EntryWriter("more_than_one_organic_waterless_forms_with_hydrate.gcd")
56                                 as writer3:
57                                     with
58                                         io.EntryWriter("one_organic_hydrates_without_waterless_form.gcd")
59                                         as writer4:
60                                             with
61                                                 io.EntryWriter("more_than_one_organic_hydrates_without_waterless_form.gcd")
62                                                 as writer5:
63                                                     with
64                                                         io.EntryWriter("one_organic_waterless_forms_without_hydrate.gcd")
65                                                         as writer6:

```

```

52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
    with
        io.EntryWriter("more_than_one_organic_waterless_forms_
_without_hydrate.gcd") as writer7:
            with
                io.EntryWriter("subset_one_organic_hydrates_witho_
ut_waterless_form.gcd") as writer8:
                    with
                        io.EntryWriter("subset_more_than_one_organic_
hydrates_without_waterless_form.gcd") as
writer9:
                            with
                                io.EntryWriter("subset_one_organic_waterl
ess_forms_without_hydrate.gcd") as
writer10:
                                    with
                                        io.EntryWriter("subset_more_than_one_
organic_waterless_forms_without_hydra
te.gcd") as writer11:
                                            for a in
                                                range(len(entry_reader1)):
                                                    #The entry formula is used
                                                    #to determine the number of
                                                    #components in each structure
                                                    #Water does not count
                                                    #toward the number of
                                                    #components for hydrate
                                                    #structures
                                                    formulas =
                                                        entry_reader1[a].formula.split(
                                                            ',')
                                                    waters = []
                                                    for b in
                                                        range(len(formulas)):
                                                            if '(H2 O1)' in
                                                                formulas[b] or '(D2
                                                                O1)' in formulas[b] or
                                                                formulas[b] == 'H2 O1'
                                                                or formulas[b] == 'D2
                                                                O1':
                                                                    waters.append(formula
                                                                        s[b])
                                                    formulas = [x for x in
                                                        formulas if x not in waters]
                                                    if len(formulas) == 1:
                                                        waters_writer.write(entry_
                                                        _reader1[a])
                                                        writer2.write(entry_reade
                                                        r2[a])
                                                    else:
                                                        writer1.write(entry_reade
                                                        r1[a])
                                                        writer3.write(entry_reade
                                                        r2[a])
#Subsets of class 2 and class 3
#are generated in this script
#The first 2,000 structures
#with 1 component and no
#disorder were written to an
#output file
#The first 2,000 structures

```

```

77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
with 2+ components and no
disorder were written to an
output file
count = 0
count1 = 0
for c in
range(len(entry_reader3)):
    formulas2 =
entry_reader3[c].formula.split(',')
waters2 = []
for d in
range(len(formulas2)):
    if '(H2 O1)' in
formulas2[d] or '(D2
O1)' in formulas2[d] or
formulas2[d] == 'H2 O1'
or formulas2[d] == 'D2
O1':
        waters2.append(formul
as2[d])
formulas2 = [x for x in
formulas2 if x not in
waters2]
if len(formulas2) == 1:
    writer4.write(entry_reade
r3[c])
    if count < 2000 and
entry_reader3[c].crystal.
has_disorder == False:
        writer8.write(entry_r
eader3[c])
        count += 1
    else:
        writer5.write(entry_reade
r3[c])
        if count1 < 2000 and
entry_reader3[c].crystal.
has_disorder == False:
            writer9.write(entry_r
eader3[c])
            count1 += 1
        count2 = 0
        count3 = 0
        for e in
range(len(entry_reader4)):
            formulas3 =
entry_reader4[e].formula.split(
',')
            if len(formulas3) == 1:
                writer6.write(entry_reade
r4[e])
                if count2 < 2000 and
entry_reader4[e].crystal.
has_disorder == False:
                    writer10.write(entry_
reader4[e])
                    count2 += 1
                else:

```

```
107 writer7.write(entry_reade  
r4[e])  
108 if count3 < 2000 and  
entry_reader4[e].crystal.  
has_disorder == False:  
109     writer11.write(entry_  
reader4[e])  
     count3 += 1  
110  
111  
112  
113  
114 if __name__ == '__main__':  
    r = Runner()  
    r.run()  
115  
116  
117
```

```

1 #30 Determine crystal system symmetry for hydrate-anhydrite pairs
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7
8 #This script records the crystal system symmetry of hydrate and waterless form
9 #structures in the hydrate-anhydrite pairs list
10
11 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
12 Classes\one_organic_hydrates_with_waterless_forms.txt"
13 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
14 Classes\more_than_one_organic_hydrates_with_waterless_forms.txt"
15 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
16 Classes\one_organic_waterless_forms_with_hydrates.txt"
17 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
18 Classes\more_than_one_organic_waterless_forms_with_hydrates.txt"
19
20 class Runner(argparse.ArgumentParser):
21
22     def __init__(self):
23         super(self.__class__, self).__init__(description=__doc__)
24         self.add_argument(
25             '-i', '--input', default=filepath1,
26             help='input database filepath')
27
28         self.add_argument(
29             '-o', '--output', default='systems_HA_subsets.gcd',
30             help='output file [systems_HA_subsets.gcd]')
31
32         self.add_argument(
33             '-m', '--maximum', default=0, type=int,
34             help='Maximum number of structures to find [all]')
35
36     args = self.parse_args()
37
38     def run(self):
39
40         crystal_reader1 = io.CrystalReader(filepath1, format='identifiers')
41         crystal_reader2 = io.CrystalReader(filepath2, format='identifiers')
42         crystal_reader3 = io.CrystalReader(filepath3, format='identifiers')
43         crystal_reader4 = io.CrystalReader(filepath4, format='identifiers')
44
45         with io.EntryWriter(self.args.output) as system_writer:
46
47             #This dictionary assigns values to the different crystal system symmetries
48             #Higher numbers are given to crystal systems with higher degrees of symmetry
49             system_dictionary = {'triclinic':1, 'monoclinic':2, 'orthorhombic':3,
50             'tetragonal':4, 'trigonal':5, 'rhombohedral':5, 'hexagonal':6, 'cubic':7}
51
52             loop = 0
53             crystal1 = crystal_reader1
54             crystal2 = crystal_reader3
55             while loop < 2:
56                 #The relationship between the degree of symmetry between the hydrate
57                 #and waterless form crystal system are compared first
58                 same_system = 0
59                 hydrate_higher = 0
60                 waterless_higher = 0
61                 #Lower crystal system symmetries (triclinic, monoclinic, and
62                 #orthorhombic) were further evaluated

```

```

60 both_triclinic = 0
61 #Htri_Wmono = hydrate is triclinic and waterless form is monoclinic
62 Htri_Wmono = 0
63 Htri_Wortho = 0
64 both_monoclinic = 0
65 Hmono_Wortho = 0
66 both_orthorhombic = 0
67 Hmono_Wtri = 0
68 Hortho_Wtri = 0
69 Hortho_Wmono = 0
70
71 for b in range(len(crystall1)):
72     if system_dictionary.get(crystall1[b].crystal_system) ==
73         system_dictionary.get(crystal2[b].crystal_system):
74             same_system += 1
75     if system_dictionary.get(crystall1[b].crystal_system) >
76         system_dictionary.get(crystal2[b].crystal_system):
77             hydrate_higher += 1
78     if system_dictionary.get(crystall1[b].crystal_system) <
79         system_dictionary.get(crystal2[b].crystal_system):
80             waterless_higher += 1
81     if system_dictionary.get(crystall1[b].crystal_system) == 1:
82         if system_dictionary.get(crystal2[b].crystal_system) == 1:
83             both_triclinic += 1
84         elif system_dictionary.get(crystal2[b].crystal_system) == 2:
85             Htri_Wmono += 1
86         elif system_dictionary.get(crystal2[b].crystal_system) == 3:
87             Htri_Wortho += 1
88     elif system_dictionary.get(crystall1[b].crystal_system) == 2:
89         if system_dictionary.get(crystal2[b].crystal_system) == 2:
90             both_monoclinic += 1
91         elif system_dictionary.get(crystal2[b].crystal_system) == 3:
92             Hmono_Wortho += 1
93     elif system_dictionary.get(crystall1[b].crystal_system) == 3:
94         if system_dictionary.get(crystal2[b].crystal_system) == 3:
95             both_orthorhombic += 1
96         if system_dictionary.get(crystal2[b].crystal_system) == 1:
97             Hmono_Wtri += 1
98         elif system_dictionary.get(crystall1[b].crystal_system) == 3:
99             Hortho_Wtri += 1
100    elif system_dictionary.get(crystal2[b].crystal_system) == 2:
101        if system_dictionary.get(crystall1[b].crystal_system) == 3:
102            Hortho_Wmono += 1
103
104 #The number of structures in each category is printed
105 print 'start of one round'
106 print same_system
107 print hydrate_higher
108 print waterless_higher
109 print both_triclinic
110 print Htri_Wmono
111 print Htri_Wortho
112 print both_monoclinic
113 print Hmono_Wortho
114 print both_orthorhombic
115 print Hmono_Wtri
116 print Hortho_Wtri
117 print Hortho_Wmono
118 print 'end of one round'
119 loop += 1
120 crystall1 = crystal_reader2
121 crystal2 = crystal_reader4
122
123
```

```
124 if __name__ == '__main__':
125     # This runs the script
126     r = Runner()
127     r.run()
128
129
130
```

```

1 #31 Determine crystal system symmetry for different classes
2
3 from __future__ import division
4 from ccdc import io, search
5 import argparse
6 import os
7 import glob
8
9 #This script determines the percentage of structures for each class (divided into 1
10 component and 2+ components) with the seven different crystal system symmetries
11
12 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
13 Classes\one_organic_hydrates_with_waterless_forms.txt"
14 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
15 Classes\more_than_one_organic_hydrates_with_waterless_forms.txt"
16 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
17 Classes\one_organic_waterless_forms_with_hydrates.txt"
18 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
19 Classes\more_than_one_organic_waterless_forms_with_hydrates.txt"
20 filepath5 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
21 Classes\one_organic_hydrates_without_waterless_form.txt"
22 filepath6 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
23 Classes\more_than_one_organic_hydrates_without_waterless_form.txt"
24 filepath7 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
25 Classes\one_organic_waterless_forms_without_hydrate.txt"
26 filepath8 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
27 Classes\more_than_one_organic_waterless_forms_without_hydrate.txt"
28 filepath9 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
29 Classes\hydrates.txt"
30 filepath10 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
31 Classes\entries_in_working_data_set.txt"
32 filepath11 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
33 Classes\chart_class1_hydrates_with_waterless_forms.txt"
34 filepath12 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
35 Classes\chart_class1_waterless_forms_with_hydrate.txt"
36 filepath13 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
37 Classes\chart_class2_hydrates_without_waterless_forms.txt"
38 filepath14 = "C:\Users\jenwe\Hydrates Manuscript\Step6 Change Composition of Three
39 Classes\chart_class3_waterless_forms_without_hydrate.txt"
40
41
42
43 class Runner(argparse.ArgumentParser):
44
45     def __init__(self):
46         super(self.__class__, self).__init__(description=__doc__)
47         self.add_argument(
48             '-i', '--input', default=filepath7,
49             help='input database file')
50         self.add_argument(
51             '-o', '--output', default='systems_four_subsets.gcd',
52             help='output file [systems_four_subsets.gcd]')
53         self.add_argument(
54             '-m', '--maximum', default=0, type=int,
55             help='Maximum number of structures to find [all]')
56
57         args = self.parse_args()
58
59         self.args = args
60         self.settings = search.Search.Settings()
61         self.settings.max_hit_structures = self.args.maximum
62
63     def run(self):
64
65         crystal_reader1 = io.CrystalReader(filepath1, format='identifiers')
66         crystal_reader2 = io.CrystalReader(filepath2, format='identifiers')

```

```

53     crystal_reader3 = io.CrystalReader(filepath3, format='identifiers')
54     crystal_reader4 = io.CrystalReader(filepath4, format='identifiers')
55     crystal_reader5 = io.CrystalReader(filepath5, format='identifiers')
56     crystal_reader6 = io.CrystalReader(filepath6, format='identifiers')
57     crystal_reader7 = io.CrystalReader(filepath7, format='identifiers')
58     crystal_reader8 = io.CrystalReader(filepath8, format='identifiers')
59     crystal_reader9 = io.CrystalReader(filepath9, format='identifiers')
60     crystal_reader10 = io.CrystalReader(filepath10, format='identifiers')
61     crystal_reader11 = io.CrystalReader(filepath11, format='identifiers')
62     crystal_reader12 = io.CrystalReader(filepath12, format='identifiers')
63     crystal_reader13 = io.CrystalReader(filepath13, format='identifiers')
64     crystal_reader14 = io.CrystalReader(filepath14, format='identifiers')
65
66     with io.EntryWriter(self.args.output) as system_writer:
67
68         duplicate_identifier_dictionary = {'GADKOL':1, 'MECYEY':2, 'OBASAN':3,
69         'OMOMAE':4, 'OTONUI':5, 'QOJWAO':6, 'TZBFZO':7, 'VOBXOZ':8, 'VUFJAI':9,
70         'VUFJAJ':10, 'WEPQOA':11, 'XADFUD':12, 'YABREY':13, 'YULCOZ':14, 'ZOVFEV':15}
71
72         loop = 0
73         crystall = crystal_reader1
74         while loop < 14:
75             triclinic = 0
76             monoclinic = 0
77             orthorhombic = 0
78             tetragonal = 0
79             trigonal = 0
80             hexagonal = 0
81             cubic = 0
82             already_found = []
83             for a in range(len(crystall)):
84                 skip = 0
85                 if loop == 0 or loop == 1 or loop == 2 or loop == 3 or loop == 10
86                 or loop == 11:
87                     if crystall[a].identifier in duplicate_identifier_dictionary:
88                         if
89                             duplicate_identifier_dictionary.get(crystall[a].identifier)
90                             not in already_found:
91
92                             already_found.append(duplicate_identifier_dictionary.get(
93                                 crystall[a].identifier))
94                         else:
95                             skip += 1
96
97                     else:
98                         if hash(crystall[a].identifier) not in already_found:
99                             already_found.append(hash(crystall[a].identifier))
100                         else:
101                             skip += 1
102
103             if skip == 0:
104                 system = crystall[a].crystal_system
105                 if system == 'triclinic':
106                     triclinic +=1
107                 elif system == 'monoclinic':
108                     monoclinic +=1
109                 elif system == 'orthorhombic':
110                     orthorhombic +=1
111                 elif system == 'tetragonal':
112                     tetragonal +=1
113                 elif system == 'trigonal' or system == 'rhombohedral':
114                     trigonal +=1
115                 elif system == 'hexagonal':
116                     hexagonal +=1
117                 else:
118                     cubic +=1
119
120             total = triclinic + monoclinic + orthorhombic + tetragonal + trigonal +
121             hexagonal + cubic

```

```

112
113     #The percentage of structures for each class is printed
114     print 'start of one round'
115     print triclinic/total
116     print monoclinic/total
117     print orthorhombic/total
118     print tetragonal/total
119     print trigonal/total
120     print hexagonal/total
121     print cubic/total
122     print 'end of one round'
123
124     loop += 1
125     if loop == 1:
126         crystall1 = crystal_reader2
127     if loop == 2:
128         crystall1 = crystal_reader3
129     if loop == 3:
130         crystall1 = crystal_reader4
131     if loop == 4:
132         crystall1 = crystal_reader5
133     if loop == 5:
134         crystall1 = crystal_reader6
135     if loop == 6:
136         crystall1 = crystal_reader7
137     if loop == 7:
138         crystall1 = crystal_reader8
139     if loop == 8:
140         crystall1 = crystal_reader9
141     if loop == 9:
142         crystall1 = crystal_reader10
143     if loop == 10:
144         crystall1 = crystal_reader11
145     if loop == 11:
146         crystall1 = crystal_reader12
147     if loop == 12:
148         crystall1 = crystal_reader13
149     if loop == 13:
150         crystall1 = crystal_reader14
151
152
153
154 if __name__ == '__main__':
155     # This runs the script
156     r = Runner()
157     r.run()
158
159
160

```

```

1 #32 Find hydrate-anhydrate pairs determined at the same temperature
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 from decimal import Decimal
8
9 #This script finds hydrate and waterless form structures in hydrate-anhydrate pairs
10 #that were collected at the same temperature
11 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
12 Classes\one_organic_hydrates_with_waterless_forms.txt"
13 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
14 Classes\more_than_one_organic_hydrates_with_waterless_forms.txt"
15 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
16 Classes\one_organic_waterless_forms_with_hydrates.txt"
17 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
18 Classes\more_than_one_organic_waterless_forms_with_hydrates.txt"
19
20 class Runner(argparse.ArgumentParser):
21
22     def __init__(self):
23         super(self.__class__, self).__init__(description=__doc__)
24         self.add_argument(
25             '-i', '--input', default=filepath1,
26             help='input database filepath')
27     )
28         self.add_argument(
29             '-o', '--output',
30             default='same_temp_one_organic_hydrates_with_waterless_forms.gcd',
31             help='output file [same_temp_one_organic_hydrates_with_waterless_forms.gcd]')
32     )
33         self.add_argument(
34             '-m', '--maximum', default=0, type=int,
35             help='Maximum number of structures to find [all]')
36
37     args = self.parse_args()
38
39     self.args = args
40     self.settings = search.Search.Settings()
41     self.settings.max_hit_structures = self.args.maximum
42
43     def run(self):
44
45         entry_reader1 = io.EntryReader(filepath1, format='identifiers')
46         entry_reader2 = io.EntryReader(filepath2, format='identifiers')
47         entry_reader3 = io.EntryReader(filepath3, format='identifiers')
48         entry_reader4 = io.EntryReader(filepath4, format='identifiers')
49
50         with io.EntryWriter(self.args.output) as same_temp_writer:
51             with
52                 io.EntryWriter("same_temp_one_organic_waterless_forms_with_hydrates.gcd")
53                     as writer1:
54                         with
55                             io.EntryWriter("same_temp_more_than_one_organic_hydrates_with_waterless_f
56 orms.gcd") as writer2:
57                                 with
58                                     io.EntryWriter("same_temp_more_than_one_organic_waterless_forms_with_
59 hydrates.gcd") as writer3:
60
61                                         loop = 0
62                                         entry4 = entry_reader1
63                                         entry5 = entry_reader3
64                                         while loop < 2:
65                                             for b in range(len(entry4)):

```

```

56     #Some structures do not report a temperature in Python
57     #and were discarded to prevent another manual sort
58     #For cases where there is a reported temperature, both
59     #the unit and number were recorded
60     if entry4[b].temperature != None and
61     entry5[b].temperature != None:
62         loop1 = 0
63         if loop == 0:
64             entry1 = entry_reader1
65         if loop == 1:
66             entry1 = entry_reader2
67         while loop1 < 2:
68             temp_text = entry1[b].temperature.split(' ')
69             if len(temp_text) == 4 and temp_text[3] == 'K':
70                 if temp_text[1].isdigit() == False:
71                     structure_temp = Decimal(temp_text[1])
72                 else:
73                     structure_temp = int(temp_text[1])
74                     structure_unit = 'K'
75             elif len(temp_text) == 3 and temp_text[2] == 'K':
76                 if temp_text[1].isdigit() == False:
77                     structure_temp = Decimal(temp_text[1])
78                 else:
79                     structure_temp = int(temp_text[1])
80                     structure_unit = 'K'
81             elif len(temp_text) == 3 and temp_text[2] ==
82                 'deg.C':
83                 if temp_text[1].isdigit() == False:
84                     structure_temp = Decimal(temp_text[1])
85                 else:
86                     structure_temp = int(temp_text[1])
87                     structure_unit = 'deg.C'
88             elif len(temp_text) == 3 and 'K' in temp_text[2]:
89                 if temp_text[2].isdigit() == False:
90                     start = temp_text[2].index('K')
91                     structure_temp =
92                         Decimal(temp_text[2][:start])
93                 else:
94                     start = temp_text[2].index('K')
95                     structure_temp =
96                         int(temp_text[2][:start])
97                     structure_unit = 'K'
98             elif len(temp_text) == 2 and temp_text[1] == 'K':
99                 if temp_text[0].isdigit() == False:
100                     structure_temp = Decimal(temp_text[0])
101                 else:
102                     structure_temp = int(temp_text[0])
103                     structure_unit = 'K'
104             elif len(temp_text) == 2 and temp_text[1] ==
105                 'deg.C':
106                 if temp_text[0].isdigit() == False:
107                     structure_temp = Decimal(temp_text[0])
108                 else:
109                     structure_temp = int(temp_text[0])
110                     structure_unit = 'deg.C'
111             elif len(temp_text) == 2 and 'K' in temp_text[1]:
112                 if temp_text[1].isdigit() == False:
113                     start = temp_text[1].index('K')
114                     structure_temp =
115                         Decimal(temp_text[1][:start])
116                     structure_unit = 'K'
117             elif len(temp_text) == 2 and 'deg.C' in
118                 temp_text[1]:

```

```

113                                     if temp_text[1].isdigit() == False:
114                                         start = temp_text[1].index('deg.C')
115                                         structure_temp =
116                                         Decimal(temp_text[1][:start])
117                                         else:
118                                             start = temp_text[1].index('deg.C')
119                                             structure_temp =
120                                             int(temp_text[1][:start])
121                                             structure_unit = 'deg.C'
122                                         if loop1 == 0:
123                                             hydrate_temp = structure_temp
124                                             hydrate_unit = structure_unit
125                                         if loop1 == 1:
126                                             waterless_temp = structure_temp
127                                             waterless_unit = structure_unit
128                                         loop1 += 1
129                                         if loop1 == 1 and loop == 0:
130                                             entry1 = entry_reader3
131                                         if loop1 == 1 and loop == 1:
132                                             entry1 = entry_reader4
133
134                                         #If the temperature for each structure is within 5
135                                         #degrees of each other, the structures are considered to
136                                         #be taken at roughly the same temperature and are
137                                         #written to output files
138                                         if hydrate_unit == waterless_unit and
139                                         abs(hydrate_temp-waterless_temp) < 5:
140                                             if loop == 0:
141                                                 same_temp_writer.write(entry4[b])
142                                                 writer1.write(entry5[b])
143                                             else:
144                                                 writer2.write(entry4[b])
145                                                 writer3.write(entry5[b])
146                                         elif hydrate_unit == 'K' and waterless_unit ==
147                                         'deg.C' and abs(hydrate_temp-(waterless_temp+273)) < 5:
148                                             if loop == 0:
149                                                 same_temp_writer.write(entry4[b])
150                                                 writer1.write(entry5[b])
151                                             else:
152                                                 writer2.write(entry4[b])
153                                                 writer3.write(entry5[b])
154                                         elif hydrate_unit == 'deg.C' and waterless_unit ==
155                                         'K' and abs((hydrate_temp+273)-waterless_temp) < 5:
156                                             if loop == 0:
157                                                 same_temp_writer.write(entry4[b])
158                                                 writer1.write(entry5[b])
159                                             else:
160                                                 writer2.write(entry4[b])
161                                                 writer3.write(entry5[b])
162
163                                         loop += 1
164                                         entry4 = entry_reader2
165                                         entry5 = entry_reader4
166
167                                         if __name__ == '__main__':
168                                             r = Runner()
169                                             r.run()

```

```

1 #33 Determine the packing fraction of hydrate-anhydrite pairs at the same temperature
2
3 from ccdc import io, search
4 import argparse
5 import os
6 import glob
7 import numpy
8 from numpy import median
9
10 #This script determines the difference in packing fraction between hydrate and
11 #waterless form structures collected at the same temperature
12
13 filepath1 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
14 Classes\same_temp_one_organic_hydrates_with_waterless_forms.txt"
15 filepath2 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
16 Classes\same_temp_more_than_one_organic_hydrates_with_waterless_forms.txt"
17 filepath3 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
18 Classes\same_temp_one_organic_waterless_forms_with_hydrates.txt"
19 filepath4 = "C:\Users\jenwe\Hydrates Manuscript\Step7 Analyze Structures in Three
20 Classes\same_temp_more_than_one_organic_waterless_forms_with_hydrates.txt"
21
22 class Runner(argparse.ArgumentParser):
23
24     def __init__(self):
25         super(self.__class__, self).__init__(description=__doc__)
26         self.add_argument(
27             '-i', '--input', default=filepath1,
28             help='input database file')
29         self.add_argument(
30             '-o', '--output', default='packing_fraction_outside_5percent_hydrate.gcd',
31             help='output file [packing_fraction_outside_5percent_hydrate.gcd]')
32
33     args = self.parse_args()
34
35     self.args = args
36     self.settings = search.Search.Settings()
37     self.settings.max_hit_structures = self.args.maximum
38
39     def run(self):
40
41         crystal_reader1 = io.CrystalReader(filepath1, format='identifiers')
42         crystal_reader2 = io.CrystalReader(filepath2, format='identifiers')
43         crystal_reader3 = io.CrystalReader(filepath3, format='identifiers')
44         crystal_reader4 = io.CrystalReader(filepath4, format='identifiers')
45
46         with io.EntryWriter(self.args.output) as packing_fraction_writer:
47             with io.EntryWriter("packing_fraction_outside_5percent_waterless_form.gcd")
48                 as writer1:
49
50                 loop = 0
51                 crystal1 = crystal_reader1
52                 crystal2 = crystal_reader3
53                 #Hydrate-anhydrite pairs whose packing fractions are within 5% of each
54                 #other were counted
55                 #The packing fraction difference was recorded at 0.5% intervals
56                 #Pairs where the hydrate has a higher packing fraction are counted
57                 #separate from those where the waterless form has a higher packing
58                 #fraction
59                 while loop < 2:
60                     Hydrate_half_percent = 0
61                     Hydrate_one_percent = 0

```

```

59     Hydrate_onehalf_percent = 0
60     Hydrate_two_percent = 0
61     Hydrate_twohalf_percent = 0
62     Hydrate_three_percent = 0
63     Hydrate_threehalf_percent = 0
64     Hydrate_four_percent = 0
65     Hydrate_fourhalf_percent = 0
66     Hydrate_five_percent = 0
67     Waterless_half_percent = 0
68     Waterless_one_percent = 0
69     Waterless_onehalf_percent = 0
70     Waterless_two_percent = 0
71     Waterless_twohalf_percent = 0
72     Waterless_three_percent = 0
73     Waterless_threehalf_percent = 0
74     Waterless_four_percent = 0
75     Waterless_fourhalf_percent = 0
76     Waterless_five_percent = 0
77 #Any pairs where one or both structure has disorder are omitted
from this analysis
78 for b in range(len(crystall1)):
79     if crystall1[b].has_disorder == False and
crystal2[b].has_disorder == False:
        crystall1[b].molecule.normalise_hydrogens()
        hydrate_fraction = crystall1[b].packing_coefficient
        crystal2[b].molecule.normalise_hydrogens()
        waterless_fraction = crystal2[b].packing_coefficient
        if hydrate_fraction - waterless_fraction > 0:
            if hydrate_fraction - waterless_fraction < 0.005:
                Hydrate_half_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.005 and
hydrate_fraction - waterless_fraction < 0.01:
                Hydrate_one_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.01 and
hydrate_fraction - waterless_fraction < 0.015:
                Hydrate_onehalf_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.015 and
hydrate_fraction - waterless_fraction < 0.02:
                Hydrate_two_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.02 and
hydrate_fraction - waterless_fraction < 0.025:
                Hydrate_twohalf_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.025 and
hydrate_fraction - waterless_fraction < 0.03:
                Hydrate_three_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.03 and
hydrate_fraction - waterless_fraction < 0.035:
                Hydrate_threehalf_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.035 and
hydrate_fraction - waterless_fraction < 0.04:
                Hydrate_four_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.04 and
hydrate_fraction - waterless_fraction < 0.045:
                Hydrate_fourhalf_percent += 1
            elif hydrate_fraction - waterless_fraction >= 0.045 and
hydrate_fraction - waterless_fraction < 0.05:
                Hydrate_five_percent += 1
        else:
            packing_fraction_writer.write(crystall1[b])
            writer1.write(crystal2[b])
    else:
        if waterless_fraction - hydrate_fraction < 0.005:
            Waterless_half_percent += 1
        elif waterless_fraction - hydrate_fraction >= 0.005 and
waterless_fraction - hydrate_fraction < 0.01:
            Waterless_one_percent += 1
        elif waterless_fraction - hydrate_fraction >= 0.01 and

```

```

114     waterless_fraction - hydrate_fraction < 0.015:
115         Waterless_onehalf_percent += 1
116     elif waterless_fraction - hydrate_fraction >= 0.015 and
117         waterless_fraction - hydrate_fraction < 0.02:
118         Waterless_two_percent += 1
119     elif waterless_fraction - hydrate_fraction >= 0.02 and
120         waterless_fraction - hydrate_fraction < 0.025:
121         Waterless_twohalf_percent += 1
122     elif waterless_fraction - hydrate_fraction >= 0.025 and
123         waterless_fraction - hydrate_fraction < 0.03:
124         Waterless_three_percent += 1
125     elif waterless_fraction - hydrate_fraction >= 0.03 and
126         waterless_fraction - hydrate_fraction < 0.035:
127         Waterless_threehalf_percent += 1
128     elif waterless_fraction - hydrate_fraction >= 0.035 and
129         waterless_fraction - hydrate_fraction < 0.04:
130         Waterless_four_percent += 1
131     elif waterless_fraction - hydrate_fraction >= 0.04 and
132         waterless_fraction - hydrate_fraction < 0.045:
133         Waterless_fourhalf_percent += 1
134     elif waterless_fraction - hydrate_fraction >= 0.045 and
135         waterless_fraction - hydrate_fraction < 0.05:
136         Waterless_five_percent += 1
137     else:
138         packing_fraction_writer.write(crystall1[b])
139         writer1.write(crystal2[b])
140
141     #The number of structures in each difference category is printed
142     print 'start of one round'
143     print "Hydrate fraction higher by 0.5% for " +
144         str(Hydrate_half_percent) + " pairs"
145     print "Hydrate fraction higher by 1.0% for " +
146         str(Hydrate_one_percent) + " pairs"
147     print "Hydrate fraction higher by 1.5% for " +
148         str(Hydrate_onehalf_percent) + " pairs"
149     print "Hydrate fraction higher by 2.0% for " +
150         str(Hydrate_two_percent) + " pairs"
151     print "Hydrate fraction higher by 2.5% for " +
152         str(Hydrate_twohalf_percent) + " pairs"
153     print "Hydrate fraction higher by 3.0% for " +
str(Hydrate_three_percent) + " pairs"
154     print "Hydrate fraction higher by 3.5% for " +
str(Hydrate_threehalf_percent) + " pairs"
155     print "Hydrate fraction higher by 4.0% for " +
str(Hydrate_four_percent) + " pairs"
156     print "Hydrate fraction higher by 4.5% for " +
str(Hydrate_fourhalf_percent) + " pairs"
157     print "Hydrate fraction higher by 5.0% for " +
str(Hydrate_five_percent) + " pairs"
158
159     print "Waterless fraction higher by 0.5% for " +
str(Waterless_half_percent) + " pairs"
160     print "Waterless fraction higher by 1.0% for " +
str(Waterless_one_percent) + " pairs"
161     print "Waterless fraction higher by 1.5% for " +
str(Waterless_onehalf_percent) + " pairs"
162     print "Waterless fraction higher by 2.0% for " +
str(Waterless_two_percent) + " pairs"
163     print "Waterless fraction higher by 2.5% for " +
str(Waterless_twohalf_percent) + " pairs"
164     print "Waterless fraction higher by 3.0% for " +
str(Waterless_three_percent) + " pairs"
165     print "Waterless fraction higher by 3.5% for " +
str(Waterless_threehalf_percent) + " pairs"
166     print "Waterless fraction higher by 4.0% for " +
str(Waterless_four_percent) + " pairs"
167     print "Waterless fraction higher by 4.5% for " +
str(Waterless_fourhalf_percent) + " pairs"

```

```
154     print "Waterless fraction higher by 5.0% for " +
155         str(Waterless_five_percent) + " pairs"
156     print 'end of one round'
157
158     loop += 1
159     if loop == 1:
160         crystal1 = crystal_reader2
161         crystal2 = crystal_reader4
162
163
164
165
166
167 if __name__ == '__main__':
168     # This runs the script
169     r = Runner()
170     r.run()
171
172
173
```