```python
#Use chemdraw API to convert cdxml to sdf
# deficiency: chemdraw has be to closed and opened everytime. This is not the case for
cdx 2 cdxml conversion
import glob
import comtypes.client as w32
import datetime



#assemble list of cdx or CDXML files

#file_list = glob.glob('D:/patent1/**/*.cdx', recursive = True)
file_list = glob.glob('D:\\Patent1\\Structures_CDXML_C1\\*.cdxml', recursive = True)



# Creates invisible ChemDraw object.
#ChemDraw = w32.CreateObject("ChemDraw.Application")

n = (len(file_list))
print(n)


current_time = datetime.datetime.now()  #check time
print (current_time) #to estimate time needed for the job using test run

i = 1
for file in file_list:
    # Creates invisible ChemDraw object.
    ChemDraw = w32.CreateObject("ChemDraw.Application")
    file_pointer = ChemDraw.Documents.Open(file)
    document = ChemDraw.Documents.Item(i)

    # Save document.
    location = 'D:\\Patent1\\Structures_CDXML_C1\\' + file[31:62] + '.sdf'
    print(location)
    document.SaveAs(location)

    # Close Document.
    document.Close()
    # Closes ChemDraw
    ChemDraw.Quit()

current_time = datetime.datetime.now() #check time
print (current_time)
```

```
1    #Create new table from parent table (Ironpython in spotfire)
2    import System
3    from Spotfire.Dxp.Data import DataType,DataTableSaveSettings
4    from Spotfire.Dxp.Data.Import import TextFileDataSource,TextDataReaderSettings
5    from System.IO import Path,StreamWriter,StringReader,StreamWriter,MemoryStream,SeekOrigin
6    from Spotfire.Dxp.Application.Visuals import VisualContent
7    from Spotfire.Dxp.Application.Visuals import CrossTablePlot
8    from Spotfire.Dxp.Data import DataValueCursor,RowSelection,IndexSet
9
10   def getVisual(visualTitle):
11       for page in Document.Pages:
12           for vis in page.Visuals:
13               if vis.Title == visualTitle:
14                   return vis.As[VisualContent]()
15
16   #Modify the line below to specify the title of the cross table visualization
17   vis = getVisual("My_Cross_Table")
18
19   stream=MemoryStream()
20   writer=StreamWriter(stream)
21   vis.ExportText(writer)
22
23   stream.Seek(0,SeekOrigin.Begin)
24
25   readerSettings=TextDataReaderSettings()
26   readerSettings.Separator="\t"
27   readerSettings.AddColumnNameRow(0)
28
29   textDataSource =TextFileDataSource(stream,readerSettings)
30
31   if Document.Data.Tables.Contains("DataFromCrossTable"):
32       Document.Data.Tables["DataFromCrossTable"].ReplaceData(textDataSource)
33   else:
34       newTable = Document.Data.Tables.Add("DataFromCrossTable", textDataSource)
35       tableSettings = DataTableSaveSettings (newTable, False, False)
36       Document.Data.SaveSettings.DataTableSettings.Add(tableSettings)
37
38   myTable = Document.Data.Tables["DataFromCrossTable"]
39
40   n = 0
41   for col in myTable.Columns:
42       n = n + 1
43       if n == 1:
44           MyCol = col.Name
45
46   LastRow = myTable.RowCount + 1
47   rowsToRemove=IndexSet(myTable.RowCount,False)
48
49   # Reference to the Column of the Table
50   dataValuesCursor=DataValueCursor.CreateFormatted(myTable.Columns[MyCol])
51
52   i = 0
53   for row in myTable.GetRows(dataValuesCursor):
54       i = i + 1
55     #Statement to remove or keep rows.
56   if i == LastRow:
57       rowsToRemove.AddIndex(row.Index)
58
59   myTable.RemoveRows(RowSelection(rowsToRemove))
60       # change date type from string to real
61   from Spotfire.Dxp.Data.Transformations import ExpressionTransformation,ColumnSelection
62   from Spotfire.Dxp.Data import *
63
64   table = Document.Data.Tables['DataFromCrossTable']
65   rowsToInclude = IndexSet(table.RowCount,True)
66   t = ExpressionTransformation()
```

```python
67  #get the list of columns
68  columns = table.Columns
69  for cc in columns:
70      ccname=cc.Name
71      #ccname='Avg(standard_value)'
72      cursor = DataValueCursor.CreateFormatted(table.Columns[ccname])
73      values=[]
74      for row in table.GetRows(rowsToInclude,cursor):
75          values.append(cursor.CurrentValue)
76          try:
77              values = [float(x) for x in values ]
78              t.ColumnReplacements.Add(
                    ccname,'real(['+ccname+'])',ColumnSelection(ccname))
79              # print ('datatypechanged')
80          except:
81              t.ColumnReplacements.Add(
                    ccname,'string(['+ccname+'])',ColumnSelection(ccname))
82              # print (ccname,'original column did not contain real data')
83              pass
84  table.AddTransformation(t)
```

```python
#combine mol files into sdf and add properties to each molecule, including the location
of its tiff format
#the code is used to construct US patent strucsture database
from rdkit import Chem
from rdkit.Chem import AllChem
#glob to get all the mol files
import glob

#assemble list of mol files

#sdfFile = open("D:\patent_sdf_new\I2023.sdf") #open sdf file in append mode
file_list = glob.glob('H:/patent1/I2023*/**/*.mol',
                      recursive = True)



print(len(file_list))
#write mol files to SD

with Chem.SDWriter("D:\patent_sdf_new\I2023.sdf") as w: #does not have to be existing
file
    for file in file_list:
        m = Chem.MolFromMolFile(file) # read mol to rdkit molecule
        if m is None: #some mol files can't be read into rdkit
#             print(file) #name of files that can't be recognized
            m = Chem.MolFromSmiles('C1=CC=CN=C1') #add random structure so compound can
            be recorded in sdf
            m.SetProp("name", file[46:77]) #add property
            location = 'D:' + file[2:78] + 'tif'
            m.SetProp("location", location) #add property
            folder= 'D:' + file[2:45]
            m.SetProp("folder", folder)   #add property
            try:
                w.write(m)
            except:
                print(file, "can't be kekulized")

        else:
#             print(file)
#             print(m)
            m.SetProp("name", file[46:77])
            location1 = 'D:' + file[2:78] + 'tif'
            m.SetProp("location", location1)
            folder1= 'D:' + file[2:45]
            m.SetProp("folder", folder1)
            try:
                w.write(m)
            except:
                print(file, "can't be kekulized")
#sdfFile.close()
w.close()
```

```sql
-- DB browser code for assemble different tables
DROP TABLE IF EXISTS aa;
CREATE TABLE aa AS SELECT

--group_concat(actp.type) AS type_g,
--group_concat(actp.relation) AS relation_g,
--group_concat(actp.value) AS value_g,
--group_concat(actp.units) AS units_g,
--group_concat(actp.text_value) AS text_value_g,
group_concat(actp.standard_type) AS standard_type_g,
group_concat(actp.standard_relation) AS standard_relation_g,
group_concat(actp.standard_value) AS standard_value_g,
group_concat(actp.standard_units) AS standard_units_g,
group_concat(actp.standard_text_value) AS standard_text_value_g,
--group_concat(actp.activity_id) AS activity_id_g,
--group_concat(actp.comments) AS comments_g,
--group_concat(actp.result_flag) AS result_flag_g,
substr(group_concat(DISTINCT di.max_phase_for_ind ORDER by di.max_phase_for_ind
DESC), 1,1) AS di_max_phase_for_ind_g,
actp.activity_id,
md.chembl_id,
md.molregno,
a.description as 'a_description',
cs.canonical_smiles,
td.pref_name as 'td_pref_name',
a.assay_organism as 'a_assay_organism',
a.assay_tissue as 'a_assay_tissue',
a.assay_type as 'a_assay_type',
a.assay_test_type as 'a_assay_test_type',

--actp.type as 'actp.type',
--actp.standard_type as 'actp.standard_type',
--actp.standard_relation,
--actp.standard_value,
--actp.standard_units,
--actp.standard_text_value,
--act.type as 'act.type',
act.pchembl_value,
act.standard_type as 'act_standard_type',
act.type as 'act_type',
act.standard_relation as 'act_standard_relation',
act.standard_value as 'act_standard_value',
act.standard_units as 'act_standard_units',
act.activity_id as 'act_activity_id',
act.record_id as 'act_record_id',
act.doc_id as 'act_doc_id',

md.first_approval,
md.oral
```

```
--ms.synonyms,
--group_concat(ms.synonyms) AS ms.synonyms_g,
--ms.syn_type,
--group_concat(ms.syn_type) AS ms.syn_type,
--group_concat(rc.company) AS rc.company_g,
--rc.company

--FROM (Select DISTINCT molregno, max_phase_for_ind from drug_indication di Where
di.max_phase_for_ind = (SELECT MAX(max_phase_for_ind) FROM drug_indication)) AS di
FROM (Select DISTINCT molregno, max_phase_for_ind from drug_indication di) AS di
--FROM (Select molregno, max_phase_for_ind from drug_indication di Where
di.max_phase_for_ind IS NOT NULL) AS di
--LEFT JOIN (Select* From activities act WHERE act.standard_value IS NOT NULL AND
act.molregno IN ('397524')) as act ON di.molregno = act.molregno --select specific
compounds
LEFT JOIN (Select* From activities act WHERE act.standard_value IS NOT NULL) as act
ON di.molregno = act.molregno
LEFT JOIN molecule_dictionary md ON act.molregno = md.molregno --AND md.molregno =
'397524'
LEFT JOIN assays a ON act.assay_id = a.assay_id
LEFT JOIN compound_structures cs ON md.molregno = cs.molregno
LEFT JOIN target_dictionary td ON a.tid = td.tid
LEFT JOIN activity_properties actp ON act.activity_id = actp.activity_id
--LEFT JOIN molecule_synonyms ms ON md.molregno = ms.molregno
--LEFT JOIN research_companies rc ON ms.res_stem_id = rc.res_stem_id
WHERE md.molecule_type = 'Small molecule' --AND md.molregno IN ('397524')
GROUP BY md.chembl_id, act.activity_id
ORDER BY act.molregno DESC, act.activity_id DESC
LIMIT 2000000;
```