

Standard Operating Procedure: Calculating R3m from Molecular Dynamic Simulation Data

Author: Kevin DeBoyace

January 2019

Abstract

This document will describe the procedure used to calculate the molecular descriptor R3m from individual molecule coordinates retrieved from a molecular dynamics (MD) simulation run in Materials Studio.

Contents

1	Introduction	2
2	Exporting the 3D Atomistic Files	2
3	Run the MATLAB Code	4
4	MATLAB code	5
4.1	GetAtomCoords.m	5
4.2	R3mCalculate_auto.m	12
4.3	AtomConnection.m	13
4.4	AtomicWeighting.m	16
4.5	EuclidDistance.m	17
4.6	InfluenceDistanceMat.m	18
4.7	MolecInfluenceMatrix.m	19
4.8	MolecMatrix.m	19

1 Introduction

This document will describe the procedure and the MATLAB code used to calculate the molecular descriptor R3m from individual molecule coordinates retrieved from a molecular dynamics (MD) simulation run in Materials Studio. It is assumed that the user has previously run a MD simulation starting with an amorphous cell built using a single molecule. This document will describe exporting the data from the simulation and applying the MATLAB code described to calculate the molecular descriptor R3m for each individual molecule in each of the selected frames from the simulation. The output is a large number of R3m values for which a distribution can be determined. This information is expected to be useful for determining if the R3m value calculated using only the SMILES file is representative of the amorphous form. Additionally, the width of the R3m distribution may help to identify molecules of interest (e.g. more flexible molecules may have a wider distribution). This data could help to explain outliers in the future. Similar methods may also be useful to explore the impact of the presence of solvent on the molecule conformation and subsequent R3m value.

2 Exporting the 3D Atomistic Files

After the MD simulation has been completed, the first step is to export several frames from the simulation. In this context, 'frame' refers to a single step of the simulation, and is a snapshot of the system. An example of a single frame of the simulation is shown in Figure 1.

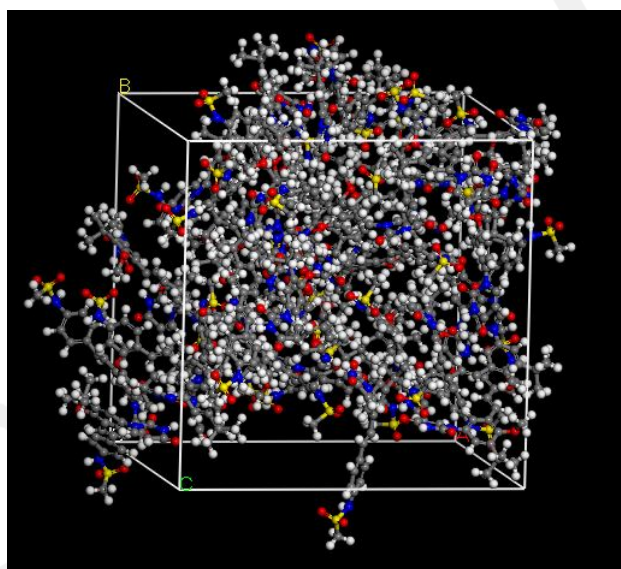


Figure 1: A frame of the MD simulation

The trajectory file which contains all of the simulation data will have a .xtd extension. Right click on the file and select 'Properties' to determine the file extension. The icon should look like the one shown in Figure 2.



Figure 2: Icon for the trajectory (.xtd) file.

Once the file is open, the user should open the 'animation' toolbar if it is not already available. Go to View > Toolbars > Animation. The toolbar should now appear, and should match Figure 3.



Figure 3: Animation toolbar

In the animation toolbar, press the 'Animation mode' button dropdown arrow. This button will look like a looping arrow. From the dropdown, select 'Options'. A new pop-up dialog will appear as shown in Figure 4.

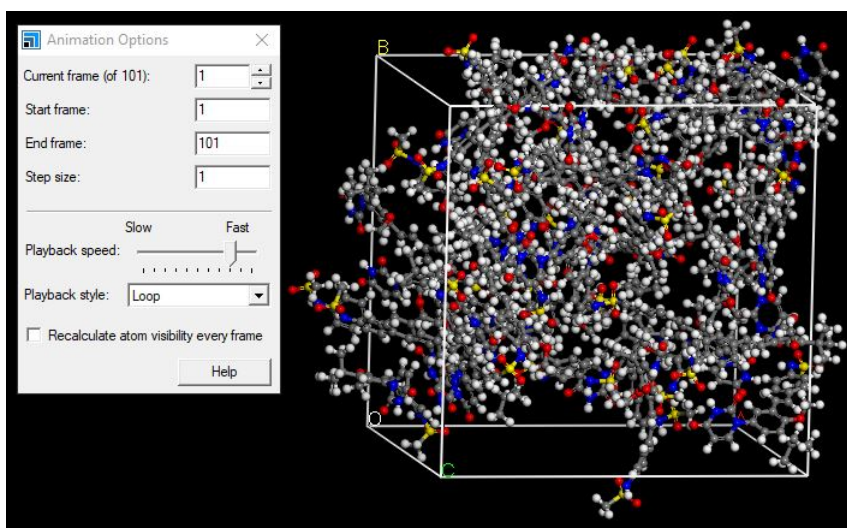


Figure 4: Animation Options dialog box shown alongside the 1st frame of the simulation.

Changing the value in the 'Current frame (of 101)' will update the position of the molecules. The total number of frames will always be 101. It is important to note that there are far more steps in the simulation than 101. As a result, to estimate the approximate location of the frame within the simulation, use percent of duration. For example, frame 40 will show a frame from the simulation at approximately 40% through the total time. The goal is to select multiple frames from the simulation once the system has equilibrated. For example, the first frame of the simulation is unlikely to be representative of the system that is being simulated because in this case the cell had a specified volume. Time is needed for the simulation to reach equilibrium, and hopefully, a more realistic cell volume. In this case, 40% of the way through the simulation corresponded with an equilibrated system. In total, 7 frames will be taken from the simulation, spaced periodically throughout the equilibrated range. In this case, frames 40, 50, 60, 70, 80, 90, and 100 will be exported. In this example, each frame contains coordinates for 40 individual molecules. Type '40' into the 'Current frame (of 100)' field and press the tab key. The selected frame should appear in the trajectory file. See Figure 5.

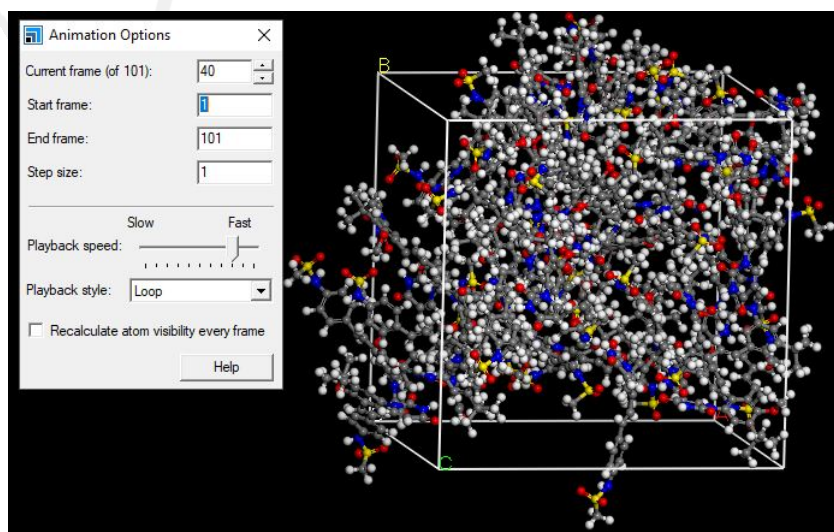


Figure 5: The 40th frame of the simulation.

The selected frame can now be exported by navigating to File > Export. Save the file as a .xsd

file in a new folder with a unique name. Repeat the export process for frames 50, 60, 70, 80, 90, and 100.

Repeat the process described above for each replicate simulation that was run, if any. For this example, the simulation had been run in triplicate. This resulted in 21 individual .xsd files.

3 Run the MATLAB Code

Once the frames (.xsd files) have been saved to the folder of your choice, the MATLAB code can be run.

IMPORTANT: The user must have [OpenBabel](#) installed for the code to run successfully. This software is called within the MATLAB code to convert from an .xyz file to an .sdf file.

The primary function is 'GetAtomCoords.m'. When run, this function will prompt the user to select all appropriate files. Hold shift or ctrl and click to select multiple files when prompted. The function will then find the appropriate information within the file (such as atom coordinates), identify individual molecules, and subsequently calculate R3m for each of the molecules. In this example, the result will be 840 values for R3m which correspond to a distribution of conformations of the molecule of interest.

To run the function, open MATLAB and type 'GetAtomCoords' or '[outputvarname] = GetAtomCoords' into the MATLAB workspace. The user will be prompted to select the files. Once the files are selected, press OK. A dialog box will appear that will show the progress of the calculation. See Figure 6. The calculation may take several minutes to complete.

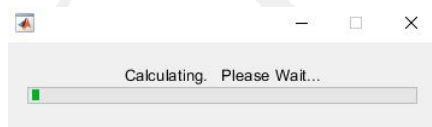


Figure 6: Once the code is running, a dialog will appear that will give the current progress of the calculation.

When the calculation is finished, a histogram plot of the distribution of R3m values will appear. This file will automatically be saved to the current directory as 'R3m_histogram.jpg' with a 600 dpi resolution. If unsure what the current directory is, type 'pwd' into the MATLAB workspace. An example of the resulting plot is shown in Figure 7. The user should change the file name before running the code again, since any file with the same name will be overwritten.

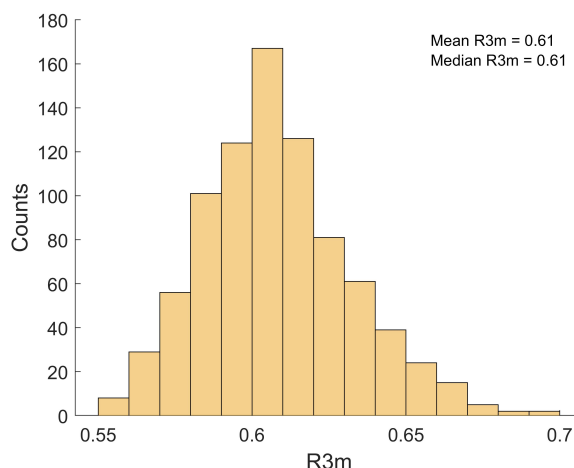


Figure 7: An example of the histogram that results from running the MATLAB code.

A matrix of R3m values will also be output to the MATLAB workspace as either 'ans' or the user

defined output variable name. It is recommended that the user rename the matrix as the drug name and save the workspace for future use. This way, the user will not need to re-run the code in the future.

A series of .xyz and .sdf files will also be output to the current working directory. These files will be all of the molecules from the final step in the code. In this case, a total of 40 of each file was output. Only the data from the last frame is output since these files are overwritten in each loop. The user can open these files in software such as [Mercury](#). Example output conformations can be seen in Figure 8.

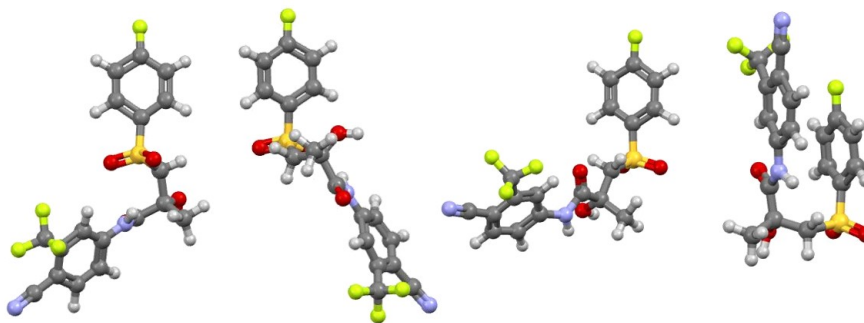


Figure 8: Example output conformations.

4 MATLAB code

The MATLAB code used to perform the calculations described above is included below. The code can be obtained online at the MathWorks file exchange via this [link](#).

4.1 GetAtomCoords.m

```

1 function [r3m_out] = GetAtomCoords()
2 %GetAtomCoords_multi - This function allows the user to select data
   exported from
3 %multiple Materials Studio files (.xsd files converted to .txt format)
   and
4 %calculate R3m values for each molecule in each frame. A distribution
   of
5 %R3m values is then plotted as histogram.
6 % This function allows the user to select text files to
7 % determine the distribution of R3m values that result from
8 % conformational differences in molecular structures as determined
   from
9 % Molecular Dynamics Simulations performed using Materials Studio.
   The
10 % user must export .xsd files from the MD run at multiple frames and
11 % then save the files in .txt format. This function can then be used.
12 %
13 % See also: R3mCalculate_auto, AtomConnection, AtomicWeighting,
14 % EuclidDistance, InfluenceDistanceMat, MolecInfluenceMatrix,
   MolecMatrix
15 %
16 % Code Author: Kevin DeBoyace
17 % Duquesne University
18 %
19 % Last Updated: January 2019
20 %
21 % NOTE: The Materials Studio data must be saved in .txt format. To
22 % accomplish this:

```

```

23 %      (1) Export the trajectory file from Materials Studio
24 %      (2) Open the file in notepad and save in .txt format
25 %      (alternatively , run a script from the command line to batch
convert
26 %      the files).
27 %      (3) Move the file to the appropriate directory (e.g. your
MATLAB
28 %      directory
29 %      (4) Run this function
30 %
31 %
32 %% IMPORT MULTIPLE FILES
33 [files_in , path_in] = uigetfile('*.txt','Select the INPUT DATA FILE(s)'
, 'Multiselect', 'on');
34
35 %% Wait bar
36 w = waitbar(0,'Initiating... ');
37
38 %% FOR 1 FILE SELECTED
39 if isstr(files_in) == 1;
40
41     file_char = ([path_in files_in]); % Import multiple files
42     filename = file_char;
43     filecontent = fileread(filename);
44
45     %% Wait bar
46     w = waitbar(0,'Wait ');
47
48     %% Set number of molecules in system
49     num_mol = 40;
50
51     %% Get Data
52     data1 = regexp(filecontent , '<Atom3d ID
= .....
', 'Match')';
53     data2 = regexp(filecontent , '<Bond ID
= .....
Match')';
54
55     % Atom ID
56     atom_id = extractBetween(data1 , '<Atom3d ID=" ', '"');
57     atom_id = str2double(atom_id); % convert to numeric
58
59     % Molecule ID
60     molecule_id = extractBetween(data1 , ' Parent=" ', '"');
61     molecule_id = str2double(molecule_id); % convert to numeric
62
63     min_mol = min(molecule_id);
64     max_mol = max(molecule_id);
65
66     % Coordinates
67     xyz = extractBetween(data1 , ' XYZ=" ', '"');
68     coords = regexp(xyz , ', ', 'split');
69     coords = vertcat(coords{:});
70     coords = str2double(coords); % convert to numeric
71     clear xyz
72
73     % Get Box limits

```

```

74     max_x = max(coords(:,1));
75     min_x = min(coords(:,1));
76     max_y = max(coords(:,2));
77     min_y = min(coords(:,2));
78     max_z = max(coords(:,3));
79     min_z = min(coords(:,3));
80
81     % Get Cell Parameters
82     data3 = regexp(filecontent, 'AVector
            =.....', 'Match')';
83     cell_param = extractBetween(data3, '"', ',0,0"');
84     cell_param = str2double(cell_param);
85
86     % Get Bond ID
87     forBond = regexp(filecontent, '<Bond ID=.....', 'Match')';
88     bond_id = extractBetween(forBond, '<Bond ID=', '');
89     clear forBond
90
91     % Get Connections
92     connections = extractBetween(data2, 'Connects="','');
93     connect = regexp(connections, ',', 'split'); %column 1 = x, column
            2 = y, column 3 = z
94     connect = vertcat(connect{:});
95     connect = str2double(connect); % convert to numeric
96     num_connect = size(connect,1)/num_mol;
97
98     % Get Atoms
99     forAtom = regexp(filecontent, 'Components=.....', 'Match')';
100    atom = extractBetween(forAtom, 'Components="','');
101    clear forAtom
102
103    connect_actual = 1:num_connect;
104
105    minMol = min(molecule_id);
106    maxMol = max(molecule_id);
107
108    % Save to .xyz format
109    for hh = minMol:maxMol;
110        waitbar(hh/maxMol, w, 'Calculating...');
111        mol1 = find(molecule_id == hh);
112        atom_id_1 = atom_id(mol1);
113        coords_1 = coords(mol1,:);
114        connect_1 = connect(connect_actual,:);
115        connect_actual = connect_actual(end):connect_actual(end)+
            num_connect;
116        atom_1 = atom(mol1);
117        for ii = 1:size(atom_id_1,1);
118            % ii = atom row of interest
119            [row, ~] = find(atom_id_1(ii) == connect_1); % row in
                connect where atom is found (column 1)
120            connected = connect_1(row,:);
121            % Get coordintes for relevant atoms
122            for jj = 1:size(connected,1);
123                for kk = 1:size(connected,2);
124                    coord_find(jj, kk) = {coords_1(find(atom_id_1==
                        connected(jj, kk)),:)};
125                end
126            end

```

```

127     % XYZ format
128     % number of elements
129     % comment line
130     % <elemente> <X> <Y> <Z>
131     % ...
132
133     % First line: num_atoms;
134     % comment line
135     % atom_1; coords_1;
136     num_atoms = size(atom_1,1);
137     atom_label = char(atom_1);
138     %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
139     % NOTE: data normalized to 1 (i.e. relative distance is the
140           same, but overall
141           % distance is no longer representative. Must convert back
142           to Angstroms
143           % HOW? —> USE CELL VOLUME
144     coords_conv = coords_1*cell_param; % Use cell parameter to
145           convert back to angstroms
146     %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
147     file_xyz = [ 'testOut', num2str(hh), '.xyz' ];
148     fileID = fopen(file_xyz, 'wt');
149     fprintf(fileID, '%d\n', num_atoms);
150     fprintf(fileID, 'NAME\n');
151     for ii = 1:size(atom_1,1);
152         fprintf(fileID, '%s\t', '%i\t%i\t%i\t\n', atom_label
153             , coords_1);
154         fprintf(fileID, '%s\t%i\t%i\t%i\t\n', atom_label(ii
155             ,:), coords_conv(ii,:));
156     end
157     fclose(fileID);
158     clear fileID
159 end
160 %% Convert from xyz to sdf to allow calculation of R3m
161 % Run Open Babel ot convert
162 % see https://openbabel.org/docs/dev/Command-line_tools/babel.
163     html
164 % -c = center coordinates
165 % -o = specify output format (e.g. -osdf —> output to sdf)
166 copyfile([ 'testOut', num2str(hh), '.xyz' ], 'testOut.xyz'); %Copy
167     file to run in Babel
168 %!obabel testOut.xyz -c -O Out.sdf
169 [yy, zz] = system('obabel testOut.xyz -c -O Out.sdf');
170 clear yy zz
171 copyfile('Out.sdf', [ 'testOut', num2str(hh), '.sdf' ]); % Copy
172     output file from Babel to new file name
173 %% Calculate R3m
174 pathname = cd;
175 file_sdf = [ 'testOut', num2str(hh), '.sdf' ];
176 filename = [pathname, '\', file_sdf]; % stitch together for full
177     path
178 [r3m] = R3mCalculate_auto(filename);
179 r3m_out((hh-minMol+1),1) = r3m; % Output r3m value to vector
180
181
182
183 end
184
185 %% FOR MULTIPLE FILES SELECTED

```



```

176 else
177     files_in = files_in';
178     for aa = 1:size(files_in,1);
179
180         file_char = ([path_in files_in{aa}]); % Import multiple files
181         filename = file_char;
182         filecontent = fileread(filename);
183
184         % Set number of molecules in system
185         num_mol = 40;
186
187         % Get Data
188         data1 = regexp(filecontent, '<Atom3d ID
            = .....', 'Match');
189         data2 = regexp(filecontent, '<Bond ID
            = .....', 'Match');
190
191         % Get Atom ID
192         atom_id = extractBetween(data1, '<Atom3d ID=', '');
193         atom_id = str2double(atom_id); % convert to numeric
194
195         % Get Molecule ID
196         molecule_id = extractBetween(data1, ' Parent=', '');
197         molecule_id = str2double(molecule_id); % convert to numeric
198
199         % Get Coordinates
200         xyz = extractBetween(data1, ' XYZ=', '');
201         coords = regexp(xyz, ',', 'split');
202         coords = vertcat(coords{:});
203         coords = str2double(coords); % convert to numeric
204         clear xyz
205
206         % Get Box limits
207         max_x = max(coords(:,1));
208         min_x = min(coords(:,1));
209         max_y = max(coords(:,2));
210         min_y = min(coords(:,2));
211         max_z = max(coords(:,3));
212         min_z = min(coords(:,3));
213
214         % Get Cell Parameters
215         data3 = regexp(filecontent, 'AVector
            = .....', 'Match');
216         cell_param = extractBetween(data3, '"', ',0,0"');
217         cell_param = str2double(cell_param);
218
219         % Get Bond ID
220         forBond = regexp(filecontent, '<Bond ID = .....', 'Match');
221         bond_id = extractBetween(forBond, '<Bond ID=', '');
222         clear forBond
223
224         % Get Connections
225         connections = extractBetween(data2, 'Connects=', '');
226         connect = regexp(connections, ',', 'split'); %column 1 = x,
            column 2 = y, column 3 = z
227         connect = vertcat(connect{:});

```

```

228     connect = str2double(connect); % convert to numeric
229     num_connect = size(connect,1)/num_mol;
230
231     %NOTE: Connections appear to refer to bond_id
232
233     % Get Atoms
234     forAtom = regexp(filecontent, 'Components=.....', 'Match');
235     atom = extractBetween(forAtom, 'Components="', '"');
236     clear forAtom
237
238     % Select only relevent connections (b/c numatoms ~= num
        connections)
239     connect_actual = 1:num_connect;
240
241     % Identify min and max molecule id numbers for iteration
242     minMol = min(molecule_id);
243     maxMol = max(molecule_id);
244     allMol = unique(molecule_id); % identify all unique values in '
        molecule_id'
245
246     %% Save to .xyz format
247     for hh = 1:size(allMol,1)
248         mol1 = find(molecule_id == allMol(hh));
249         atom_id_1 = atom_id(mol1);
250         coords_1 = coords(mol1,:);
251         connect_1 = connect(connect_actual,:);
252         connect_actual = connect_actual(end):connect_actual(end)+
            num_connect;
253         atom_1 = atom(mol1);
254
255         for ii = 1:size(atom_id_1,1);
256             [row, ~] = find(atom_id_1(ii) == connect_1); % row in
                connect where atom is found (column 1)
257             connected = connect_1(row,:);
258             % Get coordintes for relevant atoms
259             for jj = 1:size(connected,1);
260                 for kk = 1:size(connected,2);
261                     coord_find(jj, kk) = {coords_1(find(atom_id_1==
                        connected(jj, kk)), :)};
262                 end
263             end
264
265             % XYZ format:
266             % number of elements
267             % comment line
268             % <elemente> <X> <Y> <Z>
269             % ...
270
271             num_atoms = size(atom_1,1);
272             atom_label = char(atom_1);
273
274             %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
275             % NOTE: data normalized to 1 (i.e. relative distance is
                the same, but overall
276             % distance is no longer representative. Must convert
                back to
277             % Angstroms.
278             % HOW? —> use cell volume

```

```

279         coords_conv = coords_1*cell_param; % Use cell parameter
           to undo normalization and covert to Angstroms
280         %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
281         file_xyz = [ 'Out', num2str(hh), '.xyz' ];
282
283         fileID = fopen(file_xyz, 'wt');
284         fprintf(fileID, '%d\n', num_atoms);
285
286         fprintf(fileID, 'NAME\n');
287
288         for ii = 1:size(atom_1,1);
289             fprintf(fileID, '%s\t%i\t%i\t%i\t%i\n', atom_label(
                ii,:), coords_conv(ii,:));
290         end
291         fclose(fileID);
292
293         clear fileID
294
295     end
296
297     %% Convert from xyz to sdf to allow calculation of R3m
298     % Run Open Babel to convert
299     % see https://openbabel.org/docs/dev/Command-line_tools/
        babel.html
300     % -c = center coordinates
301     % -o = specify output format (e.g. -osdf --> output to sdf)
302
303     copyfile([ 'Out', num2str(hh), '.xyz' ], 'Out.xyz'); %Copy file
        to run in Babel
304
305     %! obabel Out.xyz -c -O Out.sdf
306     [yy, zz] = system('obabel Out.xyz -c -O Out.sdf'); % Use
        OpenBabel to convert from .xyz to .sdf
307     clear yy zz
308
309     copyfile('Out.sdf', [ 'Out', num2str(hh), '.sdf' ]); % Copy
        output file from Babel to new file name
310
311     %% Calculate R3m
312     pathname = cd;
313     %file taken in by function
314
315     file_sdf = [ 'Out', num2str(hh), '.sdf' ];
316     filename = [pathname, '\', file_sdf];
317     % Call other function: R3mCalculate_aut
318     [r3m] = R3mCalculate_auto(filename); % Calculate R3m for
        each molecule
319
320     % r3m_out((hh-minMol+1),aa) = r3m;
321     r3m_out(hh,aa) = r3m;
322
323     end
324
325     waitbar(aa/(size(files_in,1)), w, 'Calculating. Please Wait...');
        % update waitbar
326
327     end
328

```

```

329 end
330 close(w) % close waitbar
331
332 %% Plot a histogram of calculated R3m values
333 col = rand(1,3); % Random color
334 figure;
335 hold on
336 histogram(r3m_out, 'FaceColor', col, 'FaceAlpha', 0.5); % Plot
    histogram
337 mean_r3m = mean(mean(r3m_out)); % Calculate mean
338 median_r3m = median(median(r3m_out)); % Calculate median
339 set(gca, 'FontSize', 12);
340 % Display Mean on plot
341 text(0.7,0.95,['Mean R3m = ',num2str(round(mean_r3m,2))], 'Units', '
    Normalized', 'FontSize', 10);
342 % Display Median on plot
343 text(0.7,0.9,['Median R3m = ',num2str(round(median_r3m,2))], 'Units', '
    Normalized', 'FontSize', 10);
344 ylabel('Counts', 'FontSize', 14);
345 xlabel('R3m', 'FontSize', 14);
346
347 print(gcf, 'R3m_histogram', '-djpeg', '-r600') % Save histogram as jpg
    at 600 dpi
348
349 end

```

4.2 R3mCalculate_auto.m

```

1 function [R3m] = R3mCalculate_auto(filename)
2 %R3mCalculate - A function to calculate the R3m molecular descriptor.
3 % This function calculates the molecular matrix (M), geometry matrix
    (G),
4 % molecular influence matrix (H), and influence distance matrix (R).
    This
5 % is accomplished by calling the following subfunctions: ImportSDF,
6 % MolecMatrix, EuclidDistance, MolecInfluenceMatrix, and
7 % InfluenceDistanceMat. A GUI is called to select the file type to
    import
8 % (filetypegui). Next, the molecule is plotted using the PlotMolecule
9 % function to ensure the 3D coordinates are reasonable. Finally, the
    R3m
10 % value is calculated.
11 %
12 % See also: IMPORTSDF, MOLECMATRIX, EUCLIDDISTANCE,
    MOLECINFLUENCEMATRIX,
13 % INFLUENCEDISTANCEMAT, PLOTMOLECULE, ATOMICWEIGHTING, ATOMCONNECTION
14 %
15 % References:
16 % [1] Todeschini R, Consonni V. 2008. Handbook of molecular
17 % descriptors. ed.: John Wiley & Sons.
18 % [2] Consonni V, Todeschini R, Pavan M, Gramatica P 2002.
19 % Structure/response correlations and similarity/diversity
20 % analysis by GETAWAY descriptors. 2. Application of the
    novel
21 % 3D molecular descriptors to QSAR/QSPR studies. Journal of
22 % chemical information and computer sciences 42(3):693-705.
23 %
24 %

```

```

25 % Code Author: Kevin DeBoyace
26 % Duquesne University
27 %
28 % Date: October 2018
29
30 %% Matrices
31 % M – Molecular Matrix
32 % G – Geometry Matrix
33 % H – Molecular Influence Matrix
34 % R – Influence Distance Matrix
35
36 % Molecular Matrix
37 [ M, Atoms, numatoms, Connectivity] = MolecMatrix( filename );
38
39 % Atomic weighting
40 [ weightedmass, MolecWeight, colorset ] = AtomicWeighting( Atoms );
41 % Geometry Matrix
42 [G] = EuclidDistance(M);
43 % Molecular Influence Matrix
44 [ H, Leverage ] = MolecInfluenceMatrix( M );
45 % Influence Distance Matrix
46 [ R ] = InfluenceDistanceMat( Atoms, Leverage, G );
47
48 %% Atom connections
49 % Find atoms 1, 2 and 3 bond distances away.
50 [ connections, topology_mat ] = AtomConnection( Connectivity, Atoms );
51
52 %% R3m calculation
53 %  $R3m = \sum(i \text{ to } A-1) \sum(j>i) R*wi*wj*\delta(k;dij)$   $k = 1, 2, \dots, d$ 
54 R3m = 0;
55 for ii = 1:size(Atoms,1)-1;
56     for jj = ii+1:size(Atoms,1);
57         if topology_mat(ii,jj) == 1;
58             calc = R(ii,jj)*weightedmass(ii)*weightedmass(jj);
59             R3m = calc + R3m;
60         else
61             end
62     end
63 end
64
65 % R3m
66
67 end

```

4.3 AtomConnection.m

```

1 function [ connections, topology_mat ] = AtomConnection( Connectivity,
    Atoms )
2 %AtomConnection – A function which determines which atoms are connected
    by
3 %distances of 1, 2 and 3 bond distances.
4 % This function takes in the connectivity table (without the bond
    order
5 % column), and outputs a structure containing cell arrays, where each
6 % column corresponds to the atom number, and each row corresponds to
    the
7 % number of bond distances (e.g. col 2, row 3 represents the atoms
    which

```

```

8 % are 3 bond distance away from atom 2.)
9 %
10 % See also: ImportSDF, PlotMolecule, R3mCalculate
11 %
12 % Author: Kevin DeBoyace & Shikhar Mohan
13 % Updated: Jan 2019
14
15
16 connect = Connectivity(:,1:2);
17
18 for i = 1:size(Atoms,1)
19     B = find(connect(:,1) == i);
20     C = find(connect(:,2) == i);
21     Temp.A{1,i} = [connect(B,:);connect(C,:)];
22     Temp.A{1,i} = reshape(Temp.A{1,i},1,[]);
23     for j = 1:size(Temp.A{1,i},2);
24         if Temp.A{1,i}(j) == i
25             Temp.A{1,i}(j) = 0; %Replace where row = value with zeros (
                %Itself present in its own row)
26         else
27             end
28     end
29     %Delete zeros
30     test = Temp.A{1,i};
31     test(test==0) = [];
32     Temp.A{1,i} = test;
33 end
34 % remove empty cells
35
36 clear B C i j test
37
38
39 % Working 2 bond connection – Need to delete atom self reference
40
41 for i=1:size(Atoms,1)
42     for jj = 1:size(Temp.A{1,i},2)
43         B = find(connect(:,1) == Temp.A{1,i}(jj));
44         C = find(connect(:,2) == Temp.A{1,i}(jj));
45         for kk = 1:size(connect(B,:),2)
46             D{i, jj} = reshape([connect(B,:);connect(C,:)], 1, []);
47         end
48     end
49 end
50
51 % Place values into Temp structure
52 for i = 1:size(D,1);
53     Temp.A{2,i} = [D{i,:}];
54 end
55 clear B C D E F i ii j jj k kk
56
57 % Delete atoms which are 1 bond distance away.
58 for i = 1:size(Temp.A,2)
59     for k = 1:size(Temp.A{1,i},2)
60         [~, col] = find(Temp.A{1,i}(k) == Temp.A{2,i});
61         Temp.A{2,i}(col) = [];
62     end
63 end
64 clear i k col row

```

```

65
66 % Delete where column = atom number (where atom itself is included in
    bond
67 % connections)
68
69 for i = 1:size(Temp.A,2)
70     if isempty(Temp.A{2,i}) == 0; % NEW
71         a = find(Temp.A{2,i} == i);
72         Temp.A{2,i}(a) = [];
73     elseif isempty(Temp.A{2,i}) == 1; % NEW
74         Temp.A{2,i}(a) = []; % NEW
75
76 end
77
78
79 % 3 Bond Distances
80
81 for i=1:size(Atoms,1) % NEW
82     for jj = 1:size(Temp.A{2,i},2)
83         if isempty(Temp.A{2,i}) == 0;
84             B = find(connect(:,1) == Temp.A{2,i}(jj));
85             C = find(connect(:,2) == Temp.A{2,i}(jj));
86             for kk = 1:size(connect(B,:),2)
87                 if connect(B,kk) ~= i & connect(B,kk) ~= jj
88                     D{i, jj} = reshape([connect(B,:); connect(C,:)], 1,
                                           []);
89                 else
90                     end
91                 end
92             else
93                 end
94             end
95         end
96
97
98 % Place values into Temp structure
99 for i = 1:size(D,1);
100     Temp.A{3,i} = [D{i, :}];
101 end
102 clear B C D E F i ii j jj k kk a
103
104
105 % Delete atoms which are 1 bond distance away.
106 for i = 1:size(Temp.A,2)
107     for k = 1:size(Temp.A{1,i},2)
108         [~, col] = find(Temp.A{1,i}(k) == Temp.A{3,i});
109         Temp.A{3,i}(col) = [];
110     end
111 end
112 clear i k col row
113
114
115 % Delete atoms which are 2 bond distance away.
116 for i = 1:size(Temp.A,2)
117     for k = 1:size(Temp.A{2,i},2)
118         [~, col] = find(Temp.A{2,i}(k) == Temp.A{3,i});
119         Temp.A{3,i}(col) = [];
120

```

```

121     end
122 end
123 clear i k col row
124
125 % Delete duplicate values and sort values
126 for ii = 1:size(Temp.A,1)
127     for jj = 1:size(Temp.A, 2)
128         Temp.A{ii,jj} = unique(Temp.A{ii,jj}); % Keeps only unique
            values.
129     end
130 end
131
132 connections = Temp;
133
134
135 % Topological distance matrix
136 topology_mat = zeros(size(Connectivity,1)); %Preallocate matrix
137
138     for ii = 1:size(Atoms,1)
139         for jj = 1:size(connections.A{3,ii},2)
140             %try
141             if isempty(connections.A{3,ii}) == 0;
142                 topology_mat(ii,connections.A{3,ii}(jj)) = 1;
143                 topology_mat(connections.A{3,ii}(jj), ii) = 1;
144             else
145                 end
146             end
147         end
148     end
149
150
151 end

```

4.4 AtomicWeighting.m

```

1 function [ weightedmass , MolecWeight , colorset ] = AtomicWeighting(
    Atoms )
2 %AtomicWeigthing A function to normalize atomic masses to carbon
3 % This function is used to determine the appropriate atomic mass
4 % weighting for the calculation of R3m. In the calculation , the
    atomic
5 % mass is normalized to carbon.
6 % This function also assigns a color to each atom for plotting (see
7 % PlotMolecule) and calculates the molecular weight.
8 %
9 % Author: Kevin DeBoyace
10 % Updated: Jan 2019
11 %
12 % See also: R3mCalculate , PlotMolecule
13
14 atomicweight = zeros(size(Atoms,1),1);
15
16 % Check if molecule contains Hydrogens. If not, display warning message
    .
17 ContainHydrogen = strcmp(Atoms, 'H');
18 if any(ContainHydrogen) ~= 1;
19     h = warndlg('WARNING: This molecule is missing Hydrogens. Hydrogens
        should be added to accurately calculate molecular descriptors.

```



```

    , 'WARNING');
20     waitfor(h) %Pause code until user closes warning message.
21 else
22 end
23 % Calculate molecular weight and select atom colors for those Atoms
    which
24 % are present.
25 for b = 1:(size(Atoms,1));
26     if strcmp(Atoms(b,:), 'O ') || strcmp(Atoms(b,:), 'O') == 1;
27         atomicweight(b) = 15.9994;
28         colorset(b,:) = [1 0 0];
29     elseif strcmp(Atoms(b,:), 'C ') || strcmp(Atoms(b,:), 'C') == 1;
30         atomicweight(b) = 12.0107;
31         colorset(b,:) = [0 0 0];
32     elseif strcmp(Atoms(b,:), 'N ') || strcmp(Atoms(b,:), 'N') == 1;
33         atomicweight(b) = 14.0067;
34         colorset(b,:) = [0 0 1];
35     elseif strcmp(Atoms(b,:), 'Cl ') || strcmp(Atoms(b,:), 'Cl ') == 1;
36         atomicweight(b) = 35.453;
37         colorset(b,:) = [0 0.3 0];
38     elseif strcmp(Atoms(b,:), 'S ') || strcmp(Atoms(b,:), 'S') == 1;
39         %atomicweight(b) = 32.065;
40         atomicweight(b) = 32.066;
41         colorset(b,:) = [1 1 0];
42     elseif strcmp(Atoms(b,:), 'P ') || strcmp(Atoms(b,:), 'P') == 1;
43         atomicweight(b) = 30.973762;
44         colorset(b,:) = [1 0.5 0];
45     elseif strcmp(Atoms(b,:), 'F ') || strcmp(Atoms(b,:), 'F') == 1;
46         atomicweight(b) = 18.998;
47         colorset(b,:) = [0.5 0.8 0.4];
48     elseif strcmp(Atoms(b,:), 'H ') || strcmp(Atoms(b,:), 'H') == 1;
49         atomicweight(b) = 1.00794;
50         colorset(b,:) = [0.9 0.9 0.9];
51     end
52 end
53
54 % Throw an error if an atom is not listed in the code above. If this
    error
55 % is thrown, the atom needs to be added to the portion of the code
    above to
56 % ensure correct calculation of Molecular weight and weightedmass.
57 for b = 1:(size(Atoms,1));
58     if atomicweight(b) == 0;
59         error('Error: Atom missing from list in code')
60     end
61 end
62
63 atomicweight = atomicweight';
64 MolecWeight = sum(atomicweight);
65
66 % Masses are weighted with respect to the carbon atom : mass/ mass
    carbon
67 weightedmass = atomicweight / 12.0107;
68 weightedmass = weightedmass';
69
70 end

```

4.5 EuclidDistance.m

```
1 function [G] = EuclidDistance(M)
2 %G - Determine Euclidean Distance of atoms in a molecule
3 % G = Geometric Distance Matrix
4 % This function determines the Euclidean distance between all of the
5 % atoms in the molecule and outputs the data to the matrix G, which
6 % is
7 % known as the geometric distance matrix.
8 % See also: R3mCalculate, CheckCentroid,
9 %
10 % Author: Kevin DeBoyace
11 % Wildfong Lab
12 % Duquesne University
13 % Updated: Jan 2019
14
15 %G = pdist2(M,M,'euclidean'); % Performs same as below, but is slower
16
17 for ii = 1:size(M,1);
18     for jj = 1:size(M,1);
19         % Calculate Euclidean Distance
20         G(ii,jj) = sqrt((M(ii,1)-M(jj,1)).^2 + (M(ii,2)-M(jj,2)).^2 + (M(ii,3)-M(jj,3)).^2);
21     end
22 end
23
24 end
```

4.6 InfluenceDistanceMat.m

```
1 function [ R ] = InfluenceDistanceMat( Atoms, Leverage, G )
2 %InfluenceDistanceMat (R) - Calculation of the Influence Distance
3 %Matrix
4 % (R)
5 % This function calculates the influence/ distance matrix (R). Its
6 % initial
7 % application was for the calculation of 'R3m', but may be used to
8 % calculate any of the R-indices. The function requires input of '
9 % Atoms',
10 % 'Leverage', and 'G' (the Molecular Geometry Matrix). These values
11 % come
12 % from the functions 'ImportSDF', 'MolecInfluenceMatrix', and
13 % 'EuclidDistance', respectively.
14 %
15 % See Also: ImportSDF, MolecInfluenceMatrix, EuclidDistance,
16 % R3mCalculate
17 %
18 % Author: Kevin DeBoyace
19 % Wildfong Lab
20 % Duquesne University
21 % Updated: Jan 2019
22
23
24 for ii = 1:size(Atoms,1);
25     for jj = ii:size(Atoms,1);
26         if jj == ii
27             R(ii,jj) = 0;
28         end
29     end
30 end
```

```

23         else
24             R(ii , jj ) = (sqrt (Leverage(ii)*Leverage(jj))/G(ii , jj));
25             R(jj , ii ) = (sqrt (Leverage(jj)*Leverage(ii))/G(jj , ii));
26         end
27     end
28 end
29
30 end

```

4.7 MolecInfluenceMatrix.m

```

1 function [ H, Leverage ] = MolecInfluenceMatrix( M )
2 % Molecular Influence Matrix
3 %   This function is used to calculate the molecular influence matrix (
4 %   M)
5 %   from the 3D cartesian coordinates of a molecule.
6 %
7 % See Also: MolecMatrix, R3mCalculate
8 %
9 % Author: Kevin DeBoyace
10 %         Wildfong Lab
11 %         Duquesne University
12 % Updated: Jan 2019
13
14 H = M*pinv(M'*M)*M';
15 Leverage = diag(H);
16 end

```

4.8 MolecMatrix.m

```

1 function [ M, Atoms, numatoms, Connectivity] = MolecMatrix( filename )
2 %MolecMatrix Summary: Function to calculate the molecular matrix (M)
3 % This function extracts coordinates from a .sdf file and generates the
4 % molecular matrix (M). This function also extracts the connectivity
   table
5 % and creates a separate matrix for this information. The connectivity
6 % table contains the bond information.
7 %
8 % See Also: R3mCalculate, ImportSDF
9 %
10 % Author: Kevin DeBoyace
11 %           Wildfong Lab
12 %           Duquesne University
13 % Updated: Jan 2019
14
15 %% Read columns of data according to format string.
16 % This call is based on the structure of the file used to generate this
17 % code. If an error occurs for a different file, try regenerating the
   code
18 % from the Import Tool.
19 delimiter = ' ';
20 startRow = 2;
21 formatSpec = '%s%s%s%s%s%s%s%s%s%s%s%s [%^\n\r]';
22 fileID = fopen(filename, 'r');
23
24 dataArray = textscan(fileID , formatSpec , 'Delimiter' , delimiter , '
MultipleDelimsAsOne' , true , 'ReturnOnError' , false);
```

```

25 fclose(fileID);
26 for ii = 1:size(dataArray{1,1},1)
27     for jj = 1:size(dataArray,2)
28         dataArraynum(ii,jj) = str2double(dataArray{1,jj}{ii});
29     end
30 end
31
32 %% Extract cartesian coordinates from file and place them in a single
    matrix.
33 dataArraynum_cut = dataArraynum(:,1:3); %Take first three columns
34 datarem = rem(dataArraynum_cut,1);
35 jj = 1;
36 for ii = 1:size(dataArraynum_cut, 1);
37     if sum(datarem(ii,:)) == 0 || isnan(datarem(ii,1))
38         temp_rem(ii) = 0;
39     else
40         temp_rem(ii) = 1;
41         coords(jj,:) = dataArraynum_cut(ii,:);
42         Atoms(jj,:) = dataArray{1,4}(ii); % Extract atom names
43     end
44     jj = jj+1;
45 end
46
47 for ii = 1:size(Atoms,1)
48     temp_emp(ii) = isempty(Atoms{ii});
49     try
50         temp(ii) = str2num(Atoms{ii});
51     catch
52         temp(ii) = 0;
53     end
54 end
55 [row,col] = find(temp ~= 0 | temp_emp == 1);
56 Atoms(col) = [];
57 coords(col,:) = [];
58
59 M = coords; % Matrix of cartesian coordinates
60 clear ii jj
61
62 numatoms = size(M,1);
63
64 %% Connectivity table
65 jj = 1;
66
67 for ii = 1:size(dataArraynum_cut,1)
68     if (sum(datarem(ii,:)) == 0 && ~isnan(datarem(ii,1)) && sum(
        dataArraynum_cut(ii,1:3)~=0) == 3 %For those numbers which are
        part of connectivity or NaN —> delete
69         e(jj,:) = dataArraynum_cut(ii,:); %Build matrix with
            connectivity data
70         jj = jj + 1;
71     else
72         jj = jj;
73     end
74 end
75
76 Connectivity = e;
77
78 clear ii jj d e f %Clear Temp variables

```

79

80 **end**

S.O.P.