

# CASSANDRA

## Computational Atomistic Simulation Software At Notre Dame for Research Advances

User Manual 1.0

Written by:

Edward J. Maginn, Jindal K. Shah, Eliseo MarinRimoldi,  
Sandip Khan, Neeraj Rai, Thomas Rosch, Andrew Paluch



## Preface and Disclaimer

Cassandra is free software: you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or (at your option) any later version.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

This user manual is distributed along with the Cassandra software to aid in setting up various input files required for carrying out a Cassandra Monte Carlo simulation. Every effort is made to release the most updated and complete version of the manual when a new version of the software is released. To report any inconsistencies, errors or missing information, or to suggest improvements, send email to Edward Maginn (ed@nd.edu).

# Acknowledgements

Support for this work was provided by a grant from the National Science Foundation entitled “SI2-SSE: Development of Cassandra, a General, Efficient and Parallel Monte Carlo Multiscale Modeling Software Platform for Materials Research”, grant number ACI-1339785.

Ed Maginn would like to acknowledge financial support from Sandia National Laboratory’s Computer Science Research Institute, which enabled him to take a research leave and lay the foundation for Cassandra in collaboration with Jindal Shah, who stayed behind at Notre Dame and helped keep the group going. The hospitality of Steve Plimpton and co-workers at Sandia is gratefully acknowledged.

Finally, we would also like to thank the Center for Research Computing at Notre Dame, which provided support, encouragement, and infrastructure to help bring Cassandra to life.

People who have contributed to Cassandra through algorithm development and / or writing code (to date) include:

- Ed Maginn
- Jindal Shah
- Tom Rosch
- Neeraj Rai
- Eliseo Marin
- Sandip Khan
- Andrew Paluch

Some legacy code was used in the creation of Cassandra, and the following former students are recognized for their work:

- David Eike
- Jim Larentzos
- Craig Powers

# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	Distribution . . . . .	10
<b>2</b>	<b>Force Field</b>	<b>11</b>
2.1	Bonds . . . . .	11
2.2	Angles . . . . .	11
2.3	Dihedrals . . . . .	12
2.4	Impropers . . . . .	12
2.5	Nonbonded . . . . .	12
2.5.1	Repulsion-Dispersion Interactions . . . . .	13
2.5.2	Electrostatics . . . . .	14
2.6	Cassandra Units . . . . .	14
<b>3</b>	<b>Cassandra Basics</b>	<b>17</b>
3.1	Flow Diagram . . . . .	17
3.2	Cassandra Simulation Setup . . . . .	17

3.3	Cassandra File Preparation . . . . .	19
3.3.1	MCF File . . . . .	19
3.3.2	Input File . . . . .	19
3.3.3	Fragment Library Generation . . . . .	19
3.4	Running a Simulation . . . . .	20
3.5	Restarting a Simulation . . . . .	20
3.6	Cassandra Output Files . . . . .	20
<b>4</b>	<b>Files Required to Run Cassandra</b>	<b>23</b>
4.1	Simulation Input File . . . . .	23
4.1.1	Run Name . . . . .	23
4.1.2	Simulation Type . . . . .	24
4.1.3	Number of species . . . . .	24
4.1.4	VDW Style . . . . .	25
4.1.5	Charge Style . . . . .	27
4.1.6	Intramolecular Scaling . . . . .	28
4.1.7	Mixing Rule . . . . .	29
4.1.8	Starting Seed . . . . .	29
4.1.9	Minimum Cutoff . . . . .	30
4.1.10	Pair Energy Storage . . . . .	30
4.1.11	Molecule Files . . . . .	31
4.1.12	Simulation Box . . . . .	31

<i>CONTENTS</i>	7
4.1.13 Temperature . . . . .	32
4.1.14 Pressure . . . . .	32
4.1.15 Fugacity . . . . .	32
4.1.16 Chemical Potential . . . . .	33
4.1.17 Move Probabilities . . . . .	33
4.1.18 Start Type . . . . .	39
4.1.19 Run Type . . . . .	41
4.1.20 Frequency . . . . .	42
4.1.21 Average . . . . .	43
4.1.22 Property Output . . . . .	44
4.1.23 Fragment Files . . . . .	45
4.1.24 File Info . . . . .	46
4.1.25 CBMC parameters . . . . .	46
4.2 MCF File . . . . .	47
4.2.1 Atom Info . . . . .	48
4.2.2 Bond Info . . . . .	50
4.2.3 Angle Info . . . . .	51
4.2.4 Dihedral Info . . . . .	52
4.2.5 Fragment Info . . . . .	53
4.2.6 Fragment Connectivity . . . . .	53
<b>5 Implementing the Metropolis Acceptance Criteria</b>	<b>55</b>

5.1	Canonical Monte Carlo . . . . .	56
5.1.1	Translating a Molecule . . . . .	57
5.1.2	Rotating a Molecule . . . . .	57
5.1.3	Regrowing a Molecule . . . . .	59
5.1.4	Canonical Partition Function . . . . .	60
5.2	Isothermal-Isobaric Monte Carlo . . . . .	61
5.2.1	Scaling the Volume . . . . .	62
5.3	Grand Canonical Monte Carlo . . . . .	62
5.3.1	Inserting a Molecule with Configurational Bias Monte Carlo . . . . .	64
5.3.2	Deleting a Molecule that was Inserted via Configurational Bias Monte Carlo . . . . .	69
5.3.3	Regrowing a Molecule with Configurational Bias Monte Carlo . . . . .	71
5.4	Multiphase Systems . . . . .	73
5.4.1	Gibbs Ensemble Monte Carlo . . . . .	73
5.5	Multicomponent Systems . . . . .	73
5.5.1	Reaction Ensemble Monte Carlo . . . . .	74
5.6	Appendix . . . . .	75
5.6.1	Inserting a Molecule Randomly . . . . .	75
5.6.2	Deleting a Molecule Inserted Randomly . . . . .	76



# Chapter 1

## Introduction

Cassandra is an open source Monte Carlo package capable of simulating any number of molecules composed of rings, chains, or both. It can be used to simulate compounds such as small organic molecules, oligomers, aqueous solutions and ionic liquids. It handles a standard “Class I”-type force field having fixed bond lengths, harmonic bond angles and improper angles, a CHARMM or OPLS-style dihedral potential, a Lennard-Jones 12-6 potential and fixed partial charges. It does *not* treat flexible bond lengths. Cassandra uses OpenMP parallelization and comes with a number of scripts, utilities and examples to help with simulation setup.

Cassandra is capable of simulating systems in the following ensembles:

- Canonical (NVT)
- Isothermal-Isobaric (NPT)
- Grand canonical ( $\mu$ VT)
- Constant volume Gibbs (NVT-Gibbs)
- Constant Pressure Gibbs (NPT- Gibbs)

## 1.1 Distribution

Cassandra is distributed as a tar file `cassandra.tar`. You can unpack the distribution by running the command

```
> tar -xvf cassandra.tar
```

Upon successful unpacking of the archive file, the Cassandra directory will have a number of subdirectories. Please refer to the README file in the main Cassandra directory for a detailed information on each of the subdirectories.

**Documentation** - Contains this user guide and a document showing how molecular connectivity files are generated.

**Examples** - Contains example input files and short simulations of various systems in the above ensembles.

**MCF** - Molecular connectivity files for a number of molecules. These can be used as the basis for generating your own MCF files for molecules of interest.

**Scripts** - Useful scripts to set up simulation input files.

**Src** - Cassandra source code.

## Chapter 2

# Force Field

### 2.1 Bonds

Cassandra is designed assuming all bond lengths are fixed. If you wish to utilize a force field developed with flexible bond lengths, we recommend that you either use the nominal or “equilibrium” bond lengths of the force field as the fixed bond lengths specified for a Cassandra simulation or carry out an energy minimization of the molecule with a package that treats flexible bond lengths and utilize the bond lengths obtained from the minimization.

### 2.2 Angles

Cassandra supports two types of bond angles:

- ‘fixed’ : The angle declared as fixed is not perturbed during the course of the simulation.
- ‘harmonic’ : The bond angle energy is calculated as

$$E_{\theta} = K_{\theta}(\theta - \theta_0)^2 \tag{2.1}$$

where the user must specify  $K_{\theta}$  and  $\theta_0$ . Note that a factor of 1/2 is **not used** in the energy calculation of a bond angle. Make sure you know how the force constant is defined in any force field you use.

## 2.3 Dihedrals

Cassandra can handle four different types of dihedral angles:

- ‘OPLS’: The functional form of the dihedral potential is

$$E_\phi = a_0 + a_1 (1 + \cos(\phi)) + a_2 (1 - \cos(2\phi)) + a_3 (1 + \cos(3\phi)) \quad (2.2)$$

where  $a_0$ ,  $a_1$ ,  $a_2$  and  $a_3$  are specified by the user.

- ‘CHARMM’: The functional form of the potential is

$$E_\phi = a_0 * (1 + \cos(a_1 * \phi - \delta)) \quad (2.3)$$

where  $a_0$ ,  $a_1$  and  $\delta$  are specified by the user.

- ‘harmonic’: The dihedral potential is of the form:

$$E_\phi = K_\phi (\phi - \phi_0)^2 \quad (2.4)$$

where  $K_\phi$  and  $\phi_0$  are specified by the user.

- ‘none’: There is no dihedral potential between the given atoms.

## 2.4 Impropers

Improper energy calculations can be carried out with the following two options:

- ‘none’: The improper energy is set to zero for the improper angle.
- ‘harmonic’: The following functional form is used to calculate the energy due to an improper angle

$$E_\psi = K_\psi (\psi - \psi_0)^2 \quad (2.5)$$

where  $K_\psi$  and  $\psi_0$  are specified by the user.

## 2.5 Nonbonded

The nonbonded interactions between two atoms  $i$  and  $j$  are due to repulsion-dispersion interactions and electrostatic interactions (if any).

### 2.5.1 Repulsion-Dispersion Interactions

The repulsion-dispersion interactions can take one of the following forms:

- Lennard-Jones 12-6 potential (LJ):

$$\mathcal{V}(r_{ij}) = 4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] \quad (2.6)$$

where  $\epsilon_{ij}$  and  $\sigma_{ij}$  are the energy and size parameters set by the user. For unlike interactions, different combining rules can be used, as described elsewhere. Note that this option only evaluates the energy up to a specified cutoff distance. As described below, analytic tail corrections to the pressure and energy can be specified to account for the finite cutoff distance.

- Cut and shift potential:

$$\mathcal{V}(r_{ij}) = 4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] - 4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{cut}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{cut}} \right)^6 \right] \quad (2.7)$$

where  $\epsilon_{ij}$  and  $\sigma_{ij}$  are the energy and size parameters set by the user and  $r_{cut}$  is the cutoff distance. This option forces the potential energy to be zero at the cutoff distance. For unlike interactions, different combining rules can be used, as described elsewhere.

- Cut and switch potential:

$$\mathcal{V}(r_{ij}) = 4\epsilon_{ij} \left[ \left( \frac{\sigma_{ij}}{r_{ij}} \right)^{12} - \left( \frac{\sigma_{ij}}{r_{ij}} \right)^6 \right] f \quad (2.8)$$

The factor  $f$  takes the following values:

$$f = \begin{cases} 1.0 & r_{ij} \leq r_{on} \\ \frac{(r_{off}^2 - r_{ij}^2)(r_{off}^2 - r_{on}^2 + 2r_{ij}^2)}{(r_{off}^2 - r_{on}^2)^3} & r_{on} < r_{ij} < r_{off} \\ 0.0 & r_{ij} \geq r_{off} \end{cases} \quad (2.9)$$

where  $\epsilon_{ij}$  and  $\sigma_{ij}$  are the energy and size parameters set by the user. This option smoothly forces the potential to go to zero at a distance  $r_{off}$ , and begins altering the potential at a distance of  $r_{on}$ . Both of these parameters must be specified by the user. For unlike interactions, different combining rules can be used, as described elsewhere.

- Tail corrections: If the Lennard-Jones potential is used, standard Lennard-Jones tail corrections are used to approximate the long range dispersion interactions

### 2.5.2 Electrostatics

Electrostatic interactions are given by Coulomb's law

$$\mathcal{V}_{elec}(r_{ij}) = \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r_{ij}}. \quad (2.10)$$

where  $q_i$  and  $q_j$  are partial charges specified by the user and placed on atomic positions given by  $r_i$  and  $r_j$ . In a simulation, the electrostatic interactions are calculated using either an Ewald summation or a direct summation using the minimum image convention.

Note that the total energy that is printed out in the property file is extensive. Consequently, to obtain intensive energies, the printed energies must be divided by the total number of molecules in the system.

## 2.6 Cassandra Units

The following table provides units used in Cassandra:

Table 2.1: Cassandra units for input variables

Bond length	$l$	$\text{\AA}$
<b>Angles</b>		
Nominal bond angle	$\theta_0$	degrees
Bond angle force constant	$K_\theta$	K/rad <sup>2</sup>
<b>Dihedral angle</b>		
OPLS	$a_0, a_1, a_2, a_3$	kJ/mol
CHARMM	$a_0$	kJ/mol
	$a_1$	dimensionless
	$\delta$	degrees
harmonic	$K_\phi$	K/rad <sup>2</sup>
	$\phi_0$	degrees
<b>Improper angle</b>		
Force constant	$K_\psi$	K/rad <sup>2</sup>
	$\psi_0$	degrees
<b>Nonbonded</b>		
Energy parameter	$\epsilon/k_B$	K
Collision diameter	$\sigma$	$\text{\AA}$
Charge	$q$	e
<b>Simulation Parameters</b>		
Simulation box length		$\text{\AA}$
Volume		$\text{\AA}^3$
Distances		$\text{\AA}$
Rotational width		degrees
Temperature		K
Pressure		bar
Chemical potential		kJ/mol
Energy		kJ/mol
Fugacity		bar





## Chapter 3

# Cassandra Basics

### 3.1 Flow Diagram

A flow diagram that overviews the setup for a Cassandra simulation is displayed on figure 3.1. This diagram employs two automation scripts located in the `/Scripts/` directory: `mcfgen.py` and `library_setup.py`. These scripts are particularly useful when simulating large molecules. For details about each step, please refer to the document `MCF_tutorial.pdf`, to the README files located in the subdirectories inside the directory `/Scripts/`, and to this document.

### 3.2 Cassandra Simulation Setup

Once a system is identified, setting up a Cassandra simulation from scratch requires preparation of the following files.

- A molecular connectivity file (MCF) (\*.mcf) containing the molecular connectivity information on bonds, angles, dihedrals, impropers and whether the molecule is composed of fragments. For information on the MCF file, please refer to section 4.2.
- An input file (\*.inp) (See section 4.1)

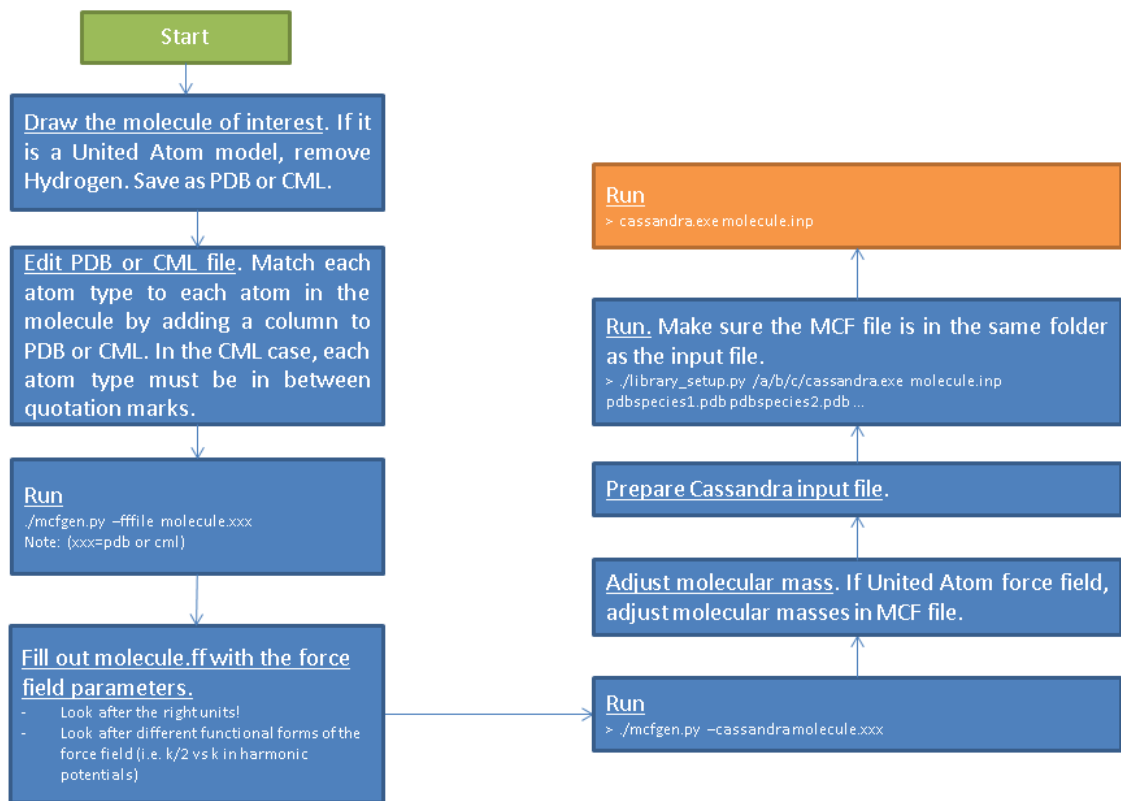


Figure 3.1: Flow diagram representing a typical setup of a Cassandra simulation

- If the molecule is composed of fragments, then a fragment library file for each of the fragments is required. For instructions on how to generate these files, please refer to the document `MCF_tutorial.pdf`.

MCF files for united-atom models of methane, isobutane, dimethylhexane, cyclohexane and diethylether are provided in the `MCF` directory. Input files for NVT, NPT, GCMC and GEMC ensembles are located in the `Examples` directory which also contains fragment library files for a number of molecules simulated in these ensembles.

## 3.3 Cassandra File Preparation

### 3.3.1 MCF File

One MCF file is required for each unique species in a simulation. A species is defined as a collection of atoms associated with each other through bonds. Thus a molecule is a species as is an ion. If you wanted to simulate sodium sulfate, you would need separate MCF files for the sodium ion and the sulfate ion. MCF files can be built “by hand” or by using the scripts provided with the code, as described in the [Documentation/MCF\\_tutorial.pdf](#). Instructions for generating an MCF file can also be found in the [Scripts/MCF\\_Generation/README](#) file. We will collect MCF files submitted to us by users and will post them on the Cassandra website <http://cassandra.nd.edu>. If you have an MCF file you would like us to post, send it to [ed@nd.edu](mailto:ed@nd.edu).

### 3.3.2 Input File

An input file is required for a Cassandra simulation. The input file specifies conditions for the simulation and various keywords required for the simulation in a given ensemble. Please refer to Chapter 4.1 for further details.

### 3.3.3 Fragment Library Generation

Cassandra makes use of reservoir sampling schemes to sample correctly and efficiently various coupled intramolecular degrees of freedom associated with branch points and rings. Details may be found in our publication [?]. The idea is to decompose the entire molecule into fragments that are either branch points or ring groups, each coupled to other fragments via a single dihedral angle. Thus, the total number of fragments of a molecule is the sum of branch points and ring groups in the molecule. The neighboring fragments are connected by two common atoms present in each of the fragments. Note that the ring group contains all the ring atoms and those directly bonded to the ring atoms. For each fragment identified, Cassandra runs a pre-simulation in the gas phase to sample the intramolecular degrees of freedom. A library of a large number of these conformations are stored for use in an actual simulation.

The gas phase library generation has been automated with the script `library_setup.py` located in the `Scripts/Frag_Library_Setup` directory. Use the following command for generating the fragment library.

```
> python path_to.Scripts/Frag_Library_Setup/library_setup.py cassandra.executable
```

```
input_filename mol1.pdb (mol1.cml) mol2.pdb (mol2.cml) ...
```

where `input_filename` is the name of the input file for the actual simulation and `mol1.pdb mol2.pdb ...` or `mol1.cml mol2.cml ...` correspond to the names of the pdb (or cml) files used to generate the MCF files. Make sure that if a file does not exist in the current working directory, its path relative to the current working directory is specified.

### 3.4 Running a Simulation

To launch a Cassandra simulation, run the following command:

```
> cassandra_executable input_filename
```

The executable will read `input_filename` and execute the instructions. Make sure that the required files (MCF, fragment library files) are located in the directories as given in the input file.

### 3.5 Restarting a Simulation

Restarting a simulation requires either a checkpoint file (\*.chk produced by Cassandra) or a configuration file obtained from xyz files generated from a previous simulation. For the set up of these simulations, a script in `Scripts/Read_Old` is provided. Detailed instructions are contained in the README file in this directory.

### 3.6 Cassandra Output Files

Cassandra generates several output files which can be used for later analysis. All have as a prefix the `Run_Name` specified in the input file. See Chapter 4.1 for details. The type of output is specified by the file name suffix. The following are generated:

- **Log file** (\*.log): Contains basic information on what the run is, timing information and reports the various parameters specified by the user. A complete copy of the input file is reproduced. Other important information includes the move acceptance rates. You can use the log file to keep track of what conditions were simulated.

- **Coordinate file** (\*.xyz): For each box in the system, a set of xyz coordinates are written out with a frequency specified by the user (Ncoordfreq). The file has as a header the number of atoms in the box. Following this, the atomic coordinates of molecule 1 of species 1 are written, then the coordinates of molecule 2 of species 1 are written, etc. After all the coordinates of the molecules of species 1 are written, the coordinates of the molecules of species 2 are written, etc. You can use this file to do all your structural analysis and post processing.
- **Checkpoint file** (\*.chk): A checkpoint file is written every Ncoordfreq steps. This can be used to restart a simulation from this point using all of the same information as the run that was used to generate the checkpoint file. To do this, you must use the checkpoint restart option (see Chapter 4.1). It will basically pick up where the simulation left off, using the same random number seed, maximum displacements, etc. This is useful in case your job crashes and you want to continue running a job. You can also use the checkpoint file to start a new simulation using the configuration of the checkpoint file as an initial configuration and the optimized maximum displacements. To do this, use the script read\_old.py. You will need to set a new random number seed if you do this. See the documentation in Chapter 4.1 for more details.
- **H-matrix file** (\*.H.box#): This file is written to every Ncoordfreq MC steps. The first line is the box volume in angstrom<sup>3</sup>. The next three lines are the box coordinates in angstrom in an H-matrix form. Since Cassandra only supports cubic boxes at the moment, this is just a diagonal and symmetric matrix, but is included here for later versions that will enable non-orthogonal boxes. After this, a blank line is written. The next line is the box number, and the final line(s) is(are) the species ID and number of molecules for that species in this box. If there are three species, there will be three lines. This output is repeated every Ncoordfreq times. This file allows you to compute the density of the box during constant pressure simulations.
- **Property file** (\*.instnt-prp#.box#): This file lists the instantaneous thermodynamic and state properties for each box. Note that you can have more than one property file (hence the# after 'prp') and more than one box (also why there is a # after 'box'). The user specifies which properties are to be written and in what order, and these are then reproduced in this file. The file is written to every Nthermofreq steps. A header is written to the first two lines to designate what each property is. You may use this file to compute thermodynamic averages.
- **Initial configuration** (\*.init\_config\_box#): If you generate your initial configuration using the make\_config command, this file will be created for each box. It contains the initial coordinates of all the species in the system. You can use this to check on whether the initial configuration is reasonable, or use it as an input to other codes. For example, the initial configuration will be generated using CBMC, so it may be a better starting configuration than if you used other methods.



## Chapter 4

# Files Required to Run Cassandra

### 4.1 Simulation Input File

This is a required file that is given as an argument to the Cassandra executable. You must generate this by hand, but you can use the input files in the Examples directory as a guide. The input file contains a number of keywords that define simulation parameters and thermodynamic state point for the simulation. A keyword is identified by a '#' while comments follow a '!'. Note that the order of the keywords in the input file is immaterial, but the format of the arguments of the keyword are important and are explained below, along with a complete listing of all keywords.

#### 4.1.1 Run Name

**# Run\_Name** - Simulation run name  
*Character\*120*

The run name is specified on the next line following the keyword. This name is used as a prefix for all the files produced by the simulation. For example,

```
# Run_Name
dee
```

Cassandra will then use **dee** as prefix for all output files created.

### 4.1.2 Simulation Type

**# Sim\_Type**

*Character\*120*

Sets the ensemble (and thus the suite of moves) of a Cassandra simulation. Currently the following ensembles are supported:

- NVT\_MC (Canonical ensemble)
- NVT\_MIN (Canonical ensemble in which minimization is carried out)
- NPT\_MC (Isothermal-isobaric ensemble)
- GCMC (Grand canonical ensemble)
- GEMC (Gibbs ensemble)
- NVT\_MC\_Fragment (Canonical ensemble simulation of a non-ring fragment)
- NVT\_MC\_Ring\_Fragment (Canonical ensemble simulation of a ring fragment)

Note that NVT\_MC\_Fragment and NVT\_MC\_Ring\_Fragment are used only for the fragment generation and are not used in the normal simulation. For example,

**# Sim\_Type**

NPT\_MC

will run an NPT MC simulation.

### 4.1.3 Number of species

**# Nbr\_Species**

*Integer*



Total number of species in the simulation. For ionic systems, each ion is counted as a separate species. For example, for a mixture of two species, use the following:

```
# Nbr_Species
2
```

#### 4.1.4 VDW Style

```
# VDW_Style
Character(i,1) Character(i,2) Real(i,3) Real(i,4)/Logical(i,4)
```

This keyword specifies the functional form of repulsion dispersion interactions to be used and if tail corrections are added for the box  $i$ . One line is required for each box. The options for  $Character(i,1)$  are “LJ” for a Lennard-Jones 12-6 potential or “None” if you wish to turn off all repulsion-dispersion interactions.  $Character(i,2)$  and  $Real(i,3)$  are specified only if  $Character(i,1)$  is set to “LJ”.  $Character(i,2)$  describes the truncation scheme used for the Lennard-Jones potential. Options are “cut”, “cut\_tail”, “cut\_switch” and “cut\_shift”. Refer to Chapter 2 for the functional forms. The other parameters  $Real(i,3)$  and  $Real(i,4)/Logical(i,4)$  depend on the selection of  $Character(i,2)$  as described below:

**cut:** This option cuts the potential off at the distance specified by  $(Real(i,3))$ . The fourth parameter is omitted.

For example, to simulate one box with a 14 Å cutoff specify the following:

```
# VDW_Style
LJ cut 14.0
```

Similarly, for a two box simulations such as used in the Gibbs ensemble where both boxes have a 14 Å cutoff, use the following:

```
# VDW_Style
LJ cut 14.0
LJ cut 14.0
```

**cut\_tail:** This options cuts the potential off at a distance corresponding to  $(Real(i,3))$  and applies analytic tail corrections to the energy and pressure. An optional fourth argument  $(Logical(i,4))$  can be set to 'TRUE' or 'true', in which case  $Real(i,3)$  is ignored and the cutoff distance is always set to half of the

simulation box length. The cutoff will change during the course of the simulation when attempting volume moves. This option is provided to enable reproduction of literature simulations that use a cut off distance of half the simulation box length, but its use is highly discouraged.

For example, to simulate one box with a 14 Å cutoff using tail corrections, specify the following:

```
# VDW_Style
LJ cut_tail 14.0
```

For a two box simulation where the first box has a 14 Å cutoff and the second one has a 20 Å cutoff, use the following:

```
# VDW_Style
LJ cut_tail 14.0
LJ cut_tail 20.0
```

**cut\_switch:** This option cuts the potential off and smoothly brings the potential to zero using a spline. The potential is cutoff and the spline turned on at a distance specified by  $Real(i,3)$  ( $r_{on}$  in Eq 2.8) and the potential goes to zero at a distance specified by  $Real(i,4)$  ( $r_{off}$  in Eq 2.8).

For example, a one box simulation using the cut\_switch option could be specified as follows:

```
# VDW_Style
LJ cut_switch 12.0 14.0
```

In this case, the Lennard-Jones potential would end at 12.0 Å and be smoothly taken to zero at 14.0 Å. Obviously,  $r_{on} < r_{off}$  or  $Real(i,3) < Real(i,4)$ .

**cut\_shift:** This option cuts the potential off at a distance specified by  $Real(i,3)$  and shifts the entire potential so that at this distance the potential is zero. The fourth parameter  $Real(i,4)/Logical(i,4)$  is omitted. The functional form of this potential is given in eq 2.7.

To perform a two box simulation with a cut\_shift option in which both boxes have a 10.5 Å cutoff, use the following:

```
# VDW_Style
LJ cut_shift 10.5
LJ cut_shift 10.5
```

**Note:** For all options, cutoff distances must be less than or equal to the shortest edge length of a simulation box.

#### 4.1.5 Charge Style

##### # Charge\_Style

*Character(i,1) Character(i,2) Real(i,3) Real(i,4)*

Cassandra allows the use of fixed partial charges on atomic centers using a Coulomb potential of the form given in Eq 2.10. To specify this option for box  $i$ , set *Character(i,1)* to “coul”. For this option, *Character(i,2)* can be set to either “Ewald” if you want to use an Ewald sum to compute Coulombic interactions or it can be set to “cut”, in which case the Coulombic interactions will be cut off and the long range interactions ignored. For the Ewald option, *Real(i,3)* is the real space cutoff distance and *Real(i,4)* specifies the accuracy of the Ewald summation. A reasonable value for the accuracy is  $10^{-5}$ . Note that the number of reciprocal vectors for the Ewald summation is determined in the code based on the accuracy parameter. For more details, see the paper by Fincham [?].

For example,

##### # Charge\_Style

```
coul Ewald 12.0 1E-5
```

will use the Ewald sum for a single box. The real space cutoff will be 12 Å and the accuracy will be  $10^{-5}$ . If you have two boxes, like in a Gibbs ensemble calculation, then you could use the following:

##### # Charge\_Style

```
coul Ewald 12.0 1E-5
```

```
coul Ewald 30.0 1E-5
```

This will use an Ewald sum for both boxes. In the first box, the real space cutoff will be 12 Å while in the second box a larger cutoff of 30 Å will be used. **Note: When performing Gibbs ensemble simulations of vapor-liquid equilibria, the vapor box is often much larger than the liquid box. In this case, you will want to use a longer real space cutoff for the larger vapor box to avoid using too many reciprocal space vectors.** Also note that the real space cutoffs must always be less than or equal to half of the shortest edge length of a simulation box.

If you do not wish to use a Coulomb potential (for example, your model has no partial charges), you still must specify `# Charge_Style`. In this case, set *Character(i,1)* to “NONE”. If “NONE” is selected for *Character(i,1)* then *Character(i,2)*, *Real(i,3)* and *Real(i,4)* are omitted.

For example,

```
# Charge_Style
NONE
```

should be used if you have no partial charges and are simulating a single box. A two box simulation with no partial charges would be specified as

```
# Charge_Style
NONE
NONE
```

**Note :** If the cutoff in *VDW\_Style* is set to half of the simulation box length, any cutoff distance specified in the *Charge\_Style* section will default to the half of the simulation box length. In the case of Ewald summation, however, the accuracy will be the same as *Real(i,4)*.

#### 4.1.6 Intramolecular Scaling

```
# Intra_Scaling
Real(i,1) Real(i,2) Real(i,3) Real(i,4)
Real(i,5) Real(i,6) Real(i,7) Real(i,8)
```

This keyword sets the intramolecular scaling for 1-2, 1-3, 1-4 and 1-N interactions within a given species. 1-2 means interactions between a given atom 1 and another atom 2 directly bonded to it, 1-3 means interactions between atom 1 and other atoms 3 separated from atom 1 by exactly two bonds, etc. The first line corresponds to the VDW scaling: *Real(i,1)* *Real(i,2)* *Real(i,3)* *Real(i,4)* apply to 1-2, 1-3, 1-4 and 1-N interactions, where N corresponds to all atoms separated from atom 1 by more than three bonds. The second line corresponds to the Coulomb scaling: *Real(i,5)* *Real(i,6)* *Real(i,7)* *Real(i,8)* apply to 1-2, 1-3, 1-4 and 1-N interactions. These lines are repeated for each species in the simulation. Note that intramolecular scaling applies to all the boxes in the simulation.

For example,

```
# Intra_Scaling
0.0 0.0 0.5 1.0
0.0 0.0 0.5 1.0
```

would turn off 1-2 and 1-3 interactions, would scale the VDW and Coulombic interactions for 1-4 atoms by 50% and would use full interactions for all other atom pairs in the species.

If you had two species in the simulation and wanted the same intramolecular scaling as above, you would specify

```
# Intra_Scaling
0.0 0.0 0.5 1.0
0.0 0.0 0.5 1.0
0.0 0.0 0.5 1.0
0.0 0.0 0.5 1.0
```

In the absence of the `# Intra_Scaling` keyword, default values of 0.0, 0.0, 0.5 and 1.0 will be used for 1-2, 1-3, 1-4 and 1-N for VDW and Coulomb interactions for all the species in the simulation. If `# Charge_Scaling` is set to “NONE”, you must still specify the `Intra_Scaling` for Coulombic interactions, but no interactions will be computed.

#### 4.1.7 Mixing Rule

```
# Mixing_Rule
Character
```

Sets the method by which Lennard-Jones interactions between unlike atoms are calculated. Acceptable options are “LB” for Lorentz-Berthelot and “geometric” for geometric. If this keyword is missing, “LB” is used as default.

#### 4.1.8 Starting Seed

```
# Seed_Info
Integer(1) Integer(2)
```

Inputs for the starting random number seeds for the simulation. Note that Cassandra uses a random

number generator proposed by L’Ecuyer [?], which takes five seeds to calculate a random number, out of which three are defined internally while two *Integer(1)* and *Integer(2)* are supplied by the user. **When a ‘checkpoint’ file is used to restart a simulation (see # Start\_Type below), the user supplied seeds will be overwritten by those present in the checkpoint file. If # Start\_Type is set to ‘read\_old’, then the seeds specified in the input file are used.**

As an example,

```
# Seed_Info
1244432 8263662
```

is an acceptable way of specifying the seeds. Note that two independent simulations can be run using the same input information if different seeds are used. If two simulations having exactly the same input information and the same seeds are run, the results will be identical.

#### 4.1.9 Minimum Cutoff

```
# Rcutoff_Low
Real(1)
```

Sets the minimum distance *Real(1)* in Å such that any MC move bringing two sites closer than this distance will be immediately rejected. It avoids numerical problems associated with random moves that happen to place atoms very close to one another such that they will have unphysically strong repulsion or attraction. This distance must be less than the intramolecular distance of all atoms in a species which are not bonded to one another. For models that use dummy sites without explicitly defining bonds between dummy and atomic sites of the molecules (for example, the TIP4P water model), it is important that the minimum distance is set to be less than the shortest distance between any two sites on the molecule. For most systems, 1 Å seems to work OK, but for models with dummy sites, a shorter value may be required.

#### 4.1.10 Pair Energy Storage

```
# Pair_Energy
Logical(1)
```

Cassandra can use a time saving feature in which the energies between molecules are stored and used during energy evaluations after a move, thereby saving a loop over all molecules. This requires more memory, but

it can be faster. The default is to not use this feature. If you wish to use this, set *Logical(1)* to 'TRUE' or 'true'.

#### 4.1.11 Molecule Files

##### **# Molecule.Files**

*Character(i,1) Integer(i,2)*

This specifies the name of the molecular connectivity file (\*.mcf) and the maximum total number of molecules of a given species specified by this MCF file. A separate line is required for each species present in the simulation. *Character(i,1)* is the name of the MCF file for species *i*. *Integer(i,2)* denotes the maximum number of molecules expected for the species.

For example

##### **# Molecule.Files**

butane.mcf 100

hexane.mcf 20

octane.mcf 5

specifies that there are three different species, and the MCF files state the names of the files where information on the three species can be found. Species 1 is butane, species 2 is hexane and species 3 is octane. There can be a maximum of 100 butane molecules, 20 hexane molecules and 5 octane molecules in the total system. The maximum number of molecules specified here will be used to allocate memory for each species, so do not use larger numbers than are needed.

#### 4.1.12 Simulation Box

##### **# Box\_Info**

*Integer(1)*

*Character(i,2)*

*Real(i,3)*

This keyword sets parameters for the simulation boxes. *Integer(1)* specifies the total number of boxes in the simulation. For now, Gibbs ensemble simulations must have two boxes. *Character(i,2)* is the shape of the  $i^{th}$  simulation box. Right now, the only supported box shape in Cassandra is cubic, so use the keyword

“CUBIC”.  $Real(i,3)$  is the length of the box edges for a cubic box in Å. Information for additional boxes is provided in an analogous fashion and is separated from the previous box by a blank line. For a two box simulation, box information is given as:

```
# Box_Info
```

```
2
```

```
CUBIC
```

```
30.0
```

```
CUBIC
```

```
60.0
```

This will construct two cubic boxes. The first will be 30 X 30 X 30 cubic Å and the second will be 60 X 60 X 60 cubic Å.

#### 4.1.13 Temperature

```
# Temperature_Info
```

```
Real(1) Real(2) ....
```

Sets the temperature in Kelvin  $Real(1) Real(2) \dots$  for simulation boxes 1, 2, ...

#### 4.1.14 Pressure

```
# Pressure_Info
```

```
Real(1) Real(2) ...
```

Specifies the pressure in bar  $Real(1) Real(2) \dots$  for simulation boxes 1, 2, ... Note that if the simulation type does not require an input pressure (for example, an NVT simulation), this command will be ignored.

#### 4.1.15 Fugacity

```
# Fugacity_Info
```

```
Real(1) Real(2) ...
```



Specifies the fugacities in bar *Real(1) Real(2) ...* of **insertable** species in the order in which they appear in the `Molecule_Files` section. The fugacity will be arbitrarily set to zero for species that cannot be swapped or exchanged with a reservoir. This option is used with grand canonical ensemble simulations. It will be ignored for all other simulation types.

#### 4.1.16 Chemical Potential

##### # Chemical\_Potential\_Info

*Real(1) Real(2) ....*

This keyword is only used for grand canonical simulations and may be used instead of fugacity. It sets the chemical potential of the insertable species in the order in which they appear in the `Molecule_Files` section. The chemical potentials will be set arbitrarily to zero for species that cannot be swapped or exchanged with a reservoir. Note that the de Broglie wavelength of the species in each box is automatically calculated and used in the acceptance rules. Units of chemical potential are kJ/mol. **Note: if you specify BOTH a fugacity AND a chemical potential, Cassandra will use whichever one appears first in the input file and will ignore the other keyword. We recommend that you do not list both keywords in any input file.**

#### 4.1.17 Move Probabilities

**# Move\_Probabilities** This section specifies the probabilities associated with different types of MC moves to be performed during the course of simulation. Please ensure that the move probabilities add up to 1.0. An error will be generated if this is not the case. All the headers are optional but an error will be produced if a move is required (for example, volume fluctuations in an NPT simulation) and the move is not specified.

##### Translation

##### # Prob\_Translation

*Real(1)*

*Real(i,j)* \*One line required for each box, and one value required for each species on each line.

*Real(1)* is the probability of performing a center of mass translation move. For each box *i*, the maximum displacement of species *j* is specified by *Real(i,j)*.

For example, if you have two species and two boxes, you would specify the translation probability as

```
# Prob_Translation
0.25
2.0 2.5
10.0 11.0
```

This will tell Cassandra to attempt center of mass translations 25% of the time. For box 1, the maximum displacement will be 2.0 Å for species 1 and 2.5 Å for species 2. For box 2, the maximum displacement for species 1 will be 10.0 Å and it will be 11.0 Å for species 2. Note that attempted moves will occur with equal probability for a give box, but attempts for a species are proportional to their mole fraction in the box.

## Rotation

### # Prob\_Rotation

*Real(1)*

*Real(i,j)* \*One line required for each box, and one value required for each species on each line.

The probability of performing a rotation move is specified by *Real(1)* while *Real(i,j)* denotes the maximum rotational width for species *j* in box *i* in degrees. If you are only simulating spherical molecules (such as Lennard-Jones particles), then do not use this keyword. If you are simulating a multi-species system where some of the species have rotational degrees of freedom and some species are spheres, then specify an appropriate value of *Real(i,j)* for the species having rotational degrees of freedom, and set *Real(i,j)* equal to zero for the spherical species. Linear molecules are a special case, where rotation is handled in Eulerian space. If you have a linear molecule such as carbon monoxide, specify any non-zero value for *Real(i,j)*. Cassandra will properly sample the rotational degrees of freedom but will not use the value set by *Real(i,j)*. Note that spherical molecules are not considered when choosing which species to perform a rotational move on.

For example, if you are simulating two species in two boxes and if the first species has rotational degrees of freedom while the second is spherical, you would specify the rotational probability as

```
# Prob_Rotation
0.25
30.0 0.0
180.0 0.0
```

This will tell Cassandra to perform rotational moves 25% of the time. For box 1, the maximum rotational width will be  $30^\circ$  for species 1 and  $0.0^\circ$  for species 2. For box 2, the maximum rotational width will be  $180^\circ$  for species 1 and  $0.0^\circ$  for species 2. Note that, since the maximum rotational width of species 2 is set to  $0^\circ$  in both boxes, no rotational moves will be attempted on species 2.

For a single box simulation with three species such that the first species has rotational degree of freedom, the second is a linear molecule and the third species is spherical, you would specify

```
# Prob_Rotation
0.25
30.0 10.0 0.0
```

This will tell Cassandra to attempt rotational move 25% of the time. The maximum rotational width for species 1 is  $30^\circ$  and that for species 2 is  $10.0^\circ$ . Since the species 2 is a linear molecule, its rotation will be attempted in Eulerian angles and Cassandra will not use this value. Since the rotational width is set equal to  $0^\circ$  for species 3, no rotational moves will be attempted for this species.

## Regrowth

```
# Prob_Regrowth
Real(1)
Real(i,2)* One for each species
```

A regrowth move consists of deleting part of the molecule randomly and then regrowing the deleted part via a configurational bias algorithm. This can result in relatively substantial conformational changes for the molecule, but the cost of this move is higher than that of a simple translation or rotation. The probability of attempting a regrowth move is specified by *Real(1)* while *Real(i,2)* specifies the relative probability of performing this move on species *i*. For monatomic species, *Real(i,2)* is set to zero. Note that the user needs to ensure that the relative probabilities add up to 1 otherwise Cassandra will display an error and quit.

For example, if you are simulating three species of which the first species is monatomic, you would specify the following:

```
# Prob_Regrowth
0.3
0.0 0.7 0.3
```

This will tell Cassandra to attempt regrowth move 30% of the time. The relative probabilities of performing regrowth moves on species 1, 2 and 3 are 0.0, 0.7 and 0.3 respectively.

## Volume

### # Prob\_Volume

*Real(1)*

*Real(2,i)* \* One line required for each box except for GEMC-NVT

Sets the probability of volume displacement moves. This flag is required for NPT, GEMC-NPT and GEMC-NVT simulations. Do not specify for any other simulation type. *Real(1)* is the relative probability of attempting a box volume change. Note that volume changes are bold and expensive moves and should be attempted infrequently. This probability should normally not exceed 0.05 and values from 0.01-0.03 are typical. *Real(2,i)* specifies the maximum volume displacement in  $\text{\AA}^3$  for box *i*. If you are simulating a two box system, a value of *Real(2,i)* is required for each box on separate lines. Note that the exception to this is for a GEMC-NVT simulation, where there are two boxes but the volume moves are coupled. In this case, only a single value of *Real(2,i)* is specified. The default is to change the box volume by random amounts up to the maximum value specified by *Real(2,i)*. For example, if you are simulating a liquid with a single box in the NPT ensemble, you would specify the following:

### # Prob\_Volume

0.02

300

This will tell Cassandra to attempt volume moves 2% of the time. The box volume would be changed by random amounts ranging from  $-300 \text{ \AA}^3$  to  $+300 \text{ \AA}^3$ . For a liquid box  $20 \text{ \AA}$  per side, this would result in a maximum box edge length change of about  $0.25 \text{ \AA}$ , which is a reasonable value. Larger volume changes should be used for vapor boxes. If you wish to perform a GEMC-NPT simulation, you might specify the following:

### # Prob\_Volume

0.02

300

5000

This will tell Cassandra to attempt volume moves 2% of the time. The first box volume (assumed here to be smaller and of higher density, such as would occur if it were the liquid box) would be changed by

random amounts ranging from  $-300 \text{ \AA}^3$  to  $+300 \text{ \AA}^3$ . The second box volume would be changed by random amounts ranging from  $-5000 \text{ \AA}^3$  to  $+5000 \text{ \AA}^3$ . As with all move probabilities, you can experiment with making bolder or more conservative moves. Note that if the `# Run_Type` is set to 'Equilibration', Cassandra will attempt to optimize the boldness of moves to achieve about 50% acceptance rates.

### Insertion and Deletion Moves

#### `# Prob_Insertion`

`Real(1)`

`Character(i,1)`

`Character(i,2)`

This flag is set only for GCMC simulations. `Real(1)` sets the probability of attempting insertion moves. If there is more than one species, each is chosen for an insertion attempt with equal probability. `Character(i,1)` and `Character(i,2)` control the manner with which the insertions are carried out for each species  $i$ . Right now, Cassandra only uses one insertion method (a reservoir sampling approach). Later versions will have other options. So for now, the only option is to set `Character(i,1)` equal to 'insertion method'. If the species can be inserted or deleted, set `Character(i,2)` equal to 'reservoir'. If the species is 'non-volatile' and should not be inserted or deleted or should stay in its original box, then set `Character(i,2)` equal to 'none'. Then whichever box that species starts in, it will remain there for the whole simulation. You must repeat these flags for each species  $i$ . For example, if you are performing a GCMC simulation with two species that can be inserted, you might specify the following

`# Prob_Insertion`

`0.1`

`insertion method`

`reservoir`

`insertion method`

`reservoir`

This will tell Cassandra to attempt insertions 10% of the time and both species will be inserted using the reservoir insertion method.

`# Prob_Insertion`

`0.1`

`insertion method`

`reservoir`

```
insertion method
none
```

This will tell Cassandra to attempt insertions 10% of the time. Only species 1 will be inserted, while species 2 will not get inserted.

### **# Prob\_Deletion**

*Real (1)*

Sets the relative probability of deletion during the course of a simulation. Each exchangeable species is randomly chosen and a deletion move is attempted on a randomly chosen molecule of this species. If a species has an insertion method ‘none’, no attempt is made to delete it. You must specify the same deletion probability as the insertion probability to satisfy microscopic reversibility. If you fail to do this, Cassandra will give an error and quit.

### **# Prob\_Swap**

*Real(1)*

*Character(i,1)*

*Character(i,2)*

This keyword is set only for a GEMC simulation to enable transfer of species between two boxes. *Real(1)* sets the relative probability of attempting transfer of a molecule from one box to the other. During the swap, the donor and receiving boxes are chosen randomly. The species chosen for transfer is selected according to its overall mole fraction which is calculated only for the species that can be exchanged between boxes. Thus, species that are “non-volatile” are not included while computing the mole fractions. A molecule is then chosen randomly for transfer.

Similar to the **# Prob\_Insertion** section, *Character(i,1)* and *Character(i,2)* describe the manner in which the swap is carried out for each species *i*. At present, the only option is to set *Character(i,1)* equal to ‘insertion method’. If the species can be swapped, set *Character(i,2)* equal to ‘reservoir’. If the species is not to be transferred between boxes, then set *Character(i,2)* to ‘none’. Then whichever box that species starts in, it will remain in that box for the whole simulation. These flags are to be repeated for each species *i*. For example, while performing a GEMC simulation for three species the first two of which are exchanged while the third is not, you might specify the following:

```
# Prob_Swap
0.1
insertion method
```

```

reservoir
insertion method
reservoir
insertion method
none

```

This will tell Cassandra to attempt swap moves 10% of the time. Attempts will be made to transfer species 1 and 2 between available boxes while molecules of species 3 will remain in the boxes they are present in at the start of the simulation.

### Flip Move

```

# Prob_Ring
Real(1) Real(2)

```

This keyword is used when flip moves are to be attempted to sample bond angles and dihedral angles in a ring fragment. For more details on this move, see our publication [?]. The relative probability of attempting a flip move is specified by *Real(1)* while the maximum angular displacement in degrees for the move is given by *Real(2)*. For example, if the flip is to be attempted 30% of the time and the maximum angular displacement for the move is 20° specify the following:

```

# Prob_Ring
0.30 20.0

```

### # Done\_Probability\_Info

This is a required keyword that marks the end of the section for specifying move probabilities. It must occur after # Move\_Probability and all the move probabilities must be specified between these two keywords. Once Cassandra reads # Done\_Probability\_info, it checks to make sure the probabilities sum to unity. If not, an error will be given.

#### 4.1.18 Start Type

```

# Start_Type
Character(1)

```

This keyword specifies whether Cassandra generates an initial configuration or uses a previously generated configuration to start a simulation. *Character(1)* takes one of the three options: ‘make\_config’, ‘checkpoint’

or ‘read\_old’ and it determines what configuration is used to start a simulation.

When ‘make\_config’ is used as the start type, Cassandra will generate an initial configuration. With this option, additional information is required on the number of molecules of each species in every box and is specified as follows:

### **make\_config**

*Integer(j,k)* \* One line for each species and one entry on each line for each box

where *Integer(j,k)* represents the number of molecules of species *j* in box *k*. Thus, for example, to generate an initial configuration for two species in two boxes such that the numbers of molecules of species 1 in box 1 and 2 are 100 and 50 respectively and those for species 2 are 75 and 25 respectively, the input file must contain the following:

```
# Start_Type
make_config
100 50
75 25
```

During the course of a simulation, Cassandra periodically generates a checkpoint file (\*.chk) containing information about the total number of translation, rotation and volume moves along with the random number seeds and the coordinates of all the molecules and their box number at the time the file is written. Cassandra provides the capability of restarting from this state point in the event that a simulation crashes or running a production simulation from an equilibrated configuration. For this purpose, in addition to the ‘checkpoint’ keyword, additional information in the form of the name of the checkpoint file *Character(2)* is required in the following format:

### **checkpoint**

*Character(2)*

For example, to continue simulations from a checkpoint file ‘methane\_vle-T148.chk’, you might specify:

```
# Start_Type
checkpoint
methane_vle-T148.chk
```

Cassandra also provides a ‘read\_old’ option to make use of just the coordinates of molecules to start a simulation. For example, a configuration generated at a lower temperature may be used to jump start a



simulation at a higher temperature. When the ‘read\_old’ option is used, additional information in the form of the file names *Character(k,3)* is required as shown below:

#### **read\_old**

*Character(k,3)* \* One line for each box

For example, to start a GEMC simulation using the configurations of the two boxes, you might specify:

```
# Start_Type
read_old
box1.readold
box2.readold
```

This will tell Cassandra to use the configurations of the two boxes stored in `box1.readold` and `box2.readold` to start a simulation. Note that configurations of the boxes can be easily extracted from the checkpoint file using the utility `read_old.py` provided in `Scripts/Read_Old`.

### 4.1.19 Run Type

#### **# Run\_Type**

*Character(1) Integer(1) Integer(2)*

This keyword is used to specify whether a given simulation is an equilibration or a production run. For an equilibration run, the maximum translational, rotational and volume widths (for an NPT or a GEMC simulation) are adjusted to achieve 50% acceptance rates. During a production run, the maximum displacement width for different moves are held constant.

Depending on the type of the simulation, *Character(1)* can be set to either “Equilibration” or “Production”. For an **Equilibration** run, *Integer(1)* denotes the number of MC steps performed for a given thermal move before the corresponding maximum displacement width is updated. *Integer(2)* is the number of MC volume moves after which the volume displacement width is updated. This number is optional if no volume moves are performed during a simulation (for example in an NVT or a GCMC simulation). When the run type is set to **Production**, the MC moves refer to the frequency at which the acceptance ratios for various moves will be computed and output to the log file. These acceptance rates should be checked to make sure proper sampling is achieved.

For an NPT equilibration run in which the widths of the thermal move are to be updated after 100 MC moves and maximum volume displacements after 10 volume moves, specify the following:

```
# Run_Type
Equilibration 100 10
```

For an NVT production run in which the acceptance ratios of various thermal moves are printed to the log file after every 250 MC steps of a given thermal move, use the following:

```
# Run_Type
Production 250
```

#### 4.1.20 Frequency

```
# Frequency_Info
freq_type Character(1)
Nthermofreq or thermofreq Integer(2)
Ncoordfreq or coordfreq Integer(3)
MCsteps or Stop Integer(4)
# Done_Frequency_Info
```

This section specifies the frequency at which thermodynamic properties and coordinates are output to a file. *Character(1)* determines the method by which the simulation is terminated and data is output. If *Character(1)* is to ‘Timed’, then the simulation stops after the specified time in minutes. The format for this option is given below:

```
freq_type Timed
thermofreq Integer(2)
coordfreq Integer(3)
Stop Integer(4)
```

With this option, thermodynamic quantities are output every *Integer(2)* minutes, coordinates are written to the disk every *Integer(3)* minutes and the total simulation time is specified in minutes by *Integer(4)*. For example, to run a simulation for 60 minutes such that thermodynamic quantities are written every minute and the coordinates are output every 10 minutes, use the following:

```
# Frequency_Info
```

```
freq_type Timed
thermofreq 1
coordfreq 10
Stop 60
# Done_Frequency_Info
```

Note that similar to `# Move_Probabilities` section, the end of the frequency section always includes the `# Done_Frequency_Info` line.

Simulations can also be run for a given number of MC steps. To enable this feature, *Character(1)* is set to 'none'. Additional information is required and is given in the following format:

```
freq_type none
Nthermofreq Integer(2)
Ncoordfreq Integer(3)
MCsteps Integer(4)
```

With this option, thermodynamic quantities are output every *Integer(2)* MC steps, coordinates are written at a frequency of *Integer(3)* MC steps and the simulation terminates after *Integer(4)* steps. Note that an MC step is defined as a single MC move, regardless of type and independent of system size.

To run a simulation of 50,000 steps such that thermodynamic quantities are printed every 100 MC steps and coordinates are output every 10,000 steps, use the following:

```
# Frequency_Info
freq_type none
Nthermofreq 100
Ncoordfreq 10000
MCsteps 50000
# Done_Frequency_Info
```

#### 4.1.21 Average

```
# Average_Info
Integer(1)
```

This section specifies how thermodynamic quantities are output. At present, Cassandra writes instant-

neous values of thermodynamic quantities at a frequency given by either *Nthermofreq* or *thermofreq* in the **# Frequency\_Info** section. *Integer(1)* is set to 1 for this purpose. Later versions of Cassandra will have the ability to output block averages as well. Thus, you will specify the following section in your input file:

```
# Average_Info
1
```

### 4.1.22 Property Output

```
# Property_Info Integer(i)
Character(i,j) * One line for each property
```

This section provides information on the properties that are output. More than one section is allowed for multiple boxes. In this case, each section is separated by a blank line. *Integer(i)* is the identity of the box for which the properties are desired. *Character(i,j)* is the property that is to be output. Each property is specified on a separate line. At present, the acceptable entries include:

```
Energy_Total: Total energy of the system (Extensive) in kJ/mol
Energy_LJ: Lennard-Jones energy of the sytem in kJ/mol
Energy_Elec: Electrostatic energy of the sytem in kJ/mol
Energy_Intra: Total intramolecular energy of the system including bonded and non-bonded interactions in kJ/mol
Enthalpy: Enthalpy of the system (Extensive) kJ/mol
Pressure: Pressure of the system in bar
Volume: Volume of the system in Å³
Nmols: Number of mols of species
Density: Density of a species in #/Å³
```

For example, if you would like total energy, volume and pressure of a one box system to be written, you may specify the following:

```
# Property_Info 1
Energy_Total
Volume
Pressure
```

For a GEMC-NVT simulation, total energy and density of all the species in box 1 and total energy, density

of all the species in box 2 along with the pressure may be output using the following format:

```
# Property_Info 1
Energy_Total
Density

# Property_Info 2
Energy_Total
Density
Pressure
```

#### 4.1.23 Fragment Files

```
# Fragment_Files
Character(i,j) Integer(i,j) * One line for each fragment i in species j
```

In this section, information about the fragment library is specified. *Character(*i,j*)* gives the location of the fragment library of fragment *i* in species *j*; *Integer(*i,j*)* is the corresponding integer id specifying the type of the fragment. This section is automatically generated by `library_setup.py`. However, if there is a need to change this section, follow the example given below.

For a simulation involving two species of which the first one contains three distinct fragments and species 2 has two identical fragments, this section might look like:

```
# Fragment_Files
frag_1_1.dat 1
frag_2_1.dat 2
frag_3_1.dat 3
frag_1_2.dat 4
frag_1_2.dat 4
```

This will tell Cassandra to use the files `frag_1_1.dat`, `frag_2_1.dat` and `frag_3_1.dat` for the three fragments of species 1. Since species 2 has two identical fragment, Cassandra will use the same fragment library `frag_1_2.dat` for these fragments.

#### 4.1.24 File Info

```
# File_Info
```

```
Character(1)
```

This section is used only while generating a fragment library. Cassandra will use the filename specified in *Character(1)* to store different conformations of the fragment being simulated. Once again, this section is automatically handled by `library_setup.py`. However, if the user wishes to modify this part, use the following template:

```
# File_Info
```

```
frag.dat
```

This will tell Cassandra to store the fragment library in the file named `frag.dat`.

#### 4.1.25 CBMC parameters

```
# CBMC_Info
```

```
kappa_ins Integer(1)
```

```
kappa_rot Integer(2)
```

```
kappa_dih Integer(3)
```

```
rcut_cbmc Real(i,4) * Number of entries equal to number of simulation boxes
```

Cassandra utilizes a configurational bias methodology based on reservoir sampling [?]. This section sets a number of parameters required for biased insertion/deletion (refer to the sections `# Prob_Insertion`, `# Prob_Deletion` and `# Prob_Swap`) and configurational regrowth (`# Prob_Regrowth` section) of molecules. For a biased insertion, a fragment is chosen at random and given a random orientation. A number of trial positions are generated for the center-of-mass of the fragment. One of the trial positions is then selected randomly based on the Boltzmann weight of the energy of the trial position. The number of trial insertion positions is given by *Integer(1)*.

Once a trial position for the insertion is chosen, rotational bias may be applied by generating a number of trial orientations. *Integer(2)* specifies the number of such trial orientations. This feature will be implemented in later versions of Cassandra and any value for *Integer(2)*, at present, is ignored. To avoid any confusion, *Integer(2)* is set to 0.

After the biased placement of the first fragment, additional fragments directly bonded to the first frag-

ment are placed. Each of these fragments undergoes a number of trial orientations with respect to the fragment to which it is added. *Integer(3)* controls the number of such orientations that are generated.

For all the trials, energy of the partially grown molecule with itself and surrounding molecules is to be calculated. For this purpose, a short cutoff is used. *Real(i,4)* specifies the cutoff distance in Å for each of the boxes in a simulation. A short cutoff is fast, but might miss some overlaps. You can experiment with this value to optimize it for your system.

For a GEMC simulation in which 12 candidate positions are generated for biased insertion/deletion, 10 trials for biased dihedral angle selection and the cutoff for biasing energy calculation is set to 5.0 Å in box 1 and 6.5 Å in box 2, this section would look like:

```
# CBMC_Info
kappa_ins 12
kappa_rot 0
kappa_dih 10
rcut_cbmc 5.0 6.5
```

## 4.2 MCF File

A Molecular Connectivity File (MCF) defines the information related to bonds, angles, dihedrals, impropers fragments and non bonded interactions for a given species. One MCF file is required for each species present in the system. The information contained in this file involves the force field parameters, atoms participating in each of the interactions and the functional form used in each potential contribution. The keywords are preceeded by a ‘#’ and comments follow a ‘!’. Similarly to the input file, the order of the keywords is not important. A complete list of the keywords is provided below.

**Note that parameters for all of the following keywords must be separated by spaces only. Do not use the tab character.**

**Note that MCF files are generated by the script `mcfgen.py` automatically. The following description is provided for the users who wish to modify the MCF file or build the MCF file on their own.**

### 4.2.1 Atom Info

#### # Atom Info

*Integer(1)*

*Integer(2) Character(3)\*6 Character(4)\*2 Real(5) Real(6) Character(7)\*20 Real(8) Real(9) Character(10)*

This keyword specifies the information for non-bonded interactions. It is a required keyword in the MCF file. If not specified, the code will abort. The inputs are specified below:

- *Integer(1)*: Total number of atoms in the species.
- *Integer(2)*: Atom index.
- *Character(3)\*6*: Atom type up to 6 characters. This string of characters should be unique for each interaction site in the system, i.e. do not use the same atom type for two atoms in the same (or different) species unless the (pseudo)atoms have the same atom types.
- *Character(4)\*2*: Atom element name up to 2 characters.
- *Real(5)*: Mass of the atom in amu. Note that for united atom models, this would be the mass of the entire pseudoatom.
- *Real(6)*: Charge on the atom.
- *Character(7)*: Specifies functional form for VDW interactions to be used in the simulation. This must match what is given for # `VDW_Style` (subsection 4.1.4) in the input file. At present only 'LJ' style is permitted.
- *Real(8)*: The energy parameter in K.
- *Real(9)*: Collision diameter ( $\sigma$ ) in Å.
- *Character(10)*: Set to 'ring' only if a given atom is part of a ring fragment. Note that a ring fragment is defined as those atoms that belong to the ring (e.g. in cyclohexane, all the six carbons) and any atom directly bonded to these ring atoms (e.g. in cyclohexane, all the hydrogens). In other words, all of the ring and exoring atoms are given the ring flag. For atoms that are not part of rings, leave this field blank.

**Note that for species with a single fragment, the branch point atom is listed as the first atom.**

For example, for a united atom pentane model:



```
# Atom_Info
5
1 C1_s1 C 15.0107 0.0 LJ 98.0 3.75
2 C2_s1 C 14.0107 0.0 LJ 46.0 3.95
3 C3_s1 C 14.0107 0.0 LJ 46.0 3.95
4 C4_s1 C 14.0107 0.0 LJ 46.0 3.95
5 C5_s1 C 15.0107 0.0 LJ 98.0 3.75
```

The number below the keyword `# Atom_Info` specifies a species with 5 interaction sites, consistent with a united atom pentane model. The first column specifies the atom ID of each of the pseudo atoms. The second and third columns provide the atom type and atom name, respectively. The fourth column represents the atomic mass of each pseudoatom. Note that the mass of `C1_s1` is 15.0107 for this united atom model, as it involves a carbon and three hydrogen atoms. The same applies for all other interaction sites. The fifth column contains the partial charges placed on each of these pseudoatoms. The sixth, seventh and eighth columns contain the repulsion-dispersion functional form, the energy parameter and the collision diameter respectively. In this case, the usual Lennard-Jones functional form is used. Note that none of these atoms used the flag 'ring', as no rings are present in this molecule.

For a molecule containing rings, for example cyclohexane:

```
# Atom_Info
6
1 C1_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
2 C2_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
3 C3_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
4 C4_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
5 C5_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
6 C6_s1 C 14.0107 0.0 LJ 52.5 3.91 ring
```

Note the flag 'ring' was appended as the last column for this cyclic molecule.

Finally, for the SPC/E water model:

```
# Atom_Info
3
1 O1_s1 O 16.00 -0.8476 LJ 78.20 3.1656
2 H2_s1 H 1.000 0.4238 LJ 0.0 0.0
```

```
3 H3_s1 H 1.000 0.4238 LJ 0.0 0.0
```

This is a molecule with a single fragment. Therefore, the branch point atom must be specified as the first atom in the list. In this case, oxygen is the branch point and thus its atom ID is 1.

### 4.2.2 Bond Info

#### # Bond\_Info

*Integer(1)*

*Integer(i,2) Integer(i,3) Integer(i,4) Character(i,5) Real(i,6) Real(i,7)*

This section provides information on the number of bonds in a molecule and atoms involved in each bond along with its type. It is a required keyword in the MCF file. If not specified, the code will abort. The inputs are specified below:

- *Integer(1)*: Total number of bonds in the species. From the next line onwards, the bonds are listed sequentially and information for each bond is included on a separate line.
- *Integer(i,2)*: Index of the  $i^{th}$  bond.
- *Integer(i,3) Integer(i,4)*: IDs of the atoms participating in the bond.
- *Character(i,5)*: Type of the bond. At present only ‘fixed’ is permitted.
- *Real(i,6)*: Specifies the bond length for a particular bond in Å.

**Note that at present, Cassandra simulations can be carried out only for fixed bond length systems.**

For example, for the water model SPC/E, the # Bond\_Info section is the following:

```
# Bond_Info
2
1 1 2 fixed 1.0
2 1 3 fixed 1.0
```

In the above example, two bonds are specified whose fixed length is set to 1.0 Å.

### 4.2.3 Angle Info

#### # Angle\_Info

*Integer(1)*

*Integer(i,2) Integer(i,3) Integer(i,4) Integer(i,5) Character(i,6) Real(i,7) Real(i,8)*

The section lists the information on the angles in the species. It is a required keyword in the MCF file. If not specified, the code will abort.

- *Integer(1)*: Number of angles in the species.
- *Integer(i,2)*: Index of the  $i^{th}$  angle.
- *Integer(i,3) Integer(i,4) Integer(i,5)*: IDs of the atoms participating in the  $i^{th}$  angle. Note that *Integer(i,4)* is the ID of the central atom.
- *Character(i,6)*: Type of the angle. Currently, Cassandra supports ‘fixed’ and ‘harmonic’ (Eq. section 2.1) angles. For the ‘fixed’ option, *Real(i,7)* is the value of the angle and *Real(i,8)* is ignored by the code if specified. In the case of ‘harmonic’ potential type, *Real(i,7)* specifies the harmonic force constant ( $K/\text{rad}^2$ ) while *Real(i,8)* is the nominal bond angle (in degrees).

For example, for a united atom pentane molecule with flexible angles, this section is the following:

```
# Angle_Info
```

```
3
```

```
1 1 2 3 harmonic 31250.0 114.0
```

```
2 2 3 4 harmonic 31250.0 114.0
```

```
3 3 4 5 harmonic 31250.0 114.0
```

In the above example, the three angles between the pseudoatoms found in the pentane model are specified. The three angles have an harmonic potential, whose force constant is equal and is set to 31250.0 K/rad<sup>2</sup>. Finally, the equilibrium angle for these angles is 114.0°.

An example for SPC/E water model with fixed angles is provided below:

```
# Angle_Info
```

```
1
```

```
1 2 1 3 fixed 109.47
```

This model has only one angle that is set to 109.47°. Note that this angle is fixed, so there is no force constant.

#### 4.2.4 Dihedral Info

##### # Dihedral\_Info

*Integer(1)*

*Integer(i,2) Integer(i,3) Integer(i,4) Integer(i,5) Integer(i,6) Character(i,7) Real(i,8) Real(i,9) Real(i,10) Real(i,11)*

This section of the MCF file lists the number of dihedral angles and associated information for a given species. It is a required keyword in the MCF file. If not specified, the code will abort.

- *Integer(1)*: Lists the number of dihedral angles.
- *Integer(i,2)*: Index of the  $i^{th}$  dihedral angle.
- *Integer(i,3)*: *Integer(i,6)* - IDs of the atoms in the  $i^{th}$  dihedral angle.
- *Character(i,7)*: Dihedral potential type. Acceptable options are ‘OPLS’, ‘CHARMM’, ‘harmonic’ and ‘none’. If ‘OPLS’ dihedral potential type is selected, then the real numbers *Real(i,8) - Real(i,11)* are the coefficients in the Fourier series (see Eq 2.2). The units are in kJ/mol. For the ‘CHARMM’ dihedral potential type, three additional parameters are specified:  $a_0, a_1$  and  $\delta$  (section 2.3). If ‘harmonic’ dihedral potential type is used, then two additional parameters,  $K_{phi}$  and  $\phi_0$  (section Eq 2.4), are specified. For the ‘none’ dihedral potential type, no additional parameters are necessary.

For example, for a united atom pentane molecule using an OPLS dihedral potential type, the dihedrals are specified as follows:

# Dihedral\_Info

2

1 1 2 3 4 OPLS 0.0 2.95188 -0.5670 6.5794

2 2 3 4 5 OPLS 0.0 2.95188 -0.5670 6.5794

In this model two dihedral angles are specified by atoms 1,2,3,4 and 2,3,4,5. This model uses an OPLS functional form and thus four parameters are provided after the OPLS flag.

### 4.2.5 Fragment Info

#### # Fragment\_Info

*Integer(1)*

*Integer(i,2) Integer(i,3) Integer(i,4) Integer(i,5) ... Integer(i,2+Integer(i,3))*

This section defines the total number of fragments in a given species. It is an optional keyword. However, if the species is composed of fragments, then this section must be specified. The inputs are specified below:

- *Integer(1)*: Total number of fragments.
- *Integer(i,2)*: Index of the  $i^{th}$  fragment.
- *Integer(i,3)*: Number of atoms in the  $i^{th}$  fragment.
- *Integer(i,4) ... Integer(i,2+integer(i,3))*: List of the atom IDs in the fragment. The first atom ID is that for the branch point atom. **Atom ordering for the remaining atoms must match the order of atoms in the fragment library files.**

For example, for a pentane united atom model:

#### # Fragment\_Info

3

1 3 2 1 3

2 3 3 2 4

3 3 4 3 5

This specifies three fragments. Each of these fragments has three atoms. The first atom specified for each of the fragments is the branch point atom.

### 4.2.6 Fragment Connectivity

#### # Fragment\_Connectivity

*Integer(1)*

*Integer(i,2) Integer(i,3) Integer(i,4)*

The section lists the fragment connectivity - which fragment is bonded to which other fragment. It is a required keyword if **Fragment\_Info** is specified.

- *Integer(1)*: total number of fragment connections.
- *Integer(i,2)*: index of the  $i^{th}$  fragment connectivity.
- *Integer(i,3) Integer(i,4)*: fragment IDs participating in the connectivity.

For example, for a pentane united atom model:

```
# Fragment_Connectivity
2
1 1 2
2 2 3
```

In this example, there are three fragments, therefore, two fragment connectivities must be specified. Note that fragment 1 is connected to fragment 2 and fragment 2 is connected to fragment 3.

## Chapter 5

# Implementing the Metropolis Acceptance Criteria

All Monte Carlo moves are implemented in Cassandra to preserve detailed balance between each pair of microstates  $m$  and  $n$

$$\Pi_{mn} \alpha_{mn} p_m = \Pi_{nm} \alpha_{nm} p_n \quad (5.1)$$

where  $\Pi_{mn}$  is the probability of accepting the move from microstate  $m$  to microstate  $n$ ,  $\alpha_{mn}$  is the probability of attempting the move that will form  $n$  from  $m$ , and  $p_m$  is the probability of  $m$  in the ensemble of interest.

In Cassandra, detailed balance is enforced via the Metropolis criterion

$$\Pi_{mn} = \min \left( 1, \frac{\alpha_{nm} p_n}{\alpha_{mn} p_m} \right) \quad (5.2)$$

The ratio in Eq. (5.2) will often involve an exponential, e.g.  $e^{-\beta\Delta U}$ . To preserve precision in the energy calculation, the acceptance probability is computed

$$\Pi_{mn} = \exp \left\{ -\max \left[ 0, \ln \left( \frac{\alpha_{nm} p_n}{\alpha_{mn} p_m} \right) \right] \right\} \quad (5.3)$$

The logarithm, defined in code as `ln_pacc`, is tested in the function `accept_or_reject()` which is defined in file `accept_or_reject.f90`. If `ln_pacc` is greater than 0 and less than a maximum numerical value,  $\Pi_{mn}$  is computed and compared to a random number.

Code 5.1: accept\_or\_reject.f90

```

47  accept = .FALSE.
48
49  IF (ln_pacc <= 0.0_DP) THEN
50
51      accept = .TRUE.
52
53  ELSE IF ( ln_pacc < max_kBT) THEN
54
55      pacc = DEXP(-ln_pacc)
56
57      IF ( rranf() <= pacc ) THEN
58
59          accept = .TRUE.
60
61      END IF
62
63  END IF

```

## 5.1 Canonical Monte Carlo

In the canonical ensemble, the number of molecules  $N$ , the volume  $V$  and temperature  $T$  are all constant. The position, orientation and conformation of a semi-flexible molecule with fixed bond-lengths containing  $M$  atoms is given by a  $2M+1$ -dimensional vector  $\mathbf{q}$ . The positions, orientations and conformations of all  $N$  molecules are denoted  $\mathbf{q}^N$ .

The probability of observing microstate  $m$  with configuration  $\mathbf{q}_m^N$  is

$$p_m = \frac{e^{-\beta U(\mathbf{q}_m^N)}}{Z(N, V, T)} d\mathbf{q}^N \quad (5.4)$$

where  $\beta$  is the inverse temperature  $1/k_B T$ ,  $U$  is the potential energy, the differential volume  $d\mathbf{q}^N$  is included to make  $p_m$  dimensionless and  $Z$  is the configurational partition function

$$Z(N, V, T) = \int e^{-\beta U(\mathbf{q}^N)} d\mathbf{q}^N. \quad (5.5)$$

The integral is over all  $N(2M + 1)$  degrees of freedom. The ratio of microstate probabilities follows from



Eq. (5.4)

$$\begin{aligned}\frac{p_m}{p_n} &= \frac{e^{-\beta U(\mathbf{q}_m^N)} d\mathbf{q}^N / Z(N, V, T)}{e^{-\beta U(\mathbf{q}_n^N)} d\mathbf{q}^N / Z(N, V, T)} \\ &= e^{\beta(U_n - U_m)} = e^{\beta \Delta U}\end{aligned}\quad (5.6)$$

The configurational partition function  $Z$  and differential volume  $d\mathbf{q}^N$  both cancel, leaving only the ratio of Boltzmann factors.

New configurations are generated by attempting moves that translate, rotate and regrow a randomly selected molecule.

### 5.1.1 Translating a Molecule

A molecule is translated by moving its center of mass in each Cartesian direction by a random amount chosen from the uniform distribution on the interval  $[-\delta r_{max}, \delta r_{max}]$ . The maximum displacement  $\delta r_{max}$  must be given in the input file. The translation move is symmetric in forward and reverse directions. That is, either microstate  $n$  can be formed from microstate  $m$  and vice versa by moving one molecule within  $\delta r_{max}$  in each Cartesian direction, or microstate  $n$  cannot be formed at all. As a result,  $\alpha_{mn} = \alpha_{nm}$ .

The acceptance probability for a translation move follows from Eq. (5.6)

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \ln \left( \frac{p_m}{p_n} \right) = \beta \Delta U \quad (5.7)$$

In Cassandra, the translation move is implemented in the subroutine `Translate` defined in `translate.f90`. The relevant lines from version 1.1 are quoted below. The variable names in the `translate.f90` code are identified with the symbols from Eq. (5.7) in Table 5.1.

Code 5.2: `translate.f90`

```
274 ln_pacc = beta(this_box) * delta_e
275 accept = accept_or_reject(ln_pacc)
```

### 5.1.2 Rotating a Molecule

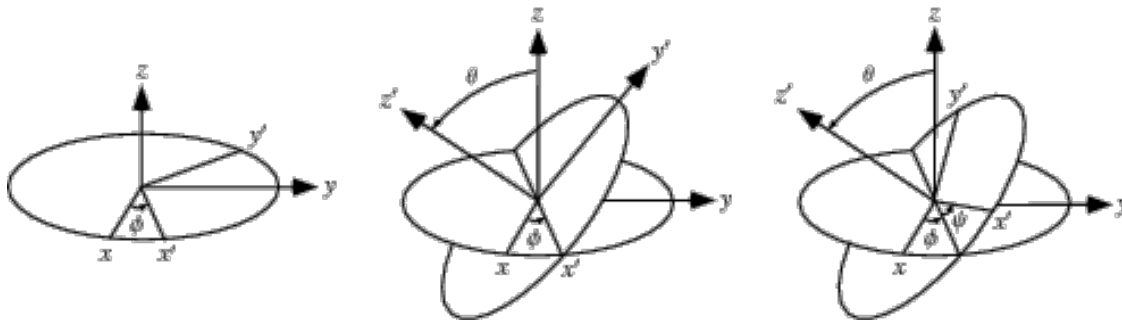
A linear molecule is rotated differently than a nonlinear molecule. A molecule is identified as linear if it is composed of 2 atoms or if all the angles are rigid with a bond angle of  $180^\circ$ . If the molecule is linear:

Table 5.1: Variable symbols and code names for translating and rotating a molecule

Symbol	Code name
$\beta$	beta(this_box)
$\Delta U$	delta_e

1. Pick three random angles:  $\phi$  on  $[-\pi, \pi]$ ,  $\cos(\theta)$  on  $[-1, 1]$ , and  $\psi$  on  $[-\pi, \pi]$ .
2. With the origin at the molecule's center of mass, rotate by  $\phi$  around  $z$ , rotate by  $\theta$  around  $x'$ , and rotate by  $\psi$  around  $z'$ , as shown in Fig. 5.1.

Even though three angles are randomly chosen, the probability of the resulting orientation is  $d\cos(\theta)d\phi/4\pi$ .

Figure 5.1: Procedure for rotating linear molecules. Image from [mathworld.wolfram.com/EulerAngles.html](http://mathworld.wolfram.com/EulerAngles.html).

If the molecule is nonlinear:

1. Randomly select an axis:  $x$ ,  $y$ , or  $z$ .
2. Choose a random angular displacement  $\delta\theta$  from  $[-\delta\theta_{max}, \delta\theta_{max}]$ .  $\delta\theta_{max}$  must be given in the input file.
3. Rotate the molecule around a vector parallel to the selected axis and through its center of mass by  $\delta\theta$ .

In either case, the rotation move is symmetric,  $\alpha_{mn} = \alpha_{nm}$ , and the acceptance criteria is given by Eq. (5.7). The rotation move is implemented in subroutine Rotate defined in rotate.f90.

Code 5.3: rotate.f90

```

261 ln_pacc = beta(this_box) * delta_e
262 accept = accept_or_reject(ln_pacc)

```

### 5.1.3 Regrowing a Molecule

Internal degrees of freedom in flexible molecules are sampled by deleting one or more fragments from the molecule and replacing the deleted fragments with conformations from a reservoir of fragment conformations. If the molecule consists of only a single fragment (e.g. water, all atom methane, united atom propane, all atom cyclohexane), the entire molecule is deleted and replaced as follows:

1. Randomly select a molecule  $i$  with uniform probability  $1/N$ , record its center-of-mass position and delete it.
2. Select a molecular conformation with Boltzmann probability  $e^{-\beta U(\mathbf{q}_{int,n}^{(i)})}/Z_{int}$ , where  $\mathbf{q}_{int,n}^{(i)}$  are the internal bond or improper angles of molecule  $i$  in microstate  $n$  and  $Z_{int}$  is the configurational partition function over internal degrees of freedom (see Eq. (5.13)).
3. Pick three random angles:  $\phi$  on  $[-\pi, \pi]$ ,  $\cos(\theta)$  on  $[-1, 1]$ , and  $\psi$  on  $[-\pi, \pi]$ . Rotate the molecule as shown in Fig. 5.1. The probability of the resulting orientation is  $d\mathbf{q}_{rot}/Z_{rot}$ , which for a nonlinear molecule is  $d\cos(\theta)d\phi d\psi/8\pi^2$ .
4. Place the molecule with the selected conformation and orientation at the same center-of-mass position as the deleted molecule.

Regrowing a monoatomic particle has no effect. Regrowing a linear molecule is the same as rotating it. The overall probability  $\alpha_{mn}$  of regrowing a molecule with the selected orientation and conformation is

$$\alpha_{mn} = \frac{1}{N} \frac{d\mathbf{q}_{rot}}{Z_{rot}} \frac{e^{-\beta U(\mathbf{q}_n^{(i)})} d\mathbf{q}_{int}}{Z_{int}} \quad (5.8)$$

where  $\mathbf{q}_n^{(i)}$  denotes the position, orientation and conformation of molecule  $i$  in microstate  $n$  and  $U(\mathbf{q}_n^{(i)})$  is the potential energy of the isolated molecule  $i$ , i.e. the intramolecular potential energy. The reverse probability  $\alpha_{nm}$  is identical except for the intramolecular potential energy  $U(\mathbf{q}_m^{(i)})$  of molecule  $i$  in microstate  $m$ . Using Eqs. (5.6) and (5.8), the acceptance criteria for the regrowth of a single fragment molecule is

$$\begin{aligned} \ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) &= \beta \left[ (U(\mathbf{q}_n^N) - U(\mathbf{q}_m^N)) - (U(\mathbf{q}_n^{(i)}) - U(\mathbf{q}_m^{(i)})) \right] \\ &= \beta \Delta U - \beta \Delta U_{int}^{(i)} = \beta \Delta U_{inter}^{(i)} \end{aligned} \quad (5.9)$$

Only the change in the intermolecular potential energy between molecule  $i$  and the other  $N - 1$  molecules contributes to the acceptance criteria. The code that implements Eq. (5.9) is shown in Code 5.7 in Section 5.3.3.

If the molecule consists of more than one fragment (e.g., all atom ethane, all atom toluene, united atom butane), a bond is cut and the severed fragments are regrown using Configurational Bias Monte Carlo (CBMC). See Section 5.3.3 for more details.

### 5.1.4 Canonical Partition Function

In Sections 5.1.1-5.1.2, the microstate probability is normalized by the configuration partition function  $Z$  because the only relevant degrees of freedom are configurational. In other ensembles, the full partition function  $Q$  appears, integrated over both configuration space  $\mathbf{q}^N$  and momenta space  $\mathbf{p}_q^N$

$$Q(N, V, T) = \frac{1}{h^{N(2M+1)} N!} \int e^{-\beta H(\mathbf{p}_q^N, \mathbf{q}^N)} d\mathbf{p}_q^N d\mathbf{q}^N \quad (5.10)$$

where the  $2M+1$  momenta  $\mathbf{p}_q$  are conjugate to the generalized coordinates  $\mathbf{q}$ . The momenta and configuration integrals are separable, and the single molecule momenta integrals are all identical.

$$\begin{aligned} Q(N, V, T) &= \frac{1}{N!} \left[ \int e^{-\beta U(\mathbf{q}^N)} d\mathbf{q}^N \right] \left[ \frac{1}{h^{2M+1}} \int e^{-\beta K(\mathbf{p}_q)} d\mathbf{p}_q \right]^N \\ &= \frac{1}{N!} Z(N, V, T) \left[ \frac{Q(1, V, T)}{Z(1, V, T)} \right]^N \end{aligned} \quad (5.11)$$

where  $Q(1, V, T)$  is the partition function of a single molecule in a box. The center of mass integrals for a single molecule are separable from the integrals over rotational and internal degrees of freedom:

$$Q(1, V, T) = Q_{com} Q_{rot+int} = V \Lambda^{-3} Q_{rot+int} \quad (5.12)$$

where  $\Lambda$  is the de Broglie wavelength of the molecule and the rotational and internal momenta integrals in  $Q_{rot+int}$  are not separable since the moments of inertia will depend on the conformation adopted by the molecule. The configurational partition function is further separable into center of mass (translational),

orientational and internal degrees of freedom:

$$Z(1, V, T) = V Z_{rot} Z_{int} \quad (5.13)$$

where the volume  $V$  is the translational partition function and  $Z_{rot}$  equals  $4\pi$  for a linear molecule and  $8\pi^2$  for a nonlinear molecule.

## 5.2 Isothermal-Isobaric Monte Carlo

In the isothermal-isobaric ensemble, the number of particles  $N$ , the pressure  $P$  and temperature  $T$  are all constant while the volume  $V$  and energy  $E$  fluctuate. The partition function is

$$\Delta(N, P, T) = \int e^{-\beta PV} Q(N, V, T) dV \quad (5.14)$$

Note that  $Q$  is dimensionless and  $\Delta$  has dimensions of volume. The probability of the system having volume  $V$  is

$$p(V) = \frac{Q(N, V, T) e^{-\beta PV}}{\Delta(N, P, T)} dV \quad (5.15)$$

The probability of observing microstate  $m$  with configuration  $\mathbf{q}_m^N$  and volume  $V_m$  is

$$\begin{aligned} p_m &= \frac{e^{-\beta U(\mathbf{q}_m^N)} d\mathbf{q}_m^N}{Z(N, V_m, T)} \frac{Q(N, V_m, T) e^{-\beta PV_m} dV}{\Delta(N, P, T)} \\ &= \frac{e^{-\beta U_m - \beta PV_m}}{\Delta(N, P, T)} \left( \frac{Q(1, V_m, T)}{Z(1, V_m, T)} d\mathbf{q}_m \right)^N dV \end{aligned} \quad (5.16)$$

where the differential volume element  $d\mathbf{q}_m^N$  has subscript  $m$  because it depends on the volume  $V_m$ . The ratio of microstate probabilities is

$$\frac{p_m}{p_n} = e^{\beta(U_n - U_m) + \beta P(V_n - V_m)} \left( \frac{d\mathbf{q}_m}{d\mathbf{q}_n} \right)^N = e^{\beta \Delta U + \beta P \Delta V} \left( \frac{d\mathbf{q}_m}{d\mathbf{q}_n} \right)^N \quad (5.17)$$

### 5.2.1 Scaling the Volume

In Cassandra, new volumes are sampled as follows:

1. Pick a random volume  $\Delta V$  with uniform probability from the interval  $[-\delta V_{max}, \delta V_{max}]$ . The trial volume is  $V + \Delta V$ .
2. Scale the box lengths uniformly.
3. Scale the center of mass of each molecule uniformly.

The probability of selecting  $\Delta V$  is the same as selecting  $-\Delta V$  which makes scaling the volume symmetric,  $\alpha_{mn} = \alpha_{nm}$ . Scaling the configurations changes the differential element  $d\mathbf{q}_m^N$  surrounding configuration  $\mathbf{q}_m^N$ . Only the molecular centers of mass change, so we can separate  $d\mathbf{q}$  into 3 center of mass coordinates  $d\mathbf{r}_{com}$  and  $2M-2$  orientational and internal coordinates  $d\mathbf{q}_{rot+int}$ . The scaled center of mass positions are held constant, making  $d\mathbf{r}_{com} = V d\mathbf{s}_{com}$ . The acceptance probability for a volume scaling move is

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \ln \left( \frac{p_m}{p_n} \right) = \beta \Delta U + \beta P \Delta V + N \ln \left( \frac{V_m}{V_n} \right) \quad (5.18)$$

The volume scaling move is implemented in subroutine Volume.Change defined in volume\_change.f90.

Code 5.4: volume\_change.f90

```

475 ln_pacc = beta(this_box) * delta_e &
476         + beta(this_box) * pressure(this_box) * delta_volume &
477         - total_molecules * DLOG(box_list(this_box)%volume/box_list_old%volume)
478 accept = accept_or_reject(ln_pacc)

```

## 5.3 Grand Canonical Monte Carlo

In the grand canonical ensemble, the chemical potential  $\mu$ , the volume  $V$  and temperature  $T$  are held constant while the number of molecules  $N$  and energy  $E$  fluctuate. The partition function is

$$\Xi(\mu, V, T) = \sum_{N=0}^{\infty} Q(N, V, T) e^{\beta \mu N} \quad (5.19)$$

Table 5.2: Variable symbols and code names for volume scaling move.

Symbol	Code name
$\beta$	beta(this_box)
$\Delta U$	delta_e
$P$	pressure(this_box)
$\Delta V$	delta_volume
$N$	total_molecules
$V_n$	box_list(this_box)%volume
$V_m$	box_list_old%volume

The probability of the system having  $N$  molecules is

$$p(N) = \frac{Q(N, V, T)e^{\beta\mu N}}{\Xi(\mu, V, T)} \quad (5.20)$$

The probability of observing microstate  $m$  with  $N_m$  molecules and configuration  $\mathbf{q}_m^{N_m}$  is

$$\begin{aligned}
p_m &= \frac{e^{-\beta U(\mathbf{q}_m^{N_m})} d\mathbf{q}^{N_m}}{Z(N_m, V, T)} \frac{Q(N_m, V, T)e^{\beta\mu N_m}}{\Xi(\mu, V, T)} \\
&= \frac{e^{-\beta U_m + \beta\mu N_m}}{\Xi(\mu, V, T)} \left[ \frac{Q(1, V, T)}{Z(1, V, T)} d\mathbf{q} \right]^{N_m}
\end{aligned} \quad (5.21)$$

Note that Eq. (5.21) does not contain the factorial  $N_m!$  that accounts for indistinguishable particles. In a simulation, particles *are* distinguishable: they are numbered and specific particles are picked for MC moves. The ratio of microstate probabilities is

$$\frac{p_m}{p_n} = e^{\beta\Delta U - \beta\mu\Delta N} \left[ \frac{Q(1, V, T)}{Z(1, V, T)} d\mathbf{q} \right]^{-\Delta N} \quad (5.22)$$

Alternatively, Eq. (5.22) can be recast to use the fugacity  $f$  instead of the chemical potential  $\mu$ . The relationship between  $\mu$  and  $f$  is

$$\mu = -k_B T \ln \left( \frac{Q(1, V, T)}{N} \right) = -k_B T \ln \left( \frac{Q(1, V, T)}{\beta f V} \right) \quad (5.23)$$

Inserting Eq. (5.23) into Eq. (5.22) yields

$$\frac{p_m}{p_n} = e^{\beta \Delta U} \left[ \frac{\beta f V}{Z(1, V, T)} d\mathbf{q} \right]^{-\Delta N} \quad (5.24)$$

Fluctuations in the number of molecules are achieved by inserting and deleting molecules. A successful insertion increases the number of molecules from  $N$  to  $N + 1$ , i.e.  $\Delta N = 1$ . A successful deletion decreases the number of molecules from  $N$  to  $N - 1$ , i.e.  $\Delta N = -1$ .

Random insertions and deletions (see Section 5.6) in the liquid phase typically have very high  $\Delta U$  due to core overlap and dangling bonds, respectively, making the probability of acceptance very low. Instead, insertions in Cassandra are attempted using Configurational Bias Monte Carlo.

### 5.3.1 Inserting a Molecule with Configurational Bias Monte Carlo

In Configurational Bias Monte Carlo (CBMC), the molecular conformation of the inserted molecule is molded to the insertion cavity. First, the molecule is parsed into fragments such that each fragment is composed of (a) a central atom and the atoms directly bonded to it (see Fig. 5.2), or (b) a ring of atoms and all the atoms directly bonded to them. Then, a position, orientation and molecular conformation of the molecule to be inserted are selected via the following steps:

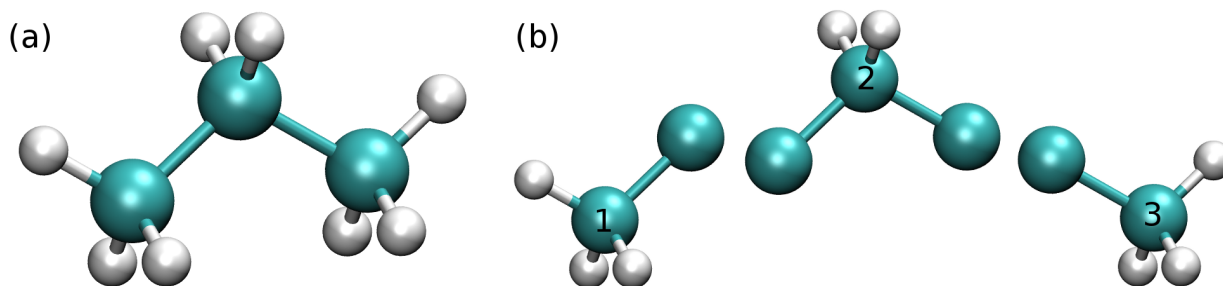


Figure 5.2: (a) An all-atom model of propane. (b) The same model as in (a), but parsed into three fragments.

1. Select the order in which each fragment of the  $(N + 1)$ th molecule will be placed. The probability of the resulting sequence is  $p_{seq}$ . (See example in Table. 5.3.)

(a) The first fragment  $i$  is chosen with uniform probability  $1/N_{frag}$ .



- (b) Subsequent fragments must be connected to a previously chosen fragment and are chosen with the uniform probability  $1/N_{cnn}$ , where the number of connections  $N_{cnn} = \sum_{ij} \delta_{ij} h_i (1 - h_j)$  is summed over all fragments  $i$  and  $j$ .  $h_i$  is 1 if fragment  $i$  has been previously chosen and 0 otherwise.  $\delta_{ij}$  is 1 if fragments  $i$  and  $j$  are connected and 0 otherwise.
2. Select a conformation for fragment  $i$  with Boltzmann probability  $e^{-\beta U(\mathbf{q}_{frag_i})} d\mathbf{q}_{frag_i} / Z_{frag_i}$ , where  $\mathbf{q}_{frag_i}$  are the internal degrees of freedom (angles and/or impropers) associated with fragment  $i$ .
  3. Select an orientation with uniform probability  $d\mathbf{q}_{rot} / Z_{rot}$ .
  4. Select a coordinate for the center of mass (COM) of fragment  $i$ :
    - (a) Select  $\kappa_{ins}$  trial coordinates  $\mathbf{r}_k$ , each with uniform probability  $d\mathbf{r}/V$ . Since one of the trial coordinates will be selected later, the individual probabilities are additive. The probability of the collection of trial coordinates is  $\kappa_{ins} d\mathbf{r}/V$ .
    - (b) Compute the change in potential energy  $\Delta U_k^{ins}$  of inserting fragment  $i$  at each position  $\mathbf{r}_k$  into configuration  $\mathbf{q}_m^N$ .
    - (c) Select one of the trial coordinates with probability  $e^{-\beta \Delta U_k^{ins}} / \sum_k e^{-\beta \Delta U_k^{ins}}$ .
  5. For each additional fragment  $j$ :
    - (a) Select a fragment conformation with Boltzmann probability  $e^{-\beta U(\mathbf{q}_{frag_j})} d\mathbf{q}_{frag_j} / Z_{frag_j}$
    - (b) Select the first of  $\kappa_{dih}$  trial dihedrals  $\phi_0$  with uniform probability from the interval  $[0, \frac{2\pi}{\kappa_{dih}})$ . Additional trial dihedrals are equally spaced around the unit circle,  $\phi_k = \phi_{k-1} + 2\pi/\kappa_{dih}$ . The probability of selecting  $\phi_0$  is  $\kappa_{dih} d\phi/2\pi$ .
    - (c) Compute the change in potential energy  $\Delta U_k^{dih}$  of attaching fragment  $j$  to the growing molecule with each dihedral  $\phi_k$ .
    - (d) Select one of the trial dihedrals with probability  $e^{-\beta \Delta U_k^{dih}} / \sum_k e^{-\beta \Delta U_k^{dih}}$ .

The overall probability  $\alpha_{mn}$  of attempting the insertion with the selected position, orientation and conformation is

Table 5.3: Possible sequences and probabilities for inserting the fragments of the all-atom model of propane shown in Fig. 5.2.

Sequence	$p_{seq}$
1 2 3	1/3
2 1 3	1/6
2 3 1	1/6
3 2 1	1/3

$$\alpha_{mn} = p_{seq} \frac{d\mathbf{q}_{rot}}{Z_{rot}} \frac{\kappa_{ins} d\mathbf{r}}{V} \frac{e^{-\beta\Delta U_k^{ins}}}{\sum_k e^{-\beta\Delta U_k^{ins}}} \times \left[ \prod_{i=1}^{N_{frag}} \frac{e^{-\beta U(\mathbf{q}_{frag_i})} d\mathbf{q}_{frag_i}}{Z_{frag_i}} \right] \left[ \prod_{j=1}^{N_{frag}-1} \frac{\kappa_{dih} d\phi}{2\pi} \frac{e^{-\beta\Delta U_k^{dih}}}{\sum_k e^{-\beta\Delta U_k^{dih}}} \right] \quad (5.25)$$

$$= p_{seq} p_{bias} \frac{e^{-\beta U(\mathbf{q}_{frag})} d\mathbf{q}}{V Z_{rot} Z_{frag} \Omega_{dih}} \quad (5.26)$$

where  $Z_{frag} = \prod_i Z_{frag_i}$  is the configurational partition function over degrees of freedom internal to each fragment,  $U(\mathbf{q}_{frag}) = \sum_i U(\mathbf{q}_{frag_i})$  is the summed potential energy of each of the (disconnected) fragments,  $\Omega_{dih} = (2\pi)^{N_{frag}-1}$  and  $p_{bias}$  is

$$p_{bias} = \frac{\kappa_{ins} e^{-\beta\Delta U_k^{ins}}}{\sum_k e^{-\beta\Delta U_k^{ins}}} \left[ \prod_{j=1}^{N_{frag}-1} \frac{\kappa_{dih} e^{-\beta\Delta U_k^{dih}}}{\sum_k e^{-\beta\Delta U_k^{dih}}} \right] \quad (5.27)$$

Note that the term  $V Z_{rot} Z_{frag} \Omega_{dih}$  in the denominator of Eq. (5.26) differs from  $Z(1, V, T) = V Z_{rot} Z_{int}$ .

In the reverse move, 1 of the  $N + 1$  particles is randomly selected for deletion. The probability  $\alpha_{nm}$  of picking the molecule we just inserted is

$$\alpha_{nm} = \frac{1}{N + 1} \quad (5.28)$$

Combining Eqs. (5.26) and (5.28) with Eq. (5.22) or Eq. (5.24) gives the acceptance probability for a CBMC insertion move

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta \left[ \Delta U - U(\mathbf{q}_{frag,n}^{(N+1)}) \right] - \beta \mu' + \ln \left( \frac{(N+1)\Lambda^3}{V} \right) + \ln(p_{seq} p_{bias}) \quad (5.29)$$

$$= \beta \left[ \Delta U - U(\mathbf{q}_{frag,n}^{(N+1)}) \right] + \ln \left( \frac{N+1}{\beta f' V} \right) + \ln(p_{seq} p_{bias}) \quad (5.30)$$

where  $\mu'$  and  $f'$  are, respectively, a shifted chemical potential and a skewed fugacity,

$$\mu' = \mu + k_B T \ln \left( Q_{rot+int} \frac{Z_{frag} \Omega_{dih}}{Z_{int}} \right) \quad (5.31)$$

$$f' = f \frac{Z_{frag} \Omega_{dih}}{Z_{int}} \quad (5.32)$$

All of the terms in Eqs. (5.31) and (5.32) are intensive. GCMC simulations using Eqs. (5.29) and (5.30) will converge to the same average density regardless of the simulation volume  $V$ . However, the values of  $\mu'$  or  $f'$  that correspond to the converged density will *not* match tabulated values of  $\mu$  or  $f$  computed from experimental data.

Note that the term  $Z^{IG}/\Omega$  from Macedonia *et al.* would be equivalent to  $Z_{int}/\Omega_{frag} \Omega_{dih}$  in the nomenclature used here. The configurational partition function of the disconnected fragments integrates over a Boltzmann factor,  $Z_{frag} = \int e^{-\beta U(\mathbf{q}_{frag})} d\mathbf{q}_{frag}$ , whereas the term  $\Omega_{frag} = \int d\mathbf{q}_{frag}$  does not.

In Cassandra, the insertion move is implemented in the subroutine Insertion in insertion.f90. The relevant lines from version 1.1 are quoted below. The variable names in the insertion.f90 code are identified with symbols in Table 5.4.

Table 5.4: Variable symbols and code names for inserting a molecule

Symbol	Code name
$\beta$	beta(this_box)
$\Delta U$	delta_e
$U(\mathbf{q}_{frag})$	E_angle + nrg_ring_frag_tot
$p_{seq}$	P_seq
$p_{bias}$	P_bias
$\mu'$	species_list(is)%chem_potential
$N$	nmols(is,this_box)
$V$	box_list(this_box)%volume
$\Lambda$	species_list(is)%de_broglie(this_box)
$f'$	species_list(is)%fugacity

Code 5.5: insertion.f90

```

441  ln_pacc = beta(this_box) * (delta_e - E_angle - nrg_ring_frag_tot)

447  ln_pacc = ln_pacc + DLOG(P_seq * P_bias) &
448      + DLOG(REAL(nmols(is,this_box)+1,DP)) &
449      - DLOG(box_list(this_box)%volume)
450
451  IF(lchempot) THEN
452      ! chemical potential is input
453      ln_pacc = ln_pacc - species_list(is)%chem_potential * beta(this_box) &
454      + 3.0_DP*DLOG(species_list(is)%de_broglie(this_box))
455  ELSE
456      ! fugacity is input
457      ln_pacc = ln_pacc - DLOG(species_list(is)%fugacity) &
458      - DLOG(beta(this_box)) &
459  END IF
460
461  accept = accept_or_reject(ln_pacc)

```

### 5.3.2 Deleting a Molecule that was Inserted via Configurational Bias Monte Carlo

The probability  $\alpha_{mn}$  of choosing a molecule to delete is

$$\alpha_{mn} = \frac{1}{N} \quad (5.33)$$

The probability of the reverse move  $\alpha_{nm}$  requires knowledge of the sequence and biasing probabilities  $p_{seq}$  and  $p_{bias}$  that would have been used to place the molecule if it was being inserted.  $p_{seq}$  and  $p_{bias}$  can be calculated using the following procedure:

1. Select the fragment order using the same procedure for inserting a molecule. The probability of the resulting sequence is  $p_{seq}$ .
2. The first fragment in the sequence is fragment  $j$ . Calculate the intramolecular potential energy of fragment  $j$ 's current conformation,  $U(\mathbf{q}_{frag_j})$ . The probability of this conformation is Boltzmann  $e^{-\beta U(\mathbf{q}_{frag_j})} d\mathbf{q}_{frag_j} / Z_{frag_j}$ .
3. The probability of the fragment's current orientation is  $d\mathbf{q}_{rot} / Z_{rot}$ .
4. Calculate the weight of the fragment's current center of mass (COM) coordinates:
  - (a) Compute the interaction potential energy  $\Delta U^{ins}$  between fragment  $j$  and the other  $N - 1$  molecules.
  - (b) Select  $\kappa_{ins} - 1$  trial coordinates  $\mathbf{r}_k$ , each with uniform probability  $d\mathbf{r} / V$ .
  - (c) Calculate the weight of the fragment's current COM amongst the trial coordinates,  $e^{-\beta \Delta U^{ins}} / \sum_k e^{-\beta \Delta U_k^{ins}}$ .
5. For each additional fragment  $j$ :
  - (a) Calculate the intramolecular potential energy of fragment  $j$ 's current conformation,  $U(\mathbf{q}_{frag_j})$ . The weight of this conformation in the Boltzmann distribution is  $e^{-\beta U(\mathbf{q}_{frag_j})} d\mathbf{q}_{frag_j} / Z_{frag_j}$ .
  - (b) Calculate the interaction potential energy  $\Delta U^{dih}$  between fragment  $j$ , on the one hand, and fragments  $i$  through  $j - 1$  and the other  $N - 1$  molecules.
  - (c) Calculate the current dihedral  $\phi_0$  of fragment  $j$ . Compute the interaction potential energy  $\Delta U_k^{dih}$  at  $\kappa_{dih} - 1$  trial dihedrals  $\phi_k = \phi_{k-1} + 2\pi / \kappa_{dih}$ .
  - (d) Compute the weight of  $\phi_0$  amongst the trial dihedrals,  $e^{-\beta \Delta U^{dih}} / \sum_k e^{-\beta \Delta U_k^{dih}}$ .

The overall probability  $\alpha_{nm}$  is

$$\alpha_{nm} = p_{seq} p_{bias} \frac{e^{-\beta U(\mathbf{q}_{frag})} d\mathbf{q}}{V Z_{rot} Z_{frag} \Omega_{dih}}. \quad (5.34)$$

The acceptance criteria for deleting a molecule inserted via CBMC is

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta \left[ \Delta U + U(\mathbf{q}_{frag,m}^{(i)}) \right] + \beta \mu' + \ln \left( \frac{V}{N \Lambda^3} \right) - \ln(p_{seq} p_{bias}) \quad (5.35)$$

$$= \beta \left[ \Delta U + U(\mathbf{q}_{frag,m}^{(i)}) \right] + \ln \left( \frac{\beta f' V}{N} \right) - \ln(p_{seq} p_{bias}) \quad (5.36)$$

In Cassandra, the deletion move is implemented in the subroutine Deletion in deletion.f90. The relevant lines are quoted below. The variable names in deletion.f90 code are identified with symbols in Table 5.5.

Code 5.6: deletion.f90

```

334  ln_pacc = beta(this_box) * (delta_e + E_angle + nrg_ring_frag_tot)

340  ln_pacc = ln_pacc - DLOG(P_seq * P_bias) &
341      - DLOG(REAL(nmols(is,this_box),DP)) &
342      + DLOG(box_list(this_box)%volume)
343
344  IF(lchempot) THEN
345      ! chemical potential is input
346      ln_pacc = ln_pacc + beta(this_box) * species_list(is)%chem_potential &
347      - 3.0_DP*DLOG(species_list(is)%de_broglie(this_box))
348  ELSE
349      ! fugacity is input
350      ln_pacc = ln_pacc + DLOG(species_list(is)%fugacity) &
351      + DLOG(beta(this_box)) &
352  END IF
353
354  accept = accept_or_reject(ln_pacc)

```

Table 5.5: Variable symbols and code names for deleting a molecule

Symbol	Code name
$\beta$	beta(this_box)
$\Delta U$	delta_e
$U(\mathbf{q}_{frag})$	E_angle + nrg_ring_frag_tot
$p_{seq}$	P_seq
$p_{bias}$	P_bias
$\mu'$	species_list(is)%chem_potential
$N$	nmols(is,this_box)
$V$	box_list(this_box)%volume
$\Lambda$	species_list(is)%de_broglie(this_box)
$f'$	species_list(is)%fugacity

### 5.3.3 Regrowing a Molecule with Configurational Bias Monte Carlo

Regrowing a molecule that has more than one fragment is a combination deletion and insertion move. Starting from microstate  $m$ :

1. Randomly select a molecule with uniform probability  $1/N$ .
2. Randomly select a bond to cut on the selected molecule with uniform probability  $1/N_{bonds}$ .
3. Delete the fragments on one side of the bond or the other with equal probability. The number of deleted fragments is  $N_{del}$ .
4. Reinsert the deleted fragments using the CBMC procedures for selecting the order of inserting the fragments, choosing a fragment conformation, and a connecting dihedral value (see Section 5.3.1).

The overall probability  $\alpha_{mn}$  of attempting to regrow the molecule with the selected conformation is

$$\begin{aligned}
 \alpha_{mn} &= \frac{p_{seq}}{NN_{bonds}} \left[ \prod_{j=1}^{N_{del}} \frac{e^{-\beta U(\mathbf{q}_{frag_j}^{(i)})} d\mathbf{q}_{frag_j}}{Z_{frag_j}} \right] \left[ \prod_{j=1}^{N_{del}} \frac{\kappa_{dih} d\phi}{2\pi} \frac{e^{-\beta \Delta U_k^{dih}}}{\sum_k e^{-\beta \Delta U_k^{dih}}} \right] \\
 &= \frac{p_{seq}}{NN_{bonds}} \frac{e^{-\beta U(\mathbf{q}_{del,n}^{(i)})} d\mathbf{q}}{Z_{del} \Omega_{del}} p_{forward}
 \end{aligned} \tag{5.37}$$

Table 5.6: Variable symbols and code names for regrowing a molecule

Symbol	Code name
$\beta$	beta(this_box)
$U(\mathbf{q}_n^N) - U(\mathbf{q}_{del,n}^{(i)})$	delta_e_n - nrg_ring_frag_forward
$U(\mathbf{q}_m^N) - U(\mathbf{q}_{del,m}^{(i)})$	delta_e_o - nrg_ring_frag_reverse
$p_{forward}$	P_forward
$p_{reverse}$	P_reverse

where  $Z_{del} = \prod_i Z_{frag_j}$  is the configurational partition function over degrees of freedom internal to the deleted fragments,  $U(\mathbf{q}_{del,n}^{(i)}) = \sum_j U(\mathbf{q}_{frag_j})$  is the summed potential energy of each deleted fragment with the conformations in microstate  $n$ ,  $\Omega_{del} = (2\pi)^{N_{del}}$  and  $p_{forward}$  is the biasing probability

$$p_{forward} = \prod_{j=1}^{N_{del}} \frac{\kappa_{dih} e^{-\beta \Delta U_k^{dih}}}{\sum_k e^{-\beta \Delta U_k^{dih}}} \quad (5.38)$$

The reverse move is identical except for the potential energy of the deleted fragments  $U(\mathbf{q}_{del,m}^{(i)})$  in microstate  $m$  and the biasing probability  $p_{reverse}$  which will depend on the values of the connecting dihedrals. Using Eqs. (5.6) and (5.37), the acceptance criteria is:

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta \left[ \left( U(\mathbf{q}_n^N) - U(\mathbf{q}_{del,n}^{(i)}) \right) - \left( U(\mathbf{q}_m^N) - U(\mathbf{q}_{del,m}^{(i)}) \right) \right] + \ln \left( \frac{p_{forward}}{p_{reverse}} \right) \quad (5.39)$$

Eq. (5.39) is implemented in subroutine cut\_N\_grow() in file cutNgrow.f90.

Code 5.7: cutNgrow.f90

```

392  ln_pacc = beta(this_box) * (delta_e_n - nrg_ring_frag_forward) &
393          - beta(this_box) * (delta_e_o - nrg_ring_frag_reverse) &
394          + DLOG(P_forward / P_reverse)
395
396  accept = accept_or_reject(ln_pacc)

```



## 5.4 Multiphase Systems

### 5.4.1 Gibbs Ensemble Monte Carlo

[SECTION INCOMPLETE]

## 5.5 Multicomponent Systems

Sections 5.1-5.4 have each been developed for pure component systems. The Monte Carlo moves and acceptance criteria for multicomponent systems are straightforward extensions of the pure component moves. The only modification needed to translate, rotate and regrow molecules is to first select a species. In these moves, a species is selected randomly in proportion to its mole fraction  $N_i/N$ . When inserting and deleting a molecule, the mole fractions of each species change. In these cases, a species in a multicomponent system is selected instead with uniform probability  $1/N_{species}$ . In either case, species selection is symmetric for both forward and reverse moves and so cancels from the acceptance criterion.

In contrast, Reaction Ensemble Monte Carlo (RxMC) fundamentally occurs in a multicomponent grand canonical ensemble because the number of molecules of all species involved in a reaction change simultaneously. In the multicomponent grand canonical ensemble, the temperature  $T$ , volume  $V$  and set of chemical potentials  $\{\mu\} = \mu_1, \dots, \mu_s$  of each species are held constant. The partition function for a system with  $s$  species is

$$\Xi(\{\mu\}, V, T) = \sum_{N_1=0}^{\infty} \dots \sum_{N_s=0}^{\infty} Q(\{N\}, V, T) e^{\beta \mu_i N_i} \quad (5.40)$$

where  $\{N\}$  is the set of  $s$  molecule numbers, and a summation is implied over the term  $\mu_i N_i$ . The probability of the system having  $\{N\}$  molecules is

$$p(\{N\}) = \frac{Q(\{N\}, V, T) e^{\beta \mu_i N_i}}{\Xi(\{\mu\}, V, T)} \quad (5.41)$$

The probability of observing microstate  $m$  with  $\{N\}_m$  molecules and configuration  $\{\mathbf{q}^N\}_m$  is

$$\begin{aligned}
p_m &= \frac{e^{-\beta U_m} \prod_{i=1}^s d\mathbf{q}_i^{N_{im}}}{Z(\{N\}_m, V, T)} \frac{Q(\{N\}_m, V, T) e^{\beta \mu_i N_{im}}}{\Xi(\{\mu\}, V, T)} \\
&= \frac{e^{-\beta U_m + \beta \mu_i N_{im}}}{\Xi(\{\mu\}, V, T)} \prod_{i=1}^s \left[ \frac{Q_i(1, V, T)}{Z_i(1, V, T)} d\mathbf{q}_i \right]^{N_{im}}
\end{aligned} \tag{5.42}$$

The ratio of microstates is

$$\frac{p_m}{p_n} = e^{\beta \Delta U - \beta \mu_i \Delta N_i} \prod_{i=1}^s \left[ \frac{Q_i(1, V, T)}{Z_i(1, V, T)} d\mathbf{q}_i \right]^{-\Delta N_i} \tag{5.43}$$

### 5.5.1 Reaction Ensemble Monte Carlo

In Reaction Ensemble Monte Carlo, the number of molecules of each species can only change by a stoichiometric amount  $\nu_i$  as a reaction occurs in the forward direction or  $-\nu_i$  as a reaction occurs in the reverse direction. Moreover, reaction equilibrium occurs when the following condition is met

$$\sum_{i=1}^s \nu_i \mu_i = 0 \tag{5.44}$$

Substituting Eq. (5.44) and  $\Delta N_i = \pm \nu_i$  into Eq. (5.43) gives

$$\frac{p_m}{p_n} = e^{\beta \Delta U} \prod_{i=1}^s \left[ \frac{Q_i(1, V, T)}{Z_i(1, V, T)} d\mathbf{q}_i \right]^{\mp \nu_i} \tag{5.45}$$

The product of partition functions can be replaced by the equilibrium constant of an ideal gas  $K_{eq}$  or the standard change of free energy  $\Delta G_{rxn}^o$

$$\prod_{i=1}^s \left( \frac{Q_i(1, V, T)}{V} \right)^{\nu_i} = K_{eq} = e^{-\beta \Delta G_{rxn}^o} (P^\circ k_B T)^{\Delta \nu} \tag{5.46}$$

where  $P^\circ$  is the standard pressure (by default, 1 atm). Substituting Eq. (5.46) and  $Z(1, V, T) = V Z_{rot} Z_{int}$  into Eq. (5.45) gives

$$\frac{p_m}{p_n} = e^{\beta\Delta U} K_{eq}^{\mp 1} \prod_{i=1}^s \left[ \frac{d\mathbf{q}_i}{Z_{rot,i} Z_{int,i}} \right]^{\mp \nu_i} \quad (5.47)$$

$$= e^{\beta\Delta U \pm \beta\Delta G_{rxn}^o} \prod_{i=1}^s \left[ \frac{d\mathbf{q}_i}{P^o k_B T Z_{rot,i} Z_{int,i}} \right]^{\mp \nu_i} \quad (5.48)$$

## 5.6 Appendix

### 5.6.1 Inserting a Molecule Randomly

To insert a molecule, a position, orientation and molecular conformation must each be selected. The probability of inserting the new molecule at a random location is  $d\mathbf{r}/V$ , where  $d\mathbf{r}$  is a Cartesian volume element of a single atom. The probability of choosing the molecule orientation is  $d\mathbf{q}_{rot}/Z_{rot}$ , which for a linear molecule is  $d\cos(\theta)d\phi/4\pi$  and for a nonlinear molecule is  $d\cos(\theta)d\phi d\psi/8\pi^2$ . The probability of the molecule conformation only depends on the remaining  $2M - 5$  internal bond angles, dihedral angles and improper angles  $\mathbf{q}_{int}$ . A thermal ensemble of configurations is Boltzmann distributed  $e^{-\beta U(\mathbf{q}_{int})}/Z_{int}$ . The overall probability  $\alpha_{mn}$  is

$$\alpha_{mn} = \frac{d\mathbf{r}}{V} \frac{d\mathbf{q}_{rot}}{Z_{rot}} \frac{e^{-\beta U(\mathbf{q}_{int,N+1,n})}}{Z_{int}} d\mathbf{q}_{int} = \frac{e^{-\beta U(\mathbf{q}_{N+1,n})}}{Z(1,V,T)} d\mathbf{q}. \quad (5.49)$$

where we have used Eq. (5.13) to recover  $Z(1,V,T)$  and recognized that only internal degrees of freedom contribute to the potential energy of the isolated  $N+1$ th molecule in microstate  $n$ ,  $U(\mathbf{q}_{N+1,n}) = U(\mathbf{q}_{int,N+1,n})$ . For a point particle with no rotational or internal degrees of freedom,  $\alpha_{mn}$  reduces to  $d\mathbf{r}/V$ . For molecules with internal flexibility, a library of configurations distributed according to  $e^{-\beta U(\mathbf{q}_{int})}/Z_{int}$  can be generated from a single molecule MC simulation. In the reverse move, 1 of the  $N+1$  particles is randomly selected for deletion. The probability  $\alpha_{nm}$  of picking the molecule we just inserted is

$$\alpha_{nm} = \frac{1}{N+1} \quad (5.50)$$

The acceptance probability for a random insertion move is

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta [\Delta U - U(\mathbf{q}_{N+1})] - \beta\mu + \ln \left( \frac{N+1}{Q(1,V,T)} \right) \quad (5.51)$$

where  $U(\mathbf{q}_{N+1})$  is the intramolecular potential energy of the inserted molecule.  $Q(1, V, T)$  is typically not known *a priori*, nor is it easily estimated. Substituting Eq. (5.12) into Eq. (5.51) and absorbing  $Q_{rot+int}$  into a shifted chemical potential  $\mu'$

$$\mu' = \mu - k_B T \ln(Q_{rot+int}) \quad (5.52)$$

gives the acceptance criteria for inserting a molecule

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta [\Delta U - U(\mathbf{q}_{N+1})] - \beta \mu' + \ln \left( \frac{(N+1)\Lambda^3}{V} \right). \quad (5.53)$$

The terms absorbed into  $\mu'$  are intensive and therefore GCMC simulations using Eq. (5.53) will converge to a specific average density. However, the value of  $\mu'$  that corresponds to the converged density will *not* match tabulated values of  $\mu$  computed from experimental data.

Substituting Eq. (5.23) into Eq. (5.51) gives

$$\ln \left( \frac{\alpha_{mn} p_m}{\alpha_{nm} p_n} \right) = \beta [\Delta U - U(\mathbf{q}_{N+1})] + \ln \left( \frac{N+1}{\beta f V} \right) \quad (5.54)$$

where no terms have been absorbed into the fugacity  $f$ . Note also that the partition function has completely been eliminated from the acceptance criteria.

### 5.6.2 Deleting a Molecule Inserted Randomly

The probability  $\alpha_{mn}$  of choosing a molecule to delete is

$$\alpha_{mn} = \frac{1}{N} \quad (5.55)$$

The probability  $\alpha_{nm}$  of inserting that molecule back in is

$$\alpha_{nm} = \frac{e^{-\beta U(\mathbf{q})}}{Z(1, V, T)} d\mathbf{q} \quad (5.56)$$

The acceptance probability for deleting a molecule inserted randomly is

$$\ln \left( \frac{\alpha_{mn}}{\alpha_{nm}} \frac{p_m}{p_n} \right) = \beta [\Delta U + U(\mathbf{q}_N)] + \beta \mu' + \ln \left( \frac{V}{N \Lambda^3} \right) \quad (5.57)$$

$$= \beta [\Delta U + U(\mathbf{q}_N)] + \ln \left( \frac{\beta f V}{N} \right) \quad (5.58)$$

Note that in  $\Delta U$  is defined differently in Eqs. (5.53) and (5.54) than in Eqs. (5.57) and (5.58). In the former, the new configuration has more molecules,  $\Delta U = U(\mathbf{q}_n^{N+1}) - U(\mathbf{q}_m^N)$ . In the latter, the new configuration has fewer molecules,  $\Delta U = U(\mathbf{q}_n^{N-1}) - U(\mathbf{q}_m^N)$ .