

# Random Number Generators à la Boltzmann

A Thesis  
Submitted for the Degree of  
MASTER OF SCIENCE (ENGINEERING)

by  
SAMARTH AGRAWAL



ENGINEERING MECHANICS UNIT  
JAWAHARLAL NEHRU CENTRE FOR ADVANCED SCIENTIFIC RESEARCH  
(A Deemed University)  
Bangalore – 560 064

JUNE 2018



*Dedicated to my family*



## DECLARATION

I hereby declare that the matter embodied in the thesis entitled “**Random Number Generators à la Boltzmann**” is the result of investigations carried out by me at the Engineering Mechanics Unit, Jawaharlal Nehru Centre for Advanced Scientific Research, Bangalore, India under the supervision of **Prof. Santosh Ansumali** and that it has not been submitted elsewhere for the award of any degree or diploma.

In keeping with the general practice in reporting scientific observations, due acknowledgment has been made whenever the work described is based on the findings of other investigators.

---

**Samarth Agrawal**



## CERTIFICATE

I hereby certify that the matter embodied in this thesis entitled “**Random Number Generators à la Boltzmann**” has been carried out by **Mr. Samarth Agrawal** at the Engineering Mechanics Unit, Jawaharlal Nehru Centre for Advanced Scientific Research, Bangalore, India under my supervision and that it has not been submitted elsewhere for the award of any degree or diploma.

---

**Prof. Santosh Ansumali**  
(Research Supervisor)





# Acknowledgements

I would like to express my sincere gratitude to Prof. Santosh Ansumali for his constant support during the course of my Masters'. He also encouraged and guided me to learn new things which has helped me professionally and personally. I would also like to thank Prof. V Kumaran, Prof. Meher P Prakash, Prof. K R Sreenivas and Prof. Ganesh Subramanian for allowing me to attend their wonderful courses and helping me during the completion of my thesis.

I would like to thank my friends — Anand, Krishnendu, Atif, Shveta, Swati, Praveen and Rhotheth for making my life at JNCASR a little easier.

I would like to thank our former administrative officer, Mr. A. Jayachandra, along with Dr. Joydeep De, Dr. Princy, employees of the library, complab, hostel and other staff for their cooperation.

I am indebted to my family who extended their support through my highs and lows throughout my years at JNCASR.



# Abstract

In the last few decades, simulations of stochastic processes have gained prominence in many fields of science and engineering. These simulations rely on Random Number generators (RNGs), routines that produce seemingly random sequence of numbers. Currently two paradigms exist for building RNGs — extracting noise from physical devices and complicated mathematical constructs based mostly on number theory. In this thesis we show that a synergy of these two paradigms, that a simulation of a stochastic process can in fact form the basis of an RNG. We illustrate this via simulating gas molecules which follow a collectively chaotic motion. This thesis also shows that various numerical schemes that can solve for hydrodynamics at a mesoscopic level generate random sequences of numbers. We propose a new algorithm with these concepts as foundation that can generate Gaussian and exponential random numbers orders-of-magnitude faster than existing methods. By employing this algorithm we simulate reaction-diffusion problem and binary gas mixtures modeled at a mesoscopic level.



# List of Figures

1.1	The central idea of the thesis. Random sequences adhering to non-uniform distributions can be generated without using complicated mathematical functions by simulating an inherently stochastic process. . . . .	2
2.1	Thermal white noise from an electric resistor can be used to generate “true” random numbers. The noise is sampled periodically. Noise greater than theoretical value corresponds to a 1 bit and a 0 bit otherwise. For the case shown, the number generated in bitwise representation is 111100111000011. Since it is equally likely for the noise to be on either side of the baseline, this produces a uniform distribution . . . . .	6
2.2	Scatter plot of consecutive numbers $(x_i, x_{i+1})$ generated in a sequence by the iterative scheme $x_{n+1} = (1093x_n + 18257) \bmod 86436$ . a) A complete picture of the plot shows considerable mesh-like pattern. b) A close-up view of the bottom left corner of the complete plot shows that sequence of numbers mainly fall in selective planes. . . . .	7
2.3	Scatter plot of consecutive numbers $(x_i, x_{i+1})$ generated in a sequence by the <code>drand()</code> function of C++. a) A complete picture of the plot lacks any discernible pattern. b) A close-up view of the bottom left corner of the complete plot shows that sequence of numbers have apparently random behaviour . . . . .	8
2.4	A geometrical interpretation of the Box-Muller method, suggests that points on a plane once transformed and plotted in a polar co-ordinate system tend to cluster around the origin. . . . .	11
2.5	Typical realization of a Brownian particle immersed in fluid obtained using Gaussian random numbers generated by the Box-Muller method. . . . .	11
2.6	Plot comparing the histogram of beta-distributed random numbers generated by the acceptance-rejection method and the theoretical probability density function. The dotted lines indicate the bounding function and the domain of the variates. . . . .	13
3.1	Sketch of the probability density functions of the specific gravity of petrol and paraffin for the hypothetical test case. . . . .	19
3.2	Plot comparing the histogram of uniformly distributed random numbers generated by a PRNG and the theoretical probability density function . . . . .	21
4.1	The geometry of collision between two particles. $\mathbf{k}$ is the vector joining the centres of the two molecules and bisecting the angle made between the pre and post-collisional velocities, $\mathbf{g}_{12}$ and $\mathbf{g}'_{12}$ . The quantity, $b$ is known as the impact parameter and is the perpendicular distance from the centre of the molecule to $\mathbf{g}_{12}$ , and $\chi$ is the angle between $\mathbf{g}_{12}$ and $\mathbf{g}'_{12}$ . . . . .	30
4.2	A computer code simulating gas dynamics is similar to Maxwell’s demon. To an observer, the values of positions and velocities act as sequence of random numbers. . . . .	34
4.3	Plot of pairwise potential $V(r)$ vs. $r$ . . . . .	35

5.1	Plots comparing the results of the deterministic solution of the chemical reaction $A \xrightarrow{k} B$ . The solid, coloured lines represent the various realizations of the naive implementation and the dashed line represents the deterministic solution. It can be seen that large number of such realization when averaged converge to the deterministic solution. These figures represent a) degradation of $A$ and b) degradation of $B$ . . . . .	47
5.2	Simulation of the reaction network presented in Eq.(5.19). The plots depict the error observed for different number of molecules present in the system for Gillespie algorithm when implemented with standard and Molecular Dice algorithms. Both methods converge with $\log(1/N)$ , with Molecular Dice being around 4 times faster.	49
5.3	Plot comparing the results of Gillespie algorithm and the deterministic solutions of Eq.(5.19). The plots depict a) progression of mole fraction of $S$ and b) progression of mole fraction of $P$ . . . . .	49
5.4	Plot of $\Delta = (n_A + 2n_{A_2} + 2n_{OA_2}) - (n_B + 2n_{B_2} + 2n_{OB_2})$ for the reaction network presented in Eq.(5.4.3). It can be seen clearly that small perturbations can lead to switching between multiple steady states. . . . .	51
5.5	Simulation of the reaction-diffusion process presented in Eq.(5.21). The plots depict a) a scatter plot of the averaged density profile of MinD protein over an oscillation cycle. A minima can be observed close to the center of the cell, indicating the site for cell division. This is in good agreement with the deterministic formulation of the process. b) A plot of the ratio of MinD protein in the left-hand 30% to right-hand 30%. The solid line is observed in the stochastic simulation while the dotted line is the deterministic solution. While the latter fails to capture the oscillations observed in experiments, the stochastic version of the model manages to capture this phenomena. . . . .	52
6.1	Sketch representing the two-step relaxation to equilibrium. . . . .	58
7.1	The three types of collisional possibilities – A-A, B-B and A-B . . . . .	65
7.2	Plot of the distribution of velocities of the light and heavy component at equilibrium for a) Model I, and b) Model II. Plot of ratio of energy at time $t$ to the initial energy ( $E(t)/E_0$ ) vs. time for individual components and the mixture for c) Model I, and d) Model II. . . . .	76
7.3	The process of gas molecules escaping through a small hole is known as effusion. The lighter particles (in this case blue) escape through the hole faster than the heavier particles, with a factor proportional to the square root of their mass ratios.	77
7.4	Model I was used to simulate a setup that could mimic Graham's law for effusion. Plot shows that results observed are in great agreement with expected behaviour, for all three cases . . . . .	78
7.5	A representative sketch of the Couette flow setup. Two walls with a separation $H$ are sheared in the opposite directions with velocity $U/2$ . . . . .	78
7.6	Plot of the concentrations after 20,000 time steps of the component a) $A$ and b) $B$ , in comparison with the analytical solution given by Eq.(7.78) . . . . .	80

# List of Tables

2.1	Speeds of two commonly used LCGs on two different computers . . . . .	8
3.1	Interpretation of $p$ -values for a test of significance . . . . .	19
3.2	Observations of a die experiment . . . . .	22
3.3	$p$ -values for different tests . . . . .	24
3.4	An output for the Birthday Test . . . . .	25
3.5	An output for the Gorilla Test . . . . .	26
3.6	Performance of drand48 and Mersenne Twister in Crush battery of tests . . . . .	27
4.1	Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by Lennard-Jones dynamics and hard-sphere system compared with the widely used MT19937 (Matsumoto & Nishimura, 1998). The Gaussian and exponential random numbers were generated using the same. . . . .	36
4.2	Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by mesoscopic methods namely – DSMC and MPCD, compared with MT19937. . . . .	38
4.3	Overview of the “Molecular Dice” algorithm . . . . .	39
4.4	Results for Crush battery of tests. The results are for all the three quantities – uniform, Gaussian and exponential with the three different methods for pair selection outlined in Sec.(4.5.1) . . . . .	40
4.5	Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by the final algorithm for all the three methods of pair selection outlined in Sec.(4.5.1), compared with MT19937. . . . .	40
5.1	Time comparisons for the three methods to calculate the mole fraction of the specie $S$ after $t = 100$ . . . . .	48
7.1	Summary of the algorithm for binary mixtures . . . . .	75
7.2	Comparison of the values of $\Pi$ between the Fokker-Planck and DVM methods, for Ne-Ar mixture at three different concentrations $C_0 = (0.1, 0.5, 0.9)$ for a range of rarefactions . . . . .	79
7.3	Comparison of the values of $\Pi$ between the Fokker-Planck and DVM methods, for He-Ar mixture at three different concentrations $C_0 = (0.1, 0.5, 0.9)$ for a range of rarefactions . . . . .	79





# Contents

<b>Abstract</b>	<b>vii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Organization of the thesis . . . . .	3
<b>2 Random Number Generators (RNGs)</b>	<b>5</b>
2.1 Introduction . . . . .	5
2.2 Hardware based Random Number Generators . . . . .	5
2.3 Pseudo Random Number Generators (PRNGs) . . . . .	6
2.3.1 Linear Congruential Generators (LCGs) . . . . .	6
2.4 Generating non-uniform distributions . . . . .	7
2.4.1 Inverse transform sampling . . . . .	8
2.4.2 Transformation methods . . . . .	9
2.4.3 Acceptance-rejection method . . . . .	10
2.5 Qualities of a good PRNG . . . . .	13
<b>3 Statistical Tests</b>	<b>15</b>
3.1 Introduction . . . . .	15
3.2 Hypothesis testing . . . . .	15
3.3 Testing for PRNGs . . . . .	20
3.4 Preliminary distribution-free tests . . . . .	20
3.4.1 Pearson's $\chi^2$ test . . . . .	20
3.4.2 Kolmogorov-Smirnov Test . . . . .	22
3.4.3 Ljung-Box Test . . . . .	23
3.5 Tests for uniform random numbers . . . . .	24
3.5.1 Serial Test . . . . .	24
3.5.2 Gap Test . . . . .	24
3.5.3 Marsaglia's tests for empirical randomness . . . . .	25
3.6 Outlook . . . . .	27
<b>4 Molecular Dice</b>	<b>29</b>
4.1 Introduction . . . . .	29
4.2 The Boltzmann equation . . . . .	29
4.2.1 Distributions at equilibrium . . . . .	31
4.3 Molecular Dynamics simulations as PRNGs . . . . .	34
4.3.1 Lennard-Jones . . . . .	34
4.3.2 Hard-sphere systems . . . . .	35
4.4 Mesoscale methods . . . . .	36
4.4.1 Direct Simulation Monte Carlo . . . . .	37
4.4.2 Multiparticle Collision Dynamics . . . . .	37
4.5 Final algorithm . . . . .	38

4.5.1	Pair selection . . . . .	38
4.6	Statistical tests . . . . .	39
4.7	Speed of RNGs . . . . .	40
4.8	Outlook . . . . .	41
<b>5</b>	<b>Chemical Reactions</b>	<b>43</b>
5.1	Introduction . . . . .	43
5.2	Deterministic formulation . . . . .	43
5.3	Stochastic formulation of chemical rate equations . . . . .	44
5.4	Simulation Algorithms . . . . .	46
5.4.1	Gillespie algorithm . . . . .	46
5.4.2	Convergence and performance of Gillespie algorithm . . . . .	48
5.4.3	Rare event sampling in biological system . . . . .	50
5.5	Reaction-Diffusion Systems . . . . .	50
5.5.1	Pattern formation in bacteria . . . . .	50
5.6	Outlook . . . . .	51
<b>6</b>	<b>Fokker-Planck model for rarefied gases</b>	<b>55</b>
6.1	Introduction . . . . .	55
6.2	Kinetic modelling of rarefied gases . . . . .	55
6.3	Quasi-equilibrium models . . . . .	57
6.4	Fokker-Planck approximation . . . . .	58
6.4.1	Transport Coefficients . . . . .	59
6.5	Numerical solution . . . . .	60
6.6	Outlook . . . . .	61
<b>7</b>	<b>Fokker-Planck model for binary mixtures</b>	<b>63</b>
7.1	Introduction . . . . .	63
7.2	Kinetic modelling of binary mixtures . . . . .	63
7.3	Quasi-equilibrium models for binary mixtures . . . . .	65
7.4	Model I: Momentum and temperature difference as the slow variable . . . . .	67
7.5	Model II: Pressure as the slow variable . . . . .	69
7.6	Transport Coefficients . . . . .	70
7.7	Numerical scheme . . . . .	73
7.7.1	Model I . . . . .	74
7.7.2	Model II . . . . .	75
7.8	Simulation results . . . . .	75
7.8.1	Graham's law for effusion . . . . .	76
7.8.2	Couette Flow . . . . .	77
7.8.3	Binary diffusion . . . . .	79
7.9	Outlook . . . . .	80
<b>8</b>	<b>Outlook</b>	<b>81</b>

# Chapter 1

## Introduction

### 1.1 Motivation

Stochastic models have gained prominence in many fields of engineering and sciences such as physics, chemistry, finance etc (Gardiner, 1985*a*) as the behaviour of many systems such as quantum particles are inherently random and hence their state is defined in probabilistic terms. This is also true for phenomena where the large amounts of information required acts as a major deterrent (Risken, 1996). Such processes require stochastic modelling where the random contributions are accounted for thereby providing an estimate of the qualitative behaviour of the system. Consider the case of a Brownian particle immersed in a fluid at rest, the particle exhibits a “zig-zag” motion in a haphazard manner. In order to ascertain the nature of the system, the position and velocities of each fluid molecule is required hence rendering it to be an infeasible approach. However, the statistics of such processes can be inferred by understanding the phenomena that drives them. For the case of Brownian motion, there are essentially two mechanisms which dictate the motion of the particle. The first being the drag faced by the particle against its velocity and the second is the constant buffeting the particle experiences from the fluid molecules (Uhlenbeck & Ornstein, 1930). The equation governing the motion of the Brownian particle immersed in a fluid at rest is

$$m \frac{d^2 \mathbf{x}}{dt^2} = -\lambda \frac{d\mathbf{x}}{dt} + \boldsymbol{\eta}(t), \quad (1.1)$$

where  $m$  and  $\mathbf{x}$  denote the mass and the position of the particle respectively,  $\lambda$  the damping coefficient associated with the viscosity of the fluid and  $\boldsymbol{\eta}(t)$  the random force component resulting from the collisions with fluid molecules. The statistics of the force component can be deduced by simple reasoning. It is argued that at a time scale when the particle has faced many collisions, the sum total of these forces would be Gaussian in nature as guaranteed by the central limit theorem with mean zero and the variance directly proportional to the temperature of the system, as it is a measure of the thermal motion of the fluid molecules (Uhlenbeck & Ornstein, 1930). In addition to this, the forces must be independent of each other in time and direction. Hence, the probability density function of the force is Gaussian with mean and variance

$$\langle \eta_\alpha(t) \rangle = 0 \quad \langle \eta_\alpha(t) \eta_\beta(t') \rangle = 2\lambda k_B T \delta_{\alpha\beta} \delta(t - t'), \quad (1.2)$$

where  $\langle . \rangle$  is the operator symbolizing an average of the quantity over many ensembles,  $\delta_{\alpha\beta}$  the Kronecker delta,  $\delta(t - t')$  the Dirac delta function and  $k_B$  is the Boltzmann constant. In order to numerically solve this equation, one would require a sequence of Gaussian random numbers which could imitate the behaviour of the random force component. This is achieved by using a class of algorithms known as Pseudo-Random Number Generators (PRNGs), which produce empirically random sequences that can be used to mimic these seemingly random contributions (Knuth, 1981). While these routines are restricted to producing sequences that are distributed uniformly, by using appropriate transformations they can be converted to sequence of Gaussian random numbers and subsequently used to simulate the motion of a Brownian particle.

Numerically solving stochastic models require random numbers from a variety of distributions such as – exponential, Poisson etc. There have been many advances in the field of PRNGs and are mostly rooted in number theory and generate uniformly distributed random numbers (Knuth, 1981). Techniques for generation of non-uniform distributions rely on complicated mathematical

transformations which are computationally expensive and thus render many large-scale scientific simulations expensive. Many problems such as whole-cell simulation have not been solved yet because of the large number of sequences required, driving the expected computational time to years, hence it is imperative to develop new methods which can mitigate this problem and open up new possibilities in the realm of large-scale scientific computation. This thesis aims to address this problem by proposing a new algorithm to generate non-uniform distributions and highlight its capabilities in the context of large-scale scientific simulations. In particular the following aspects are dealt with, in the thesis:

- Mesoscale hydrodynamic solvers as PRNGs:** The existing methods to generate random sequences either rely on the random measurements made from physical devices or complicated iterative schemes rooted in number theory. We explore whether a synergy of the two approaches — simulation of an inherently stochastic process is able to generate high quality random sequences. A visual representation of this idea is presented in Fig.(1.1). For this purpose a simple rarefied gas dynamics simulation is chosen as the probability distribution of positions and velocities of the molecules can be easily computed from the Boltzmann equation. Numerous methods exist to simulate fluid flow at a mesoscopic scale. Such solvers aim to produce desired macroscopic behaviour by considering a particular microphysical possibility. In this thesis, a new algorithm is designed along these principles which generate Gaussian and exponential random numbers orders-of-magnitude more efficient than contemporary methods.

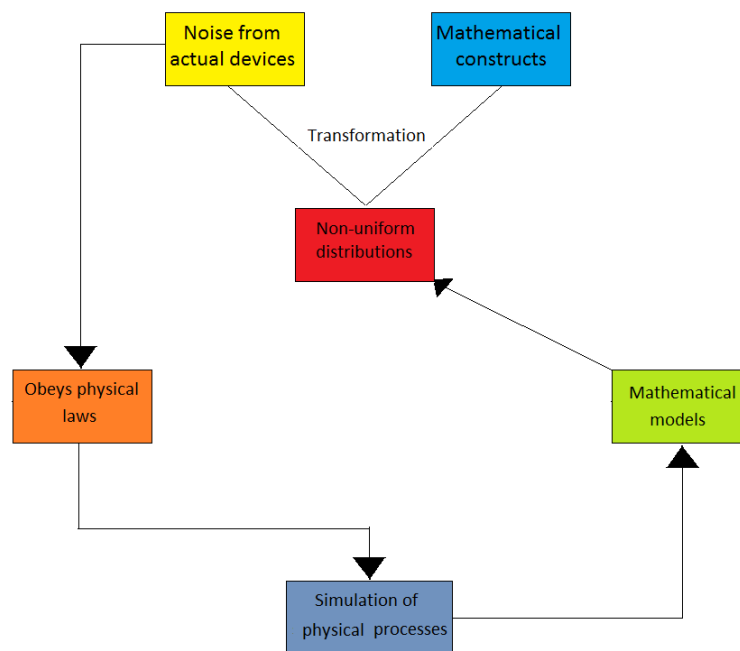


Figure 1.1: The central idea of the thesis. Random sequences adhering to non-uniform distributions can be generated without using complicated mathematical functions by simulating an inherently stochastic process.

- Stochastic simulation of chemical reactions:** For small system sizes, fluctuations and perturbations can play an important role in dictating its behaviour (Gardiner, 1985*a*). Hence, in order to incorporate these, a stochastic description of chemically reacting systems is employed. The numerical method used to solve resulting equations requires a high number of exponential and uniform random numbers (Gillespie, 1977). This problem is further exacerbated in the cases of reaction-diffusion systems where diffusion is treated as individual reaction events. The proposed algorithm is used to implement this scheme and

the resulting speedup is observed.

- **Simulation of binary mixtures:** The Fokker-Planck model for hydrodynamics has emerged as an alternative method for simulating fluid flows. Solving the resulting equations require a high number of Gaussian random numbers rendering the technique computationally expensive. We extend the existing model for binary mixtures based on the multi-relaxation scheme. This model is then benchmarked by simulating three canonical problems.

## 1.2 Organization of the thesis

The thesis is organized in the following manner:

- In Chapter 2, standard approaches to generate uniformly distributed random numbers is introduced along with various methods to generate non-uniform distributions. In addition, the characteristics of a high quality PRNG is discussed.
- In Chapter 3, the basic concepts associated with hypothesis testing and various statistical tests used to quantify and subsequently establish the quality of PRNGs is discussed.
- In Chapter 4, a new paradigm to generate random numbers is discussed. A new algorithm which could generate Gaussian and exponential random numbers orders-of-magnitude faster than contemporary algorithms is proposed.
- In Chapter 5, the basics of the stochastic formulation of chemically reactive systems along with numerical methods used to solve it is explained. Additionally, methods to solve reaction-diffusion systems are also discussed. Three problems are solved to demonstrate the efficacy of the proposed algorithm — bi-stable biochemical reaction network of proteins binding to a DNA, the Goldbeter-Koshland switch and pattern formation in bacteria.
- In Chapter 6, The Fokker-Planck approximation to the Boltzmann equation is briefly described followed by a summary of the numerical method used to solve the resulting set of equations and simulate fluid flow.
- In Chapter 7, A new model based on the Fokker-Planck equation for binary mixtures is introduced. This model is then benchmarked by solving three canonical problems — Couette flow, Graham's law for effusion and static diffusion.
- In Chapter 8, a summary of thesis is presented and further avenues resulting from this work are suggested.



# Chapter 2

## Random Number Generators (RNGs)

### 2.1 Introduction

The progression of many processes in nature such as – decay of radioactive matter, Brownian motion, queueing systems etc. cannot be predicted using simple deterministic equations because of their apparently random behaviour (Gardiner, 1985*a*). Such cases are better handled in a stochastic framework wherein the inherent randomness is accounted for and modelled accordingly. Hence, numerically simulating these models requires techniques which could mimic these random perturbations. These disturbances are incorporated in simulations in the form of stream of random numbers. There are essentially two paradigms to generate such sequences — extracting noise from hardwares exploiting the inherent fluctuations present in physical devices, and iterative schemes built on mathematical formulae which produce apparently random sequence of numbers (Knuth, 1981).

In this chapter different methods to generate random numbers used to numerically solve stochastic differential equations and different methods to generate non-uniform distributions are explained. The chapter ends with a brief discussion on the qualities of a good RNG.

### 2.2 Hardware based Random Number Generators

A simple but robust method to incorporate the randomness in systems is to use the measurements made from an actual physical device (Symul *et al.*, 2011). If the nature of noise produced from a given device can be quantified, it can be used as a stream of random numbers to numerically simulate stochastic models. For example, the electronic noise arising from thermal agitation of charge carriers inside a conductor observed in the values of the voltage, is Gaussian distributed with zero mean and variance

$$\overline{v_n^2} = 4k_B T R \Delta f, \quad (2.1)$$

where  $k_B$  is the Boltzmann constant,  $T$  the temperature,  $R$  the resistance and for a given bandwidth,  $\Delta f$  the bandwidth (Johnson, 1928; Nyquist, 1928). This is a well known result in statistical physics and referred to as the Johnson-Nyquist Noise. An example of how such observed noise can be used to generate sequence of random numbers presented in Fig.(2.1). Sequences generated using such techniques are termed as “true” random numbers and are highly reliable (Marandi *et al.*, 2012). However, the cost of providing these numbers in millions is computationally expensive owing to data transfer from the device to computers and in many cases specialized hardware is required to facilitate this transfer. In addition to this, the sequences generated using this method can be correlated or might not follow an exact distribution because of flaws in measuring instruments and techniques (Krishnan, 2015). As large-scale scientific simulations require large number of random sequences following a specified distribution such methods cannot be employed. However, this class of methods has found widespread use in the field of cryptography as it needs random sequences which cannot be predicted using statistical procedures (Jun & Kocher, 1999).

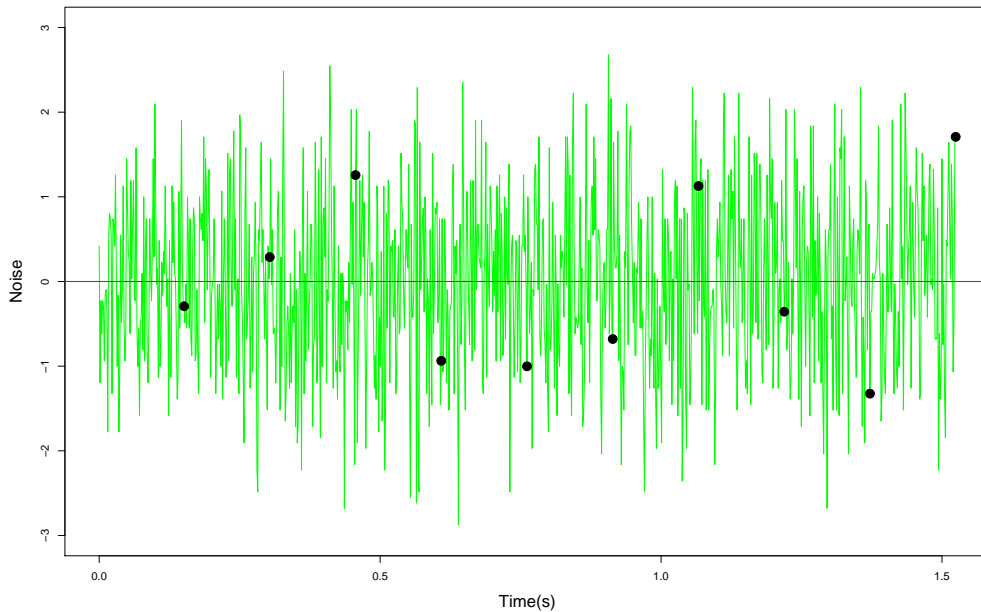


Figure 2.1: Thermal white noise from an electric resistor can be used to generate “true” random numbers. The noise is sampled periodically. Noise greater than theoretical value corresponds to a 1 bit and a 0 bit otherwise. For the case shown, the number generated in bitwise representation is 111100111000011. Since it is equally likely for the noise to be on either side of the baseline, this produces a uniform distribution

## 2.3 Pseudo Random Number Generators (PRNGs)

A computer algorithm is fundamentally restricted to produce deterministic outputs and no sequences any algorithm generates can be “truly” random. There exists, however, a class of algorithms known as Pseudo Random Number Generators (PRNGs) that can generate sequences of uniformly distributed numbers which cannot be distinguished from truly random sequences through standard statistical tests (Chorin & Hald, 2009). These sequences are considered “pseudo” random as they are the outputs of deterministic equations and would invariably start repeating after a finite period of time. The core idea of PRNGs is to define a function  $F$  such that an iterative scheme

$$x_n = F(x_0, \dots, x_{n-1}), \quad (2.2)$$

which produces an apparently random sequence of numbers (Knuth, 1981). Many such functions have been determined and have engendered different families of PRNGs such as — Xorshift (Marsaglia *et al.*, 2003), Permuted Congruential Generators (O'Neill, 2015) etc. In the following section the Linear Congruential Generator (LCG) family is discussed which would provide an insight into the framework of most PRNGs.

### 2.3.1 Linear Congruential Generators (LCGs)

A highly useful class of PRNGs is the Linear Congruential Generators (LCGs), wherein the function  $F$  is chosen in manner which yields the iterative scheme

$$x_n = (ax_{n-1} + b) \pmod{m}, \quad (2.3)$$



this scheme generates all numbers between  $[0, m)$  if and only if the following conditions are satisfied

1.  $m$  and  $b$  are co-primes, i.e, their greatest common divisor is 1.
2.  $(a - 1)$  is divisible by 4 if and only if  $m$  is.
3.  $(a - 1)$  is divisible by all prime factors of  $m$ .

these set of conditions are known as the Hull-Dobell theorem (Hull & Dobell, 1962). While these conditions guarantee that  $x_n$  attains all the values between  $[0, m)$ , it does not ensure absence of correlations and random behaviour. For example, Fig.(2.2) shows a plot of the output of an LCG which satisfies these conditions, but generates numbers which exhibit obvious pattern. In the last few decades various combinations of  $(a, b, m)$  have been determined which produce satisfactory results and have been adapted widely for scientific simulations (Knuth, 1981). As an example, a plot for the `drand()` function is shown Fig.(2.3) which clearly lacks discernible patterns as opposed to the plot in Fig.(2.2).

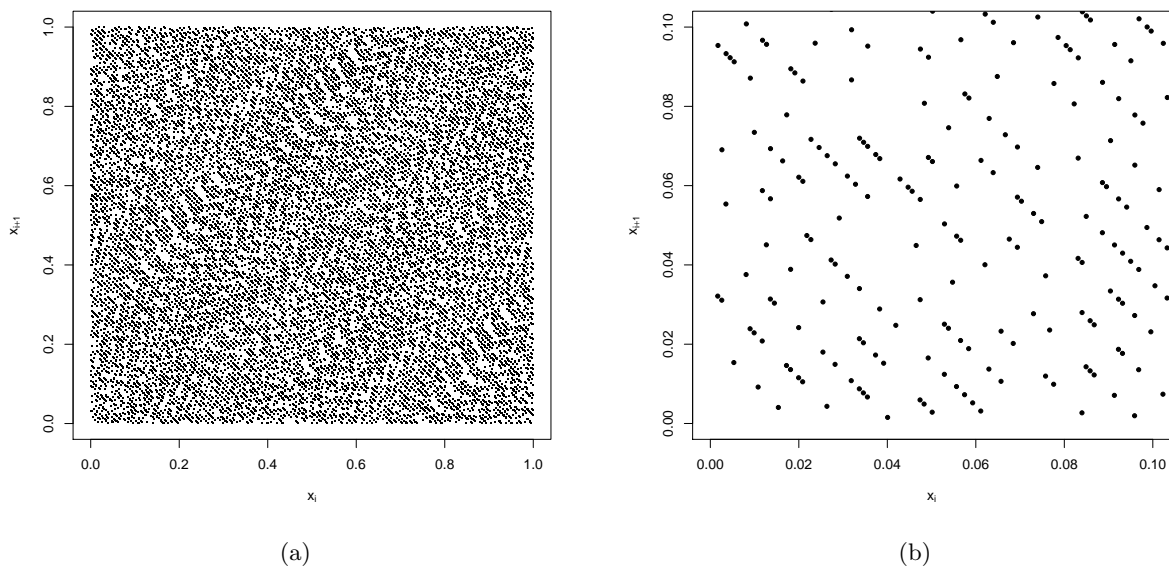


Figure 2.2: Scatter plot of consecutive numbers  $(x_i, x_{i+1})$  generated in a sequence by the iterative scheme  $x_{n+1} = (1093x_n + 18257) \bmod 86436$ . a) A complete picture of the plot shows considerable mesh-like pattern. b) A close-up view of the bottom left corner of the complete plot shows that sequence of numbers mainly fall in selective planes.

LCGs are the foundation for many canonical PRNGs such as the `drand()` family of C/C++. The simple nature of the function  $F$  associated with LCGs renders the computation to be quite fast and has been shown to hold good statistical properties for a variety of parameters. The rate of random number generation for two commonly used LCGs on different computers is tabulated in Table(2.1). The speeds of `rand()` are in integers/second while the speed of `drand48()` is in doubles/sec. The configuration of the two computers are, Computer 1 – Intel(R) Xeon(R) CPU E5-2650 v2 @ 2.60GHz and Computer 2 – Intel(R) Core(TM) i7-6800K CPU @ 3.40GHz.

## 2.4 Generating non-uniform distributions

There exists three techniques namely — inverse transform sampling, transformation methods and the acceptance-rejection method which can be used to generate non-uniformly distributed random numbers, and will be described briefly in the following sections.

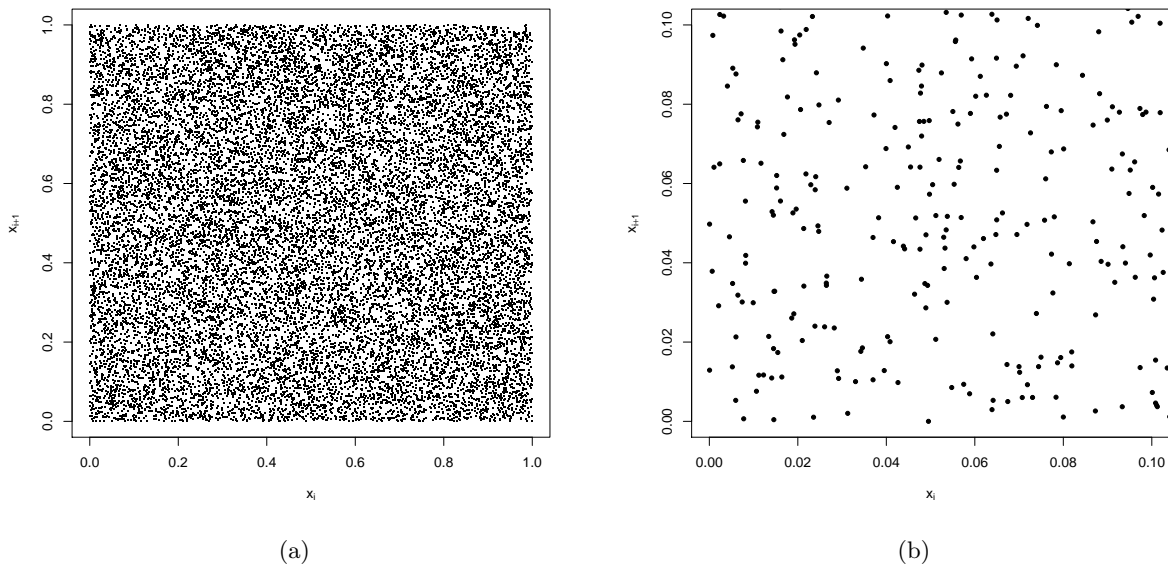


Figure 2.3: Scatter plot of consecutive numbers  $(x_i, x_{i+1})$  generated in a sequence by the `drand()` function of C++. a) A complete picture of the plot lacks any discernible pattern. b) A close-up view of the bottom left corner of the complete plot shows that sequence of numbers have apparently random behaviour

	<code>drand()</code>	<code>drand48()</code>
Computer I	$1.12 \times 10^8$	$1.32 \times 10^8$
Computer II	$1.64 \times 10^8$	$1.65 \times 10^8$

Table 2.1: Speeds of two commonly used LCGs on two different computers

### 2.4.1 Inverse transform sampling

The central idea behind inverse transform sampling is to map the uniform random numbers to the desired distribution (Chorin & Hald, 2009). This is made possible by the fact that the cumulative distribution function is monotonically increasing with range  $[0, 1)$ . Provided that a routine to generate uniformly distributed random numbers  $U$  in the interval  $[0, 1)$  exists, random numbers  $X$  with cumulative distribution function  $F_X(x)$  can be generated by using the transformation  $F_X^{-1}(U)$ , which is the inverse of the cumulative distribution function  $F_X(x)$ . Since  $F_X(x)$  is a monotonically increasing function, it must have an inverse. Consider the cumulative distribution function of the random variate  $X$

$$F_X(x) = P[X \leq x] = P[F_X^{-1}(U) \leq x], \quad (2.4)$$

which can be re-written as

$$P[F_X^{-1}(U) \leq x] = P[U \leq F_X(x)] = F_U[F_X(x)], \quad (2.5)$$

where  $F_U(u)$  is the cumulative distribution function of the uniformly distributed variate  $U$ . Since, the cumulative distribution function,  $F_U(u) = u$ , we have

$$F_U[F_X(x)] = F_X(x). \quad (2.6)$$

This proves that the transformation of uniform random numbers,  $U$  by  $F_X^{-1}(U)$  does indeed generate random streams generated with the cumulative distribution function  $F_X(x)$ . As an exam-

ple consider the exponential distribution with density function,  $f_X(x) = \lambda e^{-\lambda x}$ , with cumulative distribution function,  $F_X(x) = 1 - e^{-\lambda x}$ , then as per inverse transform sampling exponentially distributed random numbers can be generated using

$$X = -\frac{1}{\lambda} \log \frac{1}{1-U}. \quad (2.7)$$

However, this method is mostly useful to distributions whose cumulative distribution functions are explicitly invertible (Krishnan, 2015). As an example, the cumulative distribution function of the Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$  is

$$F_X(x) = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{x - \mu}{\sigma\sqrt{2}} \right) \right], \quad (2.8)$$

where erf is the error function, which does not have an analytically tractable inverse function. For such cases, this method is highly inefficient and more sophisticated techniques must be employed.

## 2.4.2 Transformation methods

An alternative to inverting the cumulative distribution function is to find a simple transformation scheme that would provide numbers generated as per the desired distribution. Different transformations produce a rich variety of statistics, for example consider the transformation

$$Y = X^2, \quad (2.9)$$

where  $X$  is distributed according to some density  $f_X(x)$  and cumulative distribution  $F_X(x)$ , then the cumulative distribution function of  $Y$  is

$$F_Y(y) = P(Y \leq y) = P(X^2 \leq y) = P(-\sqrt{y} \leq X \leq \sqrt{y}) = F_X(\sqrt{y}) - F_X(-\sqrt{y}). \quad (2.10)$$

Following this, the cumulative distribution can be differentiated to obtain the probability density function of  $Y$  as (Casella & Berger, 2002)

$$\begin{aligned} f_Y(y) &= \frac{dF_X(\sqrt{y})}{dy} - \frac{dF_X(-\sqrt{y})}{dy}, \\ &= \frac{1}{\sqrt{2y}} (f(\sqrt{y}) + f(-\sqrt{y})). \end{aligned} \quad (2.11)$$

If  $X$  is a Gaussian random variable with 0 mean and variance 1, then the density of  $Y$  is

$$f_Y(y) = \frac{1}{\sqrt{2\pi}} \frac{1}{\sqrt{y}} e^{-y/2}, \quad (2.12)$$

such a probability density function is known as chi-squared distribution with 1 degree of freedom. Similarly, the equivalent transformation for a random variable  $Y$ , distributed according to a chi-squared distribution with  $n$  degrees of freedom is

$$Y = X_1^2 + X_2^2 + \dots + X_n^2, \quad (2.13)$$

where  $\{X_1, \dots, X_n\}$  are independent Gaussian random numbers with mean 0 and variance 1. The probability density function of  $Y$  can be calculated using a  $n$ -fold convolution

$$f_n(y) = \frac{1}{2^{n/2} \Gamma(n/2)} y^{n/2-1} e^{-y/2}, \quad (2.14)$$

where  $\Gamma(\cdot)$  is the gamma function, and  $n$  is the number of degrees of freedom. Hence, sequences distributed according to some complex distribution might be easy to generate if a suitable transformation is found. An important example of such a method is presented in the next section.

### Box-Muller method

One of the earliest and most important techniques to generate Gaussian random numbers, is the Box-Muller method which transforms a pair of uniform random numbers, i.e from  $\mathbf{R}^2$  instead of  $\mathbf{R}$ , to a pair of normally distributed random numbers (Box *et al.*, 1958). The transformation is given by

$$\begin{aligned} Y_1 &= \sqrt{-2\sigma^2 \log(U_1)} \cos(2\pi U_2), \\ Y_2 &= \sqrt{-2\sigma^2 \log(U_1)} \sin(2\pi U_2), \end{aligned} \quad (2.15)$$

where  $U_1$  and  $U_2$  are uniformly distributed in  $(0, 1]$ . Geometrically, the pair  $(Y_1, Y_2)$  are points on a circle of radius  $\sqrt{-2\sigma^2 \log(U_1)}$ , with the  $\log(\cdot)$  function guaranteeing that the probability of finding a circle with larger radius decreases exponentially. The uniform distribution in angle ensures that the values for each co-ordinate average out to zero. For analyzing the distribution of the transformed variables  $Y_1$  and  $Y_2$  the Jacobian matrix is calculated as

$$\frac{\partial(u_1, u_2)}{\partial(y_1, y_2)} = -\frac{1}{2\pi} \exp\left(-\frac{1}{2}(y_1^2 + y_2^2)\right), \quad (2.16)$$

provided that  $U_1, U_2$  follow a bivariate uniform distribution,  $Y_1, Y_2$  are Gaussian distributed with zero mean and variance  $\sigma^2$ . This method transforms independent and uniformly distributed pairs of random numbers which can be represented by points on a plane, in a manner such that resulting variables are distributed along concentric circles. While the resulting variables are symmetric azimuthally, they tend to cluster around the centre thereby imitating the behaviour of Gaussian random numbers as shown in Fig.(2.4). A plot of a typical realization of Brownian motion in a 2-D domain obtained using the Box-Muller method is presented in Fig.(2.5). As expected the particle follows a “zig-zag” and haphazard trajectory.

### 2.4.3 Acceptance-rejection method

The acceptance-rejection method is particularly helpful when the probability density function  $f_X(x)$  of the desired variates is known but the respective cumulative distribution function and its inverse is not easily computable (Casella *et al.*, 2004). An outline of the method is as follows

1. A bounding function  $g(x)$  is considered such that  $g(x) \geq f_X(x)$  for all  $x$ , and a probability density function  $w_Y(x) = g(x)/c$  is constructed where  $c$  is the normalization factor given by  $c = \int_{-\infty}^{\infty} g(x)dx$ .
2. Independent and identically distributed random numbers  $y_i$  are generated as per the constructed probability density function  $w_Y(y)$ .
3. A random number uniformly distributed between  $[0, 1)$  is generated.
4. If the condition  $u_i \leq f_X(y_i)/g(y_i)$  is met then  $y_i$  is accepted to be a random number with probability density function  $f_X(x)$ , else it is rejected and steps 2 – 3 are followed until the condition is met.

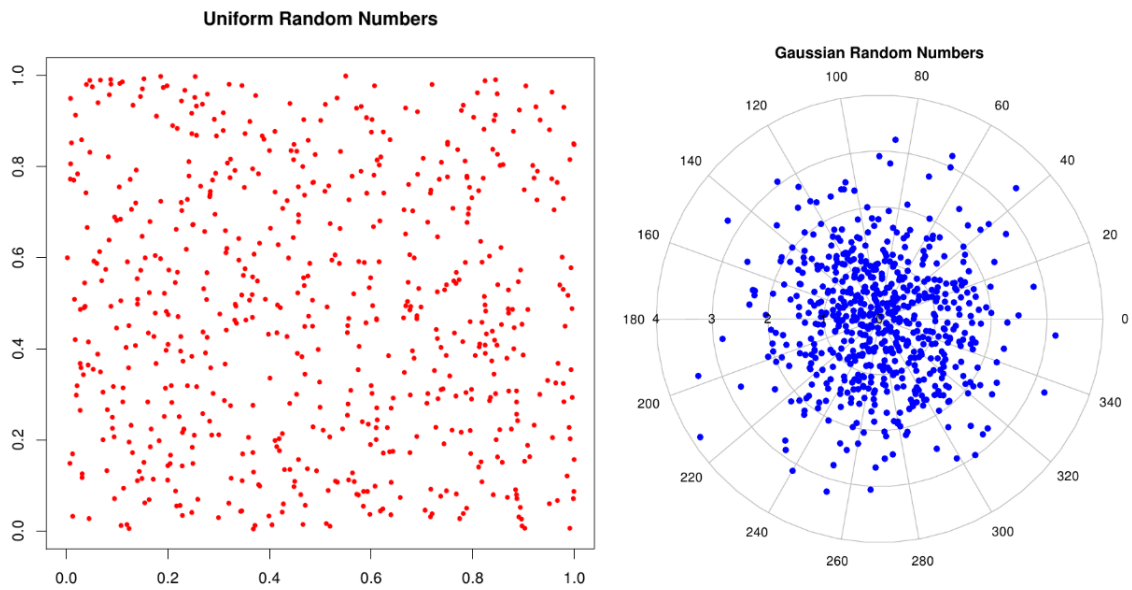


Figure 2.4: A geometrical interpretation of the Box-Muller method, suggests that points on a plane once transformed and plotted in a polar co-ordinate system tend to cluster around the origin.

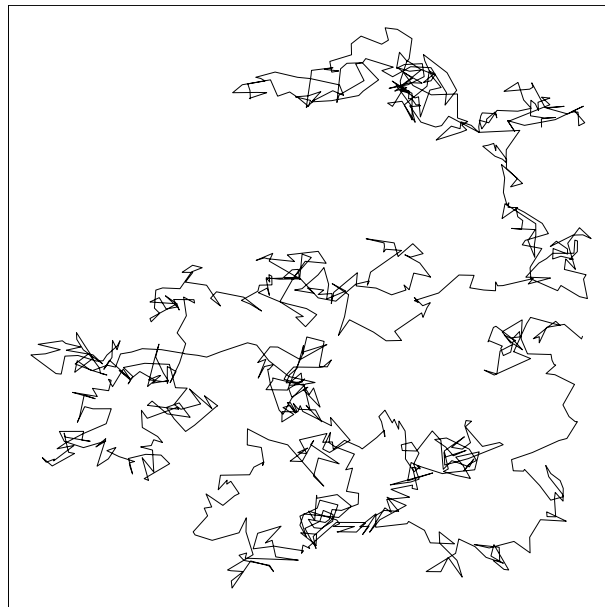


Figure 2.5: Typical realization of a Brownian particle immersed in fluid obtained using Gaussian random numbers generated by the Box-Muller method.

It can be shown that this procedure does indeed generate random sequences distributed according to  $f_X(x)$  (Krishnan, 2015). Let  $A$  be the event when the condition mentioned in the final step of the algorithm is met, then we have

$$P(X \leq x) = P(Y \leq x|A) = \frac{P(Y \leq x \cap A)}{P(A)}, \quad (2.17)$$

where  $X$  is a random variable distributed with required density function  $f_X(x)$  and  $Y$  is a random variable distributed according to the constructed density function  $w_Y(y)$ . For a given value of  $Y = y$  we have

$$P(A|Y = y) = P\left(U < \frac{f_X(y)}{g(y)}\right), \quad (2.18)$$

where  $U$  is uniformly distributed in the range  $[0, 1)$ . Similar to the inverse transform sampling method, the quantity  $P(U < f_X(y)/g(y))$  is the cumulative distribution function of the random variable  $U$  and hence we have the relation

$$P(A|Y = y) = \frac{f_X(y)}{g(y)}. \quad (2.19)$$

This can be used to calculate  $P(A)$

$$P(A) = \int_{-\infty}^{\infty} P(A|Y = y)w_Y(y)dy = \frac{1}{c}. \quad (2.20)$$

The numerator in Eq.(2.17) can now be calculated as

$$P(Y \leq x \cap A) = \int_{-\infty}^{\infty} P(A \cap Y \leq x|Y = y)w_Y(y)dy, \quad (2.21)$$

since  $y = Y$  and  $Y \leq x$ , the integrand can be simplified by considering the upper limit of integration to be  $x$  instead of  $\infty$

$$\begin{aligned} P(Y \leq x \cap A) &= \int_{-\infty}^x P(A|Y = y)w_Y(y)dy, \\ &= \int_{-\infty}^x \frac{f_X(y)}{g(y)} \frac{g(y)}{c} dy, \\ &= \frac{F_X(x)}{c}. \end{aligned} \quad (2.22)$$

Hence, we have

$$P(X \leq x) = \frac{F_X(x)/c}{1/c} = F_X(x). \quad (2.23)$$

This completes the proof for the acceptance-rejection method. An important use of this method is to generate random variates distributed according to beta distribution whose probability density function is given by

$$f_X(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1}(1-x)^{\beta-1}, \quad 0 \leq x \leq 1, \quad (2.24)$$

owing to the presence of the Gamma function, this particular density function does not yield a closed form solution for its cumulative density function and hence incapable of being inverted. The acceptance-rejection method was used to generate random numbers distributed according to the beta distribution with parameters  $\alpha = 6, \beta = 3$ . A histogram of these numbers is plotted

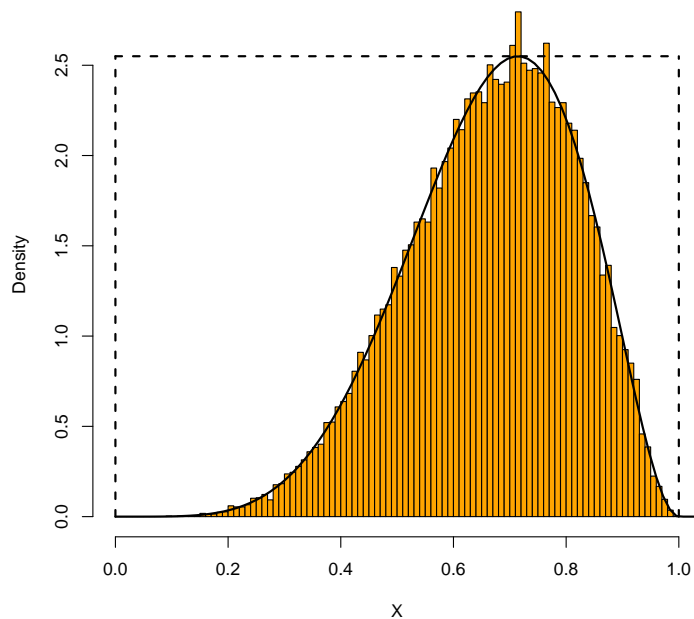


Figure 2.6: Plot comparing the histogram of beta-distributed random numbers generated by the acceptance-rejection method and the theoretical probability density function. The dotted lines indicate the bounding function and the domain of the variates.

and compared with theoretical density, presented in Fig.(2.6). As can be seen, there is good agreement between the two.

## 2.5 Qualities of a good PRNG

As we have determined, a PRNG capable of generating uniformly distributed random numbers can practically generate sequences from any specified distribution (Chorin & Hald, 2009). Large-scale scientific simulations require a large number of such streams at low computational costs for multiple realizations. The selection of the function  $F$ , which dictates the iterative scheme of the PRNG, must be such that it satisfies the following conditions

1. The sequence of numbers generated must be independently and identically distributed, must not be correlated and must satisfy as many statistical tests as possible.
2. For different realizations of the same system, one would need multiple sequences which are qualitatively similar, hence a PRNG must be able to generate significantly different sequences with relative ease.
3. The number of finite states ensures that the sequence starts repeating after a given period and given that stochastic simulations require high number of random sequences, the said period after which these sequences start repeating must be greater than the number required by the simulation.
4. Since the number of such sequences required for simulations is high, the function  $F$  must be such that it can be evaluated quickly.

In some sense, two rather contradictory requirements are imposed on  $F$ : it should be complex enough to produce a seemingly unpredictable stream of numbers but be simple enough to be

evaluated quickly. Satisfying these conditions is a well settled aspect in modern computing (Knuth, 1981). Families of PRNGs such as Xorshift, PCG, Mersenne Twister etc. have been shown to fulfill these conditions quite well and have found considerable success in the field of large-scale scientific simulations.



# Chapter 3

## Statistical Tests

### 3.1 Introduction

The random contributions to stochastic processes are incorporated in simulations by utilizing Pseudo-Random Number Generators (PRNGs). PRNGs are simple routines which produce stream of numbers which are apparently random. These sequences while being the outputs of deterministic programs closely mimic the behaviour of a predefined distribution and exhibit absence of any discernible pattern (Knuth, 1981). For scientific simulations, the requirement is that the stream of numbers produced by the PRNG should be empirically random, i.e, statistical tests should not be able to differentiate between such sequences and the ones found in nature (Chorin & Hald, 2009). Visual tests as briefly mentioned and exhibited in Chapter 2, can serve as a preliminary test to observe patterns in a generated sequence of numbers, but it is imperative to employ tests which provide a measure of the deviation of generated sequence from its expected behaviour. These tests allow us to establish the quality of PRNGs and if they are good enough to be used for scientific purposes (Knuth, 1981).

In the following sections the basic concepts and terminologies associated with statistical testing will be explained first. This is followed by descriptions of some canonical tests which can be used for random numbers from any given distribution. The chapter ends with explanations of various tests used for testing uniform random numbers, which can also be used to test random numbers from other distributions once the appropriate transformation to uniform distribution has been implemented.

### 3.2 Hypothesis testing

Statistical testing is employed when authenticity of a claim is to be studied. For example, if an oil company on an exploration project comes across oil reserves, they need to gauge the grade of the petroleum before the actual drilling as it is an expensive process. Hence for testing purposes a small batch from the reserve is extracted to characterize the reserves. Such a batch is known as sample of the population. If selected randomly and carefully the sample should be representative of the entire population, i.e, the characteristics of the sample, for example the specific gravity, should be close to that of the population itself (Freedman *et al.*, 2007). Any quantity derived from the sample is known as an estimate or a statistic. Suppose a populations comprises  $N$  elements ( $X_1, \dots, X_N$ ), with their individual statistics as

$$E[X_i] = \mu, \quad E[(X_i - \mu)^2] = \sigma^2, \quad i = 1, \dots, N, \quad (3.1)$$

where  $E[\cdot]$  is the expectation operator,  $\mu$  the population mean and  $\sigma^2$  the variance. From this population a sample of  $n$  independently and identically distributed elements such that  $n \ll N$  is drawn, thus the individual statistics of the elements are the same as the population and the mean and variance of the sample is

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n X_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2. \quad (3.2)$$

As it desirable to have the sample mirror the actual population, the statistics of the sample must converge to the true values in some limit

$$E[\hat{\theta}] = \theta, \quad (3.3)$$

where  $\hat{\theta}$  is the sample statistic and  $\theta$  the true population parameter. Estimates which satisfy this equation are termed unbiased, and denote that on an average the sample statistic converges to the population parameter (Krishnan, 2015). The expectation of the sample mean is calculated as

$$\begin{aligned} E[\hat{\mu}] &= E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] \\ &= \frac{1}{n} \cdot n\mu = \mu, \end{aligned} \quad (3.4)$$

which implies that sample mean is unbiased. The variance of the sample mean can be calculated as

$$\begin{aligned} \text{Var}(\hat{\mu}) &= E[(\hat{\mu} - \mu)^2] = E\left[\left(\frac{X_1 + \dots + X_n}{n} - \mu\right)^2\right], \\ &= \frac{1}{n^2} E[((X_1 - \mu) + \dots + (X_n - \mu))^2], \\ &= \frac{1}{n^2} E\left[\sum_{i=1}^n (X_i - \mu)^2 + \sum_{j=i+1}^n \sum_{i=1}^{n-1} (X_i - \mu)(X_j - \mu)\right], \end{aligned} \quad (3.5)$$

the quantity  $E[(X_i - \mu)(X_j - \mu)]$  for  $i \neq j$  is called correlation and effectively measures the degree of linear relationship between  $X_i$  and  $X_j$ , hence for independently distributed variables the correlation is 0. Then Eq.(3.5) reduces to

$$\begin{aligned} \text{Var}(\hat{\mu}) &= \frac{1}{n^2} E\left[\sum_{i=1}^n (X_i - \mu)^2\right], \\ &= \frac{1}{n^2} \cdot n\sigma^2 = \frac{\sigma^2}{n}. \end{aligned} \quad (3.6)$$

which indicates the fact that with an increase in sample size the sample mean better approximates the true mean of the population. The expectation of the sample variance can be calculated as

$$\begin{aligned}
E[\hat{\sigma}^2] &= E\left[\frac{1}{n}\sum_{i=1}^n(X_i - \hat{\mu})^2\right] = E\left[\frac{1}{n}\sum_{i=1}^n((X_i - \mu) - (\hat{\mu} - \mu))^2\right], \\
&= E\left[\frac{1}{n}\sum_{i=1}^n(X_i - \mu)^2 - \frac{2}{n}\sum_{i=1}^n(\hat{\mu} - \mu)(X_i - \mu) + \frac{1}{n}\sum_{i=1}^n(\hat{\mu} - \mu)^2\right], \\
&= \sigma^2 + E\left[\frac{1}{n}(\hat{\mu} - \mu)^2\sum_{i=1}^n 1\right] - E\left[\frac{2}{n}(\hat{\mu} - \mu)\sum_{i=1}^n(X_i - \mu)\right], \\
&= \sigma^2 + E[(\hat{\mu} - \mu)^2] - E\left[\frac{2}{n}(\hat{\mu} - \mu)(n\hat{\mu} - n\mu)\right], \\
&= \sigma^2 - E[(\hat{\mu} - \mu)^2], \\
&= \sigma^2 - E\left[\left(\frac{1}{n}\sum_{i=1}^n X_i - \mu\right)^2\right] = \left(1 - \frac{1}{n}\right)\sigma^2.
\end{aligned} \tag{3.7}$$

This proves that the sample variance is actually smaller than the true population variance, hence making it a biased estimator. This issue can be solved by considering the quantity

$$\bar{\sigma}^2 = \frac{1}{n-1}\sum_{i=1}^n(X_i - \hat{\mu})^2, \tag{3.8}$$

the expectation can be calculated in a manner similar to the procedure for  $\bar{\sigma}^2$

$$E[\bar{\sigma}^2] = E\left[\frac{1}{n-1}\sum_{i=1}^n(X_i - \hat{\mu})^2\right], \tag{3.9}$$

$$= \frac{n}{n-1}E\left[\frac{1}{n}\sum_{i=1}^n(X_i - \mu)^2\right], \tag{3.10}$$

$$= \frac{n}{n-1} \cdot \left(1 - \frac{1}{n}\right)\sigma^2 = \sigma^2. \tag{3.11}$$

proving that it is an unbiased estimator. Hence, it is reasonable to assume that effects of the proposed claim on the sample can be generalized to the entire population (Freedman *et al.*, 2007). It must be noted that sample statistics only approximates the true population parameters, given the sample mean and variance the probability for the true mean to lie in a certain range is

$$P(|\hat{\mu} - \mu| \geq k\bar{\sigma}) \leq \frac{1}{k^2}. \tag{3.12}$$

This is known as Chebyshev's inequality which is used for calculating the chances of obtaining a random variate  $k$  standard deviations away from the mean (Chorin & Hald, 2009). Stricter bounds can be imposed by considering that the sample mean is in fact the sum of independently and identically distributed variates and should hence tend to a Gaussian distribution. This allows one to calculate the probability of observing the true population parameter in a given range using the definition of Gaussian distribution

$$P(\hat{\mu} - k \cdot \bar{\sigma} \leq \mu \leq \hat{\mu} + k \cdot \bar{\sigma}) = F(\hat{\mu} + k \cdot \bar{\sigma}) - F(\hat{\mu} - k \cdot \bar{\sigma}) \tag{3.13}$$

where  $F(\cdot)$  is the cumulative distribution function of the Gaussian distribution function with mean  $\hat{\mu}$ , variance  $\bar{\sigma}$  and  $k$  the parameter to decided the range of the interval. This range is termed as the confidence interval and the probability associated with it is known as the confidence level

(Freedman *et al.*, 2007). For the given case, the quantitative measure that would allow one to make an inference or the qualitative conclusion is the average specific gravity of the batch extracted for testing. The average density is found within a certain confidence level and can be compared against previously available data to verify the quality of petroleum. However, it is entirely possible that the difference between the observed value and the expected value stems from “chance error”, i.e the inference made about the claim may change subject to the sample selected. Consequently a well-defined procedure, commonly known as hypothesis testing in literature, has been established to study and provide accurate conclusions from available data. Given that the probability density of specific gravity for petroleum  $f(x)$  and the probability density  $g(x)$  of some other suspected substance, say paraffin, is known then a parameter called significance level can be defined as

$$\alpha = \int_{t_{cut}}^{\infty} f(x)dx, \quad (3.14)$$

which is the total probability of observing the test statistic  $t$  to be at least  $t_{cut}$ , if the substance is petrol. Similarly, another parameter  $\beta$  can be defined as

$$\beta = \int_{-\infty}^{t_{cut}} g(x)dx, \quad (3.15)$$

which is the total probability of observing the test statistic  $t$  to be at most  $t_{cut}$ , if the substance is paraffin. The quantity  $(1 - \beta)$  is called power of the test (Cowan, 1998). Once these parameters are established, the test statistic is then calculated. In present example, say five batches are extracted from the reserves and their average specific gravity,  $t$  is found. The total probability of finding the specific gravity of petrol at least as extreme as  $t$  is

$$P = \int_t^{\infty} f(x)dx. \quad (3.16)$$

The quantity  $P$  is known as the  $p$ -value of the statistic, it is the total probability of observing a statistic at least as unlikely as  $t$ . Hence it is a good quantitative measure to facilitate the qualitative measure. Typically, it is understood the claim is rejected if the  $p$ -value is much less than the significance level. If it is close then experiments are repeated and in case the  $p$ -value is much greater than the significance level then the claim is not rejected. Essentially, the significance level is the total probability of rejecting the claim when it is true and the power of test is the total probability of not rejecting the claim it is false. The value of  $\alpha$  is set by the practitioner and depends on the field and requirement of the study. It is usually accepted that  $P$ -value close to  $\alpha = 0.05$  indicate some evidence against the claim (Freedman *et al.*, 2007). This entire procedure can then be formalized and enumerated as follows

1. **State the hypotheses:** A statement of the expected outcome, conventionally known as the null hypothesis, denoted by  $H_0$ , is stated and an alternate hypothesis, denoted by  $H_a$  different from the expectation is stated. For given example, the null and alternate hypotheses are:  
 $H_0$ :  $\rho \leq \rho_0$  and hence petroleum.  
 $H_a$ :  $\rho \geq \rho_0$  and hence paraffin.
2. **Decide the significance value:** A value for the significance level is set. The higher this value the more purity is assured in the present case, as Fig.(3.1) suggests. The value  $\alpha = 0.05$  is considered to be “statistically significant”.
3. **Characterize the expected behaviour:** The various parameters associated with the expected outcome either theoretically or experimentally. The parameter derived for present case is the probability density function of the specific gravity of high grade petroleum.

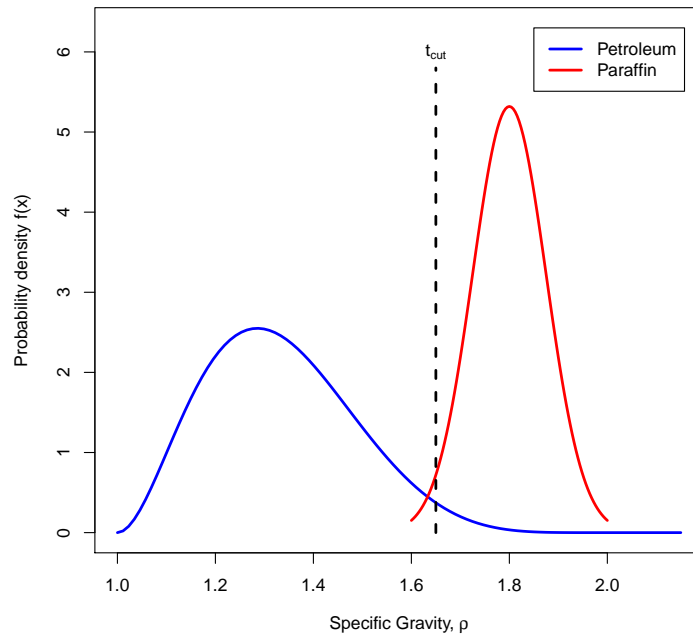


Figure 3.1: Sketch of the probability density functions of the specific gravity of petrol and paraffin for the hypothetical test case.

4. **Construct the test statistic:** From the observed values of the experiments, a test statistic is constructed. For given example, the average specific gravity of the five batches,  $t$ , is the test statistic.
5. **Calculate the  $p$ -value:** The  $p$ -value or the observed significance level, can be calculated using Eq.(3.16). An inference is made about the population using the  $p$ -value. Table(3.1) lists the interpretation of different  $p$ -values

$p$ -value	Interpretation
$\alpha > 0.1$	No evidence against null hypothesis
$0.05 < \alpha < 0.1$	Slight evidence against null hypothesis
$0.01 < \alpha < 0.05$	Moderate evidence against null hypothesis
$\alpha < 0.01$	Strong evidence against null hypothesis

Table 3.1: Interpretation of  $p$ -values for a test of significance

Tests of significance and the resulting  $p$ -values provide evidence against the null hypothesis, no statistical test actually confirms it. Hence, by convention a null hypothesis can only be rejected and not accepted. Fig.(3.1) suggests that such a procedure would lead to two kinds of errors (Cowan, 1998). The first error is to reject the null hypothesis when it is actually true, this is known as a type I error. Such an error is encountered when  $P \leq \alpha$  despite the fact that the null hypothesis is actually true, for example it is entirely possible that the specific gravity of the batch of petrol from reserves is rather high and lies somewhere on the right tail of  $f(x)$  beyond  $t_{cut}$  but the hypothesis would be rejected nonetheless. The second kind of error is encountered when the null hypothesis is not rejected while the alternate hypothesis is true, this is known as type II error. This is possible when  $P \geq \alpha$  but the actual substance is paraffin with its specific gravity lying somewhere on the left tail of  $g(x)$  and such that  $t \leq t_{cut}$ . Hence, if purity

is prioritized the value of  $\alpha$  is set higher by shifting  $t_{cut}$  to the left and if quantity is prioritized the value of  $\alpha$  is set lower by shifting  $t_{cut}$  to the right.

### 3.3 Testing for PRNGs

PRNGs generate stream of numbers which behave randomly despite being outputs of deterministic algorithms, hence it is imperative to quantify how well observed datasets compare with expected values, which help determine if they could be used for simulation purposes. Since various characteristics of random numbers can be derived theoretically, tests for individual characteristic can be constructed to assess the quality of PRNGs (Knuth, 1981). The methodology of constructing such a test is essentially the same as the procedure outlined in the previous section and can be summarized as

1. The various characteristics of a population can be determined theoretically. One simple check could to be test for the deviation from the expected mean of the population. Such statistics are usually termed as point estimators (Casella & Berger, 2002). Another class of tests are specific to the distribution being tested. For example, exponentially distributed random numbers must have the property  $E[X^n]/E[X^{n-1}] = n\lambda$ .
2. The PRNG to be tested is then used to generate a dataset which can then be compared against the theoretical prediction.
3. The generated dataset is then used to calculate the relevant the test statistic(s).
4. Once the test statistic is calculated it can be used to calculate the  $p$ -value, as previously explained.

A test can be devised corresponding to a characteristic of the theoretical distribution. An ideal PRNG would be able to satisfy as many statistical tests as possible. However, it must be noted that if a test satisfies  $n$  statistical tests, there is absolutely no guarantee that it will satisfy the  $(n + 1)^{\text{th}}$  test (Knuth, 1981). Tests of significance are inherently negative in nature, hence tests of randomness are capable of specifying major anomalies or flaws of PRNGs and do not confirm its quality, hence it is important to conduct a rather high number of tests on a PRNG before deciding on its quality.

### 3.4 Preliminary distribution-free tests

There exists a few tests which can be used for any dataset regardless of their distribution. Such tests are known as distribution-free tests and can be used as a screening process for PRNGs (Knuth, 1981). We outline some of the basic tests in the following sections.

#### 3.4.1 Pearson's $\chi^2$ test

Pearson's  $\chi^2$  test, is an important statistical tests for random numbers (Knuth, 1981). The core idea of this test is to quantify the deviation between the histogram of observed data and the theoretical probability density function (Pearson, 1900). Consider the histogram of uniformly generated random numbers presented in Fig.(3.2). Some bins have more than the expected number and others fewer, it is necessary to understand whether the deviation is within acceptable limits. For every bin, a simple statistic can be constructed as

$$t' = \frac{(\nu_i - \mu_i)}{\mu_i}, \quad (3.17)$$

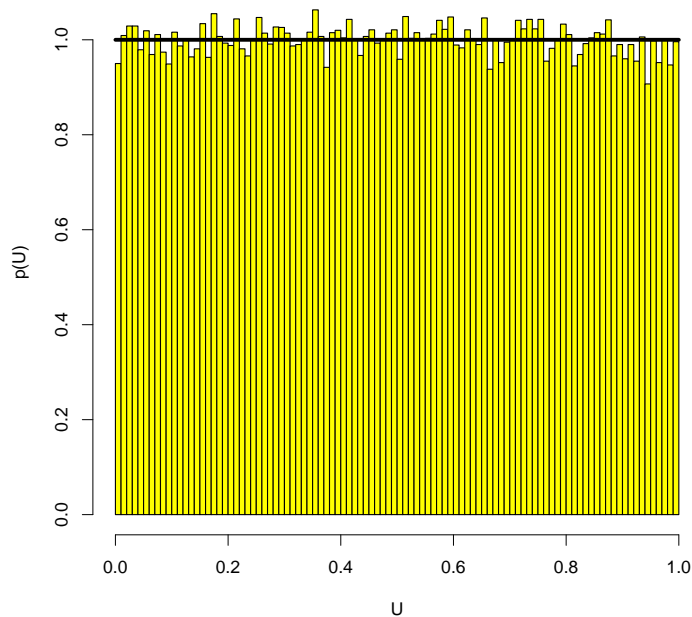


Figure 3.2: Plot comparing the histogram of uniformly distributed random numbers generated by a PRNG and the theoretical probability density function

where  $\nu_i$  is the number of data points is present in the  $i^{\text{th}}$  bin and  $\mu_i$  is the expected number for that particular bin. If the PRNG is indeed high quality, this statistic should fluctuate around zero for different samples. However it is not practical to test for such behaviour with large number of bins. An alternative is to consider the sum total of deviations from each bin. For this purpose, the deviation calculated must be in absolute terms as the sum total of all deviations might end up as zero as in the case of uniform distribution. Hence, the test statistic is constructed as

$$t = \sum_{i=1}^N \frac{(\nu_i - \mu_i)^2}{\mu_i}. \quad (3.18)$$

This is similar to the prevalent technique of calculating the mean square error in many areas of statistics and physics. We show that for the null hypothesis to hold true  $t$  must in fact follow a chi-squared distribution with  $N - 1$  degrees of freedom, which is equivalent to the sum of  $N - 1$  independently and identically distributed Gaussian random numbers. A detailed proof of this statement can be found in Appendix A. Once the test statistic  $t$  is calculated, the p-value is obtained using

$$P = \int_t^{\infty} f_{N-1}(x)dx, \quad (3.19)$$

where  $f_{N-1}$  is the chi-squared probability density function with  $(N - 1)$  degrees of freedom, as defined in the previous chapter. The  $p$ -value obtained via this procedure is a mathematical statement comparing the histogram of observed data and expected distribution. It has been noted that for this test to hold good, in general a minimum of 5 counts are required in each bin (Knuth, 1981).

This procedure can be best explained by an example. Consider that the fairness of a die is to be tested. Then following the procedure we first conduct the experiments, i.e, toss the

die 120 times. A set of sample observations are tabulated in Table(3.2). Given these values, the fairness of the die might be suspected as the face with 6 dots up considerably fewer times than expected. In order to establish the quality (fairness) of the die, a chi-squared test can be performed to quantify the deviation of the available values from expected behaviour and check if there is strong evidence against the fairness of the die.

Observation	Frequency
1	17
2	24
3	22
4	23
5	20
6	14

Table 3.2: Observations of a die experiment

The null and alternate hypotheses are stated as

$H_0$ : The die fair.  $p(X = i) = 1/6$  for  $i = 1, \dots, 6$ .

$H_a$ : The die is biased.

We decide on the significance level to be 0.05, as it is considered to be statistically significant in literature (Freedman *et al.*, 2007). The value of test statistic is now calculated

$$t = \left( \frac{(17 - 20)^2}{20} + \frac{(24 - 20)^2}{20} + \frac{(22 - 20)^2}{20} + \frac{(23 - 20)^2}{20} + \frac{(20 - 20)^2}{20} + \frac{(14 - 20)^2}{20} \right) = 3.7. \quad (3.20)$$

The  $p$ -value can then be calculated by using the definition of the chi-square distribution for  $N = 5$  degrees of freedom. This integral can be calculated using a chi-square table or a statistical package such as *R*. The value for proposed example is

$$P = \int_{3.7}^{\infty} f_5(x) dx = 0.41. \quad (3.21)$$

this denotes that there is nearly a 41% chance that an ideal die would produce results at least as extreme as the ones obtained, and since this much higher than set significance level, there is no strong evidence against the null hypothesis and the fairness of the die.

### 3.4.2 Kolmogorov-Smirnov Test

As opposed to the Pearson's chi-squared test, the Kolmogorov-Smirnov test checks for the adherence of random numbers with their respective cumulative distribution functions (Massey Jr, 1951). The empirical distribution function,  $F_n$ , of a sample of  $n$  numbers is calculated using

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[-\infty, x]}(X_i), \quad (3.22)$$

where  $I_{[-\infty, x]}(X_i)$  is the indicator function defined as

$$I_{[-\infty, x]}(X_i) = \begin{cases} 1 & X_i \leq x \\ 0 & \text{otherwise.} \end{cases}$$

The test statistic can then be calculated as:

$$D_n = \sup_x |F_n(x) - F(x)|, \quad (3.23)$$



where  $\sup_x$  is the maxima of the set of distances between the cumulative distribution function expected under the null hypothesis and the observed distribution function. It is found that if  $F(x)$  is continuous then the quantity  $\sqrt{n}D_n$  converges to the Kolmogorov distribution whose cumulative distribution is

$$F(x) = 1 + 2 \sum_{k=1}^{\infty} (-1)^k e^{-2k^2 x^2} \quad (3.24)$$

A  $p$ -value is then returned on the basis of the Kolmogorov distribution (Wang *et al.*, 2003).

### 3.4.3 Ljung-Box Test

An important quantity while testing for dependence in datasets is the correlation, defined as

$$C(X, Y) = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}, \quad (3.25)$$

where  $X$  and  $Y$  are two random variates, while  $\mu$  and  $\sigma$  denote their mean and standard deviation respectively. Correlation,  $C(X, Y)$ , measures the linear dependence between the two quantities, i.e if  $C(X, Y)$  is positive it denotes that an increase in  $X$  is likely to be accompanied with an increase in  $Y$  (Freedman, 2009). In case of perfectly linear relationship,  $C(X, Y) = 1$  ( $C(X, Y) = -1$  in case of perfect anti-correlation), and it should be as close to zero as possible if the data are completely independent. Similarly, the autocorrelation is defined as

$$R(\tau) = \frac{E[(X_\tau - \mu)(X_{t+\tau} - \mu)]}{\sigma^2}, \quad (3.26)$$

which basically defines the relationship between a random variable with itself at a given time lag  $\tau$ . This can be extremely significant for cases such as simulation of Brownian motion wherein the random force must be uncorrelated in time. Hence, for a PRNG it would be desired for the autocorrelation to be as close to zero as possible. The Ljung-Box test is a very stringent statistical test to measure the independence of data in a time series (Ljung & Box, 1978). For a series with  $n$  elements, the autocorrelation residual at lag  $i$  is defined as

$$R_i = \frac{\sum_{j=1}^{n-i} X_j \cdot X_{j+i}}{\sum_{j=1}^n X_j \cdot X_j}. \quad (3.27)$$

The test statistic is then defined as the normalized sum of square of the autocorrelation residuals for lags upto  $i$ th index

$$t = n(n+2) \sum_{i=1}^k \frac{R_i^2}{n-i}, \quad (3.28)$$

It has been shown that  $t$  must follow a chi-squared distribution with  $k$  degrees of freedom, as the sum of these residuals are equivalent to the sum of squares of independent normally distributed random numbers (Box & Pierce, 1970).

We selected  $N = 10^8$  random numbers generated from drand48 and Mersenne Twister. Table(3.3) enumerates the  $p$ -values for different methods of random number generation for all the three tests.

The  $P$ -values mentioned in Table(3.3) indicate that the preliminary tests did not find any evidence against the null hypothesis – sequences of numbers generated by drand48 and Mersenne Twister are independently and identically distributed uniform random numbers.

	Chi-Squared	Kolmogorov-Smirnov	Ljung-Box
drand48	0.12	0.76	0.25
Mersenne Twister	0.35	0.47	0.88

Table 3.3: p-values for different tests

### 3.5 Tests for uniform random numbers

Since most PRNGs are capable of generating uniformly distributed integers, most efforts to construct statistical tests have been directed towards this particular distribution (Knuth, 1981). If one is interested in testing the quality of a PRNG which generates non-uniform distribution then the numbers must be first transformed to uniform random numbers using the inverse sampling method (Thomas *et al.*, 2007). For exponential numbers:

$$U = e^{-\xi}$$

Similarly, for Gaussian random numbers we have:

$$U = \frac{1}{2} \operatorname{erfc} \left( -\frac{\eta}{\sqrt{2}} \right)$$

Tests described in the previous section are quite useful in determining a dataset's adherence to the distribution and the correlation between the numbers generated. In order to better examine the behaviour of PRNGs, tests have been designed which compare a particular characteristic of generated sequence to the theoretical value of independent and identically distributed numbers (L'Ecuyer & Simard, 2007). However, it is to be noted these tests only verify how well sequences generated from PRNGs satisfy a given aspect of the actual distribution and in no way assure or confirm that these sequences are indeed independently and identically distributed as per a specified distribution. It is desired that a given PRNG passes as many tests as possible with satisfactory  $p$ -values and not fail any test consistently.

#### 3.5.1 Serial Test

This test is performed in order to check if the PRNG generates successive pairs of uniformly distributed random integers in the range  $[0, d)$ . A sample of  $n$  pairs is generated using the PRNG with each element of pair assumed to be uniformly distributed under the null hypothesis. Each time the pair  $(X_{2i}, X_{2i+1}) = (q, r)$  the count for that particular pair  $(q, r)$  is increased. The probability of finding a given pair, provided that the generated pairs are independently and identically distributed is

$$f(X_{2i} = q, X_{2i+1} = r) = \frac{1}{d} \times \frac{1}{d} = \frac{1}{d^2}, \quad (3.29)$$

hence each pair is expected to be observed  $n/d^2$  times, i.e, the product of the probability and the total number of samples. The deviation of observed from expected values can then be used to obtain a chi-squared statistic in a manner similar to the Pearson's chi-squared test. This test can naturally be extended to triples, quadruples etc. However, this would significantly increase the sample size that would be needed, as the number of possible combinations increase with a factor of  $d$  (Knuth, 1981).

#### 3.5.2 Gap Test

The probability that an identically and independently generated random number drawn from a PRNG is contained in the interval  $[a, b)$  after  $r$  trials is

$$f_{\text{gap}}(p; r) = p(1 - p)^{(r-1)}$$

where  $p$  is the probability of finding a random number between  $[a, b)$ . The distribution  $f_{\text{gap}}(p; r)$  is called a geometric distribution and the number of trials before a number is generated in the interval  $[a, b)$  is termed as “gaps”. Although this test can be used for any distribution, for most test suites the uniform distribution is chosen with  $0 \leq a < b < 1$ . Once enough samples of  $r$  have been drawn the observed counts for each gap is compared against the theoretical distribution, which provides a chi-squared statistic.

### 3.5.3 Marsaglia’s tests for empirical randomness

Developed in the last decade, it is suite of three highly stringent tests namely — Birthday Spacings, GCD and Gorilla. These tests are designed for 32-bit uniform random integers, with each test characterizing a significantly different quality and thus providing a holistic overview of the PRNG (Marsaglia & Tsang, 2002). The birthday spacings test checks for the adherence to the distribution, the GCD test for pairwise independence and the Gorilla test examines the manner in which a sequence appears. These three tests will be briefly described in the following sections.

#### Birthday Spacings Test

For this test,  $m$  uniformly distributed integers  $U_1, \dots, U_m$ , are generated the interval  $[0, n)$  and then sorted. Then this sorted list provides  $(n - 1)$  spacings

$$S_i = U_{(i+1)} - U_{(i)} \quad (3.30)$$

where  $U_{(i)}$  are elements of the sorted list. The number of duplicates of spacings is asymptotically Poisson distributed with parameter  $\lambda = m^3/(4n)$  (Marsaglia & Tsang, 2002). This suite chooses  $m = 4096$  and  $n = 2^{32}$ , hence ending up with  $\lambda = 4$ . The program generates  $m$  uniform random integers in the range  $[0, 2^{32})$  from the PRNG to be tested. These numbers are then sorted and the spacing between them is calculated, following which the number of duplicates is enumerated, this process is repeated 5000 times. Pearson’s  $\chi^2$  test is then performed on observed data against the expected Poisson distribution and a  $p$ -value is returned. A typical output of the Birthday test for a good RNG is shown in Table(3.4).

Birthday spacings test: 4096 birthdays,  $2^{32}$  days in year

Table of Expected vs. Observed counts:

Duplicates	0	1	2	3	4	5	6	7	8	9	$\geq 10$
Expected	91.6	366.3	732.6	976.8	976.8	781.5	521.0	297.7	148.9	66.2	40.7
Observed	81	387	715	971	915	802	531	340	155	62	41
$(O-E)\hat{2}/E$	1.2	1.2	0.4	0.0	3.9	0.5	0.2	6.0	0.3	0.3	0.0

Birthday Spacings:  $\text{Sum}(O-E)\hat{2}/E = 14.023$ ,  $p = 0.828$

Table 3.4: An output for the Birthday Test

#### Gorilla Test

The idea for Gorilla Test is rooted in the monkey test which hypothesize that a monkey with a typewriter would produce absolutely random words. Streams of integers are generated in the range  $[0, 2^{26})$ , and represented in the binary format. For each bit position, sequences of 26 bits (1’s and 0’s) are considered to be a “word”. The total number of possibles words is then  $2^{26}$ . The test generates  $(2^{26} + 25)$  numbers and computes the missing sequences. The number of

these missing sequences,  $x$ , is approximately normally distributed with mean  $\mu = 24687971$  and variance standard deviation  $\sigma = 4170$ . Hence,  $\Phi((x - \mu)/\sigma)$  is uniformly distributed in  $[0, 1)$ , where  $\Phi()$  is the cumulative normal distribution function. This quantity provides the  $p$ -value for the test for each bit position. Once the  $p$ -value of each bit position is obtained, the 32 values are then subjected to a Kolmogorov-Smirnov test with the presumption that they are distributed uniformly between  $[0, 1)$ , which provides a final  $p$ -value for the entire test. A typical output of a Gorilla test for a good RNG is shown in Table(3.5):

Gorilla test for $2^{26}$ bits, positions 0 to 31:								
Note: lengthy test—for example, 20 minutes for 850MHz PC								
Bits 0 to 7— — — >	0.845	0.941	0.574	0.671	0.355	0.850	0.297	0.703
Bits 8 to 15— — — >	0.679	0.720	0.045	0.295	0.906	0.580	0.109	0.481
Bits 16 to 23— — — >	0.112	0.228	0.378	0.630	0.065	0.889	0.436	0.458
Bits 24 to 31— — — >	0.413	0.667	0.431	0.841	0.471	0.732	0.211	0.039
KS test for the above 32 p values: 0.065								

Table 3.5: An output for the Gorilla Test

### GCD Test

The GCD test relies on the statistical implications of Euler's GCD test used to determine the greatest common divisor of two integers. The procedure for calculating the GCD, is best explained by an example, consider two integers,  $u = 366$  and  $v = 297$ .

$$\begin{aligned}
 366 &= 1 * 297 + 69 \\
 297 &= 4 * 69 + 21 \\
 69 &= 3 * 21 + 6 \\
 21 &= 3 * 6 + 3 \\
 6 &= 2 * 3 + 0
 \end{aligned}$$

This procedure produces in a list of two identically and independently distributed parameters:

- The number of iterations required to find the GCD.
- The GCD itself

from empirical studies it has been shown that the first variable, steps to GCD, is normally distributed with mean,  $\mu = 18.5785$  and variance  $\sigma = 3.405$ , which can be approximated by a Binomial Distribution with parameters  $n = 50$  and  $p = 0.376$ . The distribution of the GCDs is close to the theoretical limit  $P(i = GCD) = c/i^2$  where  $c = 6/\pi^2$ . The test follows the procedure for  $10^7$  such pairs. The observed data is then compared against the expected values and a  $p$ -value is returned through the familiar  $\chi^2$  test procedure. There are many statistical tests, each built on a particular quality of the theoretical distribution. Many such canonical tests are compiled in the TestU01 library (L'Ecuyer & Simard, 2007). The Crush battery of tests, which implements 96 highly stringent statistical tests (with 144 statistics), was used to test the quality of drand48 and Mersenne Twister. The results are tabulated in Table(3.6), both these generators fail the Linear complexity test consistently.

RNG	Statistics Failed	Systemic
drand48	5	LinearComp, HammingWeight
Mersenne Twister	2	LinearComp

Table 3.6: Performance of drand48 and Mersenne Twister in Crush battery of tests

## 3.6 Outlook

The results of stochastic simulations significantly depend on the quality of the RNG used, as ones with poor statistical qualities can lead to erroneous results (Marsaglia, 1968). Although no number of statistical tests can characterize the behaviour of a given PRNG in entirety, it is important for a PRNG to satisfy at least all canonical tests before they are employed for scientific simulations.



# Chapter 4

## Molecular Dice

### 4.1 Introduction

Noise inherent to a large class of stochastic processes does not follow a uniform distribution, but other specified distributions such as Gaussian, exponential and Poissonian (Gardiner, 1985*a*). These non-uniformly distributed random variables are generated on computers by applying suitable algebraic transformations on uniformly distributed numbers. These transformations often involve multiple evaluations of computationally expensive transcendental functions, as in the case of Box-Muller method to generate Gaussian random numbers or simple inversion method for exponential random numbers. While uniform random number generation requires  $\approx 20 - 30$  FLOPS, Gaussian and exponential random number generation requires  $\approx 200$  FLOPS (Thomas *et al.*, 2007). Thus, while on a typical modern processors, it is possible to generate uniform random numbers which passes most statistical tests - at the rate of  $10^8 - 10^9$ /second, one can generate Gaussian or exponential random numbers only at the rate of  $10^6 - 10^7$ /second, hence becoming a bottleneck for many algorithms.

An alternate route is to choose the function  $F$  based on a chaotic map wherein specific choices would lead to unique distributions (Kohda & Tsuneda, 1997). Successful application of this idea, if any, is scarce, as given a distribution from which random numbers are to be sampled, a chaotic map that would asymptotically converge to the desired distribution are often not readily available. A notable exception is the tent map (Yoshida *et al.*, 1983), which produces a stream of uniform random numbers. This map is known to be qualitatively similar to the discrete version of the Navier-Stokes equation and has been used to explain the emergence of chaotic turbulent motion from a fully deterministic evolution equation (Frisch, 1995). This suggests that a discrete model of fluid motion can be a source of stochastic dynamics. On the other hand, Boltzmann showed that a stochastic description emerges at the mesoscopic level from the microscopic deterministic motion of particles (Chapman & Cowling, 1970). This symbiotic relationship between deterministic and stochastic descriptions suggests that they are entwined in the non-linear evolution of fluid motion. This suggests that fluid motion can act as a source of randomness, hence simulations of fluid flow at mesoscopic scale wherein the Boltzmann equation is solved using particle dynamics should be capable of generating stream of random numbers. The distribution of these sequences can be gauged from the Boltzmann equation.

In this chapter, the basics of kinetic theory are explained first, followed by details of the distributions of various quantities at equilibrium. Various numerical methods used to solve the Boltzmann equation are then used as model PRNGs. Finally, a new algorithm is proposed which attempts to maximize the rate of generation of random numbers while maintaining high statistical quality.

### 4.2 The Boltzmann equation

The kinetic theory of gases developed by Boltzmann and Maxwell relies on describing the time-evolution of probability distributions (Cercignani, 1998). The entity central to this theory is the distribution function,  $f(\mathbf{x}, \mathbf{v}, t)$ , which is the probability density of finding a particle in  $(\mathbf{x} + d\mathbf{x})$ , possessing velocity  $(\mathbf{v} + d\mathbf{v})$  at a given time  $t$  (Chapman & Cowling, 1970). The Boltzmann Equation describes the evolution of the distribution function in time while accounting for the effects of binary collisions between particles:

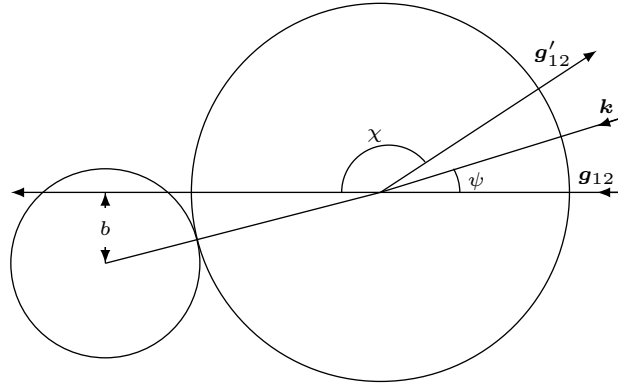


Figure 4.1: The geometry of collision between two particles.  $\mathbf{k}$  is the vector joining the centres of the two molecules and bisecting the angle made between the pre and post-collisional velocities,  $\mathbf{g}_{12}$  and  $\mathbf{g}'_{12}$ . The quantity,  $b$  is known as the impact parameter and is the perpendicular distance from the centre of the molecule to  $\mathbf{g}_{12}$ , and  $\chi$  is the angle between  $\mathbf{g}_{12}$  and  $\mathbf{g}'_{12}$

$$\frac{\partial f}{\partial t} + c_\alpha \frac{\partial f}{\partial x_\alpha} + \frac{F_\alpha}{m} \frac{\partial f}{\partial c_\alpha} = \Omega^{\text{BM}} \quad (4.1)$$

where  $\Omega^{\text{BM}}$  is the change in distribution function arising from random collisions between two particles. The kinetics of collision can be understood by considering two particles with velocities  $\mathbf{c}$  and  $\mathbf{c}_1$  which undergo an elastic collision, then their post-collisional velocities are

$$\mathbf{c}' = \mathbf{G} - \mathbf{g}_{21}' \quad \mathbf{c}'_1 = \mathbf{G} + \mathbf{g}_{21}' \quad (4.2)$$

where  $\mathbf{G}$  is the centre-of-mass velocity and  $\mathbf{g}_{21}'$  is the relative velocity post-collision. For conservation of energy it is necessary that the magnitude of  $\mathbf{g}_{21}'$  is the same as  $\mathbf{g}_{21}$ . The direction of the the relative velocity then provides a complete picture of the binary collision dynamics. Fig.(4.1) shows the geometry of a typical collision. The Boltzmann collision kernel for a dilute gas is

$$\Omega^B = \int (f' f'_1 - f f_1) |\mathbf{c} - \mathbf{c}_1| \alpha_{12} d\mathbf{e} d\mathbf{c}_1, \quad (4.3)$$

where  $\mathbf{e}'$  is the unit vector of the relative velocity post-collision (hence  $d\mathbf{e}'$  being a differential solid angle) and  $\alpha_{12}$  is a positive scalar function given by  $b|\partial b/\partial \chi|/\sin \chi$ . The derivation of the Boltzmann equation is based on the assumption of molecular chaos (also known as Stosszahlansatz), which hypothesizes that the velocities of colliding particles are statistically independent. Particles undergoing ballistic motion collide with other particles which in turn alters their course, and hence imparts randomness to the system.

Boltzmann equation satisfies the conservation equations and recovers the Navier-Stokes equation in the limit of low Knudsen numbers. Hence, it is an important tool which relates the random motion of molecules to fluid flow. In addition to preserving the conservation laws and retaining hydrodynamical behaviour, the Boltzmann equation also extends the second law of thermodynamics, known as the  $H$ -theorem for non-equilibrium situations which is also instrumental in defining the behaviour of gas particles at equilibrium (Chapman & Cowling, 1970). The non-equilibrium generalization of the entropy is known as the  $H$ -function, and is defined as

$$H = \int f \ln f d\mathbf{c}. \quad (4.4)$$



It can be shown that the evolution of the  $H$ -function takes the form:

$$\frac{dH}{dt} = \frac{1}{4} \int \int \int \ln \frac{f f_1}{f' f'_1} (f' f'_1 - f f_1) |\mathbf{c} - \mathbf{c}_1| \alpha d\mathbf{e} d\mathbf{c} d\mathbf{c}_1. \quad (4.5)$$

Hence, for the system to reach a state of equilibrium, we must have:

$$\ln f + \ln f_1 = \ln f' + \ln f'_1, \quad (4.6)$$

this relation indicates that  $\ln f$  is a conserved quantity, i.e, it remains unaffected by collisions and hence must be a linear combination of the collisional invariants - mass, momentum and energy

$$\ln f^{\text{MB}} = \alpha + \boldsymbol{\beta} \cdot m\mathbf{c} + \gamma \cdot \frac{m\mathbf{c}^2}{2}, \quad (4.7)$$

where  $f^{\text{MB}}$  is the distribution function at equilibrium. The various coefficients ( $\alpha$ ,  $\boldsymbol{\beta}$  and  $\gamma$ ) and therefore the distribution  $f^{\text{MB}}$  at equilibrium is calculated by considering the constraint that the moments of  $f^{\text{MB}}$  must retain the local density, momentum and energy. In the following section, we calculate the probability distributions for various quantities.

### 4.2.1 Distributions at equilibrium

For the simplest of systems, for example a rarefied gas in a periodic box with no external force offers a plethora of distributions.

#### Velocity distribution

For the distribution of velocities and energy Eq.(4.6) is used, which signifies the fact that the quantity  $\ln f$  remains unaffected by collisions. The expression in Eq.(4.6) can be rewritten as

$$\ln f^{\text{MB}} = \ln \alpha^{(0)} - \gamma \cdot \frac{1}{2} m \left[ \left( \frac{c_x - \beta_x}{\gamma} \right)^2 + \left( \frac{c_y - \beta_y}{\gamma} \right)^2 + \left( \frac{c_z - \beta_z}{\gamma} \right)^2 \right], \quad (4.8)$$

where  $\alpha^{(0)}$  is a constant such that

$$\ln \alpha^{(0)} = \alpha - \frac{m}{2} \frac{\beta_x^2 + \beta_y^2 + \beta_z^2}{\gamma}. \quad (4.9)$$

Using Eq.(4.8) the  $f^{\text{eq}}$  can be written as

$$f^{\text{MB}} = \alpha^{(0)} \cdot e^{-\gamma \frac{1}{2} m \mathbf{c}'^2}, \quad (4.10)$$

where  $\boldsymbol{\xi} = \mathbf{c} - \boldsymbol{\beta}/\gamma$ . To find these constants,  $f^{\text{MB}}$  is integrated over the velocity space and equated with the respective quantities. Firstly, the integral of  $f^{\text{MB}}$  over the velocity space must yield the number density

$$n = \int f^{\text{MB}} d\mathbf{c} = \alpha^{(0)} \left( \frac{2\pi}{m\gamma} \right)^{3/2}. \quad (4.11)$$

The first moment of  $f^{\text{MB}}$  provides the momentum and results in the relation

$$\rho \mathbf{u} = \int m\mathbf{c} f^{\text{MB}} d\mathbf{c} = \rho \boldsymbol{\beta}/\gamma, \quad (4.12)$$

where  $\mathbf{u}$  is the mean velocity. This relation, simplifies the expression for  $\boldsymbol{\xi} = \mathbf{c} - \mathbf{u}$ , which is defined as the peculiar velocity and provides the definition for temperature as

$$\frac{D}{2}nk_B T = \int \frac{m}{2}\xi^2 f^{\text{MB}} d\mathbf{c} = \frac{D}{2\gamma}, \quad (4.13)$$

where  $k_B$  is the Boltzmann constant,  $T$  the local temperature and  $D$  the number of dimensions. These relations then determine the form of  $f^{\text{MB}}$  to be

$$f^{\text{MB}} = \left( \frac{m}{2\pi k_B T} \right)^{\frac{D}{2}} e^{-mc^2/2k_B T}. \quad (4.14)$$

This is known as the Maxwell-Boltzmann distribution and conveys the idea that velocity of particles in each direction are Gaussian distributed with the local velocity as the mean and the local temperature as the variance.

### Energy distribution

With the expression for velocity distribution at hand, the distribution for energy can be calculated in any number of dimensions. For example, the energy in three dimensions is

$$E = m(c_x^2 + c_y^2 + c_z^2)/2, \quad (4.15)$$

the distribution for  $E$  can be calculated by expressing the Maxwell-Boltzmann distribution in spherical coordinates

$$f^{\text{MB}}(r, \theta, \phi) = \left( \frac{m}{2\pi k_B T} \right)^{3/2} r^2 \exp\left(-\frac{mr^2}{2k_B T}\right). \quad (4.16)$$

Since  $r^2 = 2E/m$ , an expression for the probability density for  $E$  can be found simply by integrating over  $\theta$  and  $\phi$  and is

$$p(E) = \left( \frac{1}{k_B T} \right)^{3/2} \sqrt{\frac{2}{\pi}} \left( \sqrt{E} \exp\left(-\frac{E}{k_B T}\right) \right). \quad (4.17)$$

Similarly, the distribution of energy can be calculated for any number of dimensions. A special case of interest is two dimensions, where the speed follows a Rayleigh distribution and the energy is exponentially distributed. The energy and speed in two dimensions are defined as

$$\begin{aligned} v &= \sqrt{c_x^2 + c_y^2}, \\ E &= mv^2/2. \end{aligned} \quad (4.18)$$

The Maxwell-Boltzmann distribution for two dimensions, in the polar co-ordinates is

$$f^{\text{MB}}(v, \theta) = \left( \frac{m}{2\pi k_B T} \right) \exp\left(-\frac{mv^2}{2k_B T}\right). \quad (4.19)$$

Integrating over  $\theta$  to obtain an expression for distribution of  $v$ , we have

$$p(v) = \left( \frac{m}{k_B T} \right) \exp\left(-\frac{mv^2}{2k_B T}\right). \quad (4.20)$$

It can be readily seen that this is equivalent to the Rayleigh distribution which has the form

$$p(x) = \frac{x}{\sigma^2} \exp\left(-\frac{x}{2\sigma^2}\right). \quad (4.21)$$

Hence the speed in two dimensions is a Rayleigh distributed with scale parameter  $\sigma = \sqrt{k_B T/m}$ . The probability distribution of  $E$  can be computed by making a change of variables in Eq.(4.20)

$$p(E) = \frac{1}{k_B T} \exp - \left( \frac{E}{k_B T} \right). \quad (4.22)$$

Hence, energy in two dimension is exponentially distributed with  $\lambda = 1/k_B T$ . Proceeding in a similar manner, many different distributions can be obtained simply by considering different number of dimensions.

### Position distribution

At equilibrium under zero external force and periodic boundary conditions the density of the system must be uniform, i.e,

$$\rho(x) = \frac{N}{L^3}, \quad (4.23)$$

where  $N$  is the number of particles in the system and  $L$  is the length of the periodic box. This is a mathematical statement of the fact that at equilibrium a particle is equally likely to be found anywhere in the domain of interest. Hence, the positions of the particles in a system must follow a uniform distribution in  $[0, L)$ .

### Poisson distribution

It can be shown that the number density, which is the number of particles in a sub-cell, follows a Poisson distribution. Suppose in a domain of volume  $V$  a sub-cell of volume  $\delta V$  is considered, then the probability of finding a particle in this sub-cell is  $p = \delta V/V$ , as the distribution in space is uniform. Then, given  $N$  particles, the probability of finding  $n$  particles inside this sub-cell is given by the binomial distribution

$$P(n) = \binom{N}{n} p^n (1-p)^{N-n}. \quad (4.24)$$

If  $\delta V \ll V$ , then  $p \ll 1$  and if  $N$  is large, then as per Poisson's limit theorem  $P(n)$  can be approximated to (Casella & Berger, 2002)

$$P(n) \approx \frac{\lambda^n e^{-\lambda}}{n!}, \quad (4.25)$$

where  $\lambda = Np$  is finite. Hence, even the simplest of systems have physical quantities that follow uniform, Gaussian, exponential and Poisson distributions making it a rich source of randomness.

### Molecular Dice hypothesis

We propose that any model of dilute gas dynamics (whether operating at hardware or software level) can be used as source of randomness. Hence, this method can be regarded as a synergy of the two existing methods - mathematical constructs and utilizing randomness of devices; a mathematical model of an actual physical process. Therefore a computer code that simulates dilute gas dynamics is similar to Maxwell's demon, an omniscient being holding the knowledge of each particle's position and velocity. To an observer not privy to these details, these state variables will seem devoid of any discernible pattern and hence can be used as stream of random numbers. A visual description of our ideas is presented in Fig.(4.2). A multitude of numerical methods to solve the Boltzmann equation or its approximation exist. In the following sections, we explore if these particle based numerical methods are capable of producing high quality random numbers.

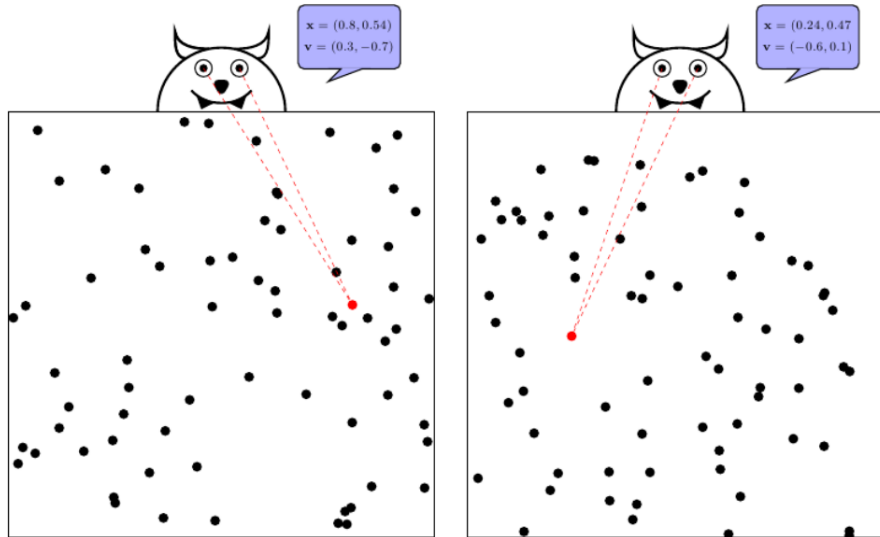


Figure 4.2: A computer code simulating gas dynamics is similar to Maxwell’s demon. To an observer, the values of positions and velocities act as sequence of random numbers.

### 4.3 Molecular Dynamics simulations as PRNGs

Particle based methods such as — molecular dynamics with Lennard-Jones potential and event-driven molecular dynamics solve the  $N$ -body problem of classical mechanics, wherein  $N$  interacting particles follow Hamiltonian dynamics. It has been shown that in the dilute gas limit (known as the Boltzmann-Grad limit), it is equivalent to solving the Boltzmann equation. Since Boltzmann dynamics assumes molecular chaos, i.e, the velocities of interacting particles are statistically independent, it is expected that for simulations in the low density limit, the particles’ positions, velocities and energies should provide stream of random numbers. In the following sections, some regular methods used to simulate gas dynamics are briefly described.

#### 4.3.1 Lennard-Jones

The Lennard-Jones model is a simple model of molecular interaction, wherein, the potential of each pair is calculated using

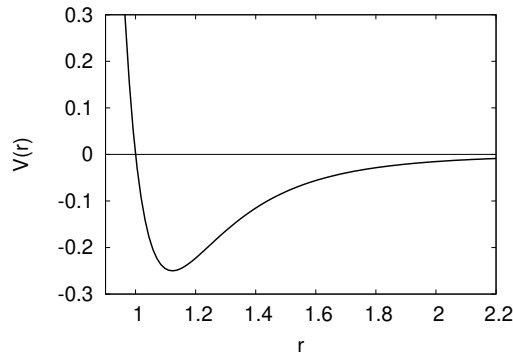
$$V(r) = 4\epsilon \left[ \left( \frac{\sigma}{r} \right)^{12} - \left( \frac{\sigma}{r} \right)^6 \right], \quad (4.26)$$

where  $\epsilon$  is the energy scale of the system and  $\sigma$  the length scale. The equation can be expressed in the dimensionless units as

$$V(r) = 4 \left[ \left( \frac{1}{r} \right)^{12} - \left( \frac{1}{r} \right)^6 \right], \quad (4.27)$$

where  $r$  is the distance between the two particles. This signifies that their interaction is strongly repulsive when they’re close, it becomes attractive as  $r$  gradually increases and eventually becomes zero as the two particles are far apart (Jones, 1924). A plot of the potential is presented in Fig.(4.3). The pairwise force is calculated by considering the gradient of the potential function,  $V(r)$ , and is found to be

$$\mathbf{f}_{ij} = 48 \left[ \left( \frac{1}{r_{ij}} \right)^{14} - \frac{1}{2} \left( \frac{1}{r_{ij}} \right)^8 \right] \mathbf{r}_{ij}, \quad (4.28)$$

Figure 4.3: Plot of pairwise potential  $V(r)$  vs.  $r$ 

where  $\mathbf{r}_{ij}$  is the position vector separating particles  $i$  and  $j$ . Since, the force is close to zero, for particles that are far away, cut-off distance was decided beyond which the interparticle force is set to zero. Once the total force exerted on each particle is calculated, their positions and velocities are updated using the Leapfrog-scheme (Rapaport, 2004)

$$\begin{aligned}\dot{x}\left(t + \frac{h}{2}\right) &= \dot{x}\left(t - \frac{h}{2}\right) + h\ddot{x}(t), \\ x(t+h) &= x(t) + h\dot{x}\left(t + \frac{h}{2}\right).\end{aligned}\tag{4.29}$$

where  $x(t)$  is the position of a given particles at time  $t$ ,  $\dot{x}(t)$  the velocity of the particle,  $\ddot{x}(t)$  is the acceleration experienced by each particle at any given instant owing to the sum total of inter-particle force, and  $h$  is a small time step of  $\mathcal{O}(10^{-2})$  (Rapaport, 2004). Such a scheme ensures that the phase space volume is conserved and the total energy remains conserved on an average. Since the pairwise potential is to be calculated in each iteration the computational complexity of the algorithm is  $\mathcal{O}(N^2)$ .

The positions, velocities and the contribution of kinetic energy in two dimensions are sampled after every  $10^3$  iterations and were considered to be streams of uniform, Gaussian and exponential numbers. Consistent with the hypothesis, these streams of numbers managed to satisfy Marsaglia's difficult-to-pass tests for randomness (Marsaglia & Tsang, 2002). This numerical experiment suggests that collective chaotic behaviour emerging from the Boltzmann picture of gases provides an alternate conceptual framework to analyse and create apparent randomness on computers. We now investigate whether simpler numerical methods are capable of generating random sequences.

### 4.3.2 Hard-sphere systems

An alternative to Lennard-Jones framework is the hard-sphere system, an event driven system wherein the motion of every particle is tracked and the velocity of a pair of particles are updated when they are just touching each other (Rapaport, 2004). Essentially, it is a simplified version of the Lennard-Jones system wherein the repulsive soft-sphere potential is replaced by a step potential. The positions and velocities of the particles are initialized randomly and possible collisions between all pairs are considered. The time between the  $i$ th and  $j$ th particle can be calculated by

$$\delta t = -\frac{b + \sqrt{d}}{\mathbf{v}_{ij} \cdot \mathbf{v}_{ij}},\tag{4.30}$$

where  $b = \mathbf{v}_{ij} \cdot \mathbf{r}_{ij}$  and  $d = b^2 - v_{ij}^2(r_{ij}^2 - \sigma^2)$ , with  $\mathbf{v}_{ij} = \mathbf{v}_i - \mathbf{v}_j$  and  $\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$  and  $\sigma$  is the diameter of the particles. The velocities of the colliding pair are updated as

$$\begin{aligned}\mathbf{v}_i &= \mathbf{v}_i - \frac{b}{\sigma^2} \mathbf{r}_{ij}, \\ \mathbf{v}_j &= \mathbf{v}_j + \frac{b}{\sigma^2} \mathbf{r}_{ij}.\end{aligned}\tag{4.31}$$

This ensures that momentum and energy are conserved. Unlike the Lennard-Jones model, the hard-sphere system needs to update only the possible collisions for the pair of particles that have undergone collision, and is therefore an  $\mathcal{O}(N)$  algorithm. Hence, the hard-sphere system turns out to be around ten times faster than the Lennard-Jones model, for a given number of particles. The positions, velocities and energies when sampled after a few hundred iterations, produce stream of random numbers corresponding to their expected distributions. These streams of numbers have been validated with statistical tests. While both Lennard-Jones model and hard-sphere system produce stream of random numbers, the rate of generation is only a few thousand per second making them unsuitable for large-scale scientific computations. The rate of generation of random numbers from both methods is listed in Table(4.1). These numerical experiments show that the concepts of randomness associated with Boltzmann dynamics are in fact consistent with statistical inference techniques used in computer science.

Quantity	Method	Speed (doubles/second)
Uniform	MT19937	$1.8 \times 10^8$
	Lennard-Jones	$0.01 \times 10^6$
	Hard-sphere system	$0.07 \times 10^6$
Gaussian	Box-Muller	$1.0 \times 10^7$
	Lennard-Jones	$0.01 \times 10^6$
	Hard-sphere system	$0.02 \times 10^6$
Exponential	Inverse sampling	$1.2 \times 10^7$
	Lennard-Jones	$0.005 \times 10^6$
	Hard-sphere system	$0.01 \times 10^6$

Table 4.1: Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by Lennard-Jones dynamics and hard-sphere system compared with the widely used MT19937 (Matsumoto & Nishimura, 1998). The Gaussian and exponential random numbers were generated using the same.

## 4.4 Mesoscale methods

The fact that the hard-sphere model generates random sequences at a higher rate than Lennard-Jones indicates that coarse-grained modelling, which operate at larger time scales are better suited to generating uncorrelated numbers. Hence, mesoscopic methods for simulating gas flows seem an attractive alternative. The key idea of such methods is a “top-bottom approach” (Succi *et al.*, 2002), wherein, the macroscopic behaviour of the system is preserved while considering one of many microphysical possibilities. This framework has led to the emergence of hydrodynamic solvers such as — Direct Simulation Monte Carlo (DSMC) (Bird, 1978), Multiparticle Collision Dynamics (MPCD) (Kapral, 2008), Dissipative Particle Dynamics (DPD) (Hoogerbrugge & Koelman, 1992) and Lattice Boltzmann Methods (Succi, 2001). The Lattice Boltzmann methods are grid-based techniques which discretize the velocity space to a set of finite discrete velocities, consequentially this is not suitable for generating random numbers. Dissipative Particle Dynamics

on the other hand is a particle-based method which considers pairwise interactions in addition to stochastic forces, hence this algorithm requires large number of computations as well as random sequences making it unsuitable for our purpose. In the following subsections, we discuss DSMC and MPCD methods and their performance as PRNGs.

#### 4.4.1 Direct Simulation Monte Carlo

The DSMC is a numerical tool which can be used to solve the Boltzmann equation for any value of Knudsen number. The central idea behind this algorithm is to emulate the Boltzmann collisional operator. Particles are sorted into cells on the basis of their positions and it is assumed that any pair of particles in the same cell can collide with each other and exchange momentum and energy (Bird, 1978). Particles in a cell are chosen via uniformly distributed random numbers and their relative velocity is modified such that the magnitude remains unchanged, and the two other parameters of collisions are taken to be

$$\phi = \cos^{-1}(1 - 2U_1) \quad , \quad \theta = 2\pi U_2, \quad (4.32)$$

where  $U_1$  and  $U_2$  are uniformly distributed random numbers, following which the components of relative velocity,  $\mathbf{v}_{ij}$  is modified as:

$$\begin{aligned} (v_{ij})_x &= v_{ij} \sin \phi \cos \theta, \\ (v_{ij})_y &= v_{ij} \sin \phi \sin \theta, \\ (v_{ij})_z &= v_{ij} \cos \phi. \end{aligned} \quad (4.33)$$

Every cell undergoes a fixed number of collisions which determines various transport coefficients of the system and parameters such as the Knudsen number. Once the required number of collisions is satisfied, the particles are streamed with their new velocities and the entire procedure is repeated again. This method is quite reliable and can be used for the entire range of rarefaction and even capture shock structures. Since the objective is to generate random numbers, the resolution of the grid can be coarser and the number of collisions required for each particle is much lower as compared to when solving for actual flows.

#### 4.4.2 Multiparticle Collision Dynamics

MPCD is a particle based mesoscopic method used to simulate complex fluids (Kapral, 2008). As opposed to DSMC, this method does not account for inter-particle collisions and focuses on collective behaviour of particles. Once the particles are sorted into cells, their velocities are updated by

$$\mathbf{v}' = \mathbf{v} + \mathbf{R}(\mathbf{v} - \mathbf{u}), \quad (4.34)$$

wherein,  $\mathbf{u}$  is the mean velocity and  $\mathbf{R}$  is the stochastic rotation matrix, associated with the cell the particle occupies. Since interparticle interactions are not taken into account, the MPCD algorithm has a lower order of computations than DSMC and hence faster. However, the nature of algorithm requires it to be used for complex system which display collective behaviour such as — colloidal solutions and polymers.

It was found that the aforementioned algorithms did indeed satisfy Marsaglia's difficult-to-pass tests. The rate of generation of various distribution using DSMC and MPCD are tabulated in Table(4.2). It is seen that while DSMC generated Gaussian and exponential random numbers at the same speed as standard methods, and uniformly distributed numbers much slower than state-of-the-art PRNGs, MPCD accelerated Gaussian random number generation by a factor of 7 and exponential random number generation by a factor of 3, although uniform random number

generation is around the same speed as existing PRNGs. It was also found that when subjected to a more stringent set of tests namely the Crush battery of tests in TESTU01 library (L'Ecuyer & Simard, 2007), it failed many statistics, indicating that sequences generated by MPCD are not high quality. The mesoscale methods are a significant improvement over deterministic solvers, however they prove to be only a slight upgrade over standard methods. Therefore, it is imperative to design a new algorithm rooted in the new framework capable of generating random numbers at higher speeds.

Quantity	Method	Speed (doubles/second)
Uniform	MT19937	$1.8 \times 10^8$
	DSMC	$2.5 \times 10^7$
	MPCD	$7.3 \times 10^7$
Gaussian	Box-Muller	$1.0 \times 10^7$
	DSMC	$3.3 \times 10^7$
	MPCD	$7.4 \times 10^7$
Exponential	Inverse sampling	$1.2 \times 10^7$
	DSMC	$1.6 \times 10^7$
	MPCD	$4.0 \times 10^7$

Table 4.2: Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by mesoscopic methods namely – DSMC and MPCD, compared with MT19937.

## 4.5 Final algorithm

Keeping in tandem with the “top-bottom approach”, we choose a system whose microphysics would be computational friendly and provides uncorrelated random numbers. We adopt an approach similar to the DSMC and borrow concepts from MPCD. The final version of the algorithm chooses a pair of particles and updates their velocities using a simple collision rule as defined by:

$$\mathbf{v}'_i = \frac{\mathbf{v}_i + \mathbf{v}_j}{2} - \mathbf{R} \frac{\mathbf{v}_i - \mathbf{v}_j}{2}, \quad \mathbf{v}'_j = \frac{\mathbf{v}_i + \mathbf{v}_j}{2} + \mathbf{R} \frac{\mathbf{v}_i - \mathbf{v}_j}{2}. \quad (4.35)$$

where  $\mathbf{R}$  is a stochastic rotation matrix. Once the velocities of the pair of particles are modified, their positions are updated with periodic boundary conditions by:

$$\mathbf{x}'_i = \mathbf{x}_i + \mathbf{v}'_i \cdot \delta t, \quad \mathbf{x}'_j = \mathbf{x}_j + \mathbf{v}'_j \cdot \delta t. \quad (4.36)$$

An explicit description of the final algorithm is provided in Table(4.3).

A key feature of this algorithm is that the uniform random numbers needed for instantiating the stochastic rotation matrix and for pair selection can be directly sampled from the positions of the particles

### 4.5.1 Pair selection

The method to select pairs of molecules for collisions is central to hydrodynamic solvers such as DSMC and even for proposed algorithm. We found the following methods of pair selection to provide satisfactory results:



Table 4.3: Overview of the “Molecular Dice” algorithm

- 
- 
1.  $N$  particles are initialized with uniformly distributed positions and normally distributed velocities in  $D$  dimensions.
  2. A stochastic rotation matrix,  $R$ , with its entries being trigonometric transformations of uniform random numbers is chosen.
  3. A pair of particles is chosen and then selected and their velocities updated using Eq.(4.35).
  4. The velocities of both the particles are returned as Gaussian random numbers.
  5. If uniform random numbers are required then the positions of the pair of particles are returned and are updated using Eq.(4.36) with periodic boundary conditions.
  6. The contribution of kinetic energy of both particles in 2-dimensions is returned as exponential random numbers.
  7. Steps 3-6 are repeated for  $W$  iterations, after which step 2 is executed once to re-initialize the stochastic rotation matrix.
- 
- 
1. Classical method: Two integers  $(i, j)$  are chosen in the range  $[0, N - 1]$ , from a uniform distribution such that  $i \neq j$ . The particles with these indices are then considered to be the colliding particles. We found that for this method  $W = N/2$  provided satisfactory results.
  2. LCG-style method: The Hull-Dobell theorem guarantees that the iterative scheme  $X_{n+1} = (aX_n + b) \bmod N$  provides all integers in the range  $[0, N - 1]$  once for  $N$  iterations under a certain set of conditions. For  $N = 2^s (s > 2)$ , the conditions are quite simplified —  $b$  must be an odd integer and that  $(a - 1)$  must be a multiple of 4 in the range  $[0, N - 1]$ . The choice of  $a$ ,  $b$  and  $X_0$  is made using uniformly distributed numbers and which satisfy the aforementioned conditions. These choices are replenished after every  $W$  iterations to avoid engendering pattern. Choosing  $W = N/2$  ensures that every particle takes part in the collision process and clears most statistical tests.
  3. Offset-and-jump method: The index for first particle is offset every iteration by a uniformly distributed integer in the range  $[0, N/W)$  and the index for second particle is chosen by adding a jump to it which is uniformly distributed between  $[1, N - 1]$ . If the index of either particle exceeds  $N$ , then modulus of the value with  $N$  is taken. The choices for offset and jump is renewed after every  $W$  iterations. We found that  $W = N/4$  provided the most satisfactory results.

In the following sections the results and the speed of the proposed algorithm for all three pair selection methods is presented.

## 4.6 Statistical tests

As discussed in the previous chapter, the quality of a PRNG is established via the standard statistical tests as outlined in the previous chapter. In addition to the tests mentioned in Chapter 3, the Crush battery of tests contained in TestU01 suite (L’Ecuyer & Simard, 2007) was also used for testing empirical randomness. It is a set of 96 highly stringent statistical tests providing a  $p$ -value for 144 statistics for 32-bit uniform integers. To test Gaussian and exponential random numbers, they are first transformed to uniform integers as explained in Chapter 3. The proposed algorithm was tested with all three methods of pair selection and for all the three quantities (uniform, Gaussian and exponential), for 100 different seeds. The tests which failed consistently, i.e, generated  $p$ -values outside the interval  $[10^{-3}, 1 - 10^{-3})$  were deemed to be systemic failures. Table(4.4) presents the results.

Pair selection method	Quantity	Systemic failures	Test(s) failed
Classical	Exponential	1	Gap
	Gaussian	2	Gap, CollisionOver
	Uniform	0	—
LCG-syle	Exponential	2	WeightDistrib, HammingWeight2
	Gaussian	1	WeightDistrib
	Uniform	0	—
Offset and jump	Exponential	1	Gap
	Gaussian	1	Gap
	Uniform	0	—

Table 4.4: Results for Crush battery of tests. The results are for all the three quantities – uniform, Gaussian and exponential with the three different methods for pair selection outlined in Sec.(4.5.1)

As can be seen, the quality of random sequences generated by proposed algorithm is at least at par with established and widely used routines such as MT19937.

## 4.7 Speed of RNGs

The speed of generation of the proposed algorithm was calculated using a program to calculate the mean of  $10^8$  numbers drawn from the generator and the time spent for the same. Such a methodology is known sequential testing and is an appropriate method to determine the speed of generators to be used for large-scale scientific computations. The speeds for all algorithms were tested on a Intel(R) Core(TM) i7-6800K CPU @ 3.40GHz machine. The results are presented in Table(4.5).

Quantity	Method	Speed (doubles/second)
Uniform	MT19937	$1.8 \times 10^8$
	Classical	$7.1 \times 10^7$
	LCG-style	$8.5 \times 10^7$
	Offset and jump	$1.0 \times 10^8$
Gaussian	Box-Muller	$1.0 \times 10^7$
	Classical	$1.8 \times 10^8$
	LCG-style	$2.0 \times 10^8$
	Offset and jump	$3.6 \times 10^8$
Exponential	Inverse sampling	$1.2 \times 10^7$
	Classical	$8.5 \times 10^7$
	LCG-style	$1.6 \times 10^8$
	Offset and jump	$1.9 \times 10^8$

Table 4.5: Rates of generation of uniform, Gaussian and exponential random numbers (in doubles/sec) by the final algorithm for all the three methods of pair selection outlined in Sec.(4.5.1), compared with MT19937.

The values in Table(4.5) suggests that proposed algorithm while providing high quality sequences is a good upgrade over standard methods in terms of speed.

## 4.8 Outlook

Existing methods to generate random numbers include extracting noise from devices or utilizing complex formulas rooted in number theory. It can be shown that a unification of the two methods – simulation of a stochastic process such as the motion of molecules produces sequences of number which are apparently random. An algorithm was designed based on existing hydrodynamics solvers capable of generating Gaussian and exponential random numbers at a much higher rate. As generation of non-uniform random numbers have often acted as a bottleneck (Thomas *et al.*, 2007), for example taking upto 90% of the total computational time for a simple Brownian motion simulation, we expect that this algorithm to generate Gaussian and exponential random numbers can speed-up simulations of many large-scale problems. This opens up the possibility of tackling problems such as whole-cell simulations which have been restricted due to the computational cost involved.



# Chapter 5

## Chemical Reactions

### 5.1 Introduction

Chemical reactions are pervasive in natural phenomena and central to the dynamics of many systems such as gene networks (Becskei & Serrano, 2000), combustion (Kraft & Wagner, 2003) etc. Simulations of such systems for long-time scales is important in engineering and scientific applications (Espenson, 1995). Typically chemical reactions have been modelled using the deterministic rate law of mass action, which describes the change in concentration of various components of a system in time (Érdi & Tóth, 1989). This method of simulating chemical reactions produces satisfactory results for large system sizes, however the discrete nature of molecules and fluctuations can play an important role in governing the behaviour for small systems (Srivastava *et al.*, 2002; Turner *et al.*, 2004). Such dynamics are often encountered in biochemical networks such as the switching between lysis and lysogeny phases in  $\lambda$ -phage which is driven from small perturbations (Arkin *et al.*, 1998). Thus, for such cases a stochastic formulation which accounts for fluctuations about the mean behaviour and the discrete nature of molecules provides a better insight into the dynamics of small systems. The stochastic simulation algorithms for chemical reactions have in fact been found to produce accurate results faster than deterministic solutions for some cases (Kraft & Wagner, 2003). Although simple, this algorithm requires a large number of exponential and uniform random numbers, the former particularly acting as a major detriment. The “Molecular Dice” algorithm introduced in the last chapter produces both exponential and uniform random numbers in a single iteration and is therefore ideally suited to this algorithm.

In this chapter, the deterministic formulation for modelling chemical kinetics is explained first followed by the stochastic formulation, following which various numerical methods used to solve the latter are discussed. A demonstration of the usefulness of the “Molecular Dice” algorithm presented in the previous chapter in context of these numerical methods is presented via simulation of a bi-stable biochemical reaction network. The chapter ends with a discussion on reaction-diffusion processes and a simulation of pattern formation in *E.Coli* bacteria which highlights the difference between stochastic and deterministic formulations.

### 5.2 Deterministic formulation

The deterministic formulation of chemical kinetics is known as the rate law of mass action and states that for a given reaction the rate of change of the concentration of various species is directly proportional to the concentration of the participating reactants, for example given a simple first order reaction



the rate law of mass action states that:

$$\begin{aligned} \frac{dC_A}{dt} &= -kC_A, \\ \frac{dC_B}{dt} &= kC_A, \end{aligned} \quad (5.2)$$

where  $C_i$  denotes the concentration of  $i$ th species and  $k$  is the rate constant associated with the reaction. The explicit solution of these differential equations is

$$\begin{aligned} C_A(t) &= C_A(0)e^{-kt}, \\ C_B(t) &= C_A(0)\left(1 - e^{-kt}\right). \end{aligned} \quad (5.3)$$

While this particular example does yield a deterministic solution, most chemically reactive systems of the form

$$\frac{d\mathbf{x}}{dt} = f(\mathbf{x}), \quad (5.4)$$

where  $f(\mathbf{x})$  denotes some function of the components of the vector  $\mathbf{x}$ , result in highly complex non-linear ordinary differential equations, which can only be solved using numerical techniques. To solve this system numerically, the Taylor expansion of this system is written as

$$\mathbf{x}(t + \delta t) = \mathbf{x} + \delta t \cdot \frac{d\mathbf{x}}{dt} + O(\delta t^2). \quad (5.5)$$

There are two simple methods to evaluate  $d\mathbf{x}/dt$  from Eq.(5.4). Either the value of the derivative is considered at  $t$  itself or  $t + \delta t$ . The former leads to a very simple set of discrete equations, using which the value of  $\mathbf{x}$  can be updated at every time step using

$$\mathbf{x}(t + \delta t) = \mathbf{x} + f(\mathbf{x}(t))\delta t. \quad (5.6)$$

Such a technique is called Euler's forward method. The second option, that is to evaluate  $d\mathbf{x}/dt$  at  $t + \delta t$ , leads to

$$\mathbf{x}(t + \delta t) = \mathbf{x} + f(\mathbf{x}(t + \delta t))\delta t. \quad (5.7)$$

In this case, as opposed to Eq.(5.6),  $f(\mathbf{x}(t + \delta t))$  cannot be evaluated directly as  $\mathbf{x}(t + \delta t)$  is unknown at time  $t$ . This implicit approach is known as Euler's backward method. Although the latter is more cumbersome, it allows for large time steps and thus is generally preferred for cases with stiff non-linear terms (Kraft & Wagner, 2003). More advanced backward solvers exist and are routinely used for solving stiff systems.

A deterministic formulation effectively does not account for the fluctuations, hence making it suitable only for large system sizes ( $N \rightarrow \infty$ ). In order to capture the correct qualitative behaviour of the system it is imperative to consider these complications and employ a stochastic formulation as opposed to a deterministic one (Gillespie, 1977).

### 5.3 Stochastic formulation of chemical rate equations

The stochastic formulation is considered to be a more detailed description of reaction dynamics as it manages to incorporate perturbations about the deterministic mean thereby providing a better insight to the behaviour of the system (Oppenheim *et al.*, 1969). For the stochastic approach, the fundamental quantity of interest is the probability of finding a given number of molecules of a species at a particular time and is denoted by  $P(\mathbf{x}, t)$ . Here  $\mathbf{x}$  is a vector with components  $\mathbf{x} \equiv (x_1, x_2, \dots, x_N)$  where  $x_i$  is the number of molecules of the  $i$ th species in the system. The key idea here is that chemical reactions are considered to be Markov processes (Gardiner, 1985a), i.e, the progression of the system only depends on the current state and does not take its history into account. In an infinitesimal time interval, there are several paths through which the system can arrive at or depart from a given state  $\mathbf{x}$ , the evolution equation for  $P(\mathbf{x}, t)$  using detailed balance is

$$\frac{\partial P(\mathbf{x}, t)}{\partial t} = \sum (\text{gain}) - \sum (\text{loss}), \quad (5.8)$$

where the summations are over all the paths that either lead to or away from  $\mathbf{x}$ . The following assumptions are made in order to describe chemical kinetics as a stochastic process

1. The process of chemical reaction itself is modelled as a birth-death process and it is assumed that the probability of a particular reaction being triggered in a infinitesimal time  $dt$  is  $kdt$ .
2. A given control system is considered to be homogenous, i.e, all possible combinations of reactants are capable of reacting with each other. For example, the probability of a first order reaction  $A + B \xrightarrow{k} \phi$  being triggered in time  $dt$  is  $kN_A N_B dt$ , where  $(N_A, N_B)$  are the number of  $(A, B)$  molecules present in the given control system.
3. It is assumed that the molecular chaos hypothesis introduced in Chapter 3 holds, using which it can be argued that the joint probability density of a transition is a product of the individual transition probabilities. In mathematical terms, can be stated as:

$$P(x_1 \pm a_1; \dots; x_N \pm a_N) = \prod_{i=1}^N P(x_i \pm a_i).$$

These assumptions lead to the Chemical master equation, a statement of the rate of change probability densities of different chemical species in the system (Oppenheim *et al.*, 1969). Its derivation can be explained by the example considered in previous section. In an infinitesimal time  $dt$ , the following transition probability law hold

$$P(N_A \rightarrow N_A - 1; N_B \rightarrow N_B + 1) = kN_A dt, \quad (5.9)$$

as per the assumptions made it is the only path in the system that brings change. The rate of change of the joint probability density is then

$$\frac{P(N_A, N_B, t + dt) - P(N_A, N_B, t)}{dt} = k(N_A + 1) \times P(N_A + 1, N_A - 1, t) - kN_A \times P(N_A, N_B, t). \quad (5.10)$$

In the limit  $dt \rightarrow 0$ , the time evolution of the probability is obtained. Similarly for a general case of  $N$  components interacting chemically through  $R$  reaction channels the chemical master equation is

$$\frac{\partial}{\partial t} P(\mathbf{x}, t) = \sum_{i=1}^R [a_i(\mathbf{x} - \nu_i) P(\mathbf{x} - \nu_i, t) - a_i(\mathbf{x}) P(\mathbf{x}, t)] \quad (5.11)$$

where  $a_i(\mathbf{x})$  known as propensity function is the total probability per unit time that the reaction  $i$  occurs in the system. It is the product of total number of available combinations of the reactants and the rate constant associated with the reaction. Another quantity of relevance is the vector,  $\nu_i$ , whose components  $\nu_{ji}$  signify the change in the number of molecules of species,  $j$ , brought about by one  $i$  reaction. It is evident that Eq.(5.11) is detailed version of Eq.(5.8), where the gain term is given by  $\sum_{i=1}^R [a_i(\mathbf{x} - \nu_i) P(\mathbf{x} - \nu_i, t)]$  which is total probability that the system makes the transition  $(\mathbf{x} - \nu_i \rightarrow \mathbf{x})$  and the loss term is  $a_i(\mathbf{x}) P(\mathbf{x}, t)$  which brings the system away from the state  $\mathbf{x}$ . (Gillespie, 2000). For the normalization of the probability density function to hold, we must have

$$\sum P(\mathbf{x}) = 1, \quad (5.12)$$

where the sum is over all possible values  $(0, \infty)$  for all components. Similarly, the mean behaviour of the system can be calculated using

$$\mu = \sum \mathbf{x}P(\mathbf{x}). \quad (5.13)$$

It has been shown that this formulation of chemical kinetics converges to the rate law of mass action in the thermodynamic limit  $(N \rightarrow \infty, V \rightarrow \infty)$ , a necessary condition for it to be considered canonical (Kurtz, 1972; Oppenheim *et al.*, 1969). While accurate, the chemical master equation is more often than not too complex to be solved analytically, hence use of efficient numerical techniques must be employed.

## 5.4 Simulation Algorithms

A naive algorithm that can be utilized to produce trajectories of the chemical master equation would be to generate a random number and check whether it satisfies the first assumption of the stochastic formulation. For the simple single reaction system introduced in Eq.(5.1), a  $\Delta t = 10^{-3}k$  can be selected which would signify  $\lambda = A(t)k\Delta t$  to be the probability of observing a reaction event. A uniform random number,  $U$ , between  $[0, 1)$  is generated, then

$$P(U \leq \lambda) = \lambda = A(t)k\Delta t \quad (5.14)$$

when satisfied, a reaction event has occurred in accordance with the first assumption made for the stochastic formulation and the number of  $A$  molecules is decreased by one. Realizations of this algorithm are presented in Fig.(5.1). While this algorithm does solve the chemical master equation exactly and produces correct realizations, it is highly inefficient as most of the time is spent generating random numbers while observing no changes in the system (Erban *et al.*, 2007).

### 5.4.1 Gillespie algorithm

The Gillespie algorithm also known as the next reaction method, is an efficient method to generate statistically correct possible realizations of the chemical master equation (Gillespie, 1976). The naive approach discussed earlier has long waiting-times, i.e, most iterations do not fulfill the criterion and are hence rejected bringing about no change in the system. As opposed to the naive scheme, the Gillespie algorithm treats time intervals between two successive reaction events as a random variable and aims to quantify its distribution. This allows one to build a scheme where every Monte-Carlo move is accepted thereby eliminating the computing time wasted for observing an event, the primary drawback in the former implementation. For this purpose, a new quantity known as *reaction probability density function*,  $P(\tau, \mu)$ , is defined as probability that next reaction occurs in the interval  $[t + \tau, t + \tau + d\tau]$ , and is the  $R_\mu$  reaction. It is assumed that the time  $d\tau$  is small enough that only a single reaction occurs during the interval  $[t + \tau, t + \tau + d\tau)$  (Gillespie, 1976). The probability that no reaction event occurs in the time interval  $[t, t + \tau)$ , denoted by  $P_0(\tau)$ , can be calculated by considering that no reaction events occur in  $K$  infinitesimally small partitions,  $\epsilon = \tau/K$  and accounting for the assumption that each event (or lack of) in a given subinterval is statistically independent, then  $P_0(\tau)$  is

$$P_0(\tau) = [1 - \sum a_\mu(\mathbf{x})\epsilon]^K. \quad (5.15)$$

In the limit  $\epsilon \rightarrow 0, K \rightarrow \infty$ , we have



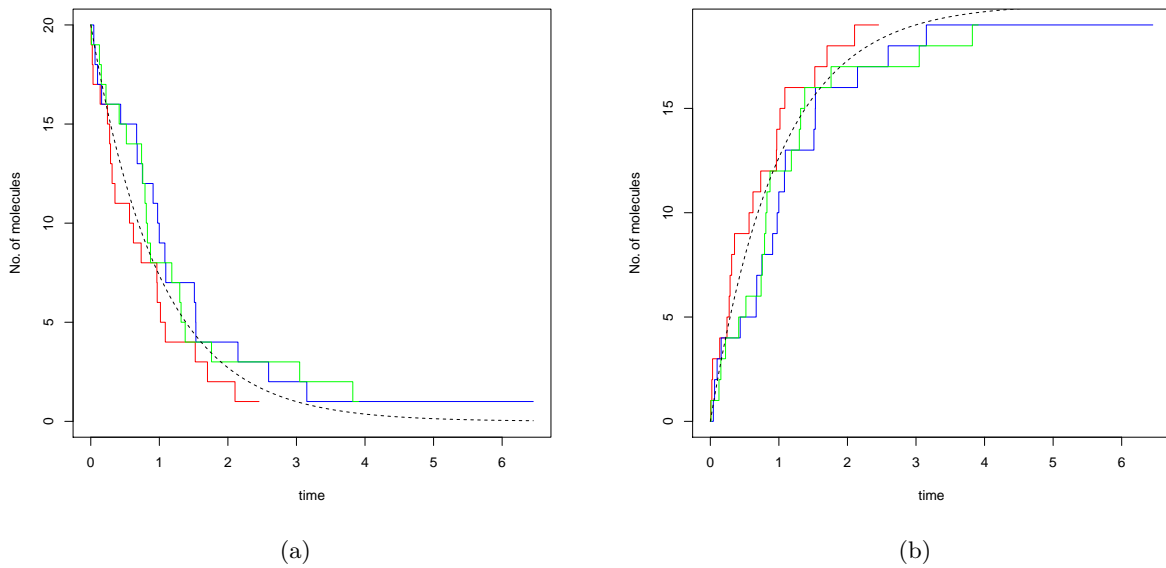


Figure 5.1: Plots comparing the results of the deterministic solution of the chemical reaction  $A \xrightarrow{k} B$ . The solid, coloured lines represent the various realizations of the naive implementation and the dashed line represents the deterministic solution. It can be seen that large number of such realization when averaged converge to the deterministic solution. These figures represent a) degradation of  $A$  and b) degradation of  $B$ .

$$P_0(\tau) = e^{-\sum a_\mu \tau} \sum a_\mu. \quad (5.16)$$

A numerical scheme can then be constructed by calculating the waiting time distribution at every time  $t$  and then generating a random variate distributed according to  $P_0$ , following which a reaction is chosen to be triggered in accordance with their respective propensities. The algorithm can be summarized as follows

1. Calculate the propensity function  $a_\mu$ , which is the total probability of observing a particular reaction event and compute their sum  $a_0 = \sum a_\mu$
2. Generate an exponential random number with the parameter  $a_0$ . Using inversion method,  $\tau$  can be simply generated as:

$$\tau = \frac{1}{a_0} \log \left( \frac{1}{U_1} \right) \quad (5.17)$$

where  $U$  is a uniform random number between  $(0, 1)$ .

3. Generate a uniform random number,  $U_2$ , and find an  $j$  such that:

$$U_2 \geq \frac{1}{a_0} \sum_{i=1}^{j-1} a_i \quad \text{and} \quad U_2 < \frac{1}{a_0} \sum_{i=1}^j a_i \quad (5.18)$$

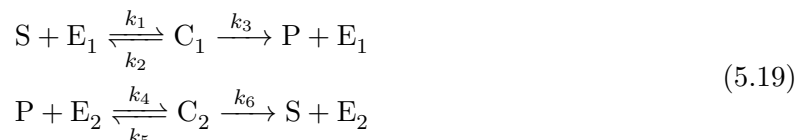
this  $j$  signifies the reaction that has occurred in  $(t + \tau, t + \tau + d\tau)$ . This is essentially choosing the reaction to be triggered based on their propensities, as  $\tau$  only indicates when the next reaction event takes place and does not specify the reaction itself.

4. The number of molecules of each species are updated as per the reaction chosen, and steps 1 – 3 are repeated until desired time is achieved.

While this algorithm eliminates a major drawback of naive implementation discussed before by calculating the waiting time distribution and is quite efficient for homogeneous cases, it spends roughly 85% time on random number generation indicating a major scope for improvement.

### 5.4.2 Convergence and performance of Gillespie algorithm

To study the convergence and performance of the algorithm, a prototype reaction the Goldbeter-Koshland switch is simulated (Goldbeter & Koshland, 1981). It is a reaction system comprising six reactions and six species



The network is simulated with initial conditions

$$\begin{bmatrix} X_S(0) \\ X_{E_1}(0) \\ X_{C_1}(0) \\ X_P(0) \\ X_{E_2}(0) \\ X_{C_2}(0) \end{bmatrix} = \begin{bmatrix} 0.275 \\ 0.25 \\ 0.075 \\ 0.075 \\ 0.25 \\ 0.075 \end{bmatrix}$$

where  $X_i$  denote the different mole fractions of the species. The values of the rate constants were  $k_1 = 0.05$ ,  $k_2 = 0.1$ ,  $k_3 = 0.1$ ,  $k_4 = 0.01$ ,  $k_5 = 0.1$ , and  $k_6 = 0.1$ . The comparison between the stochastic and deterministic solutions is shown in Fig.(5.3). The computational performance of the Gillespie algorithm with the two methods of random number generation are presented in Fig.(5.2). It is evident that the Gillespie algorithm using the standard methods to generate random numbers or with Molecular Dice converge in the same manner, with the latter proving to be faster by a factor of almost 4.

The total computational time spent for Gillespie algorithm using standard methods and the ‘‘Molecular Dice’’ algorithm are listed in Table(5.1).

Method	Error	Std. Deviation	Time(s)
Gillespie - standard	$4.2 \times 10^{-4}$	$6.2 \times 10^{-3}$	83.78
Gillespie - Molecular Dice	$6.4 \times 10^{-4}$	$5.9 \times 10^{-3}$	19.56

Table 5.1: Time comparisons for the three methods to calculate the mole fraction of the specie  $S$  after  $t = 100$ .

While the stochastic formulation does indeed converge to the deterministic solution, the computational time taken to implement this method is higher than the deterministic numerical method. However, for large reaction networks, simple numerical schemes as the one employed in present case might not be useful as the resulting equations are highly non-linear and stiff (Kraft & Wagner, 2003; Gillespie, 2007). In such cases, the Gillespie algorithm might indeed prove to be computationally cheaper than deterministic numerical schemes. The Gillespie algorithm provides a reliable measure for large system sizes, however it is highly useful for small system sizes wherein small perturbations drive the system from one stable state to another. An example is presented in the following section to highlight this issue.

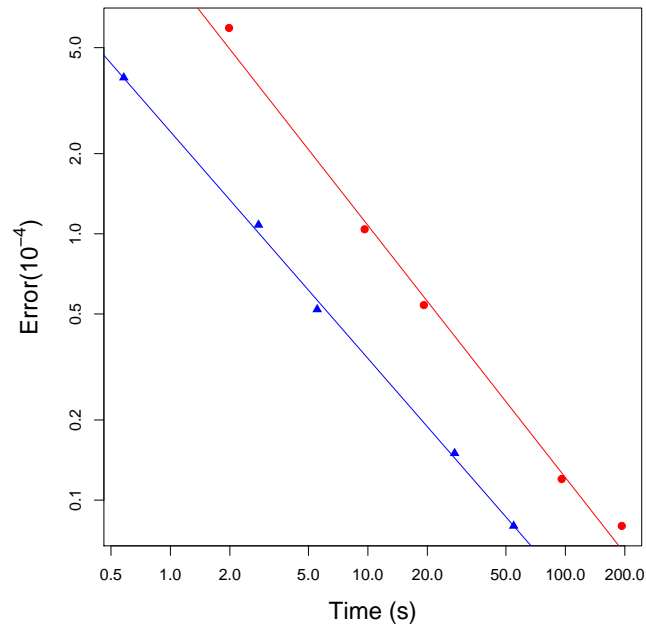


Figure 5.2: Simulation of the reaction network presented in Eq.(5.19). The plots depict the error observed for different number of molecules present in the system for Gillespie algorithm when implemented with standard and Molecular Dice algorithms. Both methods converge with  $\log(1/N)$ , with Molecular Dice being around 4 times faster.

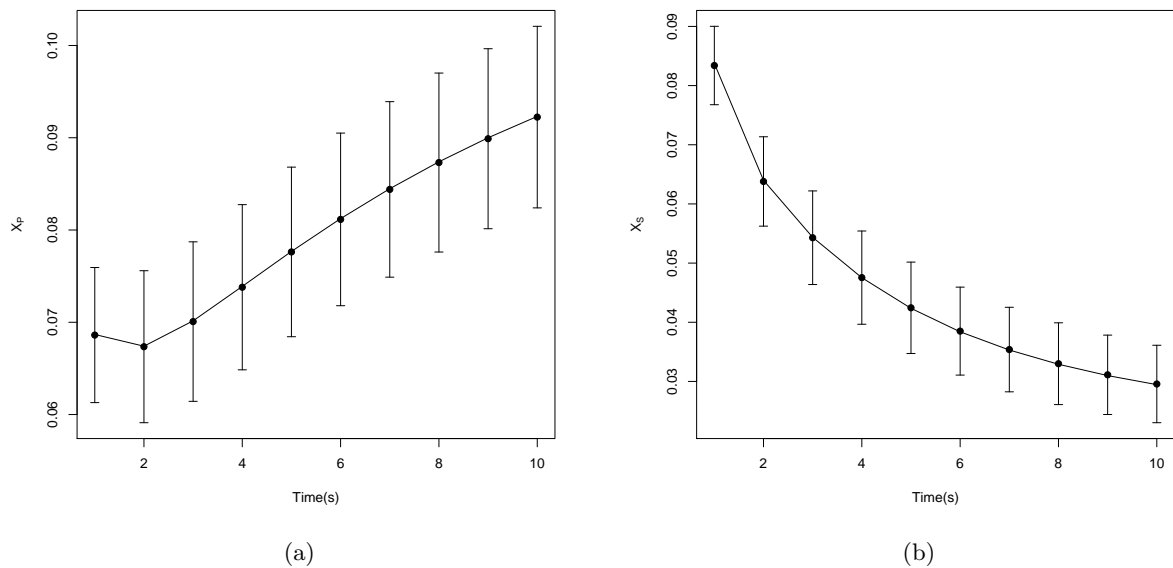
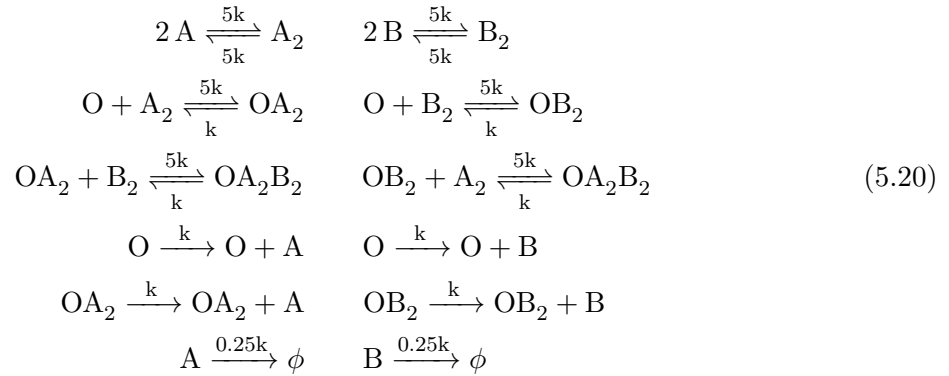


Figure 5.3: Plot comparing the results of Gillespie algorithm and the deterministic solutions of Eq.(5.19). The plots depict a) progression of mole fraction of S and b) progression of mole fraction of P

### 5.4.3 Rare event sampling in biological system

Bi-stable chemical networks are a mainstay in cellular biology, as they are central to cell fate determination (Elf & Ehrenberg, 2004), such as the lysis-lysogeny switch in  $\lambda$ -phage (Arkin *et al.*, 1998). We consider a simplified model, wherein, two proteins namely  $A$  and  $B$  try and bind with the DNA  $O$ , the reaction network is given by:



The double negative feedback nature of the network arising from the fact that both  $A$  and  $B$  inhibit each other's production after binding with the DNA  $O$  renders it to have two stable states (Allen *et al.*, 2005). There essentially two stable states of the system – one where number of  $A$  molecules is much larger than  $B$  and the other where number of  $B$  molecules is much larger than  $A$ . The switching between the two states is almost instantaneous and is driven by fluctuations. Hence for such systems, stochastic simulation algorithms are employed as they manage to capture rich and complex behaviour as opposed to solving equations arising out of rate law of mass action. The time interval between switching events is exponentially distributed,  $p(t) = k_{AB}e^{-k_{AB}t}$ . The associated parameter,  $k_{AB}$ , was calculated using the proposed algorithm and was found to be  $4.19 \times 10^{-5}$ , which is in good agreement with reported data (Allen *et al.*, 2005). Progression of  $\Delta = (n_A + 2n_{A_2} + 2n_{OA_2}) - (n_B + 2n_{B_2} + 2n_{OB_2})$  for a typical simulation is shown in Fig.(5.4). For  $10^9$  reaction events, the time spent on random number generation via traditional methods was found to be 94 seconds and 12 seconds for proposed method, thereby speeding up the entire algorithm by a factor 4.

## 5.5 Reaction-Diffusion Systems

Systems considered in the chapter were idealized to be homogenous. However, it has been observed that spatial heterogeneity imparts an extra layer of complexity and is in fact the basis for many important phenomena such as — morphogenesis (Turing, 1952), chemotaxis (Thar & Kühl, 2003) etc. The stochastic formulation can be tweaked to account for diffusion of molecules in space.

The entire domain is divided into sub-cells wherein all molecules are considered to be close enough to react with each other. The diffusion is incorporated in the system by treating movement of molecules between neighbouring cells as reaction events. The rate constant of this “pseudo-reaction” is given by  $k = D/h^2$ , where  $D$  is the diffusion constant and  $h$  the length parameter associated with the subcells (Baras & Mansour, 1997; Gardiner, 1985*a*).

### 5.5.1 Pattern formation in bacteria

Cell division in bacteria such as *E. Coli* is facilitated by a system known as the MinCDE system of proteins. These proteins while transferring between the cytoplasm and cytoplasmic membrane diffuse over the length of the bacteria with their concentrations being minimum around the

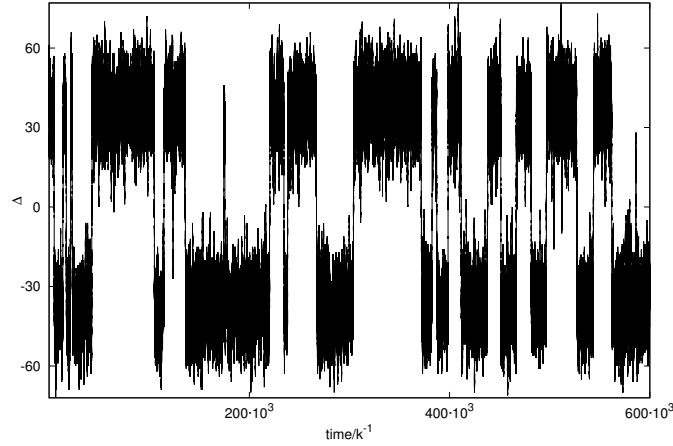


Figure 5.4: Plot of  $\Delta = (n_A + 2n_{A_2} + 2n_{OA_2}) - (n_B + 2n_{B_2} + 2n_{OB_2})$  for the reaction network presented in Eq.(5.4.3). It can be seen clearly that small perturbations can lead to switching between multiple steady states.

center of the cell which indicates the site for cell division (Howard *et al.*, 2001). However, the number of protein copies in a cell is low (around 3000) which results in significant fluctuations and drive the oscillations of the *min* proteins across the cell while the deterministic formulation is unable to capture these effects owing to suppression of noise which drives this system (Howard & Rutenberg, 2003). The rate law of mass action coupled with diffusion in space is used to describe the concentrations of various components – MinD in cytoplasm ( $C_D$ ), MinD on the cytoplasmic membrane ( $C_d$ ), MinE in the cytoplasm ( $C_E$ ) and MinE in the cytoplasmic membrane ( $C_e$ )

$$\begin{aligned}
 \frac{\partial C_d}{\partial t} &= \frac{\sigma_1 C_D}{1 + \sigma'_1 C_e} - \sigma_2 C_e C_d \\
 \frac{\partial C_D}{\partial t} &= D_D \frac{\partial^2 C_D}{\partial x^2} - \frac{\sigma_1 C_D}{1 + \sigma'_1 C_e} + \sigma_2 C_e C_d \\
 \frac{\partial C_e}{\partial t} &= \frac{\sigma_4 C_E}{1 + \sigma'_4 C_D} - \sigma_3 C_D C_E \\
 \frac{\partial C_E}{\partial t} &= D_E \frac{\partial^2 C_E}{\partial x^2} - \frac{\sigma_4 C_E}{1 + \sigma'_4 C_D} + \sigma_3 C_D C_E
 \end{aligned} \tag{5.21}$$

and the various  $\sigma_i$  are the associated constants of the various reactions. Similar to the deterministic formulation, the concentration of MinD proteins is close to the center of the cell. A scatter plot of the concentration of MinD over an oscillation cycle is plotted in Fig.(5.5a). Additionally, the ratio of this protein on either side of the center has been observed to oscillate periodically, a phenomenon that is observed in actual experiments also. The oscillation of MinD protein in time across space is presented in Fig.(5.5b). Hence, through this example it can be seen that the stochastic formulation of spatio-temporal phenomena provides a more detailed description as compared to their deterministic counterparts.

## 5.6 Outlook

Simulation of reactive systems and reaction-diffusion systems is an important problem in computational physics. The deterministic formulation provides an accurate description of the average behaviour of large systems where fluctuations about the mean are insignificant, however for small systems these perturbations can prove to be highly important and subsequently lead to

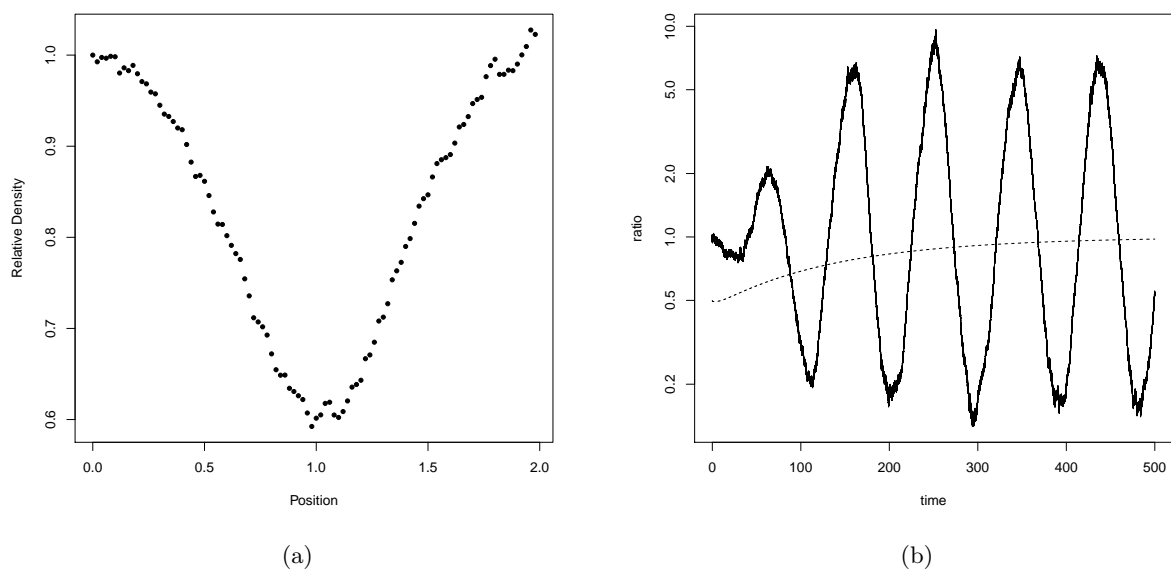


Figure 5.5: Simulation of the reaction-diffusion process presented in Eq.(5.21). The plots depict a) a scatter plot of the averaged density profile of MinD protein over an oscillation cycle. A minima can be observed close to the center of the cell, indicating the site for cell division. This is in good agreement with the deterministic formulation of the process. b) A plot of the ratio of MinD protein in the left-hand 30% to right-hand 30%. The solid line is observed in the stochastic simulation while the dotted line is the deterministic solution. While the latter fails to capture the oscillations observed in experiments, the stochastic version of the model manages to capture this phenomena.

highly complex characteristics such as random switching between multiple stable states of the system. This is also the case for reaction-diffusion systems where spatial heterogeneity plays an important role in the dynamics of the system such as pattern formation in bacteria. The diffusion of various components can be handled in a stochastic framework, by treating diffusion of molecules across neighbouring cells as reaction events. However, this method has major drawbacks as the time scales for chemical reactions is much lower than diffusion, hence the entire system progresses at the time scale of the former, hence a major fraction of the computational time is spent simulating diffusion. Since stochastic numerical schemes to simulate chemical reactions require large number of random sequences, the cost of computation for such algorithms increases drastically making it unfeasible for most real-world problems.





# Chapter 6

## Fokker-Planck model for rarefied gases

### 6.1 Introduction

The simulation of fluid flows is an important problem and finds applications in many areas of science and engineering. Typically, these simulations are achieved via numerically solving the Navier-Stokes-Fourier (NSF) equations, which are mathematical statements of momentum and energy balance of fluid flow. While these numerical methods have achieved considerable success, they are only applicable to flows where the mean free path ( $\lambda$ ) is considerably smaller than the characteristic length of the flow ( $L$ ). The ratio of  $\lambda$  and  $L$  is known as the Knudsen number (Kn) and is an important dimensionless parameter which characterizes the flow (Chapman & Cowling, 1970). Low Kn numbers ( $\leq 0.01$ ) indicate the validity of continuum hypothesis, a regime where the NSF equations are applicable. However, for setups such as fuel cells and shale gas transport the characteristic length is of the order of micrometers and the continuum approximation breaks down. For such cases, the dynamics of the fluid flow is well described by the Boltzmann equation for dilute gases (introduced in Chapter 4). This equation is difficult to solve analytically, and therefore numerical methods such as the Direct Simulation Monte Carlo (DSMC) are employed (Bird, 1978). Although DSMC provides accurate solutions to the Boltzmann equation, the computational time for low Kn cases tend to be high due to large statistical fluctuations (Bird, 1994). Numerical methods such as the Lattice Boltzmann (LB) utilize simplified collision models and discrete velocity models resulting in higher computational efficiency (Succi, 2001). While the lower-order velocity models have been shown to produce correct results for  $\text{Kn} < 0.1$ , the higher-order models provide correct results for  $\text{Kn} < 0.25$  for isothermal setups. The last decade has seen a revived interest in the Fokker-Planck approximation to Boltzmann equation and equivalent Langevin dynamics, for computational reasons. It provides results as accurate as DSMC while proving to be computationally cheaper, hence posing as a viable alternative in the context of mesoscopic simulation methods (Singh *et al.*, 2016).

In this chapter, the basics of kinetic modelling of rarefied gases is explained first, followed by a brief outline of the Fokker-Planck approximation. The chapter ends with a description of the numerical method used to solve this model and the algorithm overview.

### 6.2 Kinetic modelling of rarefied gases

Kinetic theory of gases models the molecular motion and is based on statistical description in terms of the the distribution function,  $f$ , where  $f(\mathbf{x}, \mathbf{c}, t)d\mathbf{x}d\mathbf{c}$  is the probability of finding a particle with position in the range  $(\mathbf{x}, \mathbf{x} + d\mathbf{x})$ , possessing a velocity in the range  $(\mathbf{c}, \mathbf{c} + d\mathbf{c})$ . The relevant macroscopic quantities can then be found by taking the appropriate moment of the distribution function. The lower order moments are defined as

$$n = \langle 1, f \rangle, \quad \rho = mn, \quad \rho \mathbf{u} = \langle m\mathbf{c}, f \rangle, \quad E = \langle mc^2/2, f \rangle, \quad (6.1)$$

where  $n$  is the number density,  $m$  mass of the particle,  $\rho$  the mass density,  $\rho \mathbf{u}$  is the mass flux or momentum,  $E$  the energy density and the  $\langle \cdot, \cdot \rangle$  is an operator denoting

$$\langle \phi_1, \phi_2 \rangle = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \phi_1 \phi_2 dc_x dc_y dc_z. \quad (6.2)$$

The energy density  $E$  can be seen as a sum of two parts – energy arising from the bulk motion of the fluid and the random motion of particles stemming from thermal energy. The energy density is defined as

$$E = \frac{\rho u^2}{2} + \frac{Dnk_B T}{2}, \quad (6.3)$$

where  $k_B$  is the Boltzmann constant and  $T$  the temperature. Using this definition, the temperature can then be simply defined as

$$\frac{Dnk_B T}{2} = \langle \xi^2 / 2, f \rangle \quad (6.4)$$

where  $\xi = \mathbf{c} - \mathbf{u}$  is the peculiar velocity. The Boltzmann equation introduced in Chapter 4, describes the evolution of the distribution function in time and space for the dilute limit

$$\partial_t f + \partial_{c_\alpha} f = \Omega^B, \quad (6.5)$$

where  $\Omega^B$  is the Boltzmann collisional operator quantifying the change in the distribution function from binary collisions between particles. By taking appropriate moments and integrating over the velocity space, the dynamics of various macroscopic quantities can be derived (Chapman & Cowling, 1970). Some salient features of the Boltzmann equation are

1. Conservation of density, momentum and energy as no changes happen due to binary collision, that is

$$\langle \Omega^B, \{m, m\mathbf{c}, mc^2/2\} \rangle = \{0, \mathbf{0}, 0\}, \quad (6.6)$$

using this result and calculating appropriate moments of the Boltzmann equation, the conservation laws are obtained as

$$\begin{aligned} \partial_t \rho + \partial_\alpha j_\alpha &= 0, \\ \partial_t j_\alpha + \partial_\beta (\rho u_\alpha u_\beta + p \delta_{\alpha\beta}) + \partial_\beta \sigma_{\alpha\beta} &= 0, \\ \partial_t E + \partial_\alpha ((E + p)u_\alpha + \sigma_{\alpha\gamma} u_\gamma) + \partial_\alpha q_\alpha &= 0, \end{aligned} \quad (6.7)$$

where  $p$  is the pressure,  $\sigma_{\alpha\beta}$  the stress and  $q_\alpha$  the heat flux defined in kinetic terms as

$$p = nk_B T, \quad \sigma_{\alpha\beta} = \overline{\langle \xi_\alpha \xi_\beta, f \rangle}, \quad q_\alpha = \left\langle \frac{\xi^2}{2} \xi_\alpha, f \right\rangle, \quad (6.8)$$

with  $\overline{A_{\alpha\beta}}$  indicating the traceless part of the tensor.

2. The system reaches a statistical steady state  $\Omega^B = 0$ , i.e, the collisions do not bring about any change to the distribution function when,  $f$  attains a Maxwell-Boltzmann of the form

$$f^{\text{MB}} = \rho \left( \frac{m}{2\pi k_B T} \right)^{3/2} \exp \left( -\frac{m}{2k_B T} (c - u)^2 \right). \quad (6.9)$$

3. The Boltzmann equation extends the idea of entropy to non-equilibrium situations. This is highlighted from the evolution of the  $H$ -function

$$H = \int d\mathbf{c}(f \ln f - f), \quad (6.10)$$

which is essentially the non-equilibrium generalization of entropy. The evolution of this quantity is given by the equation

$$\partial_t H + \partial_\alpha J_\alpha = -\sigma^S \quad (6.11)$$

where  $J_\alpha$  is the entropy flux term. The Boltzmann collisional operator holds the property

$$\sigma^S = \langle \Omega^B, \ln f \rangle \leq 0 \quad (6.12)$$

which ensures that entropy production is greater than 0 and hence the Boltzmann equation for rarefied gases is in accordance with the laws of thermodynamics. It is also noted that entropy production is zero at equilibrium, that is when  $f = f^{\text{MB}}$ .

The Boltzmann equation while accurate is highly complex and thus analysis is often difficult. In order to mitigate this problem, the Boltzmann collisional operator is approximated in manner such that all the aforementioned properties (conservation laws, zero of collision and  $H$ -theorem) are preserved. A highly popular example of such an approximation is the BGK approximation (Bhatnagar *et al.*, 1954a)

$$\Omega^{\text{BGK}} = \frac{1}{\tau_{\text{BGK}}} (f^{\text{MB}} - f), \quad (6.13)$$

where  $\tau_{\text{BGK}}$  is the mean free time and  $f^{\text{MB}}$  the Maxwell-Boltzmann distribution defined previously. It can be seen that such a collisional term maintains all the canonical properties of a valid kinetic model. With this approximation, the Boltzmann equation can be solved using numerical methods such as the Lattice Boltzmann (LB), which rely on discretization of the velocity space to a finite set of discrete velocities (Succi, 2001).

### 6.3 Quasi-equilibrium models

A major drawback of BGK approximation is that it is incapable of incorporating the different time scales at which the higher order moments relax to their respective equilibrium values. The quasi-equilibrium modelling approach ameliorates this problem by accounting for these different time scales. In this approach, the moments are categorised as

$$M = (M^{\text{slow}}, M^{\text{quasi-slow}}, M^{\text{fast}})$$

where  $M^{\text{slow}}$  represents the conserved quantities – mass, momentum and energy which are not perturbed from their equilibrium values,  $M^{\text{quasi-slow}}$  represents the set of higher-order moments which equilibrate fast and  $M^{\text{fast}}$  is the set of higher order moments which relax to their equilibrium values slowly. Hence, the quasi-equilibrium distribution represents the state of the system once the quantity relaxing faster has assumed its equilibrium value. The dynamics of the system is modelled as a two-step process, where the distribution function first approaches the quasi-equilibrium, the distribution function that is constrained by conservation of quasi-conserved variables and then proceeds to the Maxwell-Boltzmann distribution (Levermore, 1996). This idea is visually represented in Fig.(6.1). The quasi-equilibrium following such dynamics is

$$\Omega^{\text{BGK}} = \frac{1}{\tau_1} \left( f^*(M^{\text{slow}}, M^{\text{quasi-slow}}) - f \right) + \frac{1}{\tau_2} \left( f^{\text{MB}}(M^{\text{slow}}) - f^*(M^{\text{slow}}, M^{\text{quasi-slow}}) \right) \quad (6.14)$$

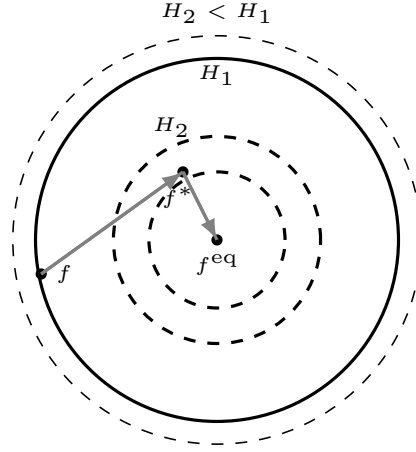


Figure 6.1: Sketch representing the two-step relaxation to equilibrium.

where  $\tau_1$  and  $\tau_2$  are different time scales of the system such that ( $\tau_1 < \tau_2$ ) and  $f^*$  is the quasi-equilibrium distribution function.

A simple example is to consider  $f^*(\rho, u_\alpha, \theta_{\alpha\beta})$  given by

$$f^* = n^{(i)} \left( \frac{m^{(i)}}{2\pi|\theta_{\alpha\beta}|} \right)^{\frac{3}{2}} \exp \left( \frac{-m^{(i)}\xi_\alpha\theta_{\alpha\beta}^{-1}\xi_\beta}{2} \right), \quad (6.15)$$

where  $\theta_{\alpha\beta} = \langle m\xi_\alpha\xi_\beta, f \rangle / n$ . This quasi-equilibrium distribution function corresponds to the case where the pressure attains its equilibrium value faster than than the heat flux and other higher-order moments. This case would correspond to systems with high Prandtl Numbers (Levermore, 1996).

## 6.4 Fokker-Planck approximation

The Fokker-Planck model approximates the Boltzmann collisional operator as (Lebowitz *et al.*, 1960a)

$$\Omega^{FP} = \frac{1}{\tau} \partial_{c_\alpha} \left( \xi_\alpha f + \frac{k_B T}{m} \frac{\partial f}{\partial c_\alpha} \right), \quad (6.16)$$

which is essentially the diffusion dynamics in velocity space, with  $\xi_\alpha$  acting as the drift constant,  $k_B T/m$  as the diffusion constant and  $\tau^{-1}$  as the friction constant. It is quite evident that this particular form of the collision kernel ensures correct form of conservation laws and that zero of collision is indeed the Maxwell-Boltzmann distribution. For this model, the entropy production is

$$\sigma^{(FP)} = -\langle \Omega^{FP}, f \rangle = \frac{1}{\tau} \int d\mathbf{c} \left( \xi_\alpha + \frac{1}{f} \frac{k_B T}{m} \frac{\partial f}{\partial c_\alpha} \right) \frac{\partial f}{\partial c_\alpha}, \quad (6.17)$$

which is positive as (Singh & Ansumali, 2015a)

$$\sigma^{(FP)} = -\frac{3\rho}{\tau} + \frac{k_B T}{m\tau} \int d\xi \frac{1}{f} \frac{\partial f}{\partial \xi_\alpha} \frac{\partial f}{\partial \xi_\alpha} = \frac{k_B T}{m\tau} \int d\xi f \left( \frac{\partial \ln \left( \frac{f}{f^{MB}} \right)}{\partial \xi_\alpha} \right)^2 \geq 0, \quad (6.18)$$

this is so because the integrand in the last term is completely positive by the virtue of  $f$  being a positive quantity which in turn is multiplied to a perfect square. To show this, the following identity was used

$$\begin{aligned} \int d\xi f \left( \frac{\partial \ln \left( \frac{f}{f^{\text{MB}}} \right)}{\partial \xi_\alpha} \right)^2 &= \int d\xi \left[ \frac{1}{f} \frac{\partial f}{\partial \xi_\alpha} \frac{\partial f}{\partial \xi_\alpha} + f \frac{\partial \ln f^{\text{MB}}}{\partial \xi_\alpha} \frac{\partial \ln f^{\text{MB}}}{\partial \xi_\alpha} - 2 \frac{\partial f}{\partial \xi_\alpha} \frac{\partial \ln f^{\text{MB}}}{\partial \xi_\alpha} \right] \\ &= \int d\xi \left[ \frac{1}{f} \frac{\partial f}{\partial \xi_\alpha} \frac{\partial f}{\partial \xi_\alpha} + f \frac{m^2 \xi^2}{k_B^2 T^2} + 2 \frac{\partial f}{\partial \xi_\alpha} \frac{m \xi_\alpha}{k_B T} \right] \\ &= -\frac{3 \rho m}{k_B T} + \int d\xi \frac{1}{f} \frac{\partial f}{\partial \xi_\alpha} \frac{\partial f}{\partial \xi_\alpha}. \end{aligned} \quad (6.19)$$

This proves that this approximation is consistent with the  $H$ -theorem. Hence, the Fokker-Planck approximation satisfies all the properties of Boltzmann collisional operator and is hence considered to be a valid kinetic description of gaseous flow.

### 6.4.1 Transport Coefficients

The stress evolution equation obtained by taking the appropriate moment of Eq.(6.16) is

$$\partial_t \sigma_{\alpha\beta} + \partial_\gamma (\sigma_{\alpha\beta} u_\gamma) + 2p \overline{\partial_\alpha u_\beta} + 2 \overline{\sigma_{\alpha\gamma} \partial_\gamma u_\beta} + \partial_\gamma Q_{\alpha\beta\gamma} + \frac{4}{D+2} \overline{\partial_\alpha q_\beta} = -\frac{2}{\tau} \sigma_{\alpha\beta}, \quad (6.20)$$

and similarly the evolution equation of heat flux is

$$\begin{aligned} \partial_t q_\alpha + \partial_\beta \left( q_\alpha u_\beta + \frac{R_{\alpha\beta}}{2} + \frac{R' \delta_{\alpha\beta}}{2D} \right) + \frac{(D+2)}{2} p \partial_\alpha \frac{p}{\rho} + \frac{2}{D+2} (q_\gamma \partial_\alpha u_\gamma + q_\alpha \partial_\beta u_\beta) \\ - \frac{\sigma_{\alpha\beta} \partial_\beta p}{\rho} + \frac{D+4}{D+2} q_\beta \partial_\beta u_\alpha + Q_{\alpha\beta\gamma} \partial_\beta u_\gamma - \frac{(D+2)p}{2\rho} \partial_\beta \sigma_{\alpha\beta} - \frac{\sigma_{\alpha\kappa} \partial_\beta \sigma_{\kappa\beta}}{\rho} = -\frac{3}{\tau} q_\alpha, \end{aligned} \quad (6.21)$$

where the higher order moments are defined as

$$Q_{\alpha\beta\gamma} = \int d\mathbf{c} f \overline{\xi_\alpha \xi_\beta \xi_\gamma}, \quad R' = \int d\mathbf{c} f \xi^2 \xi^2 - 15 \frac{p^2}{\rho}, \quad R_{\alpha\beta} = \int d\mathbf{c} f \xi^2 \overline{\xi_\alpha \xi_\beta}. \quad (6.22)$$

Here,  $D$  denotes the number of dimensions. This model yields different relaxation times for stress and heat flux evolution as compared to the BGK model, where the relaxation rates are

$$\langle \Omega^{\text{BGK}}, \xi_\alpha \xi_\beta \rangle = -\frac{1}{\tau} \sigma_{\alpha\beta} \quad \langle \Omega^{\text{BGK}}, \xi_\alpha \xi^2 / 2 \rangle = -\frac{1}{\tau} q_\alpha. \quad (6.23)$$

Hence the Prandtl number associated with BGK model is 1 while it is 3/2 for the Fokker-Planck model. However, the latter can be tuned to adjust the Prandtl number (Singh & Ansumali, 2015a).

These evolution equations form a moment chain where the evolution equation of a given quantity includes terms of a higher order. To derive the transport coefficients, the standard Chapman-Enskog methodology is used, wherein the non-equilibrium distribution function  $f$  is expanded about  $f^{\text{MB}}$  as a perturbation series with  $\tau$  being the smallness parameter, and the time derivative is also expressed as a perturbation series

$$\begin{aligned}\partial_t \phi &= \partial_t^{(0)} \phi + \tau \partial_t^{(1)} \phi + \tau^2 \partial_t^{(2)} \phi + \dots \\ f &= f^{MB}(\rho, \mathbf{u}, T) + \tau f^{(1)} + \tau^2 f^{(2)} + \dots\end{aligned}\quad (6.24)$$

The higher orders of the distribution function must satisfy the constraint

$$\left\langle f^{(n)}, \{m, m\mathbf{c}, mc^2/2\} \right\rangle = \{0, \mathbf{0}, 0\} \quad (n \geq 1), \quad (6.25)$$

this ensures that the conservation laws for mass, momentum and energy are satisfied at the zeroth order. The expressions can for the time derivative can be found from the conservation laws by considering that the distribution function is a function of conserved variables (Liboff, 2003). Other higher order moments, in terms of this expansion can be written as

$$\begin{aligned}\sigma_{\alpha\beta} &= \tau \sigma_{\alpha\beta}^{(1)} + \tau^2 \sigma_{\alpha\beta}^{(2)} + \dots \\ q_\alpha &= \tau q_\alpha^{(1)} + \tau^2 q_\alpha^{(2)} + \dots\end{aligned}\quad (6.26)$$

Similarly, other higher-order moments are written as

$$\begin{aligned}Q_{\alpha\beta\gamma} &= \tau Q_{\alpha\beta\gamma}^{(1)} + \tau^2 Q_{\alpha\beta}^{(2)} + \dots \\ R_{\alpha\beta} &= \tau R_{\alpha\beta}^{(1)} + \tau^2 R_{\alpha\beta}^{(2)} + \dots\end{aligned}\quad (6.27)$$

Using this series in the respective evolution equations and considering terms upto first order which corresponds to the hydrodynamic limit, the following set of relations are found

$$\sigma_{\alpha\beta} = -p\tau \overline{\partial_\alpha u_\beta}, \quad q_\alpha = -\tau \frac{D+2}{6} p \partial_\alpha \frac{p}{\rho} \quad (6.28)$$

then the viscosity ( $\mu$ ) and thermal conductivity ( $\kappa$ ) are found to be

$$\mu = \frac{p\tau}{2}, \quad \kappa = \tau \left( \frac{D+2}{6} p \right). \quad (6.29)$$

Using this expression for  $\mu$  and that  $\mu = \rho \bar{v} \lambda / 2$ , where  $\bar{v}$  is the root mean velocity, the expression for Kn can then be formulated as

$$Kn = \frac{\tau}{L} \sqrt{\frac{k_B T_0}{2m}}, \quad (6.30)$$

where  $L$  is the characteristic length associated with the flow, and  $T_0$  the characteristic temperature of the system.

## 6.5 Numerical solution

The solution to the Fokker-Planck model is achieved using the equivalent Langevin equations (Risken, 1996; Singh *et al.*, 2016)

$$\begin{aligned}\frac{dx_\alpha}{dt} &= c_\alpha \\ \frac{dc_\alpha}{dt} &= -\frac{\xi_\alpha}{\tau} + \sqrt{\frac{2k_B T}{m}} \frac{dW_\alpha}{dt},\end{aligned}\quad (6.31)$$

where  $dW_\alpha$  is the standard Weiner process, a random force with Gaussian distribution as encountered in the case of Brownian motion. These sets of equations can be discretized using the stochastic Verlet scheme (Ladd, 2009)

$$\begin{aligned} x_\alpha^{(1)} &= x_\alpha + \frac{c_\alpha(t)}{2} \Delta t \\ c_\alpha(t + \Delta t) &= c_\alpha(t) - \frac{\vartheta}{1 + \vartheta/2} (c_\alpha(t) - U_\alpha) + \frac{\sqrt{2\mathcal{D}\vartheta}}{1 + \vartheta/2} \phi_t \\ x_\alpha(t + \Delta t) &= x_\alpha^{(1)} + \frac{c_\alpha(t + \Delta t)}{2} \Delta t \end{aligned} \quad (6.32)$$

where  $\vartheta = \Delta t/\tau$  and  $\phi_t$  are normally distributed random variables.

The resulting algorithm is a particle based method, where the computational domain is divided into cells such that the length of each cell is of the order of mean free path and each of these are populated with particles whose positions and velocities are distributed in accordance with the initial conditions. The positions and velocities are updated as per Eqs. (6.32), and the relevant macroscopic quantities are calculated for each cell, these steps are iterated over until desired simulation time is achieved. It is evident that this algorithm has a computational complexity of  $O(N)$ , and like Brownian motion codes majority of the computational time is in fact spent on generating Gaussian random numbers.

## 6.6 Outlook

The Fokker-Planck approximation to the Boltzmann equation provides an efficient alternative to numerical methods such as the Lattice Boltzmann and Direct Simulation Monte Carlo to solve for hydrodynamics at mesoscale. Since the major bottleneck for this algorithm is the generation of Gaussian random numbers and given that an algorithm for producing Gaussian random numbers at a high rate exists, this algorithm can indeed prove to be highly efficient for simulating fluid flow. As mentioned previously, a major drawback of Gillespie algorithm for reaction-diffusion processes is the manner in which it handles diffusion. Being a particle based method, intuitively the Fokker-Planck model can be coupled with Gillespie algorithm to simulate reaction-diffusion processes. However, the current formulations of Fokker-Planck approximations are limited to simulating single components and hence must be extended to correctly model multicomponent gas mixtures.





# Chapter 7

## Fokker-Planck model for binary mixtures

### 7.1 Introduction

Kinetic modelling of gases for boundary value problems pertaining to engineering needs have found success for single component case. However, techniques dealing with gas mixtures haven't yet attained the same level of sophistication. Analytical techniques based on classical Boltzmann equation are difficult to implement for mixtures beyond stationary linearized problems. Standard approaches based on molecular dynamics such as direct simulation Monte Carlo (DSMC) while highly useful for simulating flows at large Knudsen numbers, becomes computationally expensive as one approaches the continuum limit. Since high dimensionality and complexity of the Boltzmann collision kernel render analytical and numerical solutions difficult and expensive, it is imperative to approximate this kernel. In this regard, the Bhatnagar - Gross - Krook (Bhatnagar *et al.*, 1954*b*) approximation (Lebowitz *et al.*, 1960*b*) has found considerable success. Building on this approximation, efforts have been made to kinetically model binary mixtures (Sirovich, 1962; Morse, 1964; Hamel, 1965; Sirovich, 1966; Goldman & Sirovich, 1967). These models have been shown to converge to the Navier-Stokes equation and Stefan-Maxwell diffusion equation in the hydrodynamic limit and uphold both the  $H$ -theorem and indifferenciability principle. These models such as the BGK-model are numerically solved via a class of methods known as Lattice Boltzmann methods (LBM), which have evolved as an attractive approach to simulate gas flows. Although a highly efficient method, it is limited to simulating flows in the low Knudsen number regime and also faces issues for system with sharp density gradients.

An alternative route based on Fokker-Planck approximation of the collision kernel was recently explored, and was shown to produce satisfactory results (Singh & Ansumali, 2015*b*). As opposed to traditional grid based methods such as LBM, the solution to Fokker-Planck models is obtained by discretizing the equivalent Langevin equation in the entire phase space, hence the resulting algorithm updates the velocity of each particle via a combination of drift and diffusion terms instead of modelling binary collisions. Hence unlike DSMC, this model is  $\mathcal{O}(N)$  and hence proves to be a computationally attractive alternative for low to intermediate range of Knudsen numbers (Singh *et al.*, 2015).

In this chapter, the basics of kinetic modelling of binary mixtures is explained first followed by a description of the quasi-equilibrium models. Following this, two models for binary mixtures based on the Fokker-Planck approximation is presented. In these sections, we prove the models correspondence with conservation laws, the  $H$ -theorem and indifferenciability principle. In the following section, we calculate the expression for the transport coefficients associated with two models. Lastly the numerical scheme to solve these models is present along with three canonical problems to demonstrate the efficiency of this formulation — Graham's Law for effusion, Couette flow and binary diffusion.

### 7.2 Kinetic modelling of binary mixtures

For binary mixtures, the basic dynamic quantity is the distribution function,  $f_i$  for ( $i = A, B$ ), where  $f_i(\mathbf{x}, \mathbf{c}_i, t) d\mathbf{c}_i d\mathbf{x}$  is the probability of finding a particle of  $i$ th type in the neighbourhood of the point  $(\mathbf{x}, \mathbf{c}_i, t)$  (Chapman & Cowling, 1970). The relevant macroscopic variables are defined

as

$$\begin{aligned}
n_i &= \langle 1, f_i \rangle, & n &= \sum_{i=A,B} n_i, \\
\rho_i &= m_i n_i, & \rho &= \sum_{i=A,B} \rho_i, \\
\rho \mathbf{u} &= \sum_{i=A,B} \langle m_i \mathbf{c}_i, f_i \rangle, \\
\frac{D}{2} n k_B T &= \sum_{i=A,B} \left\langle \frac{m_i (\mathbf{c}_i - \mathbf{u})^2}{2}, f_i \right\rangle,
\end{aligned} \tag{7.1}$$

where  $n_i$  is the component number density  $n$  is the number density,  $\rho_i$  the component mass density,  $\rho$  the mixture mass density,  $\mathbf{u}$  the mixture velocity and  $T$  the mixture temperature,  $D$  the number of dimensions and  $\langle \phi_1, \phi_2 \rangle$  operator denotes

$$\langle \phi_1(\mathbf{c}_i), \phi_2(\mathbf{c}_i) \rangle = \int_{-\infty}^{\infty} \phi_1(\mathbf{c}_i) \phi_2(\mathbf{c}_i) d\mathbf{c}_i. \tag{7.2}$$

Additionally, another two useful variables – component velocity denoted by  $\mathbf{u}_i$  and component temperature  $T_i$  are defined as

$$\begin{aligned}
\rho_i \mathbf{u}_i &= \langle m_i \mathbf{c}_i, f_i \rangle, \\
\frac{3}{2} n_i k_B T_i &= \left\langle \frac{m_i (\mathbf{c}_i - \mathbf{u}_i)^2}{2}, f_i \right\rangle.
\end{aligned} \tag{7.3}$$

The Boltzmann equation describes the time evolution of the distribution function by considering the different collisional possibilities, as schematically shown in Fig 7.1. The Boltzmann equation for binary mixtures is (Chapman & Cowling, 1970)

$$\begin{aligned}
\frac{\partial f_A}{\partial t} + c_{A\alpha} \frac{\partial f_A}{\partial x_\alpha} &= \Omega_A^{\text{BM}} = \underbrace{\Omega^{\text{BM}}(f_A, f_A)}_{\text{Self-collision}} + \underbrace{\Omega^{\text{BM}}(f_A, f_B)}_{\text{Cross-collision}}, \\
\frac{\partial f_B}{\partial t} + c_{B\alpha} \frac{\partial f_B}{\partial x_\alpha} &= \Omega_B^{\text{BM}} = \underbrace{\Omega^{\text{BM}}(f_B, f_B)}_{\text{Self-collision}} + \underbrace{\Omega^{\text{BM}}(f_B, f_A)}_{\text{Cross-collision}},
\end{aligned} \tag{7.4}$$

where the right-hand side of the equation is the change in distribution of the respective components arising from self collisions and cross collisions. Using this equation, the time evolution of the macroscopic quantities can be calculated. The Boltzmann collision kernel for binary mixtures satisfies the following constraints

1. The mass of individual species as well as the total momentum and energy of the mixture are conserved as binary collisions do not contribute any change to these quantities, this is represented as

$$\begin{aligned}
\langle \Omega_i^{\text{BM}}, m_i \rangle &= 0, \\
\sum_{i=A,B} \langle \Omega_i^{\text{BM}}, \{m_i \mathbf{c}_i, m_i c_i^2/2\} \rangle &= \{\mathbf{0}, 0\},
\end{aligned} \tag{7.5}$$

using which the conservations laws can be calculated similar to the single component case.

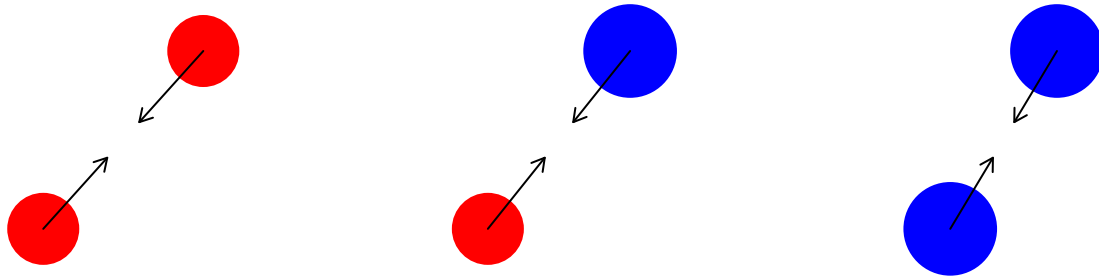


Figure 7.1: The three types of collisional possibilities – A-A, B-B and A-B

However, the component momentum and energy are not conserved as the two components exchange momentum and energy between themselves through cross-collisions (A-B type collisions). It is in fact these collisions that facilitate the relaxation of the component momentum and energy to the mixture momentum and energy (Hamel, 1965).

2. Similar to the single component case, the system reaches a state of statistical equilibrium. The distribution of a component at equilibrium is

$$f_i^{\text{MB}} = n_i \left( \frac{m_i}{2\pi k_B T} \right)^{3/2} \exp \left( -\frac{m_i}{2k_B T} (\mathbf{c}_i - \mathbf{u})^2 \right), \quad (7.6)$$

thus the distribution function at equilibrium is of the form  $f_i^{\text{MB}}(\rho_i, \mathbf{u}, T)$ . The converse for this also true, i.e, when  $\Omega_i^{\text{BM}} = 0$  then the distribution function attains the form  $f = f_i^{\text{MB}}$

3. The Boltzmann collision kernel for binary mixtures, satisfies the  $H$ -theorem, with the  $H$ -function defined as

$$H = \sum_{i=A,B} \int d\mathbf{c}_i (f_i \ln f_i - f_i), \quad (7.7)$$

which extends the idea of entropy to non-equilibrium cases.

The Boltzmann equation does not assume the continuum hypothesis thereby correctly predicts the behaviour of fluids at higher Knudsen numbers (Cercignani, 1975), but its highly complex form is a deterrent to finding its solution. Thus, approximations are made to the collisional term in order to simplify the form of the equation and devise efficient numerical schemes.

### 7.3 Quasi-equilibrium models for binary mixtures

The corresponding BGK collision kernel for a binary mixture is

$$\Omega^{\text{BGK}} = \frac{1}{\tau} (f_i^{\text{MB}}(\rho_i, \mathbf{u}, T) - f_i). \quad (7.8)$$

The fundamental drawback with such a model is that there is only a single relaxation rate for all quantities whereas for the case of a binary mixture, there are two important time scales present in the system – the rate of mass diffusion and the rate of momentum of diffusion. The dimensionless parameter that is used to characterize these time scales is known as the Schmidt number and is defined as (Bergman *et al.*, 2011)

$$\text{Sc} = \frac{\text{viscous diffusion rate}}{\text{mass diffusion rate}} = \frac{\mu}{\rho D_{AB}}, \quad (7.9)$$

where  $\mu$  is the viscosity,  $\rho$  the density and  $D_{AB}$  is the mass diffusion coefficient. In this context, quasi-equilibrium models are a simple alternative (Levermore, 1996). The collision kernel for quasi-equilibrium model is

$$\Omega_i^{\text{QE}} = \frac{1}{\tau_1}(f_i^*(M^{\text{quasi-slow}}, M^{\text{slow}}) - f_i) + \frac{1}{\tau_2}(f_i^{\text{MB}}(M^{\text{slow}}) - f_i^*(M^{\text{quasi-slow}}, M^{\text{slow}})), \quad (7.10)$$

where  $f_i^*(M^{\text{quasi-slow}}, M^{\text{slow}})$  is the quasi-equilibrium distribution function and is a function of the quasi-slow and the slow moments (Levermore, 1996). In accordance with the slow-fast dynamics that emerges from quasi-equilibrium models, two possible forms for the quasi-equilibrium distribution can be chosen – for low Sc where mass diffusion occurs at higher rate as compared to momentum diffusion and vice versa for the high Sc case. For the first case, the physically quasi-slow variables impose the following conditions on quasi-equilibrium distribution function  $f_i^*$

$$\begin{aligned} \langle m_i, f_i^* \rangle &= \rho_i, \\ \langle m_i \mathbf{c}_i, f_i^* \rangle &= \rho_i \mathbf{u}_i, \\ \left\langle m_i \frac{(\mathbf{c}_i - \mathbf{u}_i)^2}{2}, f_i^* \right\rangle &= \frac{3}{2} n_i k_B T_i. \end{aligned} \quad (7.11)$$

By minimizing the  $H$ -function as defined in Eq.(7.7) under these constraints, the form of  $f_i^*$  is (Arcidiacono *et al.*, 2006)

$$f_i^* = n_i \left( \frac{m_i}{2\pi k_B T_i} \right)^{3/2} \exp \left( -\frac{m_i (c_i - u_i)^2}{2k_B T_i} \right). \quad (7.12)$$

Similarly, for the second case where the momentum diffuses faster, the set of constraints under which the  $H$ -function is to be minimized are

$$\begin{aligned} \langle m_i, f_i^* \rangle &= \rho_i, \\ \langle m_i \mathbf{c}_i, f_i^* \rangle &= \rho_i \mathbf{u}_i, \\ \sum_{i=A,B} \langle m_i \xi_{i\alpha} \xi_{i\beta}, f_i^* \rangle &= n \theta_{\alpha\beta}, \end{aligned} \quad (7.13)$$

where  $\xi_{i\alpha} = c_{i\alpha} - u_\alpha$  and

$$\theta_{\alpha\beta} = \frac{1}{n} \sum_{i=A,B} \langle m_i \xi_{i\alpha} \xi_{i\beta}, f_i \rangle. \quad (7.14)$$

The quasi-equilibrium distribution function assumes the form (Arcidiacono *et al.*, 2006)

$$f_i^* = n_i \left( \frac{m_i}{2\pi |\theta_{\alpha\beta}|} \right)^{\frac{3}{2}} \exp \left( \frac{-m_i \xi_{i\alpha} \theta_{\alpha\beta}^{-1} \xi_{i\beta}}{2} \right), \quad (7.15)$$

where  $|\theta_{\alpha\beta}|$  is the determinant. These two distinct forms of quasi-equilibrium can be used to build two different collision kernels based on the Fokker-Planck approximation, which can solve for binary mixtures.

## 7.4 Model I: Momentum and temperature difference as the slow variable

For the form of quasi-equilibrium presented in Eq.(7.12), the Fokker-Planck approximation to the binary collisional kernel is

$$\Omega_i^{\text{FP}(1)} = \frac{1}{\tau_1} \partial_{c_{i\alpha}} \left( (c_{i\alpha} - u_{i\alpha}) f_i + \frac{k_B T_i}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right) + \frac{1}{\tau_2} \partial_{c_{i\alpha}} \left( (u_{i\alpha} - u_\alpha) f_i + \frac{k_B \Delta T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right), \quad (7.16)$$

where  $\tau_1$  and  $\tau_2$  are the relaxation times and  $\Delta T = T - T_i$ . For the model to be acceptable, it must satisfy the collisional invariants. By intergrating over the velocity space  $\mathbf{c}_i$ , it can be verified that

$$\begin{aligned} \langle \Omega_i^{\text{FP}(1)}, m_i \rangle &= 0, \\ \sum_{i=A,B} \langle \Omega_i^{\text{FP}(1)}, \{m_i \mathbf{c}_i, m_i c_i^2/2\} \rangle &= \{\mathbf{0}, 0\}. \end{aligned} \quad (7.17)$$

Using this the evolution equations for component mass, mixture momentum and energy are

$$\begin{aligned} \partial_t \rho_i + \partial_\alpha \rho_i u_{i\alpha} &= 0, \\ \partial_t \rho u_\alpha + \partial_\beta (\rho u_\alpha u_\beta + p \delta_{\alpha\beta}) + \partial_\beta \sigma_{\alpha\beta} &= 0, \\ \partial_t E + \partial_\alpha ((E + p) u_\alpha + \sigma_{\alpha\gamma} u_\gamma) + \partial_\alpha q_\alpha &= 0, \end{aligned} \quad (7.18)$$

where  $E$  is the energy,  $p$  the pressure,  $\sigma_{\alpha\beta}$  the stress and  $q_\alpha$  the heat flux, with their expressions given as

$$\begin{aligned} E &= \sum_{i=A,B} \langle m_i c_i^2/2, f_i \rangle, \\ p &= nk_B T, \\ \sigma_{\alpha\beta} &= \sum_{i=A,B} \langle m_i \overline{\xi_{i\alpha} \xi_{i\beta}}, f_i \rangle, \\ q_\alpha &= \sum_{i=A,B} \langle m_i \xi_{i\alpha} \xi_i^2/2, f_i \rangle, \end{aligned} \quad (7.19)$$

where  $\overline{A_{\alpha\beta}}$  denotes the traceless part of tensor. Hence, the collision operator  $\Omega_i^{\text{FP}(1)}$  does indeed satisfy the conservation laws. The component momentum and energy equations are as expected in the relaxation form

$$\begin{aligned} \langle \Omega_i^{\text{FP}(1)}, m_i c_{i\alpha} \rangle &= \frac{1}{\tau_2} (\rho_i u_\alpha - \rho_i u_{i\alpha}), \\ \langle \Omega_i^{\text{FP}(1)}, m_i c_i^2/2 \rangle &= \frac{1}{\tau_2} ((\rho_i u_{i\alpha} u_\alpha + Dk_B n_i T) - (\rho_i u_{i\alpha}^2 + Dk_B n_i T_i)), \end{aligned} \quad (7.20)$$

this is the case because the component velocity and temperature equilibrate to the mixture velocity and temperature. This proves that proposed model satisfies collisional invariants. In addition to the conservation laws, the model must also satisfy the  $H$ -theorem. The evolution equation for the  $H$ -function is given by

$$\partial_t H + \partial_\alpha J_\alpha^H = -\sigma^S, \quad (7.21)$$

where  $J_\alpha^H$  is the flux of the entropy and  $\sigma^S$  the entropy production term. For the model to hold  $H$ -theorem, one must have

$$\sigma^S = - \sum_{i=A,B} \left\langle \Omega_i^{FP(1)}, \ln f_i \right\rangle \geq 0. \quad (7.22)$$

For the proposed model, the expression for  $\sigma^S$  is

$$\begin{aligned} &= - \sum_{i=A,B} \int \ln f_i \left[ \frac{1}{\tau_{\text{eff}}} \partial_{c_{i\alpha}} \left( (c_{i\alpha} - u_{i\alpha}) f_i + \frac{k_B T_i}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right) \right. \\ &+ \left. \frac{1}{\tau_2} \partial_{c_{i\alpha}} \left( (c_{i\alpha} - u_\alpha) f_i + \frac{k_B T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right) \right] d\mathbf{c}_i \\ &= \sum_{i=A,B} \int \left[ \frac{1}{\tau_{\text{eff}}} \left( (c_{i\alpha} - u_{i\alpha}) f_i \frac{\partial \ln f_i}{\partial c_{i\alpha}} + \frac{k_B T_i}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial \ln f_i}{\partial c_{i\alpha}} \right) \right. \\ &+ \left. \frac{1}{\tau_2} \left( (c_{i\alpha} - u_\alpha) \frac{\partial f_i}{\partial c_{i\alpha}} + \frac{k_B T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial \ln f_i}{\partial c_{i\alpha}} \right) \right] d\mathbf{c}_i \quad (7.23) \\ &= \sum_{i=A,B} \int \left[ \frac{1}{\tau_{\text{eff}}} \left( (c_{i\alpha} - u_{i\alpha}) \frac{\partial f_i}{\partial c_{i\alpha}} + \frac{k_B T_i}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} \right) \right. \\ &+ \left. \frac{1}{\tau_2} \left( (c_{i\alpha} - u_\alpha) \frac{\partial f_i}{\partial c_{i\alpha}} + \frac{k_B T}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} \right) \right] d\mathbf{c}_i \\ &= \frac{1}{\tau_{\text{eff}}} \sum_{i=A,B} -Dn_i + \int \frac{k_B T_i}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} d\mathbf{c}_i + \frac{1}{\tau_2} \sum_{i=A,B} -Dn_i + \int \frac{k_B T}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} d\mathbf{c}_i, \end{aligned}$$

where  $\tau_{\text{eff}} = \tau_2 \tau_1 / (\tau_2 - \tau_1)$ . It can be verified that, the following equalities hold

$$\begin{aligned} \sum_{i=A,B} -Dn_i + \int \frac{k_B T_i}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} d\mathbf{c}_i &= \sum_{i=A,B} \underbrace{\int \frac{k_B T_i}{m_i} f_i \left( \frac{\partial \ln(f_i/f_i^*)}{\partial c_{i\alpha}} \right)^2 d\mathbf{c}_i}_{\text{positive}} \\ \sum_{i=A,B} -Dn_i + \int \frac{k_B T}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} d\mathbf{c}_i &= \sum_{i=A,B} \underbrace{\int \frac{k_B T}{m_i} f_i \left( \frac{\partial \ln(f_i/f_i^{\text{MB}_i})}{\partial c_{i\alpha}} \right)^2 d\mathbf{c}_i}_{\text{positive}} \end{aligned} \quad (7.24)$$

where expression of  $f_i^*$  is given Eq.(7.12). Then Eq.(7.23) suggests that

$$\sigma^S \geq 0, \quad \forall \tau_1 \leq \tau_2. \quad (7.25)$$

Therefore, proposed model satisfies the  $H$ -theorem for  $\tau_1 \leq \tau_2$ .

An important condition for  $\Omega^{FP(1)}$  to be considered valid is that at equilibrium it must follow the Maxwell-Boltzmann distribution. At equilibrium, we have

$$\begin{aligned} \left\langle \Omega_i^{FP(1)}, m_i c_{i\alpha} \right\rangle &= 0, \\ \left\langle \Omega_i^{FP(1)}, m_i c_i^2 / 2 \right\rangle &= 0, \end{aligned} \quad (7.26)$$

then as per Eq.(7.20) equilibrium  $u_{i\alpha} = u_\alpha$  and  $T_i = T$ , hence  $\Omega_i^{FP(1)} = 0$  reduces to

$$\partial_{c_{i\alpha}} \left( (c_{i\alpha} - u_\alpha) f_i + \frac{k_B T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right) = 0. \quad (7.27)$$

Integrating Eq.(7.27) with respect to the velocity space and using the fact that the distribution function and its derivatives tend to zero at infinity. We have

$$(c_{i\alpha} - u_\alpha) f_i + \frac{k_B T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} = 0. \quad (7.28)$$

Solving Eq.(7.28), we get the Maxwell-Boltzmann distribution as the solution. Additionally for the model to be consistent with the indifferentiability principle, one must be able to recover the Fokker-Planck approximation for single component case. In the case when  $\tau_1 = \tau_2 = \tau$  and  $m_A = m_B = m$ , the Fokker-Planck collision kernel for binary mixtures reduces to

$$\Omega^{\text{FP}} = \frac{1}{\tau} \partial_{c_\alpha} \left( (c_\alpha - u_\alpha) f + \frac{k_B T}{m} \frac{\partial f}{\partial c_\alpha} \right), \quad (7.29)$$

indicating that proposed model abides by indifferentiability principle.

As demonstrated, the proposed model does indeed satisfy the conservation laws,  $H$ -theorem, zero of collision and indifferentiability principle. Thus, this model is an acceptable approximation to the Boltzmann equation for binary mixtures.

## 7.5 Model II: Pressure as the slow variable

In this model, we consider the situation where the particles assume their equilibrium velocity fast, and then there is a slow transition wherein the pressure tensor relaxes to the equilibrium temperature slowly. The quasi-equilibrium distribution function for this case is mentioned in Eq.(7.15) and the corresponding collision operator is

$$\Omega_i^{\text{FP}(2)} = \frac{1}{\tau_1} \partial_{c_{i\alpha}} \left( (c_{i\alpha} - u_\alpha) f_i + \frac{\theta_{\alpha\beta}}{m_i} \frac{\partial f_i}{\partial c_{i\beta}} \right) + \frac{1}{\tau_2} \partial_{c_{i\alpha}} \left( \left( \frac{k_B T \delta_{\alpha\beta}}{m_i} - \frac{\theta_{\alpha\beta}}{m_i} \right) \frac{\partial f_i}{\partial c_{i\beta}} \right) \quad (7.30)$$

We proceed in a manner similar to the previous section. By simple integration over velocity space, it can be seen that this model satisfies the collisional invariants as

$$\begin{aligned} \langle \Omega_i^{\text{FP}(2)}, m_i \rangle &= 0, \\ \sum_{i=A,B} \langle \Omega_i^{\text{FP}(2)}, \{m_i \mathbf{c}_i, m_i c_i^2 / 2\} \rangle &= \{\mathbf{0}, 0\}, \end{aligned} \quad (7.31)$$

this results in a set of conservation laws same as the one referred in Eq.(7.18). However, one difference that is observed between the two models is relaxation of the component momentum and energy. For this case, the relaxation of component momentum and energy is

$$\begin{aligned} \langle \Omega_i^{\text{FP}(2)}, m_i c_{i\alpha} \rangle &= \frac{1}{\tau_1} (\rho_i u_\alpha - \rho_i u_{i\alpha}), \\ \langle \Omega_i^{\text{FP}(2)}, m_i c_i^2 / 2 \rangle &= \frac{1}{\tau_1} ((\rho_i u_{i\alpha} u_\alpha + D k_B n_i T) - (\rho_i u_{i\alpha}^2 + D k_B n_i T_i)). \end{aligned} \quad (7.32)$$

The component momentum relaxes with the time scale  $\tau_1$  as opposed to the first model where it relaxes with the time scale  $\tau_2$ . This is the case because this model is for high values of  $\text{Sc}$  where the momentum diffuses faster.

To test for the validity of  $H$ -theorem, we proceed in a manner similar to the previous section and find the expression for the entropy generation term

$$\sigma^S = \frac{1}{\tau_{\text{eff}}} \sum_{i=A,B} -Dn_i + \int \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\theta_{\alpha\beta}}{m_i} \frac{\partial f_i}{\partial c_{i\beta}} + \frac{1}{\tau_2} \sum_{i=A,B} -Dn_i + \int \frac{k_B T}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}}, \quad (7.33)$$

for this case, the following identities hold

$$\begin{aligned} \sum_{i=A,B} -n_i D + \int \frac{\theta_{\alpha\beta}}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\beta}} d\mathbf{c}_i &= \sum_{i=A,B} \underbrace{\int f_i \frac{\partial \ln(f_i/f_i^*)}{\partial c_{i\alpha}} \frac{\theta_{\alpha\beta}}{m_i} \frac{\partial \ln(f_i/f_i^*)}{\partial c_{i\beta}} d\mathbf{c}_i}_{\text{positive}}, \\ \sum_{i=A,B} -n_i D + \int \frac{k_B T}{m_i} \frac{1}{f_i} \frac{\partial f_i}{\partial c_{i\alpha}} \frac{\partial f_i}{\partial c_{i\alpha}} d\mathbf{c}_i &= \sum_{i=A,B} \underbrace{\int \frac{k_B T}{m_i} f_i \left( \frac{\partial \ln(f_i/f_i^*)}{\partial c_{i\alpha}} \right)^2 d\mathbf{c}_i}_{\text{positive}}, \end{aligned} \quad (7.34)$$

where the expression for  $f_i^*$  is given by Eq.(7.15). The integrand of the second term is same as the first model while the first term is of the form  $\mathbf{x}^T \mathbf{A} \mathbf{x}$  where  $\mathbf{A}$  is a positive definite matrix and  $\mathbf{x}$  is any arbitrary vector. Since,  $\theta_{\alpha\beta}$  is a positive definite matrix, we have

$$\sigma^S \geq 0, \quad \forall \tau_1 \leq \tau_2. \quad (7.35)$$

This confirms that this model for binary mixture also holds  $H$ - theorem.

Equilibrium distribution: At equilibrium  $\theta_{\alpha\beta} = nk_B T \delta_{\alpha\beta}$ , hence zero of the collision for this model reduces to

$$\partial_\alpha \left( (c_{i\alpha} - u_\alpha) f_i + \frac{k_B T}{m_i} \frac{\partial f_i}{\partial c_{i\alpha}} \right) = 0. \quad (7.36)$$

As was shown before, the solution to this equation is the Maxwell-Boltzmann distribution.

Indifferentiability principle: By substituting  $m_A = m_B$  and  $\tau_1 = \tau_2$  in Eq.(7.30) one recovers the Fokker - Planck model for single component. Hence, the indifferentiability principle is valid for this model.

## 7.6 Transport Coefficients

In order to obtain the transport coefficients, we perform the Chapman-Enskog expansion, wherein the time derivative, distribution function and other relevant variables are represented as a series with Kn acting as the smallness parameter (Chapman & Cowling, 1970). The time derivative and distribution function expressed in series form as (Arcidiacono *et al.*, 2006)

$$\begin{aligned} \partial_t &= \partial_t^{(0)} + \text{Kn} \partial_t^{(1)} + \text{Kn}^2 \partial_t^{(2)} + \dots, \\ f_i &= f_i^{\text{MB}} + \text{Kn} f_i^{(1)} + \text{Kn}^2 f_i^{(2)} + \dots, \end{aligned} \quad (7.37)$$

with the following constraints imposed on  $f_i$



$$\begin{aligned}
\langle m_i, f_i^{(n)} \rangle &= 0, \\
\sum_{i=A,B} \langle m_i c_{i\alpha}, f_i^{(n)} \rangle &= 0, \\
\sum_{i=A,B} \langle m_i c^2, f_i^{(n)} \rangle &= 0 \quad \forall \quad n \geq 1.
\end{aligned} \tag{7.38}$$

These constraints ensure that component density, mixture momentum and energy are slow moments. As stipulated,  $f_i$  can only be a function of the conserved variables, hence time derivative of  $f_i$  is (Liboff, 2003)

$$\frac{\partial f_i(\rho_i, \mathbf{u}, T)}{\partial t} = \frac{\partial f_i}{\partial \rho_i} \cdot \frac{\partial \rho_i}{\partial t} + \frac{\partial f_i}{\partial \mathbf{u}} \cdot \frac{\partial \mathbf{u}}{\partial t} + \frac{\partial f_i}{\partial T} \cdot \frac{\partial T}{\partial t}, \tag{7.39}$$

where the expression for time derivatives of the conserved variable can be calculated from the conservation laws. The higher order moments in series form are

$$\begin{aligned}
\sigma_{\alpha\beta} &= \text{Kn} \sigma_{\alpha\beta}^{(1)} + \text{Kn}^2 \sigma_{\alpha\beta}^{(2)} + \dots, \\
q_\alpha &= \text{Kn} q_\alpha^{(1)} + \text{Kn}^2 q_\alpha^{(2)} + \dots,
\end{aligned} \tag{7.40}$$

as the stress and heat flux are zero at equilibrium.

### Viscosity

The stress evolution equation for the first model is

$$\begin{aligned}
&\partial_t \sigma_{\alpha\beta} + \partial_\gamma (\sigma_{\alpha\beta} u_\gamma) + 2p \overline{\partial_\alpha u_\beta} + 2\overline{\sigma_{\alpha\gamma} \partial_\gamma u_\beta} + \partial_\gamma Q_{\alpha\beta\gamma} + \frac{4}{D+2} \overline{\partial_\alpha q_\beta} = \\
&-\frac{2}{\tau_1} \left( \sigma_{\alpha\beta} + \rho \overline{u_\alpha u_\beta} - \sum_{i=A,B} \rho_i \overline{u_{i\alpha} u_{i\beta}} \right) - \frac{2}{\tau_2} \left( \sum_{i=A,B} \rho_i \overline{u_{i\alpha} u_{i\beta}} - \rho \overline{u_\alpha u_\beta} \right),
\end{aligned} \tag{7.41}$$

where  $Q_{\alpha\beta\gamma} = \sum_{i=A,B} \langle m_i \overline{\xi_{i\alpha} \xi_{i\beta} \xi_{i\gamma}} \rangle$ . Retaining terms upto  $\mathcal{O}(\text{Kn})$ , the stress evolution equation yields

$$2p \overline{\partial_\alpha u_\beta} = -\frac{2\sigma_{\alpha\beta}^{(1)}}{\tau_1}. \tag{7.42}$$

and comparing with the Navier-Stokes law for stress tensor, we have

$$\mu = \frac{p\tau_1}{2} \tag{7.43}$$

Similarly, for the second model the right hand side of the stress evolution equation is

$$\sum_{i=A,B} \langle m_i \overline{\xi_{i\alpha} \xi_{i\beta}}, \Omega_i^{\text{FP}(2)} \rangle = -\frac{2}{\tau_2} \sigma_{\alpha\beta}. \tag{7.44}$$

Hence, the expression for viscosity for this model is

$$\mu = \frac{p\tau_2}{2}. \tag{7.45}$$

These expressions can be further used to determine the Knudsen number of the system.

### Diffusion coefficient

The diffusion coefficient can be calculated by comparison with the Stefan-Maxwell diffusion equation (Bergman *et al.*, 2011)

$$\partial_\alpha X_A = \frac{X_A X_B}{D_{AB}} \frac{V_\alpha}{m_{AB}} + (Y_A - X_A) \frac{\partial_\alpha p}{p}, \quad (7.46)$$

where  $X_i = n_i/n$  is the component mole fraction,  $Y_i = \rho_i/\rho$  is the component mass fraction and  $V_\alpha$  is the diffusive flux defined as

$$V_\alpha = m_{AB}(u_{A\alpha} - u_{B\alpha}), \quad (7.47)$$

where  $m_{AB} = (\rho_A \rho_B)/\rho$ . Diffusive flux essentially quantifies the difference between the momentum of a given component and the momentum of the mixture. The series expansion for this quantity is

$$V_\alpha = \text{Kn} V_\alpha^{(1)} + \text{Kn}^2 V_\alpha^{(2)} + \dots \quad (7.48)$$

Similar to stress and heat flux, at equilibrium the diffusive flux attains zero values as momenta of both components relax to the mixture momentum. In order to calculate the expression for  $V_\alpha$  we write the expression for individual component velocities. For the first model, we have

$$\begin{aligned} \partial_t \rho_A u_{A\alpha} + \partial_\alpha P_{A\alpha\beta} &= \frac{1}{\tau_2} (\rho u_\alpha - \rho_A u_{A\alpha}), \\ \partial_t \rho_B u_{B\alpha} + \partial_\alpha P_{B\alpha\beta} &= \frac{1}{\tau_2} (\rho u_\alpha - \rho_B u_{B\alpha}), \end{aligned} \quad (7.49)$$

where  $P_{i\alpha\beta} = \langle m_i c_{i\beta} c_{i\alpha} \rangle$  and at equilibrium attains the value  $P_{i\alpha\beta} = p_i \delta_{\alpha\beta} + \rho_i u_\alpha u_\beta$ . Subtracting one equation from another, we have

$$\begin{aligned} \frac{(u_{A\alpha} - u_{B\alpha})}{\tau_2} &= \left( \frac{\partial_t \rho_B u_{B\alpha}}{\rho_B} - \frac{\partial_t \rho_A u_{A\alpha}}{\rho_A} \right) + \left( \frac{\partial_\alpha P_{B\alpha\beta}}{\rho_B} - \frac{\partial_\alpha P_{A\alpha\beta}}{\rho_A} \right), \\ \frac{(u_{A\alpha} - u_{B\alpha})}{\tau_2} &= \left( \frac{u_{B\alpha} \partial_t \rho_B}{\rho_B} - \frac{u_{A\alpha} \partial_t \rho_A}{\rho_A} \right) + (\partial_t u_{B\alpha} - \partial_t u_{A\alpha}) + \left( \frac{\partial_\alpha P_{B\alpha\beta}}{\rho_B} - \frac{\partial_\alpha P_{A\alpha\beta}}{\rho_A} \right) \end{aligned} \quad (7.50)$$

The temporal derivatives of the component density can be eliminated using the continuity equation

$$\partial_t \rho_i = -\partial_\alpha (\rho_i u_{i\alpha}). \quad (7.51)$$

Considering quantities upto  $\mathcal{O}(\text{Kn})$ , and after some rearrangement we have

$$\begin{aligned} V_\alpha^{(1)} &= \tau_2 \frac{\rho_A \rho_B}{\rho} \left[ u_\alpha \left( \frac{\partial_\beta \rho_A u_{\beta\alpha}}{\rho_A} - \frac{\partial_\beta \rho_B u_{\beta\alpha}}{\rho_B} \right) + \left( \frac{\partial_\beta n_B k_B T}{\rho_B} - \frac{\partial_\beta n_A k_B T}{\rho_A} \right) \right. \\ &\quad \left. + \left( \frac{\partial_\beta \rho_B u_\alpha u_\beta}{\rho_B} - \frac{\partial_\beta \rho_A u_\alpha u_\beta}{\rho_A} \right) \right], \end{aligned} \quad (7.52)$$

on simplification, we obtain

$$V_\alpha^{(1)} = \tau_2 \left( \frac{\rho_A}{\rho} \partial_\alpha p_B^0 - \frac{\rho_B}{\rho} \partial_\alpha p_A^0 \right), \quad (7.53)$$

which upon simplifications leads to

$$\begin{aligned}
V_\alpha^{(1)} &= \tau_2 \left[ \frac{\rho_A}{\rho} \partial_\alpha p - \left( p \partial_\alpha \left( \frac{n_A}{n} \right) + \frac{n_A}{n} \partial_\alpha p \right) \right], \\
V_\alpha^{(1)} &= \tau_2 [Y_A \partial_\alpha p - p \partial_\alpha X_A - X_A \partial_\alpha p].
\end{aligned} \tag{7.54}$$

This provides an expression for the gradient of the molar fraction which is

$$\partial_\alpha X_A = -\frac{V_\alpha^{(1)}}{\tau_2 p} + (Y_A - X_A) \frac{\partial_\alpha p}{p}, \tag{7.55}$$

Now, comparing this with Stefan - Maxwell equation (Eq.(7.46)) we get following expression for the diffusion coefficient

$$D_{AB} = X_A X_B \frac{p}{m_{AB}} \tau_2. \tag{7.56}$$

The Schmidt number can now be computed as

$$\text{Sc} = \frac{\mu}{\rho D_{AB}} = \frac{\tau_1}{2\tau_2} \frac{m_{AB}}{X_A X_B} \frac{1}{\rho} = \frac{\tau_1}{2\tau_2} \frac{Y_A Y_B}{X_A X_B}. \tag{7.57}$$

Existence of  $H$ - theorem for this model suggests that  $\tau_1 \leq \tau_2$ , hence

$$\text{Sc} \leq \frac{Y_A Y_B}{2X_A X_B}. \tag{7.58}$$

The model has an upper limit on Schmidt number and this is in accordance with the characteristics of the quasi-equilibrium distribution. Similarly, for the second model, the Schmidt number is calculated as

$$\text{Sc} = \frac{\mu}{\rho D_{AB}} = \frac{\tau_2}{2\tau_1} \frac{m_{AB}}{X_A X_B} \frac{1}{\rho} = \frac{\tau_2}{2\tau_1} \frac{Y_A Y_B}{X_A X_B}. \tag{7.59}$$

and since the limitation  $\tau_1 \leq \tau_2$  exists, as consistent with the hypothesis there is a lower bound on the Schmidt number, which is

$$\text{Sc} \geq \frac{Y_A Y_B}{2X_A X_B}. \tag{7.60}$$

Hence, both models in conjunction can cover a large range of Schmidt numbers.

## 7.7 Numerical scheme

Analytically solving the Fokker-Planck equation for many cases is not possible, hence numerical techniques must be employed in order to reach a solution. A Fokker-Planck equation which describes the evolution of probability density function of the random variable  $\eta$ , of the form

$$\frac{dp(\eta, t)}{dt} = -\Lambda_\alpha^{(1)}(\eta, t) \frac{\partial p(\eta, t)}{\partial \eta_\alpha} + \frac{\zeta_{\alpha\beta}^{(1)}(\eta, t)}{2} \frac{\partial^2 p(\eta, t)}{\partial \eta_\alpha \partial \eta_\beta} - \Lambda_\alpha^{(2)}(\eta, t) \frac{\partial p(\eta, t)}{\partial \eta_\alpha} + \frac{\zeta_{\alpha\beta}^{(2)}(\eta, t)}{2} \frac{\partial^2 p(\eta, t)}{\partial \eta_\alpha \partial \eta_\beta}, \tag{7.61}$$

where  $\Lambda^{(i)}$  are the drift terms and  $\zeta^{(i)}$  are the diffusion coefficients. This form of Fokker-Planck equation is equivalent to the Langevin equation (Risken, 1996)

$$\dot{\eta}_\alpha = h_\alpha^{(1)}(\eta, t) + g_{\alpha\beta}^{(1)}(\eta, t) \Gamma_\beta(t) + h_\alpha^{(2)}(\eta, t) + g_{\alpha\beta}^{(2)}(\eta, t) \Gamma'_\beta(t), \tag{7.62}$$

where  $\mathbf{h}^{(i)}$  are the drift terms,  $\mathbf{g}^{(i)}$  the diffusion coefficients and  $\mathbf{\Gamma}, \mathbf{\Gamma}'$  are Gaussian distributed random numbers which hold the following properties

$$\langle \Gamma_\alpha(t) \rangle = 0, \quad \langle \Gamma_\alpha(t) \Gamma_\beta(t') \rangle = \delta(t - t') \delta_{\alpha\beta}. \quad (7.63)$$

Under certain conditions, which are satisfied by the proposed models the following relations hold (Risken, 1996)

$$\begin{aligned} \Lambda_\alpha^{(1)} &= h_\alpha^{(1)}(\xi, t), & \zeta_{\alpha\beta}^{(1)} &= g_{\alpha\gamma}^{(1)} g_{\gamma\beta}^{(1)} \\ \Lambda_\alpha^{(2)} &= h_\alpha^{(2)}(\xi, t), & \zeta_{\alpha\beta}^{(2)} &= g_{\alpha\gamma}^{(2)} g_{\gamma\beta}^{(2)} \end{aligned} \quad (7.64)$$

The central idea is that solution to Fokker-Planck equation is approximated by considering an ensemble of trajectories generated by the Langevin dynamics. In this case, a large number of particles have their position and velocities updated using Eq.(7.62). We now discuss the

### 7.7.1 Model I

For the first model the equivalent Langevin equations are

$$\begin{aligned} \frac{dx_\alpha}{dt} &= c_{i\alpha} \\ \frac{dc_{i\alpha}}{dt} &= - \left( \frac{1}{\tau_{\text{eff}}} \right) (c_{i\alpha} - u_{i\alpha}) - \frac{1}{\tau_2} (c_{i\alpha} - u_\alpha) + \sqrt{\frac{2k_B T_i}{m_i}} dW_\alpha + \sqrt{\frac{2k_B T}{m_i}} dW'_\alpha, \end{aligned} \quad (7.65)$$

where  $dW_\alpha$  and  $dW'_\alpha$  denote random forces with following statistics

$$\langle dW_\alpha \rangle = 0, \quad \langle dW'_\alpha \rangle = 0, \quad \langle dW_\alpha dW'_\alpha \rangle = 0. \quad (7.66)$$

More specifically,  $dW = W(t + \Delta t) - W(t)$  is the standard Weiner process, where  $W(t)$  is a rapidly changing random force with mean and variance as (Gardiner, 1985b)

$$\langle dW_\alpha(t) \rangle = 0, \quad \langle dW_\alpha dW_\beta \rangle = dt \delta_{\alpha\beta}. \quad (7.67)$$

Thus, the detailed binary collision description is approximated by a random collision with a heat bath in the model.

These Langevin equations can be solved efficiently using the the stochastic version of the Verlet algorithm. For the present model the discretization scheme is (Singh & Ansumali, 2015a)

$$\begin{aligned} x_\alpha^{(1)} &= x_\alpha(t) + \frac{1}{2} c_{i\alpha}(t) \Delta t, \\ c_{i\alpha}(t + \Delta t) &= c_{i\alpha}(t) - \left( \frac{\vartheta_1}{1 + \vartheta_1/2} \right) (c_{i\alpha}(t) - u_{i\alpha}) - \left( \frac{\vartheta_2}{1 + \vartheta_2/2} \right) (c_{i\alpha}(t) - u_\alpha) \\ &\quad + \frac{\sqrt{2\mathcal{D}_i^{(1)}\vartheta_1}}{1 + \vartheta_1/2} \phi_\alpha + \frac{\sqrt{2\mathcal{D}_i^{(2)}\vartheta_2}}{1 + \vartheta_2/2} \phi'_\alpha, \\ x_\alpha(t + \Delta t) &= x_\alpha^{(1)} + \frac{1}{2} c_{i\alpha}(t + \Delta t) \Delta t, \end{aligned} \quad (7.68)$$

where  $\vartheta_1 = \Delta t / \tau_{\text{eff}}$ ,  $\vartheta_2 = \Delta t / \tau_2$  and  $\phi_\alpha, \phi'_\alpha$  are Gaussian random numbers with mean zero and variance one and  $\mathcal{D}_i^{(1)}$  and  $\mathcal{D}_i^{(2)}$  are  $k_B T_i / m_i$  and  $k_B T / m_i$  respectively. This scheme works efficiently for small time step such that  $\max\{\vartheta_1, \vartheta_2\} \leq 0.001$ .

### 7.7.2 Model II

The formulation for this model remains largely unchanged and the equivalent Langevin equations are

$$\begin{aligned}\frac{dx_\alpha}{dt} &= c_{i\alpha}, \\ \frac{dc_{i\alpha}}{dt} &= -\left(\frac{1}{\tau_{\text{eff}}}\right)(c_{i\alpha} - u_\alpha) - \frac{1}{\tau_2}(c_{i\alpha} - u_\alpha) + \sqrt{2}\theta'_{i\alpha\beta}dW_\beta + \sqrt{\frac{2k_B T}{m_i}}dW'_\alpha,\end{aligned}\tag{7.69}$$

where  $\theta'_{i\alpha\gamma}\theta'_{i\gamma\beta} = \theta_{i\alpha\beta}/m_i$ . The discretization scheme for this model is

$$\begin{aligned}x_\alpha^{(1)} &= x_\alpha(t) + \frac{1}{2}c_{i\alpha}(t)\Delta t, \\ c_{i\alpha}(t + \Delta t) &= c_{i\alpha}(t) - \left(\frac{\vartheta_1}{1 + \vartheta_1/2}\right)(c_{i\alpha}(t) - u_\alpha) - \left(\frac{\vartheta_2}{1 + \vartheta_2/2}\right)(c_{i\alpha}(t) - u_\alpha) \\ &\quad + \frac{\sqrt{2\vartheta_1}\theta'_{i\alpha\beta}}{1 + \vartheta_1/2}\phi_\beta + \frac{\sqrt{2\mathcal{D}_i^{(2)}\vartheta_2}}{1 + \vartheta_2/2}\phi'_\alpha, \\ x_\alpha(t + \Delta t) &= x_\alpha^{(1)} + \frac{1}{2}c_{i\alpha}(t + \Delta t)\Delta t.\end{aligned}\tag{7.70}$$

The expression for  $\theta'_{\alpha\beta}$  can be obtained by using Cholesky decomposition of  $\theta_{\alpha\beta}/m_i$ .

In order to validate the numerical scheme, we started with a mixture with  $m_B/m_A = 2$  with  $N = 10^5$  particles in a single periodic box. For Model I, the velocities of the lighter particles were initialized uniformly in the range  $[0, 1)$  and the heavier particles in the range  $[0, 2)$ . For Model II, the velocities of lighter particle were initialized with a Gaussian distribution with mean 4 and variance 10, and the heavier particles were Gaussian distributed with mean 1 and variance 1. The plots of energy of the two components and the mixture with time over averaged an ensemble of 15 trajectories and the distribution of velocities in the x-direction, for both cases are shown in Fig.(7.2). A detailed description of the algorithm is provided in Table(7.1).

Table 7.1: Summary of the algorithm for binary mixtures

- 
1. The computational domain is divided into cells with its length parameter of the order of the mean free path.
  2. Each cell is then populated with particles in accordance with the initial conditions and the local equilibrium.
  3. The relevant variables – component density, velocity and temperature as well as the mixture density, velocity and temperature is calculated for each cell.
  4. The positions and velocities of all particles are updated using either Eq.(7.68) or Eq.(7.70).
  5. Once the updates are implemented, the particles are sorted into the cells based on their positions.
  6. Steps 3-5 are repeated until desired simulation time is achieved.
- 

## 7.8 Simulation results

In this section, we explain and present the results for three benchmark problems – Graham's law for effusion, Couette flow and binary diffusion. The first model was tested for these problems.

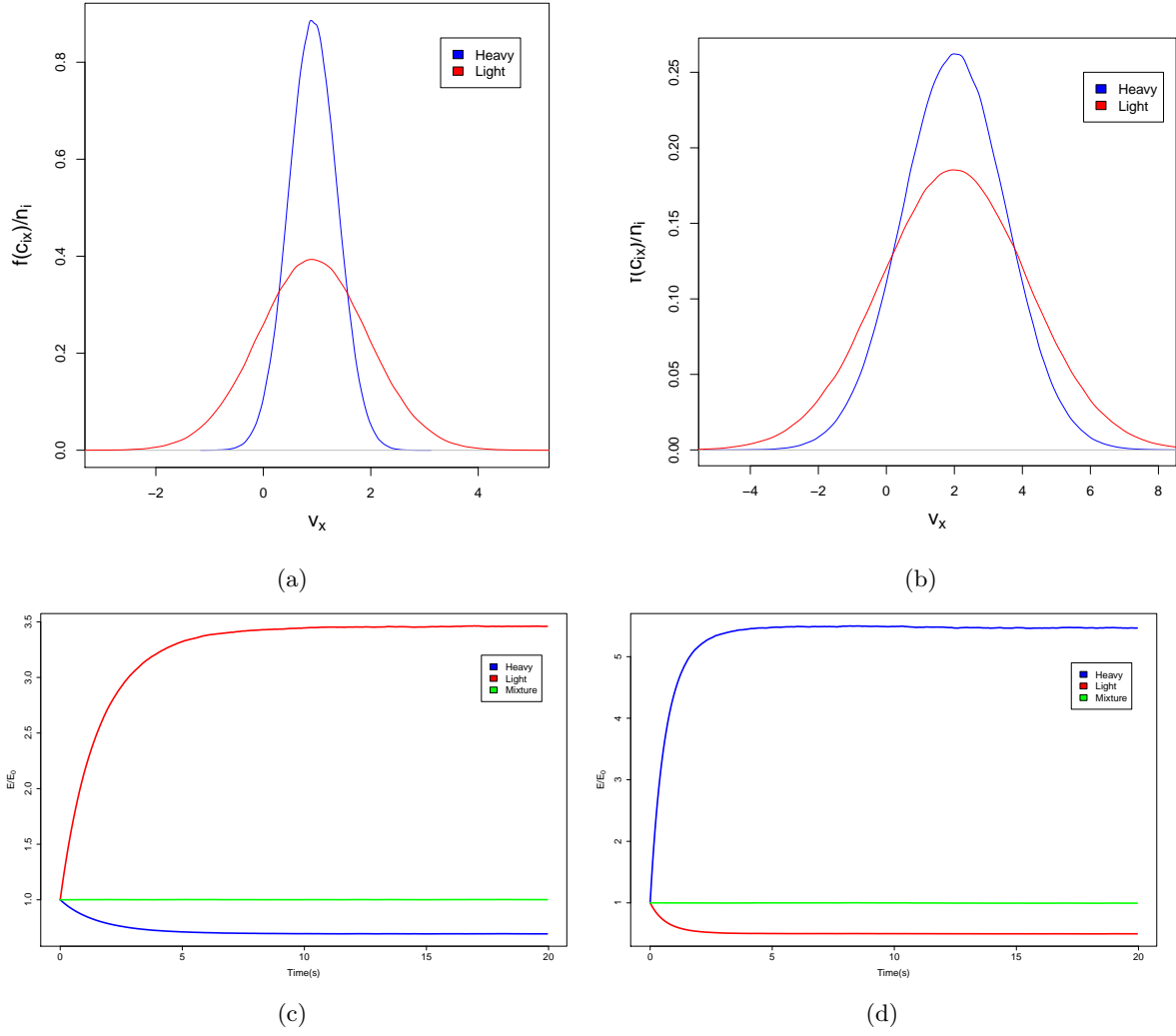


Figure 7.2: Plot of the distribution of velocities of the light and heavy component at equilibrium for a) Model I, and b) Model II. Plot of ratio of energy at time  $t$  to the initial energy ( $E(t)/E_0$ ) vs. time for individual components and the mixture for c) Model I, and d) Model II.

### 7.8.1 Graham's law for effusion

Effusion is a process wherein gas molecules escape through a small hole. The length parameter of this hole is much smaller than the mean free path of the gas, i.e.  $d \ll \lambda_{\text{mfp}}$ . A sketch of the process has been shown in Fig.(7.3). The number flux of the gas through this small hole is

$$\Phi_i = \langle c_{iz}, f(\mathbf{c}_i) \rangle, \quad (7.71)$$

where  $\Phi_i$  is the number flux and  $c_{iz}$  the molecular velocity in the direction perpendicular to the plane of the hole. By intergrating over velocity space, facilitated by a shift to the spherical co-ordinate system, the expression of  $\Phi_i$  is

$$\Phi_i = \frac{P}{\sqrt{2\pi m_i k_B T}}, \quad (7.72)$$

where  $P$  is the pressure and  $T$  the temperature of the gas. Then, for a well-mixed binary mixture the ratio of the fluxes is (Mason & Kronstadt, 1967)

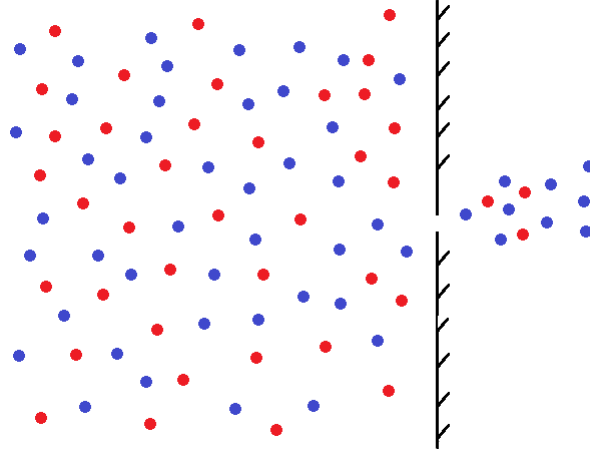


Figure 7.3: The process of gas molecules escaping through a small hole is known as effusion. The lighter particles (in this case blue) escape through the hole faster than the heavier particles, with a factor proportional to the square root of their mass ratios.

$$\frac{\Phi_A}{\Phi_B} = \sqrt{\frac{m_B}{m_A}}. \quad (7.73)$$

We simulated this system for three mass ratios  $m_B/m_A = 4$ ,  $m_B/m_A = 16$  and  $m_B/m_A = 100$ . The boundary conditions in the transverse directions were taken to be periodic while maintaining constant pressure in the system. The results have been plotted in Fig.(7.4). As can be seen, the simulations are in excellent agreement with the analytical solution.

## 7.8.2 Couette Flow

Couette flow is an important problem in fluid mechanics. The setup of the problem is simple, fluid between two plates is sheared in opposite directions with equal magnitudes, a sketch of the problem is shown in Fig.(7.5). In order to validate the model, we calculate the global stress tensor defined as (Sharipov *et al.*, 2004)

$$\Pi = -\frac{v_0}{2UP_0}P_{xy}. \quad (7.74)$$

This quantity is calculated in the entire range of rarefaction parameter, which is essentially the inverse of the Knudsen number and is defined as

$$\delta = \frac{HP_0}{\mu v_0}, \quad (7.75)$$

where  $\mu$  is the mixture viscosity and  $v_0$  the characteristic molecular velocity of the mixture defined as

$$v_0 = \sqrt{\frac{2k_B T_0}{m_0}}, \quad (7.76)$$

where  $m_0 = C_0 m_A + (1 - C_0) m_B$ , with  $C_0$  being the concentration of the lighter component. We simulated the system from two mixtures Neon-Argon (Ne-Ar) and Helium-Argon (He-Ar), for rarefaction parameters ranging from  $[0.01, 40]$  for three different concentrations - (0.1, 0.5, 0.9). The results for (Ne-Ar) are tabulated in Table(7.2) and results for (He-Ar) are tabulated in Table(7.3). Both sets of results were found to be in good agreement with the results obtained via the Discrete Velocity Method (DVM) (Sharipov *et al.*, 2004). This proves that proposed

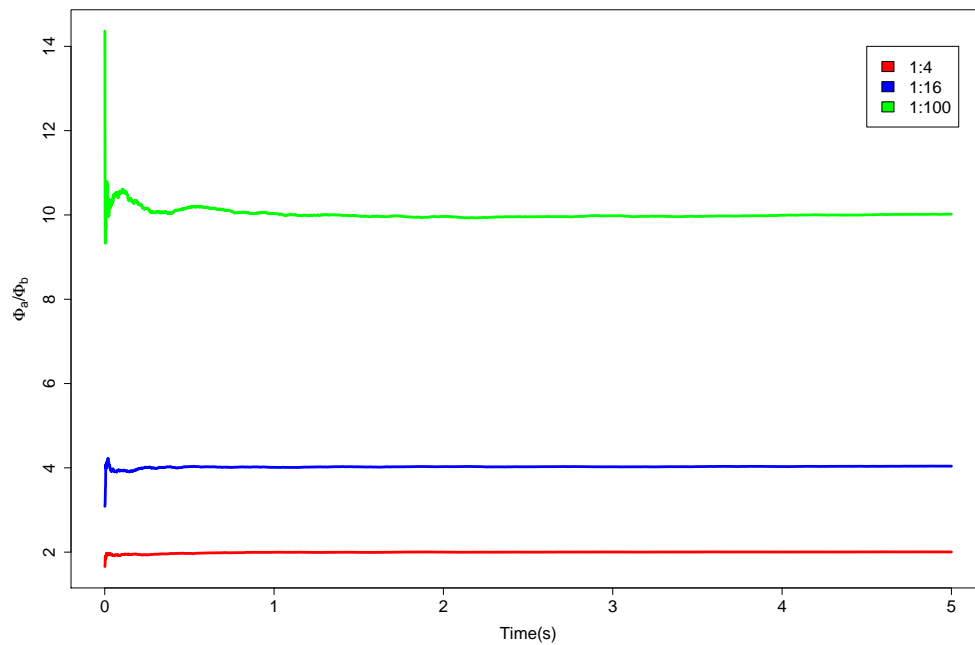


Figure 7.4: Model I was used to simulate a setup that could mimic Graham's law for effusion. Plot shows that results observed are in great agreement with expected behaviour, for all three cases

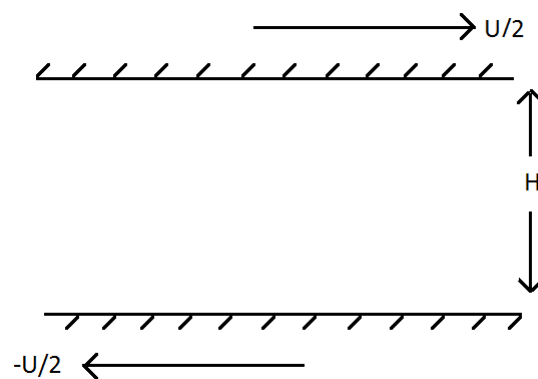


Figure 7.5: A representative sketch of the Couette flow setup. Two walls with a separation  $H$  are sheared in the opposite directions with velocity  $U/2$ .



method is indeed capable of simulating flows in a wide range of Knudsen numbers.

$\delta$	II values for Ne-Ar mixture					
	Fokker-Planck			DVM		
	$C_0 = 0.1$	0.5	0.9	0.1	0.5	0.9
0.01	0.2671	0.2731	0.2751	0.2786	0.2757	0.2778
0.1	0.2533	0.2504	0.2522	0.2601	0.2576	0.2594
1.0	0.1655	0.1636	0.1651	0.1689	0.1657	0.1685
10.0	0.0413	0.0416	0.0415	0.0415	0.0414	0.0415
40.0	0.0119	0.0119	0.0119	0.0119	0.0119	0.0119

Table 7.2: Comparison of the values of  $\Pi$  between the Fokker-Planck and DVM methods, for Ne-Ar mixture at three different concentrations  $C_0 = (0.1, 0.5, 0.9)$  for a range of rarefactions

$\delta$	II values for He-Ar mixture					
	Fokker-Planck			DVM		
	$C_0 = 0.1$	0.5	0.9	0.1	0.5	0.9
0.01	0.2704	0.2467	0.2443	0.2732	0.2484	0.2471
0.1	0.2479	0.2240	0.2232	0.2555	0.2335	0.2324
1.0	0.1617	0.1447	0.1472	0.1668	0.1566	0.1562
10.0	0.0418	0.0387	0.0398	0.0414	0.0407	0.0406
40.0	0.0121	0.0120	0.0119	0.0119	0.0118	0.0118

Table 7.3: Comparison of the values of  $\Pi$  between the Fokker-Planck and DVM methods, for He-Ar mixture at three different concentrations  $C_0 = (0.1, 0.5, 0.9)$  for a range of rarefactions

### 7.8.3 Binary diffusion

The profile of the mixture in this setup is determined by the step function

$$\begin{aligned} X_A &= 90\%, & X_B &= 10\% & \text{if } x < 0 \\ X_B &= 10\%, & X_A &= 90\% & \text{if } x \geq 0 \end{aligned} \quad (7.77)$$

where the mass ratio of the components was chosen to be  $m_B/m_A = 5$ . The step function is used instead of a smooth profile as it is a more severe check for the numerical scheme. Under the assumption that at infinity, the initial concentrations remains unchanged, this problem yields the analytical solution (Bergman *et al.*, 2011)

$$X_i = \left[ \frac{1}{2} + \frac{\Delta X_i}{2} \operatorname{erf} \left( \frac{x}{\sqrt{4D_{AB}t}} \right) \right] \quad (7.78)$$

where  $D_{AB}$  is the diffusion constant. The simulation was done for 20,000 time steps and the plots for both the components compared against their respective analytical solutions are plotted in Fig.(7.6). The simulation results were very close to the analytical solution. This exercise proves that the value of  $D_{AB}$  set by the numerical scheme is accurate.

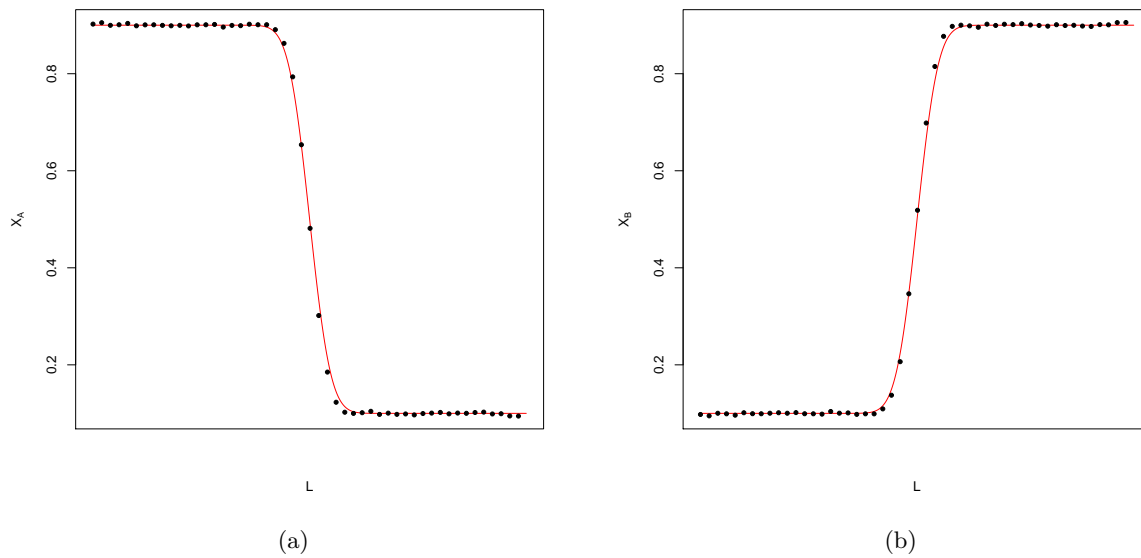


Figure 7.6: Plot of the concentrations after 20,000 time steps of the component a) A and b) B, in comparison with the analytical solution given by Eq.(7.78)

## 7.9 Outlook

We developed two new Fokker-Planck approximations to the Boltzmann equation, based on quasi-equilibrium models. These models were subjected to numerical experiments and it was determined that the algorithm is capable of simulating flow for a wide range of Knudsen numbers and diffusion coefficients. Since the algorithm relies on generation of Gaussian random numbers, a significant speedup was observed by employing use of the “Molecular Dice” algorithm, and hence proves to be an efficient alternative for solving binary mixtures.

The extension of the existing Fokker-Planck model to binary mixtures, is an indication that it could also be extended to solve for mixtures with many components. Future work is to extend this model to multi-component mixtures and possibly couple it with the Gillespie algorithm to design an efficient reaction-diffusion solver.

# Chapter 8

## Outlook

Numerical solutions of many stochastic models utilize non-uniformly distributed random numbers which are computationally expensive as they require evaluations of functions such as  $\log()$ ,  $\sin()$  etc. In this thesis, a new framework to generate random numbers is presented based on the inherent randomness present in Boltzmann dynamics. The main hypothesis is that a computer routine which simulates a stochastic process such as the motion of gas molecules is capable of generating random sequences on computers and can therefore form the basis for a new class of PRNGs. We tested our hypothesis by employing standard methods to simulate rarefied gases – Lennard Jones model, hard-sphere model, DSMC and MPCD. We considered the stream of numbers generated by these simulations as random sequences and found that these stream of numbers managed to satisfy statistical tests used to check for empirical randomness, thereby confirming our hypothesis that simulations of stochastic processes can indeed be used to create apparent randomness on computers. However, the rate of generation of random sequences using these methods was quite low and thus unsuitable for use in context of large-scale scientific computations. Based on these ideas a new algorithm was formulated capable of producing high quality Gaussian and exponential random numbers offering higher rate of exponential and Gaussian random number generation. Using canonical statistical tests, we found the quality of random sequences generated by proposed algorithm was at par with widely used PRNGs such as – drand48 and Mersenne Twister, while generating Gaussian and exponential random numbers 25 and 15 times faster than standard methods, respectively. Hence, proposed algorithm can be used to mitigate the problem of spending large fraction of computational time on random number generation for stochastic simulations.

Chemical reactions are modelled deterministically by the rate law of mass action, which provides accurate results for large systems. However for small systems fluctuations can have major impact on its behaviour and hence stochastic models must be used for such phenomena. The Gillespie algorithm is a numerical scheme to solve the stochastic models of chemical reactions and relies on generation of an exponential and a uniform random number each iteration while spending roughly 85 – 90% on the same. Since proposed “Molecular Dice” algorithm generates both exponential and uniform random numbers in the same pass, it is well suited to the Gillespie method. To test the efficiency of proposed algorithm we simulated two biochemical networks – a bi-stable system wherein two proteins attempt to bind to a DNA and another of a reaction-diffusion system used to study pattern formation inside bacteria. We found that proposed “Molecular Dice” method produced accurate results and was around 4 times faster than standard methods.

The Fokker-Planck approximation to the Boltzmann equation, has emerged as an efficient alternative to simulate flows in a wide range of Knudsen numbers. We extended this model to solve for binary mixtures, similar to quasi-equilibrium based models. Two models were proposed – one for low Schmidt number another for high Schmidt numbers. We tested proposed model with three benchmark problems – Graham’s law for effusion, Couette flow and binary diffusion and found that results for all three problems were accurate. Since, the solution to Fokker-Planck methodology involves use of Gaussian random numbers, the “Molecular Dice” algorithm adds to the computational efficiency of the Fokker-Planck approach to solve for hydrodynamics.

Thus, the proposed approach of generating random numbers using hydrodynamic solvers can indeed prove to be beneficial for large-scale scientific simulations which have hitherto been difficult owing to the huge number of random streams required and the associated computational cost involved. In addition, we also extended the Fokker-Planck model for hydrodynamics to solve

for binary mixtures. We expect that this model can act as base to devise a scheme for multi-component mixtures, which can then be coupled with the Gillespie algorithm to form an efficient advection-diffusion-reaction solver.

# Appendix A: Proof of chi-squared test

The proof of Pearson's chi-squared test is completed by building a histogram of an arbitrary distribution. We consider  $N$  bins  $B_1, \dots, B_N$  in which  $r$  balls  $X_1, \dots, X_r$  are thrown with probabilities

$$P(X_i \in B_j) = p_j, \quad (8.1)$$

such that  $\sum_{j=1}^N p_j = 1$ . As can be seen, this is qualitatively the same as generating random sequences and binning them to plot a histogram. The number of balls in the  $i$ th bin is then found to be

$$\nu_i = \sum_{j=1}^r I(X_j \in B_i), \quad (8.2)$$

where  $I(X_j \in B_i)$  is the indicator function defined as

$$I(X_j \in B_i) = \begin{cases} 1 & X_j \in B_i \\ 0 & X_j \notin B_i. \end{cases} \quad (8.3)$$

The sum of number of balls in each bin is equal to the total number of balls and that it is impossible for one ball to be in two bins at the same time. The constraints on the system can then be expressed as

$$\sum_{i=1}^N \nu_i = r \quad \text{and} \quad I(X_i \in B_j)I(X_i \in B_k) = 0 \quad \text{for} \quad j \neq k. \quad (8.4)$$

Hence, the indicator function must follow a Bernoulli distribution with its mean and variance being

$$E[I(X_i \in B_j)] = p_j \quad \text{Var}(I(X_i \in B_j)) = p_j(1 - p_j). \quad (8.5)$$

The sum of indicator function provides the number of balls in a particular bin and its statistics can then be used to determine the behaviour of observed number of balls in a given bin. It can be seen that, as per the central limit theorem

$$\frac{\sum_{i=1}^r I(X_i \in B_j) - rE[I(X_i \in B_j)]}{\sqrt{r\text{Var}(I(X_i \in B_j))}} \rightarrow N(0, 1), \quad (8.6)$$

where  $N(0, 1)$  is shorthand for standard Gaussian distribution. This can be simplified to

$$Z_j = \frac{\nu_j - rp_j}{\sqrt{rp_j}} \rightarrow \sqrt{1 - p_j}N(0, 1) = N(0, 1 - p_j), \quad (8.7)$$

hence, the observed number of balls in each bin is indeed a Gaussian random number. However, owing to the constraints  $I(X_i \in B_j)I(X_i \in B_k) = 0$  for  $j \neq k$ , the quantities  $Z_j$  and  $Z_k$  are not independent. The covariance matrix can be derived and then diagonalized for translating  $Z_j$  to independent Gaussian random numbers. The cross-correlations are found to be

$$E[Z_j Z_k] = -\sqrt{p_j p_k}.$$

Hence, the structure of the covariance matrix associated with the  $Z$  vector is

$$\Sigma = \text{Cov}(Z) = \begin{bmatrix} 1 - p_1 & -\sqrt{p_1 p_2} & \cdots \\ -\sqrt{p_1 p_2} & 1 - p_2 & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}. \quad (8.8)$$

it can be readily observed that  $\Sigma = \mathbf{I} - \mathbf{p}\mathbf{p}^T$ , where  $\mathbf{p}$  is the column vector of the quantities  $(\sqrt{p_1}, \sqrt{p_2}, \dots)$ . Since  $\Sigma$  has this particular form its determinant can be calculated easily using Sylvester's determinant theorem and thereby making it simple to calculate the eigenvalues of this matrix. These eigenvalues can then be used to diagonalize the covariance matrix. The solution of the characteristic equation of the covariance matrix  $\Sigma$  is found to be

$$\text{Det}(\Sigma - \lambda \mathbf{I}) = (1 - \lambda)^{n-1} \lambda = 0, \quad (8.9)$$

the covariance matrix has  $(n - 1)$  eigenvalues that are 1 and a single eigenvalue equal to 0. Therefore, a matrix  $A$  exists such that

$$A\Sigma A^T = \begin{bmatrix} 0 & 0 & 0 & \cdots \\ 0 & 1 & 0 & \cdots \\ 0 & 0 & 1 & \cdots \\ \vdots & \vdots & \vdots & \ddots \end{bmatrix}. \quad (8.10)$$

A new vector  $X = AZ$  is defined, it can be shown that  $X$  is Gaussian distributed with a covariance matrix  $A\Sigma A^T$ . Given that one of the eigenvalues is 0, the form of  $X$  is  $(0, X_1, \dots, X_{N-1})$  where  $X_i$  are independent and identically distributed normal random variables. Since the norm of the vector doesn't change under rotation, we infer

$$\sum_{i=1}^{N-1} X_i^2 = \sum_{i=1}^N Z_i^2 = \sum_{i=1}^N \frac{(\nu_i - \mu_i)^2}{\mu_i}, \quad (8.11)$$

hence proving the original claim that the test statistic,  $t = \sum_{i=1}^N (\nu_i - \mu_i)^2 / \mu_i$  is equal to the sum of the square of  $(N - 1)$  independent Gaussian random numbers.

# References

- ALLEN, R. J., WARREN, P. B. & TEN WOLDE, P. R. 2005 Sampling rare switching events in biochemical networks. *Phys. Rev. Lett.* **94** (1), 018104.
- ARCIDIACONO, S., MANTZARAS, J., ANSUMALI, S., KARLIN, I. V., FROUZAKIS, C. & BOULOUSCHOS, K. B. 2006 Simulation of binary mixtures with the lattice boltzmann method. *Phys. Rev. E* **74**, 056707.
- ARKIN, A., ROSS, J. & MCADAMS, H. H. 1998 Stochastic kinetic analysis of developmental pathway bifurcation in phage  $\lambda$ -infected escherichia coli cells. *Genetics* **149** (4), 1633–1648.
- BARAS, F. & MANSOUR, M. M. 1997 Microscopic simulations of chemical instabilities. *Advances in Chemical Physics, Volume 100* pp. 393–474.
- BECSKEI, A. & SERRANO, L. 2000 Engineering stability in gene networks by autoregulation. *Nature* **405** (6786), 590.
- BERGMAN, T. L., INCROPERA, F. P., DEWITT, D. P. & LAVINE, A. S. 2011 *Fundamentals of heat and mass transfer*. John Wiley & Sons.
- BHATNAGAR, P. L., GROSS, E. P. & KROOK, M. 1954*a* A model for collision processes in gases. i. small amplitude processes in charged and neutral one-component systems. *Physical review* **94** (3), 511.
- BHATNAGAR, P. L., GROSS, E. P. & KROOK, M. 1954*b* A model for collision processes in gases. i. small amplitude processes in charged and neutral one-component systems. *Phy. Rev.* **94**, 511–525.
- BIRD, G. 1978 Monte carlo simulation of gas flows. *Ann. Rev. Fluid Mech.* **10** (1), 11–31.
- BIRD, G. 1994 *Molecular gas dynamics and the direct simulations of the gases*. Oxford University Press.
- BOX, G. E., MULLER, M. E. *et al.* 1958 A note on the generation of random normal deviates. *The annals of mathematical statistics* **29** (2), 610–611.
- BOX, G. E. & PIERCE, D. A. 1970 Distribution of residual autocorrelations in autoregressive-integrated moving average time series models. *Journal of the American statistical Association* **65** (332), 1509–1526.
- CASELLA, G. & BERGER, R. L. 2002 *Statistical inference*, , vol. 2. Duxbury Pacific Grove, CA.

- CASELLA, G., ROBERT, C. P. & WELLS, M. T. 2004 Generalized accept-reject sampling schemes. *Lecture Notes-Monograph Series* pp. 342–347.
- CERCIGNANI, C. 1975 *Theory and application of the Boltzmann equation*. Scottish Academy Press, Edinburgh.
- CERCIGNANI, C. 1998 *Ludwig Boltzmann: the man who trusted atoms*. Oxford University Press.
- CHAPMAN, S. & COWLING, T. G. 1970 *The Mathematical Theory of Non-Uniform Gases: An Account of the Kinetic Theory of Viscosity, Thermal Conduction and Diffusion in Gases*. Cambridge university press.
- CHORIN, A. J. & HALD, O. H. 2009 *Stochastic tools in mathematics and science*, , vol. 3. Springer.
- COWAN, G. 1998 *Statistical data analysis*. Oxford university press.
- ELF, J. & EHRENBERG, M. 2004 Spontaneous separation of bi-stable biochemical systems into spatial domains of opposite phases. *Syst. Biol.* **1** (2), 230–236.
- ERBAN, R., CHAPMAN, J. & MAINI, P. 2007 A practical guide to stochastic simulations of reaction-diffusion processes. *arXiv preprint arXiv:0704.1908* .
- ÉRDI, P. & TÓTH, J. 1989 *Mathematical models of chemical reactions: theory and applications of deterministic and stochastic models*. Manchester University Press.
- ESPENSON, J. H. 1995 *Chemical kinetics and reaction mechanisms*, , vol. 102. Citeseer.
- FREEDMAN, D., PISANI, R. & PURVES, R. 2007 Statistics (international student edition). *Pisani, R. Purves.*—4th edition.—NY, USA: WW Norton & Company **720**.
- FREEDMAN, D. A. 2009 *Statistical models: theory and practice*. cambridge university press.
- FRISCH, U. 1995 Turbulence.
- GARDINER, C. W. 1985a *Handbook of stochastic methods*, , vol. 3. Springer Berlin.
- GARDINER, C. W. 1985b *Stochastic methods*. Springer-Verlag, Berlin–Heidelberg–New York–Tokyo.
- GILLESPIE, D. T. 1976 A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of computational physics* **22** (4), 403–434.
- GILLESPIE, D. T. 1977 Exact stochastic simulation of coupled chemical reactions. *J. Phys. Chem.* **81** (25), 2340–2361.
- GILLESPIE, D. T. 2000 The chemical langevin equation. *The Journal of Chemical Physics* **113** (1), 297–306.
- GILLESPIE, D. T. 2007 Stochastic simulation of chemical kinetics. *Annu. Rev. Phys. Chem.* **58**, 35–55.



- 
- GOLDBETER, A. & KOSHLAND, D. E. 1981 An amplified sensitivity arising from covalent modification in biological systems. *Proceedings of the National Academy of Sciences* **78** (11), 6840–6844.
- GOLDMAN, E. & SIROVICH, L. 1967 Equations for gas mixtures. *Phys. Fluids* **10**, 1928.
- HAMEL, B. B. 1965 Kinetic model for binary gas mixtures. *Phys. Fluids* **8**, 418 – 425.
- HOOPERBRUGGE, P. & KOELMAN, J. 1992 Simulating microscopic hydrodynamic phenomena with dissipative particle dynamics. *EPL (Europhysics Letters)* **19** (3), 155.
- HOWARD, M. & RUTENBERG, A. D. 2003 Pattern formation inside bacteria: fluctuations due to the low copy number of proteins. *Physical Review Letters* **90** (12), 128102.
- HOWARD, M., RUTENBERG, A. D. & DE VET, S. 2001 Dynamic compartmentalization of bacteria: accurate division in e. coli. *Physical review letters* **87** (27), 278102.
- HULL, T. E. & DOBELL, A. R. 1962 Random number generators. *SIAM review* **4** (3), 230–254.
- JOHNSON, J. B. 1928 Thermal agitation of electricity in conductors. *Physical review* **32** (1), 97.
- JONES, J. E. 1924 On the determination of molecular fields.ii. from the equation of state of a gas. In *Proc. R. Soc. Lond. A*, , vol. 106, pp. 463–477. The Royal Society.
- JUN, B. & KOCHER, P. 1999 The intel random number generator. *Cryptography Research Inc. white paper* .
- KAPRAL, R. 2008 Multiparticle collision dynamics: simulation of complex systems on mesoscales. *Adv. Chem. Phys.* **140**, 89.
- KNUTH, D. E. 1981 *Seminumerical algorithms, volume 2 of The Art of Computer Programming*. Addison-Wesley, Reading, MA, second edition.
- KOHDA, T. & TSUNEDA, A. 1997 Statistics of chaotic binary sequences. *IEEE Trans. Inform. Theory* **43** (1), 104–112.
- KRAFT, M. & WAGNER, W. 2003 Numerical study of a stochastic particle method for homogeneous gas-phase reactions. *Computers & Mathematics with Applications* **45** (1-3), 329–349.
- KRISHNAN, V. 2015 *Probability and random processes*. John Wiley & Sons.
- KURTZ, T. G. 1972 The relationship between stochastic and deterministic models for chemical reactions. *The Journal of Chemical Physics* **57** (7), 2976–2978.
- LADD, A. 2009 Numerical methods for molecular and continuum dynamics. *Lectures at the 3rd Warsaw School of Statistical Physics, Kazimierz, Poland* .
- LEBOWITZ, J., FRISCH, H. & HELFAND, E. 1960a Nonequilibrium distribution functions in a fluid. *The Physics of Fluids* **3** (3), 325–338.

- LEBOWITZ, J. L., FRISCH, H. L. & HELFAND, E. 1960*b* Nonlinear distribution function in a fluid. *Phys. Fluids* **3**, 325–338.
- L’ECUYER, P. & SIMARD, R. 2007 Testu01: Ac library for empirical testing of random number generators. *ACM Transactions on Mathematical Software (TOMS)* **33** (4), 22.
- LEVERMORE, C. D. 1996 Moment closure hierarchies for kinetic theories. *Journal of statistical Physics* **83** (5-6), 1021–1065.
- LIBOFF, R. L. 2003 *Kinetic theory: classical, quantum, and relativistic descriptions*. Springer Science & Business Media.
- LJUNG, G. M. & BOX, G. E. 1978 On a measure of lack of fit in time series models. *Biometrika* **65** (2), 297–303.
- MARANDI, A., LEINDECKER, N. C., VODOPYANOV, K. L. & BYER, R. L. 2012 All-optical quantum random bit generation from intrinsically binary phase of parametric oscillators. *Optics express* **20** (17), 19322–19330.
- MARSAGLIA, G. 1968 Random numbers fall mainly in the planes. *Natl. Acad. Sci. Proc.* **61** (1), 25–28.
- MARSAGLIA, G. & TSANG 2002 Some difficult-to-pass tests of randomness. *Journal of Statistical Software* **7** (3), 1–9.
- MARSAGLIA, G. *et al.* 2003 Xorshift rngs. *Journal of Statistical Software* **8** (14), 1–6.
- MASON, E. & KRONSTADT, B. 1967 Graham’s laws of diffusion and effusion. *Journal of Chemical Education* **44** (12), 740.
- MASSEY JR, F. J. 1951 The kolmogorov-smirnov test for goodness of fit. *Journal of the American statistical Association* **46** (253), 68–78.
- MATSUMOTO, M. & NISHIMURA, T. 1998 Mersenne twister: a 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* **8** (1), 3–30.
- MORSE, T. F. 1964 Kinetic model equations for a gas mixture. *Phys. Fluids* **7**, 2012 – 2013.
- NYQUIST, H. 1928 Thermal agitation of electric charge in conductors. *Physical review* **32** (1), 110.
- OPPENHEIM, I., SHULER, K. & WEISS, G. 1969 Stochastic and deterministic formulation of chemical rate equations. *The Journal of Chemical Physics* **50** (1), 460–466.
- ONEILL, M. E. 2015 Pcg: A family of simple fast space-efficient statistically good algorithms for random number generation. *ACM Transactions on Mathematical Software* .

- 
- PEARSON, K. 1900 X. on the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **50** (302), 157–175.
- RAPAPORT, D. C. 2004 *The Art of Molecular Dynamics Simulation*. Cambridge university press.
- RISKEN, H. 1996 Fokker-planck equation. In *The Fokker-Planck Equation*, pp. 63–95. Springer.
- SHARIPOV, F., CUMIN, L. M. G. & KALEMPA, D. 2004 Plane couette flow of binary gaseous mixture in the whole range of the knudsen number. *Eur. J. Mech. B Fluids* **23** (6), 899–906.
- SINGH, S. & ANSUMALI, S. 2015a Fokker-planck model of hydrodynamics. *Physical Review E* **91** (3), 033303.
- SINGH, S., THANTANAPALLY, C. & ANSUMALI, S. 2016 Gaseous microflow modeling using the fokker-planck equation. *Physical Review E* **94** (6), 063307.
- SINGH, S. K. & ANSUMALI, S. 2015b Fokker - planck model of hydrodynamics. *Phy. Rev. E* **91**(3), 033303.
- SINGH, S. K., THANTANAPALLY, C. & ANSUMALI, S. 2015 Gaseous microflow modeling using the fokker-planck equation. *Phy. Rev. E* **94**(6), 063307.
- SIROVICH, L. 1962 Kinetic modeling of gas mixtures. *Phys. Fluids* **5**, 908.
- SIROVICH, L. 1966 Mixtures of maxwell molecules. *Phys. Fluids* **9**, 2323.
- SRIVASTAVA, R., YOU, L., SUMMERS, J. & YIN, J. 2002 Stochastic vs. deterministic modeling of intracellular viral kinetics. *Journal of theoretical biology* **218** (3), 309–321.
- SUCCI, S. 2001 *The lattice Boltzmann equation: for fluid dynamics and beyond*. Oxford university press.
- SUCCI, S., KARLIN, I. V. & CHEN, H. 2002 Colloquium: Role of the h theorem in lattice boltzmann hydrodynamic simulations. *Rev. Mod. Phys.* **74** (4), 1203.
- SYMUL, T., ASSAD, S. & LAM, P. K. 2011 Real time demonstration of high bitrate quantum random number generation with coherent laser light. *Applied Physics Letters* **98** (23), 231103.
- THAR, R. & KÜHL, M. 2003 Bacteria are not too small for spatial sensing of chemical gradients: an experimental evidence. *Proceedings of the National Academy of Sciences* **100** (10), 5748–5753.
- THOMAS, D. B., LUK, W., LEONG, P. H. & VILLASENOR, J. D. 2007 Gaussian random number generators. *ACM Comput. Surv.* **39** (4), 11.
- TURING, A. M. 1952 The chemical basis of morphogenesis. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* **237** (641), 37–72.

- TURNER, T. E., SCHNELL, S. & BURRAGE, K. 2004 Stochastic approaches for modelling in vivo reactions. *Computational biology and chemistry* **28** (3), 165–178.
- UHLENBECK, G. E. & ORNSTEIN, L. S. 1930 On the theory of the brownian motion. *Phys. Rev.* **36** (5), 823.
- WANG, J., TSANG, W. W. & MARSAGLIA, G. 2003 Evaluating kolmogorov's distribution. *Journal of Statistical Software* **8** (18).
- YOSHIDA, T., MORI, H. & SHIGEMATSU, H. 1983 Analytic study of chaos of the tent map: band structures, power spectra, and critical behaviors. *J. Stat. Phys.* **31** (2), 279–308.