



**KTH Computer Science  
and Communication**

# **Complexity and Error Analysis of Numerical Methods for Wireless Channels, SDE, Random Variables, and Quantum Mechanics**

HÅKON ANDREAS HOEL

Doctoral Thesis  
Stockholm, Sweden, 2012

TRITA-CSC-A 2012:06  
ISSN-1653-5723  
ISRN KTH/CSC/A-12/06-SE  
ISBN 978-91-7501-350-3

School of Computer Science and Communication  
KTH  
SE-100 44 Stockholm  
SWEDEN

Akademisk avhandling som med tillstånd av Kungl Tekniska högskolan framlägges till offentlig granskning för avläggande av Doktorgradsexamen onsdagen den 30 maj 2012 kl 10.15 i sal F3, Lindstedtsvägen 26, Kungliga Tekniska högskolan, Stockholm.

© Håkon Andreas Hoel, 2012

Tryck: E-print, [www.eprint.se](http://www.eprint.se)

## Abstract

This thesis consists of four papers considering different aspects of stochastic process modeling, error analysis, and minimization of computational cost.

In Paper I, we construct a Multipath Fading Channel (MFC) model for wireless channels with noise introduced through scatterers flipping on and off. By coarse graining the MFC model a Gaussian process channel model is developed. Complexity and accuracy comparisons of the models are conducted.

In Paper II, we generalize a multilevel Forward Euler Monte Carlo method introduced by Giles [16] for the approximation of expected values depending on solutions of Itô stochastic differential equations. Giles work [16] proposed and analyzed a Forward Euler Multilevel Monte Carlo (MLMC) method based on realizations on a hierarchy of uniform time discretizations and a coarse graining based control variates idea to reduce the computational cost required by a standard single level Forward Euler Monte Carlo method. This work is an extension of Giles' MLMC method from uniform to adaptive time grids. It has the same improvement in computational cost and is applicable to a larger set of problems.

In paper III, we consider the problem to estimate the mean of a random variable by a sequential stopping rule Monte Carlo method. The performance of a typical second moment based sequential stopping rule is shown to be unreliable both by numerical examples and by analytical arguments. Based on analysis and approximation of error bounds we construct a higher moment based stopping rule which performs more reliably.

In paper IV, Born-Oppenheimer dynamics is shown to provide an accurate approximation of time-independent Schrödinger observables for a molecular system with an electron spectral gap, in the limit of large ratio of nuclei and electron masses, without assuming that the nuclei are localized to vanishing domains. The derivation, based on a Hamiltonian system interpretation of the Schrödinger equation and stability of the corresponding hitting time Hamilton-Jacobi equation for non ergodic dynamics, bypasses the usual separation of nuclei and electron wave functions, includes caustic states and gives a different perspective on the Born-Oppenheimer approximation, Schrödinger Hamiltonian systems and numerical simulation in molecular dynamics modeling at constant energy.



# Preface

This thesis consists of an introduction and the four following papers:

**Paper I** Håkon Hoel and Henrik Nyberg. *Gaussian Coarse Graining of a Master Equation Extension of Clarke's Model*, TRITA-NA 2012:5. Submitted to Advances in Applied Probability.

The author contributed to all sections of the paper.

**Paper II** Håkon Hoel, Erik von Schwerin, Anders Szepessy, and Raúl Tempone. *Implementation and Analysis of an Adaptive Multilevel Monte Carlo Algorithm*, TRITA-NA 2012:6. Shorter version published in Numerical Analysis of Multiscale Computations, Lecture Notes in CSE, Springer, cf. [20]. Full version submitted to Monte Carlo Methods and Applications.

The author contributed to the analysis and the description of the algorithms.

**Paper III** Christian Bayer, Håkon Hoel, Erik von Schwerin, and Raúl Tempone. *On Non-Asymptotic Optimal Stopping Criteria in Monte Carlo Simulations*, TRITA-NA 2012:7. In preparation.

The author contributed to all sections of the paper.

**Paper IV** Christian Bayer, Håkon Hoel, Ashraful Kadir, Petr Plecháč, Mattias Sandberg, Anders Szepessy, and Raúl Tempone. *How Accurate is Molecular Dynamics?* TRITA-NA 2012:8.

Earlier version available on Arxiv, cf. [2], full version to be submitted.

The author contributed to the numerical examples and to a small degree to the analysis.



# Acknowledgments

I would like to thank my supervisor Anders Szepessy for providing interesting and challenging projects. My collaborators Christian Bayer, Ashraful Kadir, Henrik Nyberg, Petr Plecháč, Mattias Sandberg, Erik von Schwerin, and Raúl Tempone. A special thanks goes to Raúl Tempone for being a great host during my visits at his university and to Christian Bayer for much needed help on sampling problems. My colleagues Jesper Karlsson, Jonas Kiessling, Love Lindholm, and Georgios Zouraris for discussions during the weekly meetings of professor Szepessy's research group.

I would also like to thank all my colleagues at the numerical analysis department at KTH, especially my kind friends Murtazo Nazarov, Jelena Popovic, and Sara Zahedi.

This work is funded by Center for Industrial and Applied Mathematics (CIAM) and King Abdullah University of Science and Technology, Saudi Arabia.





# Contents

Contents	vii
<b>I Introductory Chapters</b>	<b>1</b>
1 Introduction	3
2 Probability background	7
2.1 Random Variables and Probability Measures . . . . .	7
2.2 Convergence of Random Variables and Limit Theorems for Sampling . . . . .	11
3 Monte Carlo Methods	13
3.1 Examples and Properties . . . . .	13
3.2 Variance Reduction . . . . .	15
3.3 Sequential Stopping Rules . . . . .	17
4 Adaptive Weak Approximations for SDE	19
4.1 Stochastic Differential Equations (SDE) . . . . .	19
4.2 The Euler Method for Itô SDE . . . . .	22
4.3 Adaptive Weak Solution Approximation . . . . .	28
5 Wireless Channel Modeling	31
5.1 The Multipath Fading Channel . . . . .	31
5.2 From MFC Models to Gaussian Processes . . . . .	36
6 Classical and Quantum Mechanics	39
6.1 Classical Mechanics . . . . .	39
6.2 Quantum Mechanics . . . . .	43
Bibliography	51
<b>II Included Papers</b>	<b>55</b>



## Part I

# Introductory Chapters



# Chapter 1

## Introduction

If a system does not always produce the same output from a given initial state, we call it a non-deterministic system. Considering a non-deterministic system whose uncertainty is described in form of a density  $P$ , mean value quantities on the form

$$\int g(x)dP(x) \tag{1.0.1}$$

are often sought. For example:

- Let  $X$  be a random variable modeling the number of active users in a network with the with  $P(n)$  denoting the probability for having  $n \in \mathbb{N}$  active users (within a fixed time interval). Then (1.0.1) with  $g(n) = n$  represents the mean number of active users in the network.
- A mechanical system with  $N$  particles positioned at  $q \in \mathbb{R}^{3N}$ , with momentum  $p \in \mathbb{R}^{3N}$ , potential function  $V(q)$ , and  $P(q, p)$  denoting the probability for the particles having the phase space configuration  $(q, p)$ . Then the average potential energy of the system  $g(q) = V(q)$  is a quantity of interest.

In some settings, e.g., if the density  $P$  is not explicitly known or if it takes a subtle form, integrals on the form (1.0.1) have to be approximated. A typical way of approximating (1.0.1) is by Monte Carlo sampling

$$\sum_{i=1}^M \frac{g(X_i)}{M},$$

where the realizations  $X_i$  are generated according to or approximately according to the density  $P$ . In this thesis, we will for four different non-deterministic systems study questions on how to efficiently generate realizations that are approximately distributed according to a density  $P$  and give estimates on weak approximation errors of type

$$\left| \int g(x)dP(x) - \sum_{i=1}^M \frac{g(X_i)}{M} \right| \tag{1.0.2}$$

in terms of the computational cost. Let us be more specific.

## Paper I

Wireless channel models model the received version of a signal transmitted wirelessly from a transmitter to a receiver. For industries developing wireless transmission equipment, such models are used to estimate the performance of equipment and software by simulations instead of more costly real world tests. But for such models to be of interest, they must be accurate and computationally efficient with respect to running time.

One of the most popular models of today, the Multipath Fading Channel model (MFC), is based on approximating the signal from superpositioning a finite number of contributing wave paths yielding an output signal on the form

$$Y_t = \frac{1}{\sqrt{M}} \sum_{k=1}^M a_k e^{i\theta_k(t)}.$$

The superpositioning of wave paths is computationally costly, and, consequently, it is also costly to generate output signal realizations using MFC models. In Paper I, we propose a new MFC model with scatterers flipping on and off adding noise to the total signal, and by coarse graining this MFC model we derive a Gaussian process wireless channel model. Computational cost estimates in Paper I indicate that signal realizations are generated more efficiently by using a Gaussian process algorithm than by using an MFC algorithm.

## Paper II

Stochastic Differential Equations (SDE) are non-deterministic processes whose future evolution is described by a probability distribution, as opposed to the deterministic evolution of an ordinary differential equation. Complex phenomena which might seem non-deterministic, such as stock market evolution, are frequently modeled by stochastic processes, cf. [5, 24].

In the second paper, we develop an adaptive Forward Euler Monte Carlo algorithm which for any sufficiently well behaved function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  approximates the expected value

$$\mathbb{E}[g(X_T)] = \int g(X_T(\omega)) dP(\omega), \quad (1.0.3)$$

where  $X_T$  is the solution of an Itô SDE. The algorithm we have developed constructs numerical realizations of the SDE on adaptive time grids using the adaptive Forward Euler method. To obtain variance reduction, realizations are constructed on different tolerance levels and the expected value (1.0.3) is approximated from the numerical realizations by the Monte Carlo method.

The multilevel Monte Carlo method, first introduced by Giles in [16], showed the improvement of computational cost of approximating  $\mathbb{E}[g(X_T)]$  with accuracy  $\mathcal{O}(\text{TOL})$  from  $\mathcal{O}(\text{TOL}^{-3})$  by a single level method to  $\mathcal{O}(\text{TOL}^{-2} \log(\text{TOL})^2)$  by the multilevel method when the realizations  $X_T(\omega)$  of the given SDE problem were generated numerically on uniform time grids. In Paper II we extend Giles' multilevel method to the setting of adaptive time grid SDE realizations relying on the adaptive weak approximation methods for SDE developed by Szepessy et al. [30, 25, 26]. Our extended multilevel method is applicable to a larger set of SDE problems and, when comparable, it is shown to have the same computational cost as Giles' method has.

### Paper III

Given i.i.d. random variables  $X_1, X_2, \dots$  the typical way of approximating their expected value  $\mu = \mathbb{E}[X]$  using  $M$  samples is the sample average

$$\bar{X}_M := \sum_{i=1}^M \frac{X_i}{M}.$$

In Paper III, we consider the objective of choosing  $M$  sufficiently large so that

$$P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta, \quad (1.0.4)$$

for small, fixed accuracy-confidence constants  $\text{TOL} > 0$  and  $\delta > 0$ . Clearly,  $P(|\bar{X}_M - \mu| > \text{TOL})$  decreases as  $M$  increases, but at the same time the cost of computing  $\bar{X}_M$  increases. From an application and cost point of view it is therefore of interest to have theory giving sharp bounds on the number of samples  $M$  needed to fulfill (1.0.4). For some settings this exists. For example, if  $\mathbb{E}[|X|^\infty] < C$ , it is possible to derive good theoretical upper bounds for  $M$ , but in the general case when no or little information of the distribution is given, however, little theory is known and the typical way of estimating  $\mathbb{E}[X]$  is by a sequential stopping rule; sequentially increasing the number of samples  $M$  until the sampled moments fulfill a stopping criterion. In Paper III, we show that the “intuitive” stopping criterion

$$2 \left( 1 - \Phi \left( \frac{\sqrt{M} \text{TOL}}{\bar{\sigma}_M} \right) \right) < \delta,$$

where

$$\Phi(z) := \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-x^2/2} dx$$

gives a stopping rule that performs unreliably when sampling heavy-tailed r.v. From approximations of error bounds we construct a new stopping criterion based on second, third, and fourth order sample moments which according to numerical experiments performs more reliably and is only slightly more costly than the stopping rule with the “intuitive” stopping criterion.

### Paper IV

Molecular dynamics is a computational method to study molecular systems in materials science, chemistry and molecular biology. The simulations are used, for example, in designing and understanding new materials or for determining biochemical reactions in drug design. The wide popularity of molecular dynamics simulations relies on the fact that in many cases it agrees very well with experiments. Indeed, given experimental data it is easy to verify correctness of the method by comparing with experiments at certain parameter regimes. However, if we want the simulation to predict something that has no comparing experiment, we need a mathematical estimate of the accuracy of the computation. In the case of molecular systems with few particles such studies are made by directly solving the Schrödinger equation. A fundamental and still open question in classical molecular dynamics simulations is how to verify the accuracy computationally, i.e., when the solution of the Schrödinger equation is not a computational alternative.

The aim of this paper is to derive qualitative error estimates for molecular dynamics and present new mathematical methods which could be used also for a more demanding

quantitative accuracy estimation, without solving the Schrödinger equation. That is, let  $\Phi$  be a solution of the time-independent Schrödinger equation

$$\left( -\frac{1}{2}M^{-1} \sum_{n=1}^N \Delta_{X^n} + \mathcal{V} \right) \Phi = E\Phi,$$

where  $\mathcal{V}$  is a given potential operator,  $E \in \mathbb{R}$  is energy and  $M$  is a the mass constant, and let  $\mathbf{X}(t)$  be a molecular dynamics path with total energy also equal to  $E$ . Then, our object of study is the approximation error of

$$\int_{\mathbb{R}^{3(N+n)}} g(\mathbf{X}) \Phi(\mathbf{X}, \mathbf{x})^* \Phi(\mathbf{X}, \mathbf{x}) d\mathbf{X} d\mathbf{x} - \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(\mathbf{X}(t)) dt \quad (1.0.5)$$

in terms of  $M$  as  $M \rightarrow \infty$ , for any observable  $g(\mathbf{X})$ .

Having molecular dynamics error estimates opens, for instance, the possibility of systematically evaluating which density functionals or empirical force fields are good approximations and under what conditions the approximation properties hold. Computations with such error estimates could also give improved understanding when quantum effects are important and when they are not, in particular in cases when the Schrödinger equation is too computationally complex to solve.

### Outline of the Introductory Chapters

In Chapter 2, we give a short description of random variables, probability distributions, and limit theorems in probability theory. In Chapter 3, we present topics on Monte Carlo methods such as variance reduction and dimension independent convergence rate. In Chapter 4, we give a short introduction to numerical methods and weak solution approximation for SDE. Chapter 5 presents statistical wave path modeling of wireless channels, and in Chapter 6, we give a short outline on Classical and Quantum Mechanics.



## Chapter 2

# Probability background

*Although the Events of Games, which Fortune solely governs, are uncertain, yet it may be certainly determin'd, how much one is more ready to lose than gain...*

*It is impossible for a Die, with such determin'd force and direction, not to fall on such determin'd side, only I don't know the force and direction which makes it fall on such determin'd side, and therefore I call it Chance, which is nothing but the want of art.*

—John Arbuthnot, *Of the Laws of Chance*<sup>1</sup>.

### 2.1 Random Variables and Probability Measures

During renaissance times, loose guidelines involving reason, intuition, and observations were applied when assessing the uncertainty of evidence material in court, for discussing betting strategies in terms of odds, and for setting maritime insurance premiums. A first mathematical treatment of probability can be traced back to letters between Piere de Fermat and Blaise Pascal in the year 1654 where they discuss, among other things, the division of stakes in fair gambling games. Let us therefore, in the spirit of gambling, develop our first random variable (r.v.) as a model of the game Heads or Tails. Heads or Tails is the game of predicting which side will face up when a coin is flipped. The game is quite successfully modeled by

$$\text{Outcome} = \begin{cases} \text{Heads,} & \text{with probability } p \\ \text{Tails,} & \text{with probability } (1 - p), \end{cases} \quad (2.1.1)$$

where  $p \in [0, 1]$  models the coin's proclivity towards Heads. (Typically with  $p = 1/2$  when the coin is of fair shaped, or alternatively configured from experiments or measurements. For example, tossing the coin  $N$  times and setting  $p = \#\text{Heads}/N$ .)

It is often possible to motivate both deterministic and stochastic models for the problem one is studying. In the case of Heads or Tails, for example, one might argue that if all input parameters—initial orientation, velocity and spin for the coin; air density; material and geometrical data for the coin and the ground a.s.f.—is known prior to the coin flip, the outcome of the coin flip ought to be deterministically predictable. But if the input data is uncertain and/or the deterministic model for the outcome is too complicated, a stochastic

---

<sup>1</sup>English translation (with additions) of *De Ratiociniis in Ludo Aleae* by Frans van Schooten and Christiaan Huygens (foreword, at least) in 1656.

model might be preferable in comparison to a deterministic one. Furthermore, the choice of model should depend on the problem you are facing. For example, if you wish to estimate how likely ten consecutive Heads throws are, the stochastic model straightforwardly gives you the estimate  $p^{10}$ , while if you wish to estimate a more specific property of the coin throws, for example the precise path of each coin flip, a deterministic model of the coin flips might be needed.

To present and analyze more complicated r.v. and stochastic models we introduce some terminology for probability spaces and r.v.

**Definition 2.1.1 (Probability space,  $\sigma$ -algebra, and Probability Measure)** A probability space is a measure space triple  $(\Omega, \mathfrak{F}, P)$  with the sample space  $\Omega$ , the event space  $\mathfrak{F}$ , and the probability measure  $P : \mathfrak{F} \rightarrow [0, 1]$ . The sample space is the set of outcomes, with an outcome signifying the result of single execution of the model. An event denotes the union of one or more outcomes, and the event space is the set of events. The event space  $\mathfrak{F}$  is a  $\sigma$ -algebra, i.e., it is a collection of subsets of  $\Omega$  fulfilling the following:

- (i)  $\mathfrak{F}$  is non-empty.
- (ii) Closed under complement: If  $A \in \mathfrak{F}$ , then  $A^C \in \mathfrak{F}$ , with  $A^C$  denoting the complement.
- (iii) Closed under countable unions: If  $A_i \in \mathfrak{F}$  for  $A_0, A_1, \dots$ , then  $\cup_i A_i \in \mathfrak{F}$ .

The probability measure acts on the measurable space  $(\Omega, \mathfrak{F})$  fulfilling the following criteria:

- (i) Non-negative:  $P(A) \geq 0$ .
- (ii) Countable additivity: For all collections a collection  $\{A_i\}_i$  contained in  $\mathfrak{F}$  and pairwise disjoint,

$$P(\cup_i A_i) = \sum_i P(A_i).$$

- (iii) Probability measure:  $P(\Omega) = 1$ .

A r.v. is a mapping defined on a probability space.

**Definition 2.1.2 (Random Variable)** Given the probability space  $(\Omega, \mathfrak{F}, P)$  and an arbitrary measurable space  $(\Gamma, \mathfrak{G})$ , then  $X : \Omega \rightarrow \Gamma$  is said to be a random variable/vector if the map from  $(\Omega, \mathfrak{F})$  to  $(\Gamma, \mathfrak{G})$  is measurable, i.e.,

$$\{\omega : X(\omega) \leq G\} \in \mathfrak{F}, \quad \forall G \in \mathfrak{G}.$$

A frequently encountered type of r.v. is the real-valued r.v. mapping from  $(\Omega, \mathfrak{F})$  to  $(\mathbb{R}, \mathfrak{B})$  where  $\mathfrak{B}$  denoting the Borel  $\sigma$ -algebra. A real-valued r.v.  $X$  has a non-decreasing, right continuous *Cumulative Distribution Function (CDF)*  $F(x) := P(X \leq x)$  defined for  $x \in \mathbb{R}$  and a *Probability Density Function (PDF)*  $f(x) = P(X = x)$  which is connected to CDF by

$$F(x) = \int_{-\infty}^x dF(s) = \int_{-\infty}^x f(s) ds.$$

Often, r.v. are indirectly defined by their PDF, and we further note that moments such as the *expected value*

$$\mu = E[X] := \int_{\mathbb{R}} xf(x) dx$$

and the *variance*

$$\sigma^2 = \text{Var}(X) := E[|X - E[X]|^2]$$

are values often used both to describe and to analyze r.v.

### Random Variables—Examples

We now include examples of some of the r.v. which appear in the papers of this thesis.

**Example 2.1.3 (Heads or Tails a.k.a. Bernoulli distribution)** *Based on the model (2.1.1) let us define the real-valued r.v.*

$$X(\omega) = \begin{cases} 1, & \text{if } \omega = \{\text{Heads}\} \\ 0, & \text{if } \omega = \{\text{Tails}\} \end{cases} \quad (2.1.2)$$

It has the probability space

$$\Omega = \{\text{Heads}, \text{Tails}\}, \quad \mathfrak{F} = \{\emptyset, \{\text{Heads}\}, \{\text{Tails}\}, \Omega\},$$

and  $P(X = 1) = p = P(X = \{\text{Tails}\})$ ,

the PDF

$$f(x) = \frac{\delta_x + \delta_{1-x}}{2},$$

with  $\delta_x$  denoting the Dirac delta distribution,  $\mu = p$ , and  $\sigma^2 = p(1-p)$ . (Strictly speaking, the function  $f$  is called a probability mass function when it only attains non-zero values at a countable number of  $x \in \mathbb{R}$ , not a PDF.)

**Example 2.1.4 (The Poisson distribution)** *The Poisson distribution models the likelihood for a given number of identical independent events occurring within a fixed time interval. For example, the number of telephone calls in a call system or the number of wave scatterers becoming active in a wireless multi path fading channel, cf. Paper I. The PDF is*

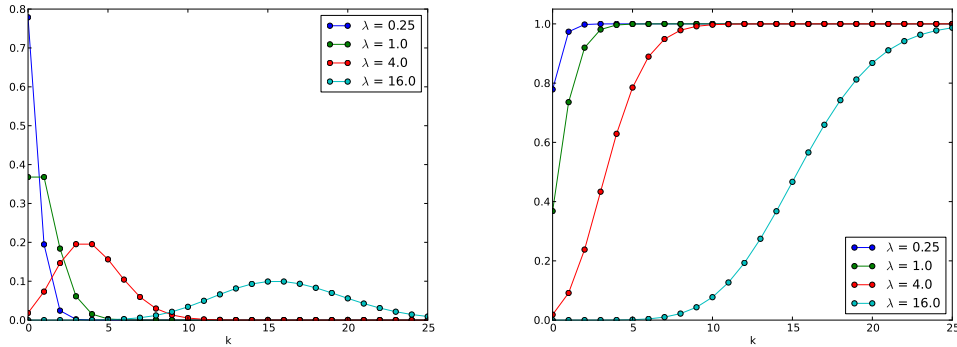


Figure 2.1: The Poisson PDF and CDF (both only defined at the dots for  $k \in \mathbb{N}_0$ —interpolating lines are included to illustrate the transitions).

given by

$$f(k) = \frac{\lambda^k e^{-\lambda}}{k!} \text{ for } k \in \mathbb{N}_0,$$

where  $\lambda > 0$  is a positive parameter equaling the expected number of events on a given interval. That is,  $E[X] = \lambda$  and (it turns out that) also  $\text{Var}(X) = \lambda$ .

**Example 2.1.5 (Normal distribution)** *The normal distribution is a very important distribution with a tremendous amount of applications. It is innately connected to the limit distribution of the scaled mean of samples of independent, identically distributed (i.i.d.) r.v. (see Section 2.2), and many real life uncertainty estimates depend on this distribution. Considering a normal r.v.  $X$  with mean  $\mu \in \mathbb{R}$  and variance  $\sigma^2 \in \mathbb{R}_+$ , the PDF is given by*

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

We write  $X \sim \mathcal{N}(\mu, \sigma^2)$ . We further note that the CDF of a standard normal r.v.  $X \sim \mathcal{N}(0, 1)$ , which often is used in estimates of sampling error, is in this thesis represented by  $\Phi(x)$ .

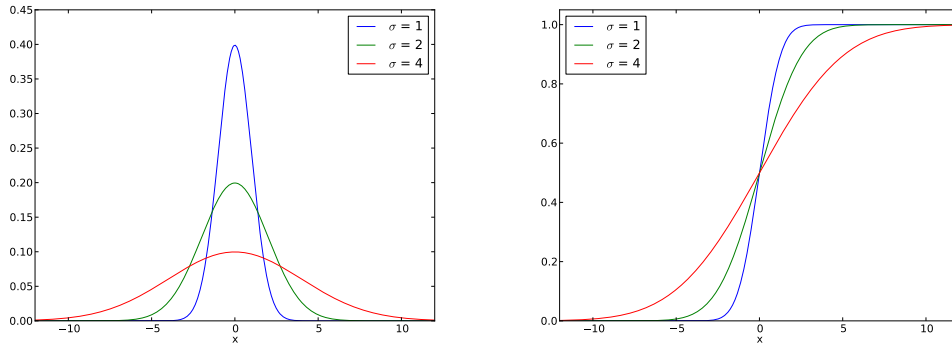


Figure 2.2: The PDFs and CDFs for normal distributions with  $\mu = 0$ .

Generalizations of the normal r.v. to  $\mathbb{R}^n$  and  $\mathbb{C}^n$  exist. In  $\mathbb{R}^n$ , the multivariate normal with input parameters  $\mu \in \mathbb{R}^n$  and symmetric, positive definite covariance matrix  $K \in \mathbb{R}^{n \times n}$  has the PDF

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{n/2} \det(K)^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T K^{-1}(\mathbf{x} - \mu)\right).$$

**Example 2.1.6 (Pareto distribution)** *The Pareto-distribution has the PDF*

$$f(x) = \begin{cases} \alpha x_m^\alpha x^{-(\alpha+1)} & \text{if } x \geq x_m \\ 0 & \text{else,} \end{cases} \quad (2.1.3)$$

where  $\alpha, x_m \in \mathbb{R}_+$  are respectively the shape and the scale parameter. The moments of  $E[X^n]$  for the Pareto r.v. only exists for  $n < \alpha$  and, supposing  $\alpha > 2$ , its mean and variance are given by

$$\mu = \frac{\alpha x_m}{\alpha - 1} \text{ and } \sigma^2 = \frac{x_m^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)}.$$

The Pareto-distribution is a so called heavy-tailed distribution, meaning that its tail is not exponentially bounded, i.e.

$$\lim_{x \rightarrow \infty} e^{\lambda x} P(X \geq x) = \infty, \quad \forall \lambda \in \mathbb{R}_+.$$

The Italian economist Vilfredo Pareto introduced the Pareto distribution as model for income distribution. The pareto index  $\alpha$  models the breadth of income; the larger the pareto index, the smaller the proportion of high-income people.

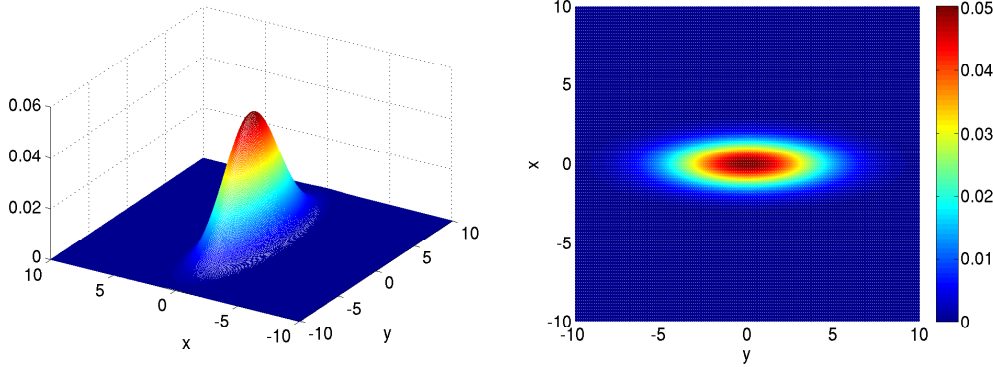


Figure 2.3: The PDF of the 2D multivariate normal distribution with  $\mu = (0,0)$  and the covariance matrix  $K = \text{diag}(1,10)$ .

## 2.2 Convergence of Random Variables and Limit Theorems for Sampling

In this section we review relevant convergence definitions and results from probability theory. For simplicity, the results are presented for sequences of scalar, real-valued r.v.

**Definition 2.2.1 (Convergence in distribution)** A sequence of real-valued r.v.  $X_0, X_1, \dots$  converges in distribution to the r.v.  $X$  with the CDF  $F(x)$  if

$$\lim_{n \rightarrow \infty} P(X_n \leq x) = F(x)$$

at all continuity points  $x \in \mathbb{R}$  of  $F$ .

**Definition 2.2.2 (Convergence in probability)** A sequence of real-valued r.v.  $X_0, X_1, \dots$  converges in probability to the r.v.  $X$  if for any  $\epsilon > 0$ ,

$$\lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0.$$

We write  $X_n \rightarrow X$  in probability.

**Definition 2.2.3 (Almost sure convergence)** A sequence of real-valued r.v.  $X_0, X_1, \dots$  converges almost surely to the r.v.  $X$  if

$$P(\lim_{n \rightarrow \infty} X_n - X = 0) = 1.$$

We write  $\lim_{n \rightarrow \infty} X_n = X$  a.s.

The strength relations between the convergence notions are

$$\text{Almost sure conv.} \implies \text{Conv. in probability} \implies \text{Conv. in distribution.}$$

Given a sequence of independent identically distributed (i.i.d.) r.v.  $\{X_i\}_{i=1}^n$ , a question of interest is if and how fast the average  $S_n/n$ , where  $S_n := \sum_{i=1}^n X_i$  converges to  $\mu = E[X_1]$  as the number of samples  $M$  increases. Before reviewing results regarding this question, we first recall the definition of independent r.v.

**Definition 2.2.4 (Independent random variables)** Two r.v.  $X_1$  and  $X_2$  mapping from  $(\Omega, \mathfrak{F})$  to  $(\mathbb{R}, \mathfrak{B})$  are independent if for all  $A_1, A_2 \in \mathfrak{B}$

$$P(X_1 \in A_1, X_2 \in A_2) = P(X_1 \in A_1)P(X_2 \in A_2).$$

In particular, independence implies that  $E[X_1 X_2] = E[X_1] E[X_2]$ . The simplest asymptotic convergence results for  $S_n/n$  is the weak law of large numbers.

**Theorem 2.2.5 (Weak law of large numbers)** Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with  $E[|X_1|] < \infty$ . Then  $S_n/n \rightarrow \mu$  in probability as  $M \rightarrow \infty$ .

The result in the general form stated above follows from the Dominated Convergence Theorem, but let us prove it in the setting when  $\text{Var}(X_1) < \infty$  using Chebycheff's inequality and independence: For any  $\epsilon > 0$ ,

$$\begin{aligned} P(|S_n/n - \mu| > \epsilon) &\stackrel{\text{Chebycheff's ineq.}}{\leq} E\left[\frac{|S_n/n - \mu|^2}{\epsilon^2}\right] \\ &\leq \sum_{i=1}^M \frac{\text{Var}(X_i)}{M^2 \epsilon^2} \rightarrow 0 \text{ as } M \rightarrow \infty. \end{aligned}$$

The Strong law of large numbers is a, not surprisingly, a slightly stronger result telling us that in the limit  $n \rightarrow \infty$ , the sample average  $S_n/n$  equals  $\mu$  on a measure 1 set (i.e., a.s.).

**Theorem 2.2.6 (Strong law of large numbers)** Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with  $E[|X_1|] < \infty$ . Then  $S_n/n \rightarrow \mu$  a.s. as  $n \rightarrow \infty$ .

The weak and strong laws of large numbers describes the asymptotic pointwise convergence of the sample mean  $S_n/n$ . The Central Limit Theorem (CLT) gives you additional information on the asymptotic distribution of  $S_n/\sqrt{n}$ .

**Theorem 2.2.7 (The Central Limit Theorem)** Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with (as before)  $\mu = E[X_1]$  and  $\sigma^2 = \text{Var}(X_1) < \infty$ . Then

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - \mu}{\sigma\sqrt{n}} \leq x\right) = \Phi(x),$$

where  $\Phi(x)$  is the CDF of a standard normal distributed r.v. cf. Example 2.1.5.

An explicit  $n^{-1/2}$  convergence rate for the distributional convergence in the CLT was derived independently by Berry and Esseen in the 1940s through bounding multi-convoluted characteristic functions, cf. [4, 14].

**Theorem 2.2.8 (Berry-Esseen)** Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with (as before)  $\mu = E[X_1]$  and  $\sigma^2 = \text{Var}(X_1)$  and  $\rho = E[|X_1 - \mu|^3] < \infty$ . Then

$$\sup_{x \in \mathbb{R}} |F_n(x) - \Phi(x)| \leq C \frac{\rho}{\sqrt{n}\sigma^3}$$

where

$$F_n(x) := P\left(\frac{S_n}{\sigma\sqrt{n}} \leq x\right).$$

For multi-dimensional and other extensions of the above presented convergence results, see [12, 3].

## Chapter 3

# Monte Carlo Methods

*The infinite we shall do right away. The finite may take a little longer.*

—Stan Ulam

Monte Carlo methods are a class of algorithms relying on repeated sampling of r.v. to compute the quantities of interest. A simple example would be the sample average

$$\bar{g}_M := \sum_{i=1}^M \frac{g(X_i)}{M}, \quad (3.0.1)$$

for some sequence of i.i.d. r.v.  $X_1, X_2, \dots$  in  $\mathbb{R}^n$ ,  $g : \mathbb{R}^n \rightarrow \mathbb{R}$  and the number of samples  $M$  predetermined or increased until the variance of  $\bar{g}_M$  is sufficiently small. MC methods are generally easy to implement, even for seemingly complicated problems, and their convergence rate is  $\mathcal{O}(M^{-1/2})$ , independent of the dimension of the problem. MC methods have become popular among statisticians and scientists in applied fields which apply the methods in estimations such as medicine efficiency, election predictions, and finance.

The field of MC methods grows by the day, so a clear cut definition of the term and summary of the subject becomes increasingly difficult. Here, we restrict ourselves to an informal presentation of the three central themes for MC methods through examples: how many samples  $M$  are needed to estimate the quantity of interest reliably, what variance reduction is, and the dimension independent convergence rate of MC methods.

### 3.1 Examples and Properties

When approximating a quantity of interest by the MC estimate (3.0.1), a central performance question is how many samples  $M$  are needed to make the approximation sufficiently accurate. To construct a reliable MC algorithm, it is important to choose  $M$  sufficiently large, but to make the algorithm efficient, it is important that it does not use more samples than needed to meet the accuracy demand;  $M$  should also be low as possible. In this section we will demonstrate that the MC convergence rate is  $\sigma M^{-1/2}$ , and use this fact to construct a reasonable choice for  $M$  in the setting of sampling Bernoulli r.v.

The vote of a citizen participating in a two party election might be modeled as a Bernoulli r.v. The American presidential election could be assigned the r.v.

$$X(\omega) = \begin{cases} 1, & \text{if } \omega = \{\text{Republican}\} \\ 0, & \text{if } \omega = \{\text{Democrat}\}. \end{cases} \quad (3.1.1)$$

The quantity of interest is  $\mu = P(X = 1)$ ; if  $\mu < 0.5$ , the Democrats win, and if  $\mu > 0.5$ , the Republicans win. Suppose we seek to predict the outcome of the election  $\mu$  through sampling i.i.d. voters  $X_1, X_2, \dots$  by the sample average

$$\bar{X}_M := \sum_{i=1}^M \frac{X_i}{M},$$

with the accuracy-confidence constraint

$$P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta, \quad (3.1.2)$$

where we call TOL is the accuracy and  $\delta$  the confidence. A reasonable question to raise is how large does  $M$  have to be—how many voters must we sample—to achieve (3.1.2). For large  $M$ , the CLT motivates the approximation

$$P(|\bar{X}_M - \mu| > \text{TOL}) = P\left(\frac{|\sum_{i=1}^M X_i - \mu|}{\sqrt{M}\sigma} > \frac{\sqrt{M}\text{TOL}}{\sigma}\right) \approx 2\left(1 - \Phi\left(\frac{\sqrt{M}\text{TOL}}{\sigma}\right)\right).$$

implies that to ensure that

$$2\left(1 - \Phi\left(\frac{\sqrt{M}\text{TOL}}{\sigma}\right)\right) \leq \delta,$$

we should use

$$M = \frac{\sigma^2 (\Phi^{-1}(1 - \delta/2))^2}{\text{TOL}^2} \quad (3.1.3)$$

vote samples. In general, the variance  $\sigma^2$  is unknown, making the above expression for  $M$  incomplete, but for Bernoulli r.v.  $\sigma^2 = \mu(1 - \mu) \leq 1/4$ , so we may conservatively choose  $M$  according to

$$M = \frac{(\Phi^{-1}(1 - \delta/2))^2}{4\text{TOL}^2}.$$

So, if for example given the demands of accuracy  $\text{TOL} = 0.02$  and confidence  $\delta = 0.05$ , we would, according to our derivations, have to sample at least  $M = 2401$  i.i.d. voters.

The MC convergence rate  $\sigma M^{-1/2}$  in the norm  $\sqrt{\mathbb{E}[|\bar{X}_M - \mu|^2]}$  follows from noting that for i.i.d. samples

$$\mathbb{E}\left[\frac{|\sum_{i=1}^M X_i - \mu|^2}{M^2}\right] = \frac{\mathbb{E}[|X - \mu|^2]}{M} = \frac{\sigma^2}{M}. \quad (3.1.4)$$

In euclidean norm, the convergence rate is often represented on the r.v. form

$$\frac{\bar{X}_M - \mu}{\sqrt{M}} \sim \sigma\chi,$$

where  $\chi$  denotes a standard normal r.v. and one should keep in mind that this is an asymptotic representation.



### Dimension Independent Convergence Rate

To illustrate the dimension independent convergence rate of MC methods we consider an example of multivariate integration

$$I = \int_A g(x) dx \quad (3.1.5)$$

where  $A$  is a compact subset of  $\mathbb{R}^n$  and  $g : A \rightarrow \mathbb{R}$ . To approximate the value of  $I$  by the MC method we generate i.i.d. r.v.  $X_1, X_2, \dots$  uniformly distributed in  $A$  and set

$$\bar{I}_M = \sum_{i=1}^M \frac{g(X_i)}{M}. \quad (3.1.6)$$

By construction we have that  $\mathbb{E}[\bar{I}_M] = I$  and by the reasoning of (3.1.4),

$$\sqrt{\mathbb{E}[|\bar{I}_M - I|^2]} = \sqrt{\text{Var}(g(X))} M^{-1/2},$$

equaling the expected convergence rate. In comparison, a cubature<sup>1</sup> approximation of (3.1.5) of order  $k$  will have the convergence rate  $\mathcal{O}(M^{-k/d})$  if we assume it resolves each dimension using  $M^{1/d}$  grid points. So when  $k/d < 1/2$ , the MC method will asymptotically outperform the cubature method.

To visualize the dimensional independent convergence we consider the problem of computing the volume of the unit ball in  $\mathbb{R}^6$ . This problem is equivalent to computing the integral (3.1.5) with

$$g(x) = \begin{cases} 1, & \text{if } |x| \leq 1, \\ 0, & \text{otherwise,} \end{cases}$$

and  $A = [-1, 1]^n$ . Figure 3.1 shows the performance of the MC algorithm compared to the first order cubature method

$$\tilde{I}_M := \sum_{i,j,k,l,m,n=1}^{\lceil M^{1/6} \rceil} g(x_i, x_j, x_k, x_l, x_m, x_n) h^6, \quad (3.1.7)$$

where  $\lceil x \rceil := \min\{n \in \mathbb{Z} \mid n \geq x\}$ ,  $h = 2/\lceil M^{1/6} \rceil$ , and  $x_i = -1 + (i - 1/2)h$  for  $i = 1, 2, \dots, \lceil M^{1/6} \rceil$ .

### 3.2 Variance Reduction

In Section 3.1, we observed that the convergence rate for the MC method is  $\sigma M^{-1/2}$  and that the number of samples needed to meet the accuracy-confidence constraint (3.1.2) was (approximately) given by

$$M = \frac{\sigma^2 (\Phi^{-1}(1 - \delta/2))^2}{\text{TOL}^2}.$$

Generally, the cost of computing an MC estimate is

$$\mathcal{O}(M) = \mathcal{O}\left(\frac{\sigma^2 (\Phi^{-1}(1 - \delta/2))^2}{\text{TOL}^2}\right),$$

---

<sup>1</sup>Cubature is the name for numerical integration in dimensions higher than 1.

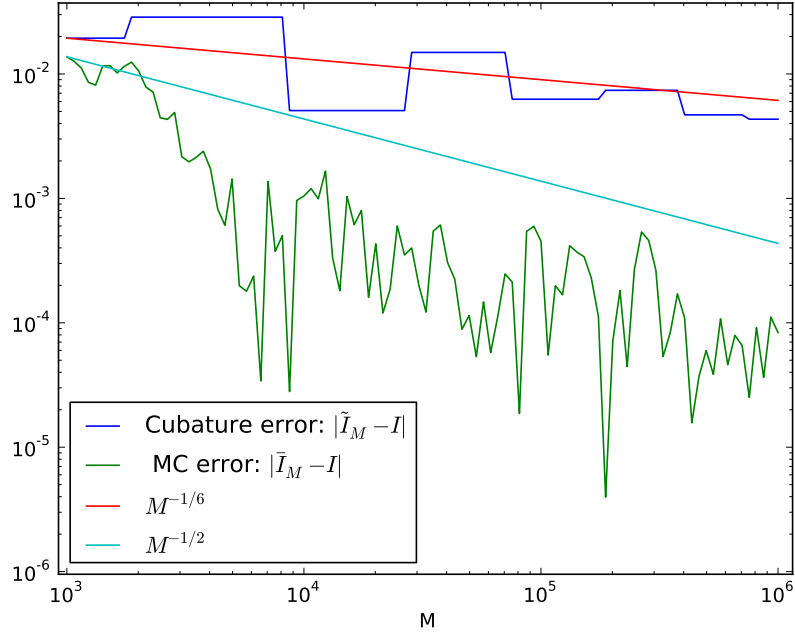


Figure 3.1: Comparison of convergence rates for approximating the volume of the unit ball in  $\mathbb{R}^6$ ,  $I = \pi^3/6$ , when using either the MC algorithm (3.1.6) or the cubature (3.1.7). Due to the high dimensionality of the problem, the MC algorithm outperforms the first order cubature method. The observed convergence rates for the experiment fit quite well with theoretically predicted rates.

so if it somehow is possible to reduce the variance  $\sigma^2$  of the samples used in the estimate, the computational cost of the MC estimate will also be reduced. Variance reduction is a vast subject for MC methods motivated from the goal of reducing computational cost. Here, we will restrict ourselves to describing the variance reduction technique named control variates which we will return to in Chapter 4. For a nice overview of other interesting variance techniques such as antithetic variables and importance sampling, we refer to [7].

### Control Variates

Let  $X$  be a given r.v. which we seek to approximate the expected value  $E[X]$  by a Monte Carlo method, say

$$\bar{X}_M = \sum_{i=1}^M \frac{X_i}{M}.$$

Suppose that  $Y$  is another r.v. for which the expected value  $E[Y]$  is known. The constructed random variable  $Z = X + E[Y] - Y$  will then have the same expected value as  $X$ , so  $E[X]$  can also be approximated by applying the Monte Carlo method on the r.v.  $Z$ , say

$$\bar{Z}_M = \sum_{i=1}^M \frac{Z_i}{M}.$$

If the variance is reduced,  $\text{Var}(Z) = \text{Var}(X - Y) < \text{Var}(X)$ , we reason from the introduction of this section that from a cost perspective, it is more efficient to approximate  $E[X]$  by sampling  $Z_i$  r.v. than by sampling  $X_i$  r.v.

### 3.3 Sequential Stopping Rules

Given no prior distributional information of samples used in an MC estimate, determining the number of samples  $M$  required to meet an accuracy-confidence constraint is very difficult; some distributional properties of the samples, such as their variance, generally has to be inferred to choose  $M$  sensibly. However daunting it might seem to perform reliable MC estimates in settings with no prior distributional sample information, such settings are frequently encountered and therefore they deserve some attention.

A sequential stopping rule MC algorithm progressively gathers distributional information on the samples through sampling higher moments and iteratively increasing the number of samples used in the estimate until a stop criterion is met. For example, Algorithm 1 represents a sample variance based stopping rule.

---

#### Algorithm 1 Sample Variance Based Stopping Rule

---

**Input:** Initial number of samples  $M_0$ , accuracy TOL, confidence  $1 - \delta$ , and the CDF of the standard normal distributed r.v.  $\Phi(x)$ .

**Output:**  $\bar{X}_M$ .

Set  $k = 0$ , generate  $M_k$  samples  $\{X_i\}_{i=1}^{M_k}$  and compute the sample variance

$$\bar{\sigma}_{M_k}^2 := \frac{1}{M_k - 1} \sum_{i=1}^{M_k} (X_i - \bar{X}_{M_k})^2. \quad (3.3.1)$$

**while**  $2(1 - \Phi(\sqrt{M_k} \text{TOL} / \bar{\sigma}_{M_k})) > \delta$  **do**

    Set  $k = k + 1$  and  $M_k = 2M_{k-1}$ .

    Generate a batch of  $M_k$  i.i.d. samples  $\{X_i\}_{i=1}^{M_k}$ .

    Compute the sample variance  $\bar{\sigma}_{M_k}^2$  as given in (3.3.1).

**end while**

Set  $M = M_k$ , generate samples  $\{X_i\}_{i=1}^M$  and compute the output sample mean  $\bar{X}_M$ .

---

Chow and Robbins proved that for r.v. with bounded second moments, sample variance based stopping rules are asymptotically reliable, cf. [9]. That is, supposing the generalized sample average estimator  $\bar{\sigma}_M^2$  used is positive and asymptotically consistent in the sense  $\bar{\sigma}_M \rightarrow \sigma$  almost surely as  $\text{TOL} \rightarrow 0$ , then for any confidence  $\delta \in (0, 1)$ , the stopping criterion

$$\frac{\bar{\sigma}_M^2}{M} \leq \frac{\text{TOL}^2}{(\Phi^{-1}(1 - \delta/2))^2} \quad (3.3.2)$$

implies the asymptotic fulfillment of the accuracy-confidence constraint

$$\lim_{\text{TOL} \rightarrow 0} P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta.$$

For practical purposes however, it is the reliability of sequential stopping rules for a fixed accuracy-confidence combination  $\text{TOL}, \delta \in (0, 1)$  which is of interest; the non-asymptotic regime. To the best of our knowledge, this is largely an open problem. For a fixed number of

samples  $M$ , an upper bound of the convergence rate in the non-asymptotic regime is given by the Berry-Esseen Theorem 2.2.8 which presents the rate  $\mathcal{O}(M^{-1/2})$ . But, the leading order constant of the Berry-Esseen bound contains the factor

$$\frac{\mathbb{E}[|X - \mu|^3]}{\sigma^2},$$

which in the settings we are confined to would have to be sampled/estimated, giving a bound estimate instead of an explicit bound.

In Paper III of this thesis, we study stopping rules' performance in the non-asymptotic regime. There we show by examples that for certain heavy-tailed r.v. and fixed TOL,  $\delta \in (0, 1)$ , the sample variance based stopping rule presented in Algorithm 1 fails to meet the accuracy-confidence constraint (3.1.2), and we construct a more reliable higher moments based stopping rule.

## Chapter 4

# Adaptive Weak Approximations for Stochastic Differential Equations

*It is not the strongest of the species that survives, nor the most intelligent that survives. It is the one that is the most adaptable to change.*

—Charles Darwin

### 4.1 Stochastic Differential Equations (SDE)

SDE is an extension of ordinary differential equations that seeks to incorporate uncertainty into a differential equation:

$$\begin{aligned}dX &= a(t, X)dt + \text{“noise”}, \\ X_0 &= x_0,\end{aligned}\tag{4.1.1}$$

In general, the noise in the differential equation will be motivated from physical considerations for the model. If the “noise” is a Wiener process, the differential equation (4.1.1) becomes an SDE which (on Itô form) may be represented as follows

$$\begin{aligned}dX &= a(t, X)dt + b(t, X)dW_t, \\ X_0 &= x_0.\end{aligned}\tag{4.1.2}$$

Here, the  $dW_t$  term is called the Wiener process increment.

#### The Wiener Process

The history of the Wiener process can be traced back to 1828 when the British botanist Robert Brown observed through a microscope that pollen grains suspended in water performed an irregular motion. The physical process for the irregular motion was, in honor of its discoverer, given the name Brownian motion. The mathematician Norbert Wiener later proposed a mathematical model for standardized Brownian motion which was to become known as the Wiener process.

The *Wiener process* is a Gaussian process<sup>1</sup>  $W : [0, T] \times \Omega \rightarrow \mathbb{R}$  which is characterized by the properties

---

<sup>1</sup>For more on Gaussian processes, see Chapter 5.

1.  $W_0 = 0$ ,
2. The function  $t \rightarrow W_t$  is almost surely continuous.
3.  $W_t$  has independent increments with  $W_t - W_s \sim \mathcal{N}(0, |t - s|)$ .

As a consequence of property 3.

$$\mathbb{E}[W_t W_s] = \min(t, s),$$

and, supposing  $t_1 \leq t_2 \leq t_3 \leq t_4$ ,

$$\mathbb{E}[(W_{t_2} - W_{t_1})(W_{t_4} - W_{t_3})] = \mathbb{E}[W_{t_2} - W_{t_1}] \mathbb{E}[W_{t_4} - W_{t_3}] = 0.$$

The Wiener process  $W_t(\omega)$  is a function of two variables. For fixed  $\omega \in \Omega$ , we call  $W(\omega) : [0, T] \rightarrow \mathbb{R}$  a *realization* or *sample path* of the Wiener process, cf. Figure 4.1. For each  $t \in \mathbb{R}_+$  fixed,  $W_t = W_t(\cdot)$  is a r.v. with the distribution (property 3)  $W_t \in \mathcal{N}(0, t)$ , that is, for any Borel set  $B \subset \mathbb{R}$ ,

$$P(W_t \in B) = \int_B \frac{\exp(-\frac{x^2}{2t})}{\sqrt{2\pi t}} dx.$$

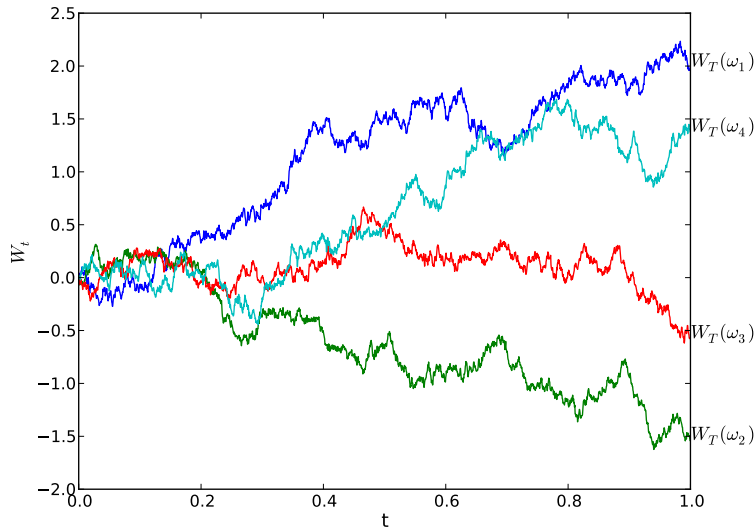


Figure 4.1: Four independent Wiener process realizations  $W_t(\omega_i)$ .

### Itô Integrals

According to the theory for ordinary differential equations, we expect solutions of the SDE (4.1.2) to be representable on integral form

$$X_t = X_0 + \int_0^t a(s, X) ds + \int_0^t b(s, X) dW_s. \quad (4.1.3)$$

However, what that is meant by the rightmost integral integrating over Wiener increments  $dW_s = W_{s+ds} - W_s$ , is not clear. Let us investigate this issue by considering integrals on the form

$$I(f) = \int_0^T f(t, \omega) dW_t.$$

The theory of Riemann integrals for deterministic functions leads us to expect that the rightmost integral could be considered as the limit sum of infinitesimal contributions

$$\int_0^t f(t, \omega) dW_s = \lim_{\Delta t \rightarrow 0} \sum_i f(t_i^*, \omega) \Delta W_i \quad (4.1.4)$$

with  $\Delta W_i := W_{t_{i+1}} - W_{t_i}$  and  $t_i^*$  any point in  $[t_i, t_{i+1}]$ . But, an example with  $f(t, \omega) := W_t(\omega)$  will illustrate that the theory for deterministic integrals is not directly extendable to integrals involving stochastic processes: Choosing the leftmost point of each interval as integration points  $t_i^* = t_i$  yields

$$I_-(f) := \lim_{\Delta t \rightarrow 0} \sum_i W_{t_i} \Delta W_i,$$

and, if instead choosing the rightmost point of each interval as integration points,  $t_i^* = t_{i+1}$ , we get

$$I_+(f) := \lim_{\Delta t \rightarrow 0} \sum_i W_{t_{i+1}} \Delta W_i.$$

From property 3. of the Wiener process, we may compute that while  $E[I_-(f)] = 0$ ,  $E[I_+(f)] = t$  on the other hand. This illustrates that for SDE it indeed does matter which point  $t_j^* \in [t_i, t_{i+1}]$  of each infinitesimal interval is used in the integral evaluation (4.1.4).

For SDE, different choices of  $t_j^*$  have given rise to different stochastic integrals, but in this thesis we will only consider *Itô integrals*. Itô integrals use the integrand points  $t_j^* = t_j$ , and is thus defined by

$$\int_0^t f(s, \omega) dW_s := \lim_{\Delta t \rightarrow 0} \sum_i f(t_i, \omega) \Delta W_i. \quad (4.1.5)$$

Itô integrals may be considered as the limiting form of the Forward Euler integrating method for SDE. In correspondence with Itô integrals, we call the SDE when written on the form (4.1.2) an Itô SDE. Note further that as an implication of having Wiener sample paths, the representation

$$X_t(\omega) = X_0 + \int_0^t a(s, X_s(\omega)) ds + \int_0^t b(s, X_s(\omega)) dW_s(\omega),$$

implies that for each fixed  $\omega$ ,  $t \rightarrow X_t(\omega)$  is a sample path solution of the SDE, and for each fixed  $t$ ,  $\omega \rightarrow X_t(\omega)$  is a r.v. See Figure 4.2 for an illustration of Itô SDE sample path solutions.

There are two types of solutions for Itô SDE; weak and strong solutions. In essence, one might say that weak solutions are sample path solutions which are “unique” in distributional sense, while strong solutions are unique in sample path sense according to the norm

$$\|X\|_{L^2(P; [0, T])} = E \left[ \int_0^T |X_t|^2 dt \right]^{1/2}.$$

See Section 4.2 for more on solution classification.

### Itô Calculus

*Itô calculus* is a handy tool for analyzing Itô SDE. Let  $X_t$  be a solution of (4.1.2) and  $g(t, X_t)$  a sufficiently smooth function. Then the Itô “laws”

$$dW_s^2 = dt, \quad dt \cdot dW_s = 0, \quad dt^2 = 0, \dots \quad (4.1.6)$$

yields following truncated Itô-Taylor expansion of  $dg$

$$\begin{aligned} dg(t, X_t) &= \partial_t g(t, X_t)dt + \partial_X g(t, X_t)dX + \frac{1}{2}\partial_{XX}g(t, X_t)dX^2 \\ &= \left( \partial_t g(t, X_t) + a(t, X_t)\partial_X g(t, X_t) + \frac{(b(t, X_t))^2}{2}\partial_{XX}g(t, X_t) \right) dt \\ &\quad + b(t, X_t)\partial_X g(t, X_t)dW_t. \end{aligned} \quad (4.1.7)$$

Here we used (4.1.2) and the Itô “laws” (4.1.6) to derive that  $dX^2 = b^2 dt + o(dt)$ .

### Solving an SDE by Using Itô Calculus

For an application of Itô calculus, let us consider the following problem of geometric Brownian motion

$$dX_t = rX_t dt + vX_t dW_t, \quad r, v \in \mathbb{R}_+. \quad (4.1.8)$$

This problem might be interpreted as model for stock price evolution with  $r$  representing a constant interest rate, and  $v$  the constant volatility rate, i.e., the stochastic evolution of the stock price. To solve this SDE, we assume  $X > 0$ , and note that then (4.1.8) is equivalent to

$$\frac{dX_t}{X_t} = rdt + v dW_t.$$

Considering the function  $g(X_t) = \log(X_t)$ , we see by (4.1.7) with  $a(X_t) = rX_t$  and  $b(X_t) = vX_t$  that

$$d \log(X_t) = \left( r - \frac{v^2}{2} \right) dt + v dW_t \implies X_t = X_0 \exp \left( \left( r - \frac{v^2}{2} \right) t + v W_t \right).$$

Figure 4.2 holds examples of sample path solutions of the SDE (4.1.8).

## 4.2 The Euler Method for Itô SDE

We recall that a solution of an Itô SDE

$$\begin{aligned} dX &= a(t, X)dt + b(t, X)dW_t, \\ X_0 &= x_0. \end{aligned}$$

can be written on the Itô integral form

$$X_t = X_0 + \int_0^t a(s, X_s) ds + \int_0^t b(s, X_s) dW_s,$$

where the rightmost integral is the limit sum of a forward Euler integration

$$\int_0^t b(s, X_s) dW_s := \lim_{\Delta t \rightarrow 0} \sum_i b(t_i, X_{t_i}) \Delta W_i.$$



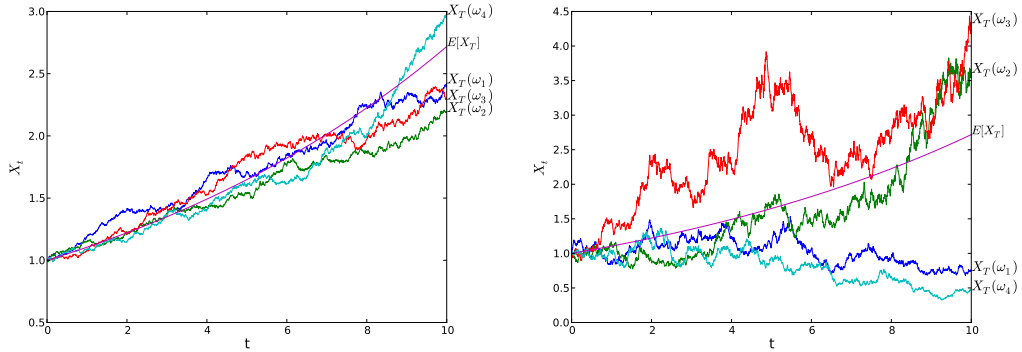


Figure 4.2: **Left plot:** Four independent sample path solutions of the SDE (4.1.8) when  $X_0 = 1$ ,  $r = 0.1$ , and  $v = 0.05$ . **Right plot:** Four independent sample path solutions of the the SDE (4.1.8) when  $X_0 = 1$ ,  $r = 0.1$ , and  $v = 0.25$ .

For numerical solutions, the (forward) Euler method generates realizations  $\bar{X}$  of the Itô SDE by the scheme

$$\begin{aligned} \bar{X}_{t_{n+1}} &= \bar{X}_{t_n} + a(t_n, \bar{X}_{t_n})\Delta t + b(t_n, \bar{X}_{t_n})\Delta W_n, \\ \bar{X}_0 &= x_0, \end{aligned} \quad (4.2.1)$$

with uniform step size  $\Delta t = T/N$ ,  $t_n = n\Delta t$ , and Wiener increments  $\Delta W_n = W(t_{n+1}) - W(t_n) \sim N(0, \Delta t)$ .

The rate of convergence of  $\bar{X}_T$  to the real solution  $X_T$  depends on the measure. The *strong convergence* is given by  $\mathbb{E}[|\bar{X}_T - X_T|]$  and is a measure for the pathwise convergence of each numerical realization  $\bar{X}_T(\omega)$  to its corresponding real solution realization  $X_T(\omega)$ . For many problems, however, pathwise convergence might be more than you seek. Instead, you might be interested in some distributional quantity of interest expressible on the form  $\mathbb{E}[g(X_t)]$  for some smooth function  $g$ . If so, the convergence criterion of interest is the *weak convergence*<sup>2</sup>

$$\sup_{g \in C_c^\infty(\mathbb{R})} |\mathbb{E}[g(X_T)] - \mathbb{E}[g(\bar{X}_T)]|,$$

where  $C_c^\infty(\mathbb{R})$  denotes the set of smooth compactly supported functions. For the Euler method on uniform grids with step size  $\Delta t$ , the strong convergence rate is  $\mathcal{O}(\Delta t^{1/2})$ , cf. [22, 8], and the weak convergence was shown by Talay and Tubaro to be  $\mathcal{O}(\Delta t)$ , cf. [31], see also cf. [8].

**Remark 4.2.1** *In this thesis we consider SDE solved numerically using the Euler method, but we remark that through truncation of Itô-Taylor expansions, one may develop higher order numerical schemes for SDE. One example is the Milstein scheme,*

$$\bar{X}_{n+1} = \bar{X}_n + a_n \Delta t + b_n \Delta W_n + \frac{b_n \partial_X b_n}{2} (\Delta W_n^2 - \Delta t),$$

*which has first order strong convergence. As for deterministic problems, however, higher order methods require stronger regularity assumptions and are generally more difficult to*

<sup>2</sup>Weak convergence is equivalent to convergence in distribution, as defined in Chapter 2.

expand to adaptive methods in higher dimensions. For more on higher order methods, we refer to [22].

### An adaptive Euler Method for Weak Solutions

For problems where either the drift or diffusion in the SDE lacks regularity, numerical solutions using adaptive time stepping may improve convergence rates, or even be your only option for obtaining reliable numerical solutions. In this subsection, we will describe the adaptive Euler method for weak solutions of SDE which was developed by Szepessy et al. in [30] in the beginning of the 2000s.

#### Objective

Consider the problem of generating solutions of the initial value SDE

$$\begin{aligned} dX &= a(t, X)dt + b(t, X)dW_t, \quad 0 \leq t \leq T, \\ X_0 &= x_0, \quad x_0 \in \mathbb{R}. \end{aligned} \quad (4.2.2)$$

with an adaptive Euler method which fulfills the weak accuracy constraint

$$|\mathbb{E}[g(X_T) - g(\bar{X}_T)]| \leq \text{TOL}_T, \quad (4.2.3)$$

where  $X_t$  and  $\bar{X}_t$  represents the explicit and numerical solution of (4.3.1), respectively, and  $\text{TOL}_T > 0$  is a given tolerance constraint. In this setting, the goal of adaptivity is to minimize  $N$ , the number of adaptive time steps needed in the grid  $0 = t_0 < t_1 < \dots < t_N = T$  such that your adaptive Euler solution  $\bar{X}_t$  fulfill (4.2.3). As expected, the adaptive Euler method is given by

$$\bar{X}_{t_{n+1}} = \bar{X}_{t_n} + a(t_n, \bar{X}_{t_n})\Delta t_n + b(t_n, \bar{X}_{t_n})\Delta W_n, \quad (4.2.4)$$

where  $\Delta t_n(\omega) = (t_{n+1} - t_n)(\omega)$  denote the  $n$ -th time step of the sample path solution  $\omega$ .

#### Grid Refinement

The adaptive method by Szepessy et al. expresses the weak discretization error by the error expansion

$$\mathbb{E}[g(X_T) - g(\bar{X}_T)] \simeq \sum_{n=0}^N \mathbb{E}[\rho(t_n, \cdot)\Delta t_n^2(\cdot)] + \text{h.o.t.} \quad (4.2.5)$$

Here,  $\rho(t_n, \omega)$  represents the error density for a given sample path solution  $\omega$  at the time  $t_n$  and the error indicators  $\rho(t_n, \omega)\Delta t_n^2(\omega)$  then represents, up to higher order terms, the error contribution from the  $n$ -th time step of the sample path  $\omega$ . The error indicators  $\rho(t_n, \omega)\Delta t_n^2(\omega)$  provide information for further refinement of the time grid. Error control at a low computational cost is obtained through the condition

$$\rho(t_n, \omega)\Delta t_n^2(\omega) \leq C \frac{\text{TOL}_T}{\mathbb{E}[N]}, \quad \forall n = 0, 1, 2, \dots, N(\omega). \quad (4.2.6)$$

Given an initial time grid, time increments, and a Wiener process numerical sample path which we represent by  $t.(\omega, 0)$ ,  $\Delta t.(\omega, 0)$ , and  $W_{t.(\omega, 0)}$ , respectively, the following adaptive mesh refinement algorithm is used by Szepessy et al. to ensure that the condition (4.2.6) is met.

1. Set  $\ell = 0$ .
2. Compute the error density  $\rho(t, \omega)$  for the numerical path represented by  $t.(\omega, \ell)$ ,  $\Delta t.(\omega, \ell)$ , and  $W_{t.(\omega, \ell)}$ , cf. (4.2.9).
3. If condition (4.2.6) is fulfilled, then stop; otherwise
4. Refine the grid  $t.(\omega, \ell)$  at all points where (4.2.6) is not fulfilled by halving the steps:

$$\Delta t_{\bar{n}}(\omega, \ell + 1) = \frac{\Delta t_n(\omega, \ell)}{2} \quad \text{and} \quad \Delta t_{\bar{n}+1}(\omega, \ell + 1) = \frac{\Delta t_n(\omega, \ell)}{2}, \quad (4.2.7)$$

where  $\bar{n} \geq n$  is defined as the natural number such that  $t_{\bar{n}}(\omega, \ell + 1) = t_n(\omega, \ell)$ . When adding a point  $t_{\bar{n}}(\omega, \ell + 1)$  to your grid by the refinement (4.2.7), add a corresponding Wiener process sample path point by Brownian bridges:

$$W_{t_{\bar{n}+1}(\omega, \ell + 1)} = \frac{W_{t_n(\omega, \ell)} + W_{t_{n+1}(\omega, \ell)}}{2} + \frac{\sqrt{\Delta t_n(\omega, \ell)}}{2} \xi, \quad \text{where} \quad \xi \sim \mathcal{N}(0, 1).$$

5. Set  $\ell = \ell + 1$  and return to 2.

### The Error Density

For notational convenience, let us write the Euler method scheme (4.2.4) as follows

$$\bar{X}_{t_{n+1}} = c(t_n, \bar{X}_{t_n}), \quad \text{where} \quad c(t_n, x) := x + a(t_n, x)\Delta t_n + b(t_n, \bar{X}_{t_n})\Delta W_n. \quad (4.2.8)$$

The error density derived in Theorem 2.2 of [30] can then be represented by

$$\begin{aligned} \rho(t_n, \cdot) = & \frac{1}{2\Delta t_n} \left( (a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})) \phi(t_{n+1}) \right. \\ & \left. + (b^2(t_{n+1}, \bar{X}_{t_{n+1}}) - b^2(t_n, \bar{X}_{t_n})) \partial_{\bar{X}_{t_{n+1}}} \phi(t_{n+1}) \right), \end{aligned} \quad (4.2.9)$$

where  $\phi$  is the solution of the discrete dual backward problem

$$\begin{aligned} \phi(t_n) &= \partial_x c(t_{n+1}, \bar{X}_{t_{n+1}}) \phi(t_{n+1}), \quad t < T \\ \phi(T) &= g'(\bar{X}_T), \end{aligned} \quad (4.2.10)$$

and  $\partial_{\bar{X}_{t_n}} \phi(t_n)$  is obtained by a longer scheme by linearising the forward problem (4.1.2), cf. [30, Theorem 2.2].

To derive the error density (4.2.9), we extend the numerical solutions of the Euler method (4.2.4) to  $t \in [t_n, t_{n+1})$  by

$$\bar{X}_t = \bar{X}_{t_n} + \int_{t_n}^t \bar{a}(s; \bar{X}) ds + \int_{t_n}^t \bar{b}(s; \bar{X}) dW_s,$$

where  $\bar{a}$  and  $\bar{b}$  are piecewise constant approximations

$$\bar{a}(s; \bar{X}) = a(t_n, \bar{X}_{t_n}) \quad \text{and} \quad \bar{b}(s; \bar{X}) = b(t_n, \bar{X}_{t_n}) \quad \text{for} \quad s \in [t_n, t_{n+1}), \quad n = 0, 1, 2, \dots$$

The expected value  $\mathbb{E}[g(X_T)]$  with  $X_T$  solving the SDE (4.1.2) is related to a Partial Differential Equation (PDE): the utility function  $u(t, x) := \mathbb{E}[g(X_T)|X_t = x]$  solves the Kolmogorov backward equation

$$\begin{aligned} \partial_t u &= - \left( a \partial_x + \frac{b^2}{2} \partial_{xx} \right) u, & (t, x) \in [0, T) \times \mathbb{R} \\ u(T, x) &= g(x), \end{aligned} \quad (4.2.11)$$

cf. [22]. By Itô calculus and the relation (4.2.11), it follows that

$$\begin{aligned} du(t, \bar{X}_t) &= \left( \partial_t u(t, \bar{X}_t) + \bar{a}(t; \bar{X}) \partial_x u(t, \bar{X}_t) + \frac{\bar{b}^2(t; \bar{X})}{2} \partial_{xx} u(t, \bar{X}_t) \right) dt \\ &\quad + \bar{b}(t; \bar{X}) \partial_x u(t, \bar{X}_t) dW_t \\ &= \left( (\bar{a}(t; \bar{X}) - a(t, \bar{X}_t)) \partial_x u(t, \bar{X}_t) + \frac{\bar{b}^2(t; \bar{X}) - b^2(t, \bar{X}_t)}{2} \partial_{xx} u(t, \bar{X}_t) \right) dt \\ &\quad + \bar{b}(t; \bar{X}) \partial_x u(t, \bar{X}_t) dW_t. \end{aligned}$$

Recalling that  $\bar{X}_0 = X_0$ , we see that  $u(0, \bar{X}_0) = \mathbb{E}[g(X_T)]$ . Hence,

$$\begin{aligned} \mathbb{E}[g(X_T)] - g(\bar{X}_T) &= - \int_0^T du(t, \bar{X}_t) \\ &= \sum_n \int_{t_n}^{t_{n+1}} (a(t, \bar{X}_t) - a(t_n, \bar{X}_{t_n})) \partial_x u(t, \bar{X}_t) + \frac{b^2(t, \bar{X}_t) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{xx} u(t, \bar{X}_t) dt \\ &\quad - \sum_n \int_{t_n}^{t_{n+1}} \bar{b}(t_n, \bar{X}_{t_n}) \partial_x u(t, \bar{X}_t) dW_t. \end{aligned} \quad (4.2.12)$$

Taking the expected value of equation (4.2.12) yields

$$\begin{aligned} \mathbb{E}[g(X_T) - g(\bar{X}_T)] &= \sum_n \int_{t_n}^{t_{n+1}} \mathbb{E}[(a(t, \bar{X}_t) - a(t_n, \bar{X}_{t_n})) \partial_x u(t, \bar{X}_t)] dt \\ &\quad + \sum_n \int_{t_n}^{t_{n+1}} \mathbb{E} \left[ \frac{b^2(t, \bar{X}_t) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{xx} u(t, \bar{X}_t) \right] dt, \end{aligned} \quad (4.2.13)$$

where we used that the expected value of the integral with  $dW_t$  increments in (4.2.12) is zero, cf. [28]. Introducing the utility function for the numerical solution  $\bar{u}(t, x) := \mathbb{E}[g(\bar{X}_T)|\bar{X}_t = x]$ , the following approximations of the terms in (4.2.13) are valid:

$$\begin{aligned} &\int_{t_n}^{t_{n+1}} \mathbb{E} [a(t, \bar{X}_t) - (a(t_n, \bar{X}_{t_n})) \partial_x u(t, \bar{X}_t)] dt \\ &= \mathbb{E} [(a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})) \partial_x \bar{u}(t_{n+1}, \bar{X}_{t_{n+1}})] \frac{\Delta t_n}{2} + \text{h.o.t.} \end{aligned}$$

and

$$\begin{aligned} &\int_{t_n}^{t_{n+1}} \mathbb{E} \left[ \frac{b^2(t, \bar{X}_t) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{xx} u(t, \bar{X}_t) \right] dt \\ &= \mathbb{E} \left[ \frac{b^2(t_{n+1}, \bar{X}_{t_{n+1}}) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{xx} \bar{u}(t_{n+1}, \bar{X}_{t_{n+1}}) \right] \frac{\Delta t_n}{2} + \text{h.o.t.} \end{aligned}$$

Further, since

$$\partial_x \bar{u}(t, x) = \frac{\partial}{\partial x} \mathbb{E}[g(\bar{X}_T) | \bar{X}_t = x] = \mathbb{E}\left[g'(\bar{X}_T) \partial_{\bar{X}_t} \bar{X}_T \mid \bar{X}_t = x\right],$$

we see that

$$\partial_x \bar{u}(t_{n+1}, \bar{X}_{t_{n+1}}) = \mathbb{E}\left[g'(\bar{X}_T) \partial_{\bar{X}_{t_{n+1}}} \bar{X}_T \mid \mathcal{F}_{t_{n+1}}\right], \quad (4.2.14)$$

where  $\mathcal{F}_t$  denotes the  $\sigma$ -algebra generated by  $\{W_s\}_{s \in [0, t]}$ . By definition (4.2.10),

$$\begin{aligned} 0 &= \sum_{\ell=n}^{N-1} (\phi(t_\ell) - \partial_x c(\bar{X}_{t_\ell}, t) \phi(t_{\ell+1})) \partial_{\bar{X}_{t_n}} \bar{X}_{t_\ell} \\ &= \sum_{\ell=n}^{N-1} \phi(t_{\ell+1}) \underbrace{\left( \partial_{\bar{X}_{t_n}} \bar{X}_{t_{\ell+1}} - \partial_x c(\bar{X}_{t_\ell}, t) \partial_{\bar{X}_{t_n}} \bar{X}_{t_\ell} \right)}_{=0 \text{ cf. (4.2.8)}} - \phi(T) \frac{\partial \bar{X}_T}{\partial \bar{X}_{t_n}} + \phi(t_n) \frac{\partial \bar{X}_{t_n}}{\partial \bar{X}_{t_n}} \\ &= -\phi(T) \frac{\partial \bar{X}_T}{\partial \bar{X}_{t_n}} + \phi(t_n). \end{aligned} \quad (4.2.15)$$

Recalling that  $\phi(T) = g'(X_T)$ , we see that  $\phi(t_n) = g'(\bar{X}_T) \partial_{\bar{X}_{t_n}} \bar{X}_T$ . Since  $a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})$  is  $\mathcal{F}_{t_{n+1}}$  measurable and by (4.2.14), we derive that

$$\begin{aligned} &\mathbb{E}[(a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})) \partial_x \bar{u}(t_{n+1}, \bar{X}_{t_{n+1}})] \\ &= \mathbb{E}[(a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})) \mathbb{E}[\phi(t_{n+1}) | \mathcal{F}_{t_{n+1}}]] \\ &= \mathbb{E}[(a(t_{n+1}, \bar{X}_{t_{n+1}}) - a(t_n, \bar{X}_{t_n})) \phi(t_{n+1})], \end{aligned} \quad (4.2.16)$$

and by a similar argument it follows that

$$\begin{aligned} &\mathbb{E}\left[\frac{b^2(t_{n+1}, \bar{X}_{t_{n+1}}) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{xx} \bar{u}(t_{n+1}, \bar{X}_{t_{n+1}})\right] \\ &= \mathbb{E}\left[\frac{b^2(t_{n+1}, \bar{X}_{t_{n+1}}) - b^2(t_n, \bar{X}_{t_n})}{2} \partial_{\bar{X}_{t_{n+1}}} \phi(t_{n+1})\right]. \end{aligned} \quad (4.2.17)$$

We conclude the derivation of the density representation (4.2.9) by inserting (4.2.16) and (4.2.17) into equality (4.2.13).

**Remark 4.2.2** *By a longer analysis, Szepessy et al. refines the error density to the following representation*

$$\begin{aligned} \rho(t_n, \omega) &= \frac{1}{2} \left( \partial_t a + a \partial_x a + \frac{b^2 \partial_{xx} a}{2} \right) \phi(t_{n+1}) \\ &\quad + \left( \partial_t \frac{b^2}{2} + a \partial_x \frac{b^2}{2} + \frac{b^2}{2} \partial_{xx} \frac{b^2}{2} + b^2 \partial_x a \right) \partial_{\bar{X}_{t_{n+1}}} \phi(t_{n+1}) \\ &\quad + \frac{b^2}{2} \partial_{xx} b^2 \frac{\partial^2 \phi(t_{n+1})}{\partial \bar{X}_{t_{n+1}}^2}, \end{aligned}$$

cf. Theorem 3.3 of [30].

**Remark 4.2.3** *The use of dual functions is standard in both in optimal control theory and in adaptive grid algorithms for ordinary and partial differential equations, cf. [1, 13].*

### 4.3 Adaptive Weak Solution Approximation

For the Itô SDE

$$\begin{aligned} dX &= a(t, X)dt + b(t, X)dW_t, \quad 0 \leq t \leq T, \\ X_0 &= x, \end{aligned} \tag{4.3.1}$$

we consider the problem of minimizing the computational cost of approximating  $\mathbb{E}[g(X(T))]$  within a given tolerance  $\text{TOL} > 0$ . Problems of this kind arise, for example, in computing option prices in mathematical finance, cf. [21] and [17].

Approximations of weak solutions are typically obtained using Monte Carlo (MC) methods. For SDE problems, it is possible to reduce the variance of the MC estimate, and thus reduce the complexity, by generating solution realizations on grids of different step size. This variance reduction technique is called Multilevel Monte Carlo (MLMC) methods and is in some ways similar to multigrid methods for PDE problems [6]. Paper II of this thesis develops an adaptive time step MLMC algorithm. As a preparation for that paper, we will here give an outline of the uniform time step MLMC algorithm. But let us first consider the single level MC method.

#### Single Level Weak Approximation

The single level MC method generates  $M$  numerical realizations  $\bar{X}_T(\omega_i)$  by the uniform time step Euler method and approximates  $\mathbb{E}[g(X_T)]$  by the MC sample average

$$\mathcal{A}(g(\bar{X}_T); M) = \sum_{i=1}^M \frac{g(\bar{X}_T(\omega_i))}{M}. \tag{4.3.2}$$

According to the problem formulation, we seek to fulfill the weak error bound  $|\mathbb{E}[g(X_T)] - \mathcal{A}(g(\bar{X}_T); M)| \leq \text{TOL}$  at minimal computational cost. For error control, the approximation error is split into two parts

$$\begin{aligned} |\mathbb{E}[g(X_T)] - \mathcal{A}(g(\bar{X}_T); M)| &\leq |\mathbb{E}[g(X_T)] - \mathbb{E}[g(\bar{X}_T)]| \\ &\quad + |\mathbb{E}[g(\bar{X}_T)] - \mathcal{A}(g(\bar{X}_T); M)| =: \mathcal{E}_T + \mathcal{E}_S, \end{aligned}$$

where  $\mathcal{E}_T$  is the time discretization error and  $\mathcal{E}_S$  the statistical error. Under appropriate smoothness assumptions on the drift  $a$  and diffusion  $b$ , we mentioned in Section 4.2 that weak convergence rate, i.e., the time discretization error, for the Euler method fulfills  $\mathcal{E}_T = \mathcal{O}(\Delta t)$ . Further, by the CLT, the asymptotic convergence rate for the statistical error is  $\mathcal{E}_S = \mathcal{O}(M^{-1/2})$ . So to ensure that  $|\mathbb{E}[g(X(T))] - \mathcal{A}(g(\bar{X}_T); M)| = \mathcal{O}(\text{TOL})$ , one should choose  $\Delta t = \mathcal{O}(\text{TOL})$  and  $M = \mathcal{O}(\text{TOL}^{-2})$ . The computational cost thus becomes  $\mathcal{O}(\text{TOL}^{-3})$  for the single level MC method.

#### The Multilevel Monte Carlo Method

The MLMC method developed by Giles in [16] expanded the single level MC method by constructing Euler method realizations of the SDE (4.3.1) on hierarchies of uniform time grids, typically with the size relations

$$\Delta t^{(\ell)} = 2^{-\ell} \Delta t^{(0)}, \quad \ell = 0, 1, 2, \dots, L.$$

Let  $\bar{X}_t^{(\ell)}$  denote an Euler method realization on a grid with uniform step size  $\Delta t^{(\ell)}$ . Then the MLMC method approximates  $E[g(X(T))]$  by the telescoping sum

$$\mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0) = \sum_{i=1}^{M_0} \frac{g(\bar{X}_T^{(0)}(\omega_{i,0}))}{M_0} + \sum_{\ell=1}^L \sum_{i=1}^{M_\ell} \frac{g(\bar{X}_T^{(\ell)}(\omega_{i,\ell})) - g(\bar{X}_T^{(\ell-1)}(\omega_{i,\ell}))}{M_\ell}, \quad (4.3.3)$$

where  $M_0$  and  $M_\ell := 2^{-\ell} M_0$ ,  $\ell = 1, \dots, L$ , represents the number of samples generated at respective grid levels. On each level  $\ell$  in the above estimator, the realization pairs  $\bar{X}_T^{(\ell)}(\omega_{i,\ell})$  and  $\bar{X}_T^{(\ell-1)}(\omega_{i,\ell})$  are generated by the same Wiener process sample path  $W_t(\omega_{i,\ell})$ , but on different temporal grids with step size  $\Delta t^{(\ell)}$  and  $\Delta t^{(\ell-1)}$ , respectively. The values of a Wiener process sample path is first computed on the coarse grid  $t_n^{(\ell-1)} = n\Delta t^{(\ell-1)}$ , let us write  $W_{t_n^{(\ell-1)}}(\omega)$ , and thereafter, the values of the sample path on the finer grid  $t_n^{(\ell)} = n\Delta t^{(\ell)}$  is computed by Brownian bridges

$$W_{t_{2n}^{(\ell)}}(\omega) = W_{t_n^{(\ell-1)}}(\omega) + \frac{W_{t_n^{(\ell-1)}}(\omega) + W_{t_{n+1}^{(\ell-1)}}(\omega)}{2} + \frac{\sqrt{\Delta t^{(\ell-1)}}}{2} \xi, \quad \text{where } \xi \sim \mathcal{N}(0,1),$$

cf. Figure 4.3.

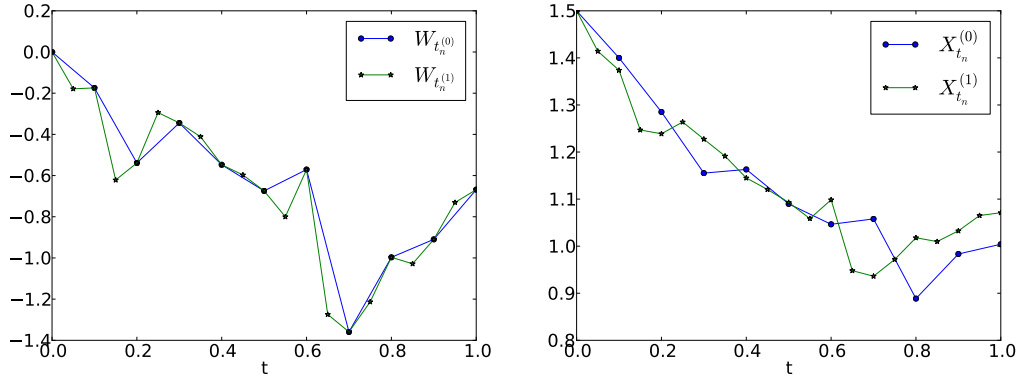


Figure 4.3: **Left plot:** A sample path  $W_t(\omega)$  plotted on the coarse grid  $W_{t_n^{(0)}}(\omega)$  (blue line) with  $\Delta t^{(0)} = 1/10$ , and its finer level pair generated by Brownian bridges,  $W_{t_n^{(1)}}(\omega)$  (green line) with  $\Delta t^{(1)} = \Delta t^{(0)}/2$ . **Right plot:** Euler method numerical solutions of the Ornstein-Uhlenbeck SDE problem  $dX_t = 2(1 - X_t)dt + 0.2dW_t$ ,  $X_0 = 3/2$ . for the Wiener process sample path plotted in the left plot.  $\bar{X}_{t_n}^{(0)}(\omega)$  (blue line) generated from using the Wiener increments from the path  $W_{t_n^{(0)}}(\omega)$  and  $\bar{X}_{t_n}^{(1)}(\omega)$  (green line) generated using Wiener increments from the path  $W_{t_n^{(1)}}(\omega)$

The MLMC method has the consistent estimator

$$E[\mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)] = E[g(\bar{X}_T^{(L)})],$$

and since sample paths are i.i.d.

$$E\left[\left(g(\bar{X}_T^{(\ell)}) - g(\bar{X}_T^{(\ell-1)})\right) \left(g(\bar{X}_T^{(m)}) - g(\bar{X}_T^{(m-1)})\right)\right] = 0, \quad \text{when } \ell \neq m.$$

The strong convergence  $\mathbb{E}\left[\left(g(\bar{X}_T^{(\ell)}) - g(X_T)\right)^2\right] = \mathcal{O}(\Delta t^{(\ell)})$ , cf. Section 4.2, then further implies that

$$\begin{aligned} \text{Var}(\mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)) &= \frac{\text{Var}\left(g\left(\bar{X}_T^{(0)}\right)\right)}{M_0} + \sum_{\ell=1}^L \frac{\text{Var}\left(g\left(\bar{X}_T^{(\ell)}\right) - g\left(\bar{X}_T^{(\ell-1)}\right)\right)}{M_\ell} \\ &= \mathcal{O}\left(M_0^{-1} + \sum_{\ell=1}^L 2^{-\ell} M_\ell^{-1}\right) = \mathcal{O}(LM_0^{-1}). \end{aligned} \quad (4.3.4)$$

Giles showed that by choosing  $L = \mathcal{O}(\log(\text{TOL}^{-1}))$  and  $M_0 = \mathcal{O}(L\text{TOL}^{-2})$ , one obtains

$$\left|\mathbb{E}[g(X_T)] - \mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)\right| \leq \mathcal{O}(\text{TOL})$$

at the computational cost

$$\mathcal{O}\left(\sum_{\ell=0}^L \frac{M_\ell}{\Delta t^{(\ell)}}\right) = \mathcal{O}\left((\log(\text{TOL}^{-1}))^2 \text{TOL}^{-2}\right).$$

This might be shown by splitting MLMC error into two contributions

$$\begin{aligned} \left|\mathbb{E}[g(X_T)] - \mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)\right| &\leq \left|\mathbb{E}\left[g(X_T) - g\left(\bar{X}_T^{(L)}\right)\right]\right| \\ &\quad + \left|\mathbb{E}\left[g\left(\bar{X}_T^{(L)}\right)\right] - \mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)\right| =: \mathcal{E}_T + \mathcal{E}_S. \end{aligned}$$

When  $L = \mathcal{O}(\log(\text{TOL}^{-1}))$  and  $M_0 = \mathcal{O}(L\text{TOL}^{-2})$ , the weak rate of convergence for the Euler method implies that  $\mathcal{E}_T = \mathcal{O}(\text{TOL})$ , and by the CLT and (4.3.4), one may derive that  $\mathcal{E}_S = \mathcal{O}(\text{TOL})$ .

With Giles' MLMC method, the computational cost is thus improved from  $\mathcal{O}(\text{TOL}^{-3})$  for the single level MC method to  $\mathcal{O}((\log(\text{TOL}^{-1}))^2 \text{TOL}^{-2})$ . The cost improvement is due to the multilevel variance reduction which is similar to control variates, cf. Section 3.2. That is, the multilevel average operator  $\mathcal{A}_{\text{MLC}}(g(\bar{X}_T); M_0)$  has lower variance than the single level average operator  $\mathcal{A}(g(\bar{X}_T); M)$ , and this implies that fewer samples have to be used to control the statistical error for the MLMC method than for the single level MC method.

**Remark 4.3.1** *The theory presented for 1-dimensional SDE problems in this chapter generalizes to the  $n$ -dimensional setting, cf. [8, 22, 16]. MC methods for weak approximations of SDE do in fact become more interesting in higher dimensions due to the dimension independent convergence rate of MC methods, cf. Section 3.1. Lower dimensional weak approximation problems are often more efficiently solved by solving the Kolmogorov backward partial differential equation relating to the utility function  $u(t, x) = \mathbb{E}[X_T | X_t = x]$ , cf. (4.2.11).*



## Chapter 5

# Wireless Channel Modeling

Wireless signal transmission is realized through electromagnetic radiation from transmitter to receiver, and the signal received is described by electromagnetic field at the receiver as a function of time. In principle, a numerical solution of Maxwell's equations with well resolved scattering boundary would yield the electrical field at the receiver and everywhere around it. But, for standard Modeling scenarios where the communication carrier frequency is of order  $10^9$ Hz and the wavelength consequently is of order centimeters, localization of the scattering boundary would have to be made to the order of centimeter accuracy to obtain an acceptable electrical field solution. Since this typically implies the need for resolving billions of boundary points, determining the electrical field through solving Maxwell's equations is considered too costly both from a boundary measurement and computational perspective.

Multipath Fading Channel (MFC) models give a more cost efficient, but less accurate model of the received signal than Maxwell's equations by superpositioning a large number of incoming signal wave paths. In Paper I of this thesis, we study an MFC model with noise introduced by scatterers flipping on and off. To prepare the reader for the contents that paper, we will here give a short introduction to MFC models and some more general signal theory concepts.

### 5.1 The Multipath Fading Channel

An MFC model approximates the output signal of a wireless channel by a superposition of a large number  $M$  of incoming signal wave paths, cf. Figure 5.1. With  $X_t$  denoting the baseband input signal, the time invariant MFC model has the baseband output representation

$$Z_{t,M} = \sum_{j=1}^M a(\alpha_j) e^{-i2\pi f_c \tau(\alpha_j, t)} X_{t-\tau(\alpha_j, t)}, \quad (5.1.1)$$

where  $f_c$  denotes the carrier frequency,  $\alpha_j$  the horizontal angle of arrival for the  $j^{\text{th}}$  wave path, and  $a(\alpha_j)$  and  $\tau(\alpha_j, t)$  its amplitude and time delay function, respectively.

The terminology baseband and passband refers to the frequency modulation of the signal in question. Let  $\mathcal{F}(\cdot)$  denote the Fourier transform. Then the *baseband input signal*  $X_t$  has the spectral representation  $\mathcal{F}(X_t)(f) \in C([-W, W])$ . The spectral representation is centered at origo, band-limited to  $[-W, W]$ , and it does not have to be symmetric w.r.t.  $f$ ; the baseband signal  $X_t$  may be complex-valued. Before transmission the baseband signal

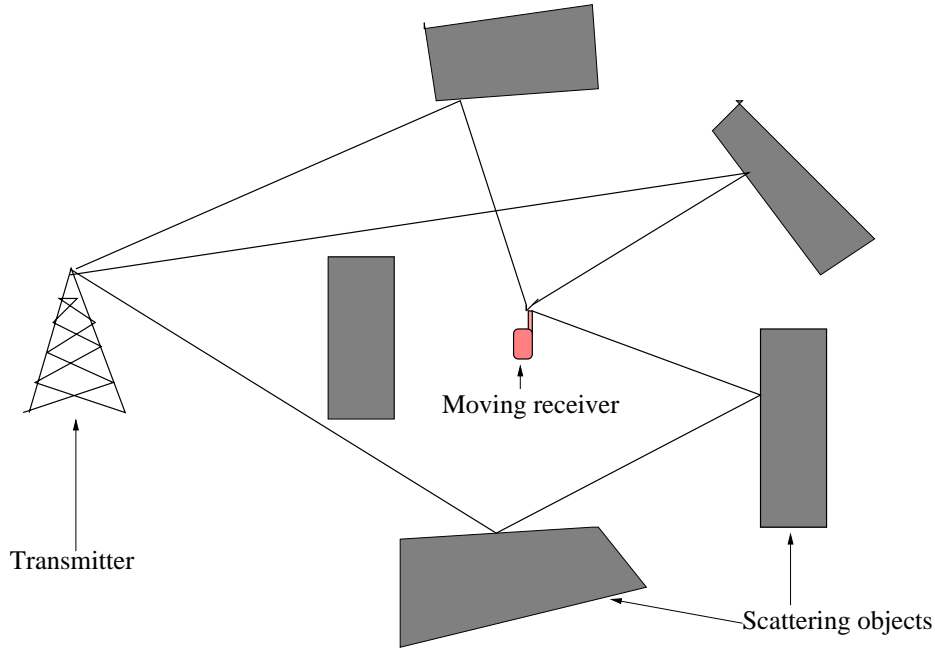


Figure 5.1: Illustration of typical scattering scenario considered for the MFC model. Lines between the transmitter, scatterers and the receiver represent different radio wave paths.

$X_t$  is modulated to the carrier frequency  $f_c$  by the procedure

$$\tilde{X}_t = \mathcal{F}^{-1} \left( \frac{\mathcal{F}(X_t)(f - f_c) + \mathcal{F}(X_t)(-f - f_c)}{\sqrt{2}} \right) (t),$$

or, in time domain,  $\tilde{X}_t = \sqrt{2} \text{Re}[X_t \exp(i2\pi f_c t)]$ . The transmitter transmits the real-valued *passband input signal*  $\tilde{X}_t$ , and the receiver receives the *passband output signal*

$$\tilde{Z}_t = \int_{\mathbb{R}} \tilde{H}(s, t) \tilde{X}_{t-s} ds \quad (5.1.2)$$

where the passband channel response for MFC channels is given by

$$\tilde{H}(s, t) := \sum_{j=1}^M a(\alpha_j) \delta(s - \tau(\alpha_j, t)), \quad (5.1.3)$$

and  $\delta(\cdot)$  denotes the Dirac delta function. To extract the information of the received signal, the passband output signal is demodulated to the *baseband output signal*

$$Z_t = \mathcal{F}^{-1} \left( \sqrt{2} \mathcal{F}(\tilde{Z}_t)(f - f_c) 1_{|f| < W} \right) (t).$$

Frequency modulation makes it possible to operate many channels simultaneously in the same physical area by transmitting at different carrier frequencies (for example, FM radio broadcasting with many active channels), but the information each channel transmits is contained and most simply analyzed as a baseband signal. In fact, by introducing the

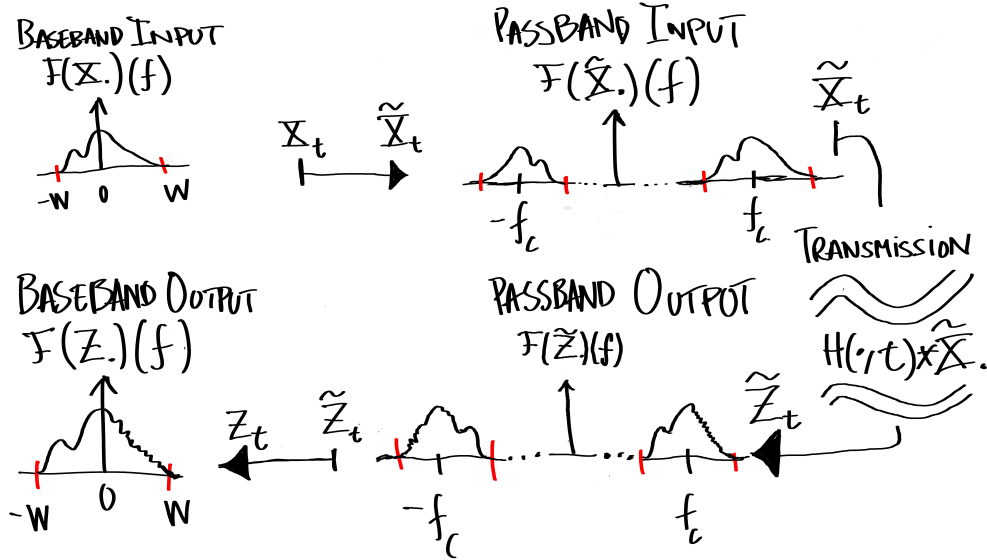


Figure 5.2: The modulation and demodulation taking place when transmitting data in a frequency modulated wireless channel. The baseband input is not drawn equal to the baseband output in the figure, i.e.,  $\mathcal{F}(Z) \neq \mathcal{F}(X)$ ; this is to visualize that some signal information is lost during transmission.

baseband channel response

$$H(s, t) := \sum_{j=1}^M a(\alpha_j) e^{-i2\pi f_c \tau(\alpha_j, t)} \delta(s - \tau(\alpha_j, t)), \quad (5.1.4)$$

it is possible to represent the output baseband as a function of the input baseband

$$Z_t = \int_{\mathbb{R}} H(s, t) X_{t-s} ds, \quad (5.1.5)$$

cf. [32]. This representation equals (5.1.1), avoids dealing with the passband, and it shows that given the baseband input signal, the the output signal is described by the baseband channel response.

In Paper I, we study the setting with baseband input signal  $X_t := 1$ . This setting gives the baseband output signal

$$Z_{t_m} = \sum_{j=1}^M a(\alpha_j) e^{-i2\pi f_c \tau(\alpha_j, t_m)},$$

which is closely related to the discrete channel response.

Although the wireless channel (5.1.1) is presented in a continuum setting, real life channels signal processing is performed on discrete time steps with all wave paths arriving at the receiver within a sample period  $\Delta t$  are averaged to become the received value for that time

sample, cf. [32]. In the transition from continuum to discrete, the Nyquist-Shannon Sampling Theorem describes how small the sample period  $\Delta t$  have to be to resolve a continuous signal  $X_t$ .

**Theorem 5.1.1 (Nyquist-Shannon Sampling Theorem [33])** *Suppose a function  $X_t$  contains no frequencies higher than  $B$  Hertz. Then it is completely determined by giving its ordinates at a series of points spaced  $1/(2B)$  seconds apart.*

**Remark 5.1.2** *For more on the ideas presented in this wireless channels, see [32] and [11].*

### Clarke's Model

Clarke's model is a famous MFC model which was motivational for our work, cf. [10]. It considers the superposition of  $M$  wave paths with amplitudes  $a = 1/\sqrt{M}$ ;

$$Z_{t,M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M e^{-i(2\pi f_c v (\cos(\alpha_m)t/c + \theta_m))}. \quad (5.1.6)$$

Here, apart from the variables already introduced in Section 5.1,  $c$  is the speed of light,  $\{\theta_m\}_{m=1}^M$  are i.i.d. initial phase shifts uniformly distributed in  $[0, 2\pi)$  and the scatterer angle of arrivals  $\{\alpha_m\}_{m=1}^M$  are distributed according to the scatterer angle density  $p(\alpha)$  which is independent from the initial phase shift distribution. The delay function is on the form  $\tau(\alpha, t) = -f_c v \cos(\alpha)t/c$  and the factor  $\partial_t \tau(\alpha, t) = -f_c v \cos(\alpha)/c$  is the Doppler shift under the assumption that the receiver moves in the direction  $(1, 0)$ . The Doppler shift describes the change of frequency of a wave for a receiver moving relative to the wave transmitter (the tone of the siren of a passing ambulance is a classical illustration of this phenomenon).

Among the functions used to analyze channel properties is the autocorrelation function and the Power Spectral Density (PSD). The *autocorrelation*  $A_M(t) := \mathbb{E}[Z_{t,M} Z_{0,M}^*]$  describes the correlation between  $Z_{t,M}$  and  $Z_{0,M}$ , and for Wide Sense Stationary (WSS) signals, the *PSD* is the Fourier transform pair of the autocorrelation function. Let us state this formally.

**Definition 5.1.3 (Wide Sense Stationary Random Process)** *A random process  $Z_t$  is WSS if there is a function  $g$  such that*

$$\mathbb{E}[Z_{\bar{t},M} Z_{t,M}^*] = g(|\bar{t} - t|) \quad \text{for any } \bar{t}, t \in \mathbb{R}.$$

**Theorem 5.1.4 (Wiener-Khintchine Theorem [11, p. 49])** *The power spectral density and autocorrelation of a wide sense stationary process are Fourier transform pairs.*

Considering the scenario with scatterer angle density  $p(\alpha) = (2\pi)^{-1}$ , Clarke noted that for his model the autocorrelation function  $A_M(t)$  converges to the zeroth-order Bessel function of the first kind,  $J_0(2\pi f_c vt/c)$ , as  $M \rightarrow \infty$ , and that its PSD then is on the form

$$S(f) = \begin{cases} \frac{c}{\pi \sqrt{(vf_c)^2 - (cf)^2}} & |f| < vf_c/c \\ 0 & |f| \geq vf_c/c, \end{cases} \quad (5.1.7)$$

cf. Figure 5.3.

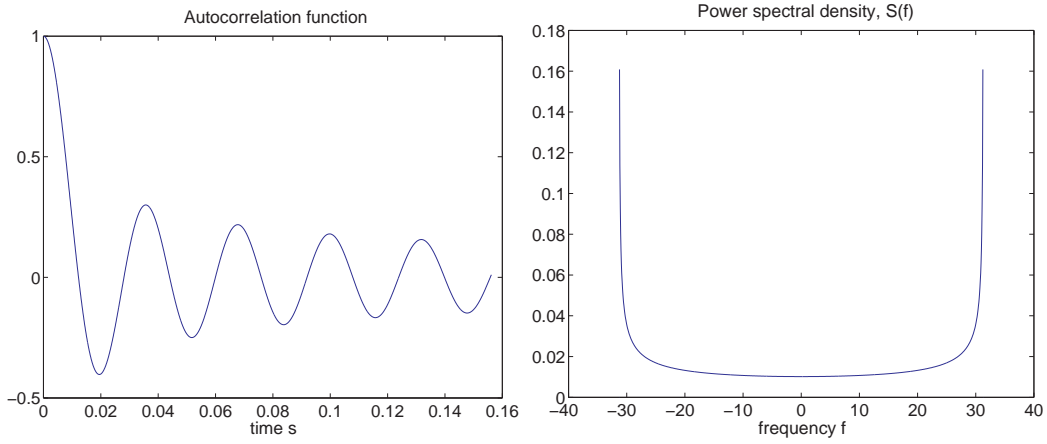


Figure 5.3: **Left plot:** The autocorrelation function for Clarke's model with azimuth density  $p(\alpha) = (2\pi)^{-1}$ ,  $v = 5m/s$  and  $f_c = 1.8775GHz$ . **Right plot:** The power spectral density of Clarke's model, often called Jakes' spectrum with the same model conditions as for the left plot.

### MFC Model with Flipping Scatterers

Due to local shadowing by moving cars, pedestrians, leaves blowing in the wind, weather conditions etc. scatterers can flip from being active to passive and vice versa. Seeking to include scattering objects in our MFC model, we introduce the amplitude function as a stochastic process  $a(\alpha, t)$  which flips on when it changes value from 0 to  $a^+(\alpha) \geq 0$  and off when it changes values oppositely, cf. Figure 5.4. It is here assumed that the mapping

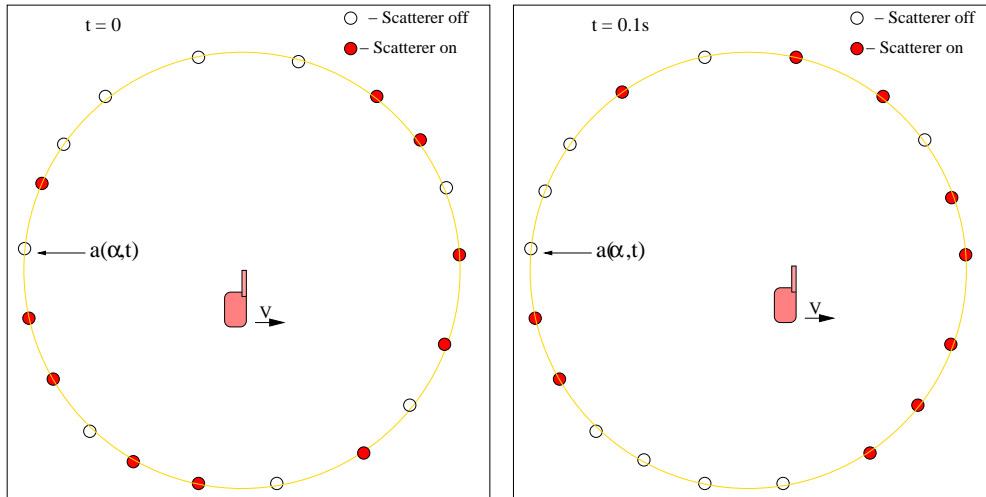


Figure 5.4: Moving receiver in the center of a circular scattering environment with scatterers flipping on and off as time goes.

$a^+ : \Omega \rightarrow \mathbb{R}_+$  is smooth; it might for example be constant or depend on the distance from scatterer to the receiver. The flip process is taken to be Poisson distributed with flip rate

constant  $C$ :

$$P(a(\alpha) \text{ flips } k \text{ times on time step } \Delta t) = \frac{(C\Delta t)^k \exp(-C\Delta t)}{k!},$$

where flips are independent from the scatterers' random initial phase shifts.

With the above defined amplitude function and the set of arrival angles  $\{\alpha_j\}_{j=1}^M$  distributed according to a scatterer density  $p(\alpha)$ , we propose the following flip process extension of Clarke's model

$$Z_{t,M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M a(\alpha_m, t) e^{-i2\pi f_c v \cos(\alpha_m) t/c + \theta_m(t)}. \quad (5.1.8)$$

Realizations of (5.1.8) is then generated by generating scatterer angles  $\{\alpha_j\}_{j=1}^M$  and i.i.d. initial phase shifts  $\{\theta_m\}_{m=1}^M$ , generating sample paths of the stochastic process amplitudes  $a(\alpha_m, t)$ , and summing contributions according to (5.1.8). The left plot of Figure 5.5 illustrates the difference between an MFC signal realization envelope generated with  $C = 0$  and one with positive flip rate. The positive flip rate gives the signal envelope more small scale temporal noise and less smoothness than what is found in the non-flipping signal envelope. The right plot of Figure 5.5 is a measured signal envelope from Ericsson Labs which shows that the measured signal has a small scale noise contribution similar to the MFC signal realization with positive flip rate.

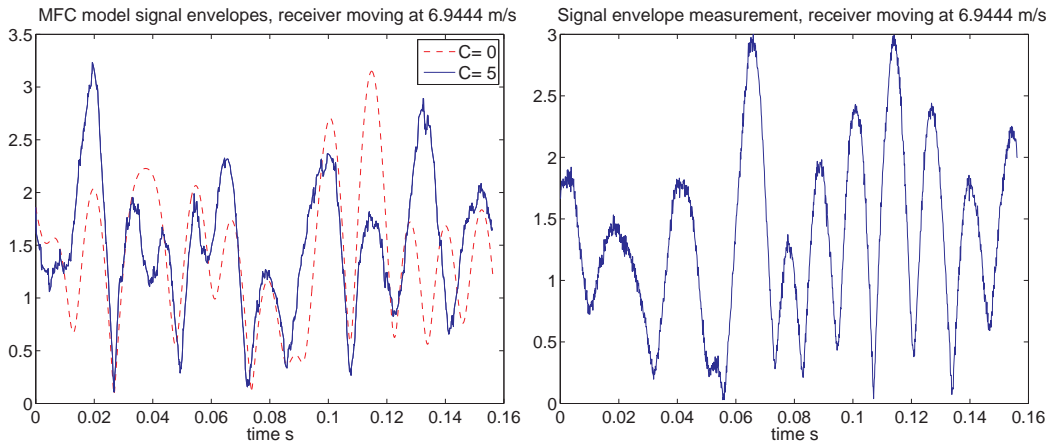


Figure 5.5: **Left plot:** Two computer generated signal realizations of the MFC model generated with same random seed initialization. Both realizations are generated with the modeling parameters  $f_c = 1.8775GHz$ ,  $v = 6.944m/s$ ,  $a^+(\alpha) = 2/\sqrt{M}$ , but the flip rate separates the realizations with the red dashed line corresponding to a realization having  $C = 0$ , and the whole line to a realization with  $C = 5$ . **Right plot:** Measured urban environment signal envelope from Ericsson Labs. The carrier frequency and receiver speed is identical to the corresponding values for the left plot.

## 5.2 From MFC Models to Gaussian Processes

In Paper I, we show that under some assumptions, the stochastic process  $Z_{t,M}$  of the MFC model (5.1.8) converges to a Gaussian process as the number of included wave paths

$M \rightarrow \infty$ . Here we will give a short description of Gaussian processes, starting with the definition.

**Definition 5.2.1 (Gaussian Process)** *A Gaussian process is a stochastic process  $\{Z_t\}_{t \in [0, T]}$ ,  $Z_t \in \mathbb{R}^n$ , for which any finite length sample vector  $\mathbf{Z} = (Z_{t_1}, Z_{t_2}, \dots, Z_{t_n})$  with  $0 \leq t_1 \leq t_2 \leq \dots \leq t_n$  is multivariate normal distributed.*

A numerical realization of a Gaussian process on a set of times  $\{t_j\}_{j=1}^N$  can be created by first computing the process' covariance matrix,

$$K_{i,j} = E[Z_{t_i} Z_{t_j}^*], \quad i, j \in \{1, 2, \dots, N\}$$

and setting

$$Z_{t_i} = \sum_{j=1}^N \sqrt{K}_{i,j} \chi_j, \quad (5.2.1)$$

where  $\sqrt{K}$  is the square root of  $K$  in the sense that it fulfills  $K = \sqrt{K} \sqrt{K}^H$  (for example the lower diagonal Cholesky factorization), and  $\chi_1, \chi_2, \dots, \chi_N$  is a set of i.i.d. standard normal distributed r.v.

**Example 5.2.2 (The Wiener Process)** *The Wiener process  $W_t$  is a Gaussian process which has the increment property  $W_{t_2} - W_{t_1} \sim N(0, |t_2 - t_1|)$  and thereby the covariance matrix*

$$K_{i,j} = E[W_{t_i} W_{t_j}] = E[W_{\min(t_i, t_j)}^2] = \min(t_i, t_j), \quad (5.2.2)$$

cf. Section 4.1. For the Wiener process, the structure of the Cholesky factorized  $\sqrt{K}$  is particularly simple:

$$\sqrt{K} = \begin{pmatrix} \sqrt{t_1} & 0 & \dots & 0 \\ \sqrt{t_1} & \sqrt{t_2 - t_1} & \ddots & 0 \\ \vdots & \vdots & \ddots & 0 \\ \sqrt{t_1} & \sqrt{t_2 - t_1} & \dots & \sqrt{t_N - t_{N-1}} \end{pmatrix}. \quad (5.2.3)$$

This is fortunate because it makes the computational cost of generating a Wiener process realization  $\{W_{t_j}\}_{j=1}^N$  to  $\mathcal{O}(N)$  which compares favorably to  $\mathcal{O}(N^2)$  for general Gaussian process realizations. The Wiener process realization cost  $\mathcal{O}(N)$  may be concluded from the scheme (5.2.1) which when  $\sqrt{K}$  is given by (5.2.3) becomes

$$W_{t_{j+1}} = W_{t_j} + \sqrt{t_{j+1} - t_j} \chi_j, \quad \chi_j \sim \mathcal{N}(0, 1).$$

Generating signal realizations from an MFC model is generally computationally costly, in particular in more realistic settings when a high number of wave paths  $M$  are included in your model. Comparatively, the analysis and computational experiments of Paper I indicate that generation of Gaussian process signal realizations are substantially less costly.





## Chapter 6

# Classical and Quantum Mechanics

For particles at human scale, motion is accurately described by classical mechanics, while for particles at atomic scale, motion is described by a quantum mechanics. In settings with many particles, quantum mechanical computations are very difficult, and, often, the best approximations one can make are by means of classical mechanics. Molecular dynamics is a classical mechanics approximation of the quantum scale motion of electrons and nuclei in molecular bindings, and in Paper IV of this thesis we study how well the motion of electrons and nuclei are approximated by molecular dynamics when the nucleus-to-electron mass ratio becomes large. In this chapter we give an outline of some classical and quantum mechanical concepts that will be useful when reading Paper IV.

### 6.1 Classical Mechanics

The motion of a single particle in a  $d$ -dimensional potential  $V : \mathbb{R}^d \rightarrow \mathbb{R}$  and with mass  $M$  is described by the force acting on the particle;  $\mathbf{F} = -\nabla V$ . This yields the equations of motion

$$\dot{\mathbf{q}} = \mathbf{v} \quad \text{and} \quad M\dot{\mathbf{v}} = -\nabla V, \quad (6.1.1)$$

where  $\mathbf{q}(t)$  and  $\mathbf{v}(t)$  represents the position and velocity of the particle, respectively. A preserved quantity during this motion is the energy  $E = M|\mathbf{v}|^2/2 + V(\mathbf{q})$  which is easily verified through differentiating:

$$\dot{E} = M\dot{\mathbf{v}} \cdot \mathbf{v} + \nabla V \cdot \dot{\mathbf{q}} \stackrel{(6.1.1)}{=} 0. \quad (6.1.2)$$

Introducing the Hamiltonian  $H(\mathbf{q}, \mathbf{p}) = |\mathbf{p}|^2/2M + V(\mathbf{q})$  with the momentum  $\mathbf{p} = M\mathbf{v}$ , it is often more convenient to represent the equations of motion on the Hamiltonian form

$$\dot{\mathbf{q}} = \nabla_{\mathbf{p}} H \quad \text{and} \quad \dot{\mathbf{p}} = -\nabla_{\mathbf{q}} H. \quad (6.1.3)$$

We also note that the Hamiltonian is preserved since by construction  $H = E$ .

For a system of  $N$  particles with the  $i$ th particle position denoted  $\mathbf{q}_i \in \mathbb{R}^d$  and with mass  $M_i$ , the force fields acting on each particle can often be divided into the external potential  $V_i : \mathbb{R}^d \rightarrow \mathbb{R}$ , and the potentials for the interaction between  $i$ th and  $j$ th particle  $V_{ij}(\mathbf{q}_i - \mathbf{q}_j)$ . Introducing the Hamiltonian

$$H(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^N \frac{|\mathbf{p}_i|^2}{2M_i} + V_i(\mathbf{q}_i) + \sum_{i < j} V_{ij}(\mathbf{q}_i - \mathbf{q}_j),$$

where  $\mathbf{q} = (\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_N) \in \mathbb{R}^{dN}$  and, similarly,  $\mathbf{p} \in \mathbb{R}^{dN}$ , the equations of motion for  $N$  particles take the familiar form

$$\dot{\mathbf{q}}_i = \nabla_{\mathbf{p}_i} H \quad \text{and} \quad \dot{\mathbf{p}}_i = -\nabla_{\mathbf{q}_i} H, \quad \text{for } i = 1, 2, \dots, N. \quad (6.1.4)$$

The term particle is used in a relative sense in our above outline: we might be dealing with stars and planets or grains of sand, conditioned that the reduction to mass points makes sense. Let us present two examples of Hamiltonian dynamics.

**Example 6.1.1 (The harmonic oscillator)** *The simple harmonic oscillator is 1-dimensional mass-spring system which when displaced from its equilibrium position experiences a restoring force  $\mathbf{F} = -Mkq$ , with  $k \in \mathbb{R}_+$  being a spring constant and  $q \in \mathbb{R}$  the displacement from equilibrium. The Hamiltonian for this dynamics becomes*

$$H(q, p) = \frac{p^2}{2M} + Mkq^2/2.$$

**Example 6.1.2 ( $N$ -body problem)** *Newton explained the observed planetary motion by a negligible external field and interplanetary forces  $\mathbf{F}_{ij} = -GM_i M_j (\mathbf{q}_i - \mathbf{q}_j) / |\mathbf{q}_i - \mathbf{q}_j|^3$ , where  $G$  is a gravitational constant. Considering the setting with  $N$  planets, Newton's interplanetary forces translate to the interaction potentials*

$$V_{ij}(\mathbf{q}_i - \mathbf{q}_j) = -\frac{GM_i M_j}{|\mathbf{q}_i - \mathbf{q}_j|}, \quad i, j \in \{1, 2, \dots, N\} \text{ and } i \neq j,$$

and the negligible external potentials to  $V_i = 0$ , for  $i = 1, 2, \dots, N$ . Consequently, the  $N$ -body equations of motion are implicitly given by (6.1.4) with the Hamiltonian

$$H(\mathbf{q}, \mathbf{p}) = \sum_{i=1}^N \frac{|\mathbf{p}_i|^2}{2M_i} - \sum_{i < j} \frac{GM_i M_j}{|\mathbf{q}_i - \mathbf{q}_j|}.$$

## Evolution of States

For a given  $N$ -particle system, the set of possible  $(\mathbf{q}, \mathbf{p})$  restricted to some subset of  $\mathbb{R}^{dN} \times \mathbb{R}^{dN}$  is referred to as the phase space, and the vector  $(\mathbf{q}(t), \mathbf{p}(t))$  at a given time  $t$  is referred to as the particle system's state at time  $t$ . Considering experiments which involve measurements, we will extend the notion of states to densities over the phase space  $\rho(\mathbf{q}, \mathbf{p})$  with the restriction that  $\rho$  is non-negative. A pure state  $(\tilde{\mathbf{q}}, \tilde{\mathbf{p}})$  may then be represented on density form by dirac functions  $\rho(\mathbf{q}, \mathbf{p}) = \delta(\mathbf{q} - \tilde{\mathbf{q}}, \mathbf{p} - \tilde{\mathbf{p}})$ . For a given state  $\rho$ , the expected value of an observable, like the total energy  $H$ , is given by

$$E_\rho[H] = \frac{\int H(\mathbf{q}, \mathbf{p}) \rho(\mathbf{q}, \mathbf{p}) d\mathbf{q} d\mathbf{p}}{\int \rho(\mathbf{q}, \mathbf{p}) d\mathbf{q} d\mathbf{p}}.$$

A convex combination of states also becomes a state. E.g., given two states  $\rho_1$  and  $\rho_2$ , the convex combination  $\rho_3 = \alpha\rho_1 + (1 - \alpha)\rho_2$  is also a state.

The motion of  $(\mathbf{q}, \mathbf{p})$  satisfies the Hamiltonian dynamics (6.1.4), and this implies that the density  $\rho(\mathbf{q}, \mathbf{p})$  is convected by the flow  $(\dot{\mathbf{q}}, \dot{\mathbf{p}})$ . The convection makes it practical to consider the density time dependent in the following sense

$$\rho(\mathbf{q}(t), \mathbf{p}(t); t) = \rho(\mathbf{q}(0), \mathbf{p}(0)). \quad (6.1.5)$$

This is often referred to as the Liouville picture, cf. [15]. For any material region  $R(t)$  of phase space convected by the flow  $(\dot{\mathbf{q}}, \dot{\mathbf{p}})$ , the Liouville picture (6.1.5) straightforwardly implies that

$$\frac{d}{dt} \int_{R(t)} \rho(\mathbf{q}, \mathbf{p}; t) d\mathbf{q} d\mathbf{p} = 0,$$

and by analysis involving time differentiation of the Jacobian determinant of  $\partial(\mathbf{q}, \mathbf{p})/\partial(\mathbf{q}(0), \mathbf{p}(0))$ , one may further derive the transport theorem

$$\frac{d}{dt} \int_{R(t)} \rho(\mathbf{q}, \mathbf{p}; t) d\mathbf{p} d\mathbf{q} = \int_{R(t)} \partial_t \rho + \nabla_{\mathbf{q}, \mathbf{p}} \cdot ((\dot{\mathbf{q}}, \dot{\mathbf{p}})\rho) d\mathbf{q} d\mathbf{p}, \quad (6.1.6)$$

cf. [27]. For any fixed region  $D$ , there is at any moment in time a coinciding material volume  $R(t)$ , and combining this property with (6.1.5) and (6.1.6), yields the following density PDE

$$\partial_t \rho + \nabla_{\mathbf{q}, \mathbf{p}} \cdot ((\dot{\mathbf{q}}, \dot{\mathbf{p}})\rho) = 0.$$

Furthermore, the Hamiltonian dynamics (6.1.4) implies that

$$\nabla_{\mathbf{q}, \mathbf{p}} \cdot (\dot{\mathbf{q}}, \dot{\mathbf{p}}) = \sum_{i=1}^N (\nabla_{\mathbf{q}_i} \cdot \nabla_{\mathbf{p}_i} - \nabla_{\mathbf{p}_i} \cdot \nabla_{\mathbf{q}_i}) H = 0, \quad (6.1.7)$$

so that we end up with the following PDE for  $\rho(\mathbf{q}, \mathbf{p}; t)$

$$\partial_t \rho + \sum_{i=1}^N \nabla_{\mathbf{p}_i} H \cdot \nabla_{\mathbf{q}_i} \rho - \nabla_{\mathbf{q}_i} H \cdot \nabla_{\mathbf{p}_i} \rho = 0 \quad (6.1.8)$$

with initial condition  $\rho(\mathbf{q}, \mathbf{p}; 0) = \rho(\mathbf{q}, \mathbf{p})$ . This is often referred to as Liouville's theorem.

### Preservation of Volume

For Hamiltonian dynamics, the phase space volume is preserved. This might be verified by considering the material region  $R(t)$  convected by the Hamiltonian flow  $(\dot{\mathbf{q}}, \dot{\mathbf{p}})$  with the evolving phase space volume

$$V(t) = \int_{R(t)} 1 d\mathbf{q} d\mathbf{p}.$$

Equation (6.1.6), in this case with  $\rho = 1$ , and (6.1.7) implies that

$$\dot{V}(t) = \int_{R(t)} \nabla_{\mathbf{q}, \mathbf{p}} \cdot (\dot{\mathbf{q}}, \dot{\mathbf{p}}) d\mathbf{q} d\mathbf{p} = 0.$$

The preservation of phase space implies that if a fluid convected by Hamiltonian dynamics is expanding in  $\mathbf{q}$ -space, then it is contracting in  $\mathbf{p}$ -space, and vice versa.

### Symplectic Numerical Methods

The word symplectic derives from the Ancient Greek *συμπλεκτικός* which is a composite of “braided” and “together”. In mathematics, it was introduced by Herman Weyl in 1939 to describe the group of unitary transformations of  $\mathbb{C}^{2n}$  that for vectors  $\xi = (\xi^q, \xi^p)$ ,  $\eta = (\eta^q, \eta^p)$  preserves the operation

$$\sum_{i=1}^n (\xi_i^p)^* \eta_i^q - (\xi_i^q)^* \eta_i^p, \quad (6.1.9)$$

where  $\xi^p, \xi^q, \eta^p, \eta^q \in \mathbb{C}^n$  are the  $p$  and  $q$  components of the vector. Introducing the matrix

$$J = \begin{bmatrix} 0 & I \\ -I & 0 \end{bmatrix},$$

with  $I$  the  $n$ -dimensional identity matrix, the preservation of the operation (6.1.9) is equivalent to the preservation of  $\xi^H J \eta$ , and restricting ourselves to  $\mathbb{R}^{2n}$ , we may with the aid of  $J$  define linear and differentiable symplectic maps as follows.

**Definition 6.1.3 (Symplectic linear map, cf. [19])** A linear map  $A : \mathbb{R}^{2n} \rightarrow \mathbb{R}^{2n}$  is called symplectic if

$$A^T J A = J$$

**Definition 6.1.4 (Symplectic differentiable map, cf. [19])** A differentiable map  $g : U \rightarrow \mathbb{R}^{2n}$ , (where  $U \subset \mathbb{R}^{2n}$  is an open set) is called symplectic if the Jacobian matrix  $g'(\mathbf{q}, \mathbf{p})$  is everywhere symplectic, i.e., if

$$g'(\mathbf{q}, \mathbf{p})^T J g'(\mathbf{q}, \mathbf{p}) = J.$$

For any Hamiltonian and fixed  $t$ , the flow mapping  $\phi_t(p(0), q(0)) = (p(t), q(t))$  is symplectic, cf. [19]. This property opens an alternative proof of the flow  $\phi_t$ : Symplecticity means that

$$\phi_t'(\mathbf{q}, \mathbf{p})^T J \phi_t'(\mathbf{q}, \mathbf{p}) = J.$$

By observing that  $\det(J) = (-1)^{2n} = 1$ ,  $\det(\phi_0'(\mathbf{q}, \mathbf{p})) = 1$  and assuming the Hamiltonian is sufficiently smooth to make  $\phi_t$  a continuous function with respect to  $t$ , we conclude that  $\det(\phi_t'(\mathbf{q}, \mathbf{p})) = 1$  for all valid times. This implies conservation of volume.

Symplecticity is also sought property for numerical integrators of the Hamiltonian dynamics.

**Definition 6.1.5 (Symplectic one step-method, cf. [19])** A numerical one step method  $(p_{n+1}, q_{n+1}) = \phi_h((p_n, q_n))$ , is called symplectic if whenever applied to a smooth Hamiltonian system,

$$\left( \frac{\partial(p_{n+1}, q_{n+1})}{\partial(p_n, q_n)} \right)^T J \left( \frac{\partial(p_{n+1}, q_{n+1})}{\partial(p_n, q_n)} \right) = J.$$

For Hamiltonian systems, symplectic numerical integrators preserve the phase space volume perfectly and generally preserve energy approximately even when integrating over long time intervals. These properties are often highly important to preserve; in the case of planetary motion, for example, the preservation of the energy ensures that planets stay on orbit. Let us illustrate by an example comparing the performance simplest one step method, the forward Euler method

$$(q_{n+1}, p_{n+1}) = (q_n + \Delta t \partial_p H(p_n, q_n), p_n - \Delta t \partial_q H(p_n, q_n)),$$

with the simplest symplectic one step method, the symplectic Euler method

$$(p_{n+1}, q_{n+1}) = (q_n + \Delta t \partial_p H(p_{n+1}, q_n), p_n - \Delta t \partial_q H(p_{n+1}, q_n)).$$

**Example 6.1.6 (Numerical solutions for the Harmonic oscillator)** We consider the Harmonic oscillator of Example 6.1.1 with  $M = 1, k = 1$  and the Hamiltonian  $H =$

$p^2/2 + q^2/2$ . For this problem, the forward Euler and symplectic Euler scheme take the following respective forms

$$\begin{aligned}(q_{n+1}, p_{n+1}) &= (q_n + \Delta t p_n, p_n - \Delta t q_n), && \text{(forward Euler),} \\ (q_{n+1}, p_{n+1}) &= (q_n + \Delta t(p_n - \Delta t q_n), p_n - \Delta t q_n), && \text{(symplectic Euler).}\end{aligned}$$

With the initial condition  $(q(0), p(0)) = (1, 0)$ , the preservation of the Hamiltonian implies that the dynamics should fulfill  $H(q, p) = 1/2$  for any  $t$ , and one may further derive that  $q(t) = \cos(t)$  is the solution of this problem. The comparison of the methods given in Figure 6.1 shows that while the volume is preserved and the energy is approximately preserved by the symplectic Euler method, even for large time steps, the forward Euler method solutions have increasing material regions and energy.

## 6.2 Quantum Mechanics

Quantum mechanics is a description of the mechanics of particles on atomic scale. Quantum mechanics differ from classical mechanics in the following sense:

- The state of a particle in classical mechanics is described by position and momentum (or, if extending notions, by a density) and the state evolves deterministically.
- The state of a particle in quantum mechanics is described by the amplitude of the complex-valued wave function  $\Phi(x, t)$  solving the time-dependent Schrödinger equation

$$i\partial_t \Phi = \left( -\frac{1}{2M} \Delta + V(x, t) \right) \Phi.$$

Furthermore, since it is impossible to obtain perfect information on the state of a particle (an implication of Heisenberg's uncertainty principle, cf. [18]), whether or not a quantum state evolves deterministically is a question never to be verified experimentally—probably...

### Derivation of Schrödinger's Equations

Around the beginning of the 1900s, various experiments, such as the slit experiment(s), cf. [15], showed that elementary particles like photons and electrons in addition to having particle properties have wave-like diffraction properties. This wave-particle duality led to the conjecture that elementary particles both have particle attributes such as mass, charge, and spin, and wave attributes such as frequency and, consequently, that there are two equivalent expressions for the total energy. Classical mechanics states that the total energy for a single particle is given by  $E = |\mathbf{p}|^2/2m + V(\mathbf{x}, t)$  with  $\mathbf{p}, \mathbf{x} \in \mathbb{R}^3$ . For waves, work by Planck, Einstein, and de Broglie led to the conjecture an elementary particle with momentum  $\mathbf{p}$  and energy  $E$  in some sense corresponds to a wave

$$\Phi(\mathbf{x}, t) = \phi(\mathbf{x}, t) e^{i(\mathbf{k} \cdot \mathbf{x} - 2\pi f t)} \quad (6.2.1)$$

with wave vector  $\mathbf{k} \in \mathbb{R}^3$  and frequency  $f$  for which the following relations hold

$$\mathbf{p} = \hbar \mathbf{k}, \quad |\mathbf{k}| = 2\pi f/c, \quad \text{and} \quad E = |\mathbf{p}|c = \hbar |\mathbf{k}|c = \hbar f. \quad (6.2.2)$$

Here  $c$  denotes the speed of light,  $h \approx 1.05 \times 10^{-27}$  erg-sec denotes Planck's constant, and  $\hbar := h/(2\pi)$ . By the relations (6.2.2), the wave function (6.2.1) may be represented as follows

$$\Phi(\mathbf{x}, t) = \phi(\mathbf{x}, t) e^{i(\mathbf{p} \cdot \mathbf{x} - Et)/\hbar}. \quad (6.2.3)$$

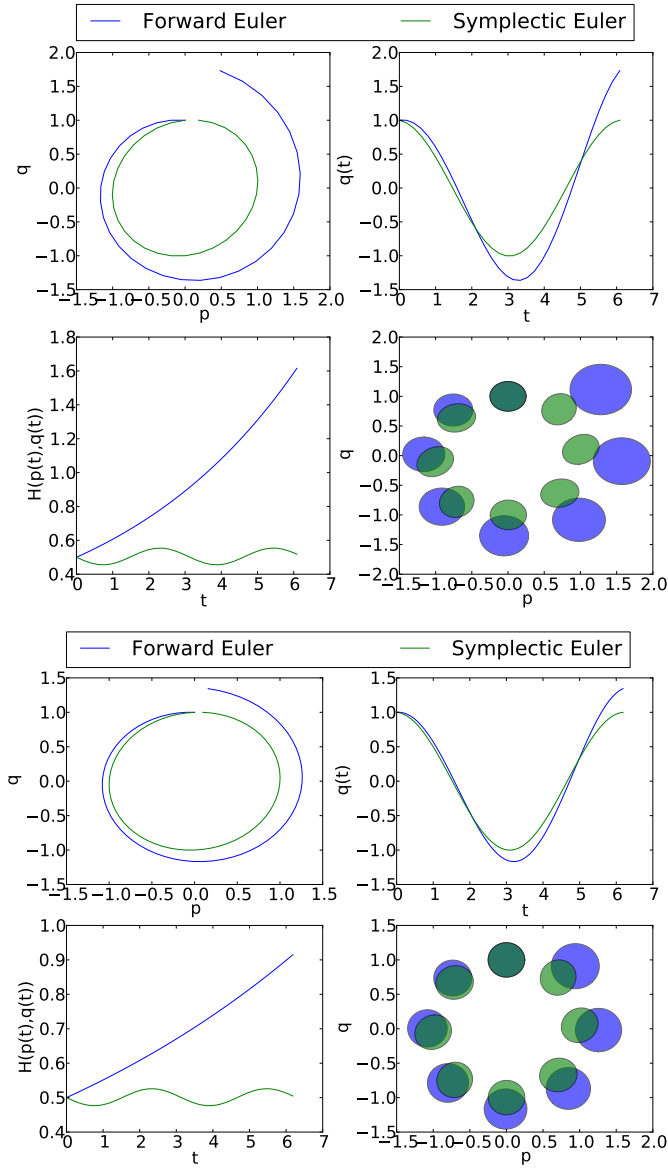


Figure 6.1: Numerical solutions of the Harmonic oscillator problem of Example 6.1.6. The top four plots are numerical solutions with  $\Delta t = 2\pi/32 \approx 0.2$  integrated on the interval  $t \in [0, 2\pi)$  and consist of: phase space solution  $(p(t), q(t))$  (upper left),  $(t, q(t))$  (upper right), the energy  $(t, H(p(t), q(t)))$  (lower left), anti-clockwise flow simulation of phase space ball material regions (lower right). The bottom four plots are analogous numerical solutions with  $\Delta t = 2\pi/64 \approx 0.1$ .

Supposing the amplitude  $\phi$  is smooth and slowly varying compared to the scales  $\mathcal{O}(f/(hc))$  and  $\mathcal{O}(E/h)$  the wave function (6.2.3) approximately fulfills the relation

$$(i\hbar\partial_t + \hbar^2 \frac{\Delta_x}{2m})\Phi = (E - \frac{|\mathbf{p}|^2}{2m})\Phi. \quad (6.2.4)$$

By coordinate transformations  $t \rightarrow t/\hbar$ ,  $x \rightarrow x/\hbar$  and replacing  $E$  with the classical mechanics total energy  $|\mathbf{p}|^2/2m + V(\mathbf{x}, t)$ , equation (6.2.4) becomes the time-dependent Schrödinger equation

$$i\partial_t\Phi = \left(-\frac{1}{2m}\Delta + V(\mathbf{x}, t)\right)\Phi.$$

### Time-Independent Schrödinger Equation

Assuming the potential is time-independent, we may argue as above using the wave ansatz (6.2.3) to observe that  $i\hbar\partial_t\Phi = E\Phi$  and to further derive the time-independent Schrödinger equation

$$\left(-\frac{1}{2m}\Delta + V(\mathbf{x})\right)\Phi = E\Phi.$$

Solutions of the time-independent Schrödinger equation describe particle motions with preserved total energy—steady state solutions. The time-independent Schrödinger equation requires no initial data  $\Phi(\cdot, 0)$ , and in this sense it is a more fundamental PDE than the time-dependent Schrödinger equation. In what follows, we will restrict ourselves to the time-independent Schrödinger equation.

### Many Particle Schrödinger Equation

Considering a molecule with  $N$  nuclei and  $n$  electrons, the time-independent Schrödinger equation takes the form

$$\left(-\sum_{i=1}^N \frac{1}{2M_i} \Delta_{\mathbf{x}_i} - \sum_{i=1}^n \frac{1}{2m} \Delta_{\mathbf{x}_i} + V(\mathbf{X}, \mathbf{x})\right)\Phi(\mathbf{X}, \mathbf{x}) = E\Phi(\mathbf{X}, \mathbf{x}),$$

where  $\mathbf{X} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N) \in \mathbb{R}^{3N}$  and  $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \in \mathbb{R}^{3n}$  denote the nuclear and electron positions, respectively,  $M_i$  denotes the mass of the  $i$ -th nucleus, and  $m \approx 9.11 \times 10^{-31}$  kg the electron mass. Supposing all nuclei have the same mass, i.e.,  $M_i = M_1 \forall i$ , and introducing the coordinate transformations  $(\mathbf{X}, \mathbf{x}) \rightarrow \sqrt{m}(\mathbf{X}, \mathbf{x})$  and the short hand notation  $\Delta_{\mathbf{X}} := \sum_{i=1}^N \Delta_{\mathbf{x}_i}$ , we obtain the time-independent Schrödinger equation on the form studied in paper IV:

$$\left(-\frac{1}{2M}\Delta_{\mathbf{X}} - \underbrace{\sum_{i=1}^n \frac{1}{2}\Delta_{\mathbf{x}_i}}_{=: \mathcal{V}(\mathbf{X}, \mathbf{x})} + V(\mathbf{X}, \mathbf{x})\right)\Phi(\mathbf{X}, \mathbf{x}) = E\Phi(\mathbf{X}, \mathbf{x}), \quad (6.2.5)$$

with solutions sought in a Hilbert subspace of  $L^2(\mathbb{R}^{3N} \times \mathbb{R}^{3n})$  with symmetry conditions based on the Pauli exclusion principle for bosons and fermions, cf. [23]. In (6.2.5),  $M = M_1/m$  represents the nucleus-to-electron mass ratio which ranges from approximately 1836 for the  $^1\text{H}$  Hydrogen isotope<sup>1</sup> (corresponding to the proton-to-electron mass ratio) to approximately  $244 \times 1836$  for the Uranium isotope  $^{244}\text{U}$ ,  $^{244}\text{U}$  being the heaviest known isotope to exist on earth in its natural form.

<sup>1</sup>Isotope refers to the number of protons and neutrons in the nucleus of a chemical element.

### Quantum States

Each solution of the time-independent Schrödinger equation with energy  $E \in \mathbb{R}$  corresponds to a quantum state for the particle at the given energy level. That is, for any given solution with an energy  $E$ , we assume the scaling

$$\iint |\Phi(\mathbf{X}, \mathbf{x})|^2 d\mathbf{X} d\mathbf{x} = 1,$$

and that  $|\Phi|^2$  is the probability distribution of the particle positions  $(\mathbf{X}, \mathbf{x})$  in the state  $\Phi$ . For any observable, i.e., self-adjoint operators  $\mathcal{A}$  on  $L^2(\mathbb{R}^{3N} \times \mathbb{R}^{3n})$ , the expected value is then given by

$$E_\Phi[\mathcal{A}] = \iint \Phi(\mathbf{X}, \mathbf{x})^* \mathcal{A}(\mathbf{X}) \Phi(\mathbf{X}, \mathbf{x}) d\mathbf{X} d\mathbf{x}.$$

In paper IV, we consider the subset of observables consisting of self-adjoint operators on the nuclei coordinates  $L^2(d\mathbf{X})$ . An example of an observable is the  $x$ -position of nucleus coordinate  $i$ :  $\mathcal{A}(\mathbf{X}) = \mathbf{X}_{i,1}$ .

For solutions of the time-independent Schrödinger equation a particular kind of linearity is observed: if  $\Phi_1$  and  $\Phi_2$  are solutions of (6.2.5) for a given energy  $E$  (so called degenerate solutions), then any linear combination of these solutions also solves (6.2.5) with the energy  $E$ . That is,

$$\left( -\sum_{i=1}^N \frac{1}{2M} \Delta_{\mathbf{x}_i} - \mathcal{V} \right) (a_1 \Phi_1 + a_2 \Phi_2) = E(a_1 \Phi_1 + a_2 \Phi_2), \quad \forall a_1, a_2 \in \mathbb{C}$$

### The WKB Ansatz

In Paper IV, we compare the density generated by a state of the Schrödinger equation to a density generated by so called Born-Oppenheimer molecular dynamics. In what follows, we will give a short description of the mentioned densities, and some of the tools used to derive them, starting with the WKB ansatz.

The WKB ansatz assumes that solutions of the time-independent Schrödinger equation (6.2.5) are on the form

$$\Phi(\mathbf{X}, \mathbf{x}) = \phi(\mathbf{X}, \mathbf{x}) e^{i\sqrt{M}\theta(\mathbf{X})}, \quad (6.2.6)$$

where the amplitude  $\phi : \mathbb{R}^{3N} \times \mathbb{R}^{3n} \rightarrow \mathbb{C}$  and the phase  $\theta : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  are smooth functions varying on a much slower scale than the mass ratio  $M$ . The WKB function is a solution of (6.2.5) provided that

$$\begin{aligned} 0 &= \left( -\frac{1}{2} \Delta_{\mathbf{X}} + \mathcal{V}(\mathbf{X}, \mathbf{x}) - E \right) \phi(\mathbf{X}, \mathbf{x}) e^{i\sqrt{M}\theta(\mathbf{X})} \\ &= \left( \underbrace{\left( \frac{1}{2} |\nabla\theta|^2 + V_0 - E \right)}_{=:I} \phi \right. \\ &\quad \left. - \underbrace{\left( \frac{1}{2M} \Delta_{\mathbf{X}} \phi + (\mathcal{V} - V_0) \phi - \frac{i}{M^{1/2}} (\nabla_{\mathbf{X}} \phi \cdot \nabla_{\mathbf{X}} \theta + \frac{1}{2} \phi \Delta_{\mathbf{X}} \theta) \right)}_{=:II} \right) e^{i\sqrt{M}\theta(\mathbf{X})}. \end{aligned} \quad (6.2.7)$$



Here the term

$$V_0(\mathbf{X}) := \frac{\langle \phi(\mathbf{X}, \cdot), \mathcal{V}(\mathbf{X}, \cdot) \phi(\mathbf{X}, \cdot) \rangle}{\langle \phi(\mathbf{X}, \cdot), \phi(\mathbf{X}, \cdot) \rangle},$$

with  $\langle \cdot, \cdot \rangle$  representing the electron coordinates inner product of complex-valued functions in  $L^2(\mathbb{R}^{3n})$ , is introduced so that setting  $I = 0$  gives a well defined limit as  $M \rightarrow \infty$  (see Paper IV, p. 9 for details). Setting  $I = 0$  in (6.2.7) gives a (Eikonal) Hamilton-Jacobi PDE for the phase with characteristics given by the Hamiltonian system

$$\dot{\mathbf{X}} = \partial_{\mathbf{P}} H_S, \quad \dot{\mathbf{P}} = -\partial_{\mathbf{X}} H_S, \quad \text{and } H_S(\mathbf{X}, \mathbf{P}) = \frac{1}{2} |P|^2 + V_0(\mathbf{X}) - E,$$

with  $\mathbf{P}(t) = \nabla_{\mathbf{X}} \theta(\mathbf{X}(t))$ . Thereafter, setting  $II = 0$  in (6.2.7), we derive in Theorem 3.1 of Paper IV that under some conditions, the phase fulfills the equation

$$\phi(\mathbf{X}(t), \mathbf{x}) = \frac{\psi(\mathbf{X}(t), \mathbf{x})}{G(\mathbf{X}(t))},$$

where  $\psi(\mathbf{X}, \mathbf{x})$  is the solution of the time-dependent Schrödinger equation

$$iM^{-1/2} \frac{d}{dt} \psi(\mathbf{X}(t), \mathbf{x}) = \left( \mathcal{V}(\mathbf{X}(t), \mathbf{x}) - V_0(\mathbf{X}(t)) \right) \psi(\mathbf{X}(t), \mathbf{x}) - \frac{G(\mathbf{X}(t))}{2M} \Delta_{\mathbf{x}} \frac{\psi(\mathbf{X}(t), \mathbf{x})}{G(\mathbf{X}(t))},$$

and  $G : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  is implicitly defined by the integrating factor

$$\frac{d}{dt} \log G(\mathbf{X}(t)) = \frac{1}{2} \Delta \theta(\mathbf{X}(t)).$$

Supposing the quantum state generated from the above equations for the phase and amplitude is well defined, the WKB ansatz (6.2.6) takes the form

$$\Phi(\mathbf{X}, \mathbf{x}) = \phi(\mathbf{X}, \mathbf{x}) e^{i\sqrt{M}\theta(\mathbf{x})} = \frac{\psi(\mathbf{X}, \mathbf{x})}{G(\mathbf{X})} e^{i\sqrt{M}\theta(\mathbf{x})}$$

and the nuclear coordinate density becomes

$$\rho_S(\mathbf{X}) = \frac{G^{-2}(\mathbf{X}) \langle \psi(\mathbf{X}, \cdot), \psi(\mathbf{X}, \cdot) \rangle}{\|G^{-1}\psi\|_{L^2(\mathbb{R}^{3N} \times \mathbb{R}^{3n})}}.$$

We next present an approximation of the density  $\rho_S$  generated by Born-Oppenheimer molecular dynamics.

### The Born-Oppenheimer Approximation

The Born-Oppenheimer approximation of solutions of the time-independent Schrödinger equation (6.2.5) consists of the following two steps:

1. Clamp the nuclear coordinate and neglect the nuclear kinetic energy in the equation (6.2.5) and solve the remaining electron coordinate equation

$$\mathcal{V}(\mathbf{X}, \cdot) \psi_0(\mathbf{X}, \cdot) = \lambda_0(\mathbf{X}) \psi_0(\mathbf{X}, \cdot). \quad (6.2.8)$$

Here  $\lambda_0(\mathbf{X})$  represents fixed eigenvalue function<sup>2</sup> of the operator  $\mathcal{V}(\mathbf{X}, \cdot)$  which we assume is spectrally separated from other eigenvalues. This equation is generally solved approximately.

---

<sup>2</sup>Often referred to as a potential energy surface.

2. Reintroduce the kinetic energy into the equation and approximately solve the nuclear coordinate equation

$$\left( \frac{1}{2M} \Delta_{\mathbf{X}} + \lambda_0(\mathbf{X}) \right) \Phi(\mathbf{X}) = E\Phi(\mathbf{X}). \quad (6.2.9)$$

The separation of nuclear and electron coordinates in the Born-Oppenheimer approximation is motivated from the assumption that the relatively speaking heavy nuclei move at a much slower speed than the light electrons (a valid assumption if nuclear and electron momenta are of comparable magnitude).

To motivate the Born-Oppenheimer molecular dynamics we first state the ansatz

$$\Phi_{\text{BO}}(\mathbf{X}, \mathbf{x}) = \phi_{\text{BO}}(\mathbf{X}, \mathbf{x}) e^{i\sqrt{M}\theta_{\text{BO}}(\mathbf{X})},$$

where, as before, we assume that the amplitude  $\phi_{\text{BO}} : \mathbb{R}^{3N} \times \mathbb{R}^{3n} \rightarrow \mathbb{C}$  and phase  $\theta_{\text{BO}} : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  are smooth functions which vary on a much slower scale than  $M$ . To conform with step 1. of the Born-Oppenheimer approximation, we assume that the electrons are in an eigenstate  $\psi_0(\mathbf{X}, \cdot)$  of (6.2.8) with eigenfunction  $\lambda_0(\mathbf{X})$  so that the amplitude takes the form  $\phi_{\text{BO}}(\mathbf{X}, \mathbf{x}) = \sqrt{\rho_{\text{BO}}}(\mathbf{X})\psi_0(\mathbf{X}, \mathbf{x})$  with unknown density function  $\rho_{\text{BO}} : \mathbb{R}^{3N} \rightarrow \mathbb{R}$ . Inserting  $\Phi_{\text{BO}}$  into equation (6.2.7), yields

$$\begin{aligned} 0 &= \left( \frac{1}{2M} \Delta_{\mathbf{X}} + \lambda_0(\mathbf{X}) - E \right) \Phi_{\text{BO}} \\ &= \left( \underbrace{\left( \frac{1}{2} |\nabla \theta_{\text{BO}}|^2 + \lambda_0 - E \right)}_{=: I} \phi_{\text{BO}} \right. \\ &\quad \left. - \frac{1}{2M} \Delta_{\mathbf{X}} \phi_{\text{BO}} - \frac{i}{M^{1/2}} (\nabla_{\mathbf{X}} \phi_{\text{BO}} \cdot \nabla \theta_{\text{BO}} + \frac{1}{2} \phi_{\text{BO}} \Delta \theta_{\text{BO}}) \right) e^{i\sqrt{M}\theta_{\text{BO}}(\mathbf{X})}. \end{aligned} \quad (6.2.10)$$

We approximately solve this equation by truncating the  $\mathcal{O}(M^{-1/2})$  terms and considering the remaining the Eikonal equation of term I:

$$\frac{1}{2} |\nabla \theta_{\text{BO}}(\mathbf{X})|^2 + \lambda_0(\mathbf{X}) - E = 0,$$

whose characteristics are given by the Hamiltonian system

$$\begin{aligned} \dot{\mathbf{X}}_i &= \partial_{\mathbf{P}_i} H_{\text{BO}}(\mathbf{X}, \mathbf{P}), & \dot{\mathbf{P}} &= -\partial_{\mathbf{X}_i} H_{\text{BO}}(\mathbf{X}, \mathbf{P}), \\ \text{where } H_{\text{BO}}(\mathbf{X}, \mathbf{P}) &:= \sum_{i=1}^N \frac{|\mathbf{P}_i|^2}{2M} + \lambda_0(\mathbf{X}), & \text{and } \mathbf{P}(t) &= \nabla_{\mathbf{X}} \theta_{\text{BO}}(\mathbf{X}(t)). \end{aligned} \quad (6.2.11)$$

To determine the density  $\rho_{\text{BO}}$  for the nuclear coordinates described by the dynamics (6.2.11), we first observe that conservation of mass implies that

$$\nabla \cdot (\rho_{\text{BO}}(\mathbf{X}) \nabla \theta_{\text{BO}}(\mathbf{X})) = 0.$$

By the conservation of mass and the dynamics (6.2.11) we derive the ordinary differential equation

$$\frac{d}{dt} \rho_{\text{BO}}(\mathbf{X}(t)) = \nabla \rho_{\text{BO}}(\mathbf{X}(t)) \cdot \nabla \theta_{\text{BO}}(\mathbf{X}(t)) = -\rho_{\text{BO}}(\mathbf{X}(t)) \Delta \theta_{\text{BO}}(\mathbf{X}(t)).$$

Introducing the integrating factor

$$\frac{d}{dt} \log G_{\text{BO}}(\mathbf{X}(t)) = \frac{1}{2} \Delta \theta_{\text{BO}}(\mathbf{X}(t)),$$

we obtain the following relation for the Born-Oppenheimer molecular dynamics density

$$\rho_{\text{BO}}(\mathbf{X}) = \frac{C}{G_{\text{BO}}^2(\mathbf{X})}.$$

The main study of Paper IV is the comparison of the densities

$$\rho_S(\mathbf{X}) = \frac{G^{-2}(\mathbf{X}) \langle \psi(\mathbf{X}, \cdot), \psi(\mathbf{X}, \cdot) \rangle}{\|G^{-1}\psi\|_{L^2(\mathbb{R}^{3N} \times \mathbb{R}^{3n})}} \quad \text{and} \quad \rho_{\text{BO}}(\mathbf{X}) = \frac{C}{G_{\text{BO}}^2(\mathbf{X})},$$

and we end this chapter by noting that Theorem 7.1 of Paper IV we prove that under some assumptions

$$\int_{\mathbb{R}^{3N}} g(X) \rho_{\text{BO}}(\mathbf{X}) d\mathbf{X} = \int_{\mathbb{R}^{3N}} g(X) \rho_S(\mathbf{X}) dX + \mathcal{O}(M^{-1+\delta}), \quad (6.2.12)$$

for any  $\delta > 0$  and  $g \in \mathcal{C}^3(\mathbb{R}^{3N})$ . The proof of (6.2.12) relies strongly on the stability of symplectic numerical methods for Hamiltonian-Jacobi equations derived by Szepessy et al. in [29].



# Bibliography

- [1] I. Babuška, A. Miller, and M. Vogelius. Adaptive methods and error estimation for elliptic problems of structural mechanics. In *Adaptive computational methods for partial differential equations (College Park, Md., 1983)*, pages 57–73. SIAM, Philadelphia, PA, 1983.
- [2] Christian Bayer, Håkon Hoel, Petr Plecháč, Anders Szepessy, and Raúl Tempone. How accurate is molecular dynamics? *arXiv:1104.0953*, 2011.
- [3] V. Bentkus. On the dependence of the Berry-Esseen bound on dimension. *J. Statist. Plann. Inference*, 113(2):385–402, 2003.
- [4] Andrew C. Berry. The accuracy of the Gaussian approximation to the sum of independent variates. *Trans. Am. Math. Soc.*, 49:122–136, 1941.
- [5] Fischer Black and Myron Scholes. The pricing of options and corporate liabilities. *Journal of Political Economy*, 81(3):637, 1973.
- [6] William L. Briggs, Van Emden Henson, and Steve F. McCormick. *A multigrid tutorial. 2nd ed.* Philadelphia, PA: SIAM, Society for Industrial and Applied Mathematics., 2000.
- [7] Russel E. Caflisch. Monte Carlo and quasi-Monte Carlo methods. *Acta Numerica*, 7(-1):1–49, 1998.
- [8] Jesper Carlsson, Kyoung-Sook Moon, Anders Szepessy, Raúl Tempone, and Georgios Zouraris. Stochastic differential equations: Models and numerics. Lecture notes used at KTH, 2010.
- [9] Y. S. Chow and Herbert Robbins. On the asymptotic theory of fixed-width sequential confidence intervals for the mean. *Ann. Math. Statist.*, 36:457–462, 1965.
- [10] R. H. Clarke. A statistical theory of mobile-radio reception. *Bell Sys. Tech.*, 47:957–1000, 1968.
- [11] Gregory Durgin. *Space-time wireless channels*. Prentice Hall Press, Upper Saddle River, NJ, USA, 2002.
- [12] Richard Durrett. *Probability: theory and examples*. Duxbury Press, Belmont, CA, second edition, 1996.
- [13] Kenneth Eriksson, Don Estep, Peter Hansbo, and Claes Johnson. Introduction to adaptive methods for differential equations. In *Acta numerica, 1995*, Acta Numer., pages 105–158. Cambridge Univ. Press, Cambridge, 1995.

- [14] Carl-Gustav Esseen. Fourier analysis of distribution functions. A mathematical study of the Laplace-Gaussian law. *Acta Math.*, 77:1–125, 1945.
- [15] L. D. Faddeev and O. A. Yakubovskii. *Lectures on quantum mechanics for mathematics students*. Translated by Harold McFaden. Student Mathematical Library 47. Providence, RI: American Mathematical Society (AMS). xii, 2009.
- [16] Michael B. Giles. Multilevel Monte Carlo path simulation. *Oper. Res.*, 56(3):607–617, 2008.
- [17] Paul Glasserman. *Monte Carlo methods in financial engineering*, volume 53 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 2004. Stochastic Modelling and Applied Probability.
- [18] Stephen J. Gustafson and Israel Michael Sigal. *Mathematical concepts of quantum mechanics*. 2nd ed. Universitext. Berlin: Springer. xiii, 2011.
- [19] Ernst Hairer, Christian Lubich, and Gerhard Wanner. *Geometric numerical integration. Structure-preserving algorithms for ordinary differential equations*. Reprint of the second 2006 ed. Springer Series in Computational Mathematics, 31. Berlin: Springer. xviii, 2010.
- [20] Håkon Hoel, Erik Schwerin, Anders Szepessy, and Raúl Tempone. Adaptive multilevel monte carlo simulation. In *Numerical Analysis of Multiscale Computations*, volume 82 of *Lecture Notes in Computational Science and Engineering*, pages 217–234. Springer Berlin Heidelberg, 2012.
- [21] E. Jouini, J. Cvitanić, and Marek Musiela, editors. *Option pricing, interest rates and risk management*. Handbooks in Mathematical Finance. Cambridge University Press, Cambridge, 2001.
- [22] Peter E. Kloeden and Eckhard Platen. *Numerical solution of stochastic differential equations*, volume 23 of *Applications of Mathematics (New York)*. Springer-Verlag, Berlin, 1992.
- [23] Elliott H. Lieb and Robert Seiringer. *The stability of matter in quantum mechanics*. Cambridge: Cambridge University Press. xv, 2009.
- [24] Robert C. Merton. Theory of rational option pricing. *The Bell Journal of Economics and Management Science*, 4(1):141–183, 1973.
- [25] Kyoung-Sook Moon, Anders Szepessy, Raúl Tempone, and Georgios E. Zouraris. Convergence rates for adaptive weak approximation of stochastic differential equations. *Stoch. Anal. Appl.*, 23(3):511–558, 2005.
- [26] Kyoung-Sook Moon, Erik von Schwerin, Anders Szepessy, and Raúl Tempone. An adaptive algorithm for ordinary, stochastic and partial differential equations. In *Recent advances in adaptive computation*, volume 383 of *Contemp. Math.*, pages 325–343. Amer. Math. Soc., Providence, RI, 2005.
- [27] John C. Neu. *Training manual on transport and fluids*. Providence, RI: American Mathematical Society (AMS), 2010.
- [28] Bernt Øksendal. *Stochastic differential equations*. Universitext. Springer-Verlag, Berlin, sixth edition, 2003. An introduction with applications.

- [29] Mattias Sandberg and Anders Szepessy. Convergence rates of symplectic Pontryagin approximations in optimal control theory. *M2AN Math. Model. Numer. Anal.*, 40(1):149–173, 2006.
- [30] Anders Szepessy, Raúl Tempone, and Georgios E. Zouraris. Adaptive weak approximation of stochastic differential equations. *Comm. Pure Appl. Math.*, 54(10):1169–1214, 2001.
- [31] Denis Talay and Luciano Tubaro. Expansion of the global error for numerical schemes solving stochastic differential equations. *Stochastic Anal. Appl.*, 8(4):483–509 (1991), 1990.
- [32] David Tse and Pramod Viswanath. *Fundamentals of wireless communication*. Cambridge University Press, New York, NY, USA, 2005.
- [33] Wikipedia. Nyquist-shannon sampling theorem. [http://en.wikipedia.org/wiki/Nyquist-Shannon\\_sampling\\_theorem](http://en.wikipedia.org/wiki/Nyquist-Shannon_sampling_theorem), 2012 (accessed April 10, 2012).





**Part II**  
**Included Papers**



# Paper I

# GAUSSIAN COARSE GRAINING OF A MASTER EQUATION EXTENSION OF CLARKE'S MODEL

HÅKON HOEL AND HENRIK NYBERG

ABSTRACT. We study the error and computational cost of generating output signal realizations for the channel model of a moving receiver in a scattering environment, as in Clarke's model, with the extension that scatterers randomly flip on and off. At micro scale, the channel is modeled by a Multipath Fading Channel (MFC) model, and by coarse graining the micro scale model we derive a macro scale Gaussian process model. Four algorithms are presented for generating stochastic signal realizations, one for the MFC model and three for the Gaussian process model. A computational cost comparison of the presented algorithms indicates that Gaussian process algorithms generate signal realizations more efficiently than the MFC algorithm does. Numerical examples of generating signal realizations in time independent and time dependent scattering environments are given, and the problem of estimating model parameters from real life signal measurements is also studied.

## 1. INTRODUCTION

We consider the Multipath Fading Channel (MFC) model with a transmitter fixed and the receiver moving with a constant speed in an urban environment with buildings obstructing the line of sight between scatterer and receiver. Incoming rays at the receiver are thus modeled as scattered off the receiver's surroundings. Looking at scenarios where the distance between transmitter and receiver is large and the majority of scattering surfaces are flat walls, the vertical angle of arrival of incoming rays at the receiver is assumed to be 0 degrees, cf. Figure 1, i.e. scatterers are assumed to lie in the horizontal plane. The receiver receives  $M$  incoming signal rays whose horizontal angle of arrival  $\{\alpha_m\}_{m=1}^M$  are distributed according to a prescribed scatterer density  $p : [0, 2\pi) \rightarrow \mathbb{R}_+$ , and the resulting baseband output signal is modeled by

$$(1) \quad Z_{t,M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M a(\alpha_m, t) e^{-i(2\pi f_c \tau(\alpha_m, t) + \theta_m(t))}.$$

Here  $f_c$  denotes the carrier frequency,  $\tau$  the delay function, and  $a(\alpha_m, t)$  and  $\theta_m(t)$  are respectively amplitude and phase shift random processes further described in Subsection 1.1. Two different delay functions are used in this paper; the full delay function depending on an explicit description of the scattering boundary, and  $\tau(\alpha, t) = -vt \cos(\alpha)/c$ , with  $v$  and  $c$  denoting the speed of the moving receiver and the speed light, respectively. The latter delay function is a first order approximation of the full delay function, cf. Appendix B.

---

2000 *Mathematics Subject Classification.* Primary 94A12; Secondary 60G15.

*Key words and phrases.* Wireless channel modeling; signal theory; master equations; Gaussian processes.

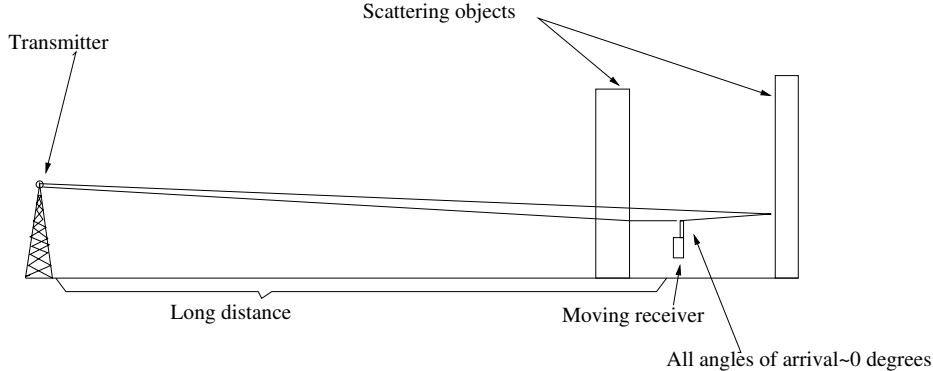


FIGURE 1. Illustration of the typical scattering environment which we wish to model the wireless channel in.

The choice  $\tau(\alpha, t) = -vt \cos(\alpha)/c$ , is common in channel modeling since it results in a Wide Sense Stationary (WSS) channel model, i.e. a model whose signal realizations  $Z_{t,M}$  have constant expected value as a function of time and the autocorrelation function  $E[Z_{t+\Delta t, M} Z_{t, M}^*]$  is a function of only  $\Delta t$ , and therefore are, relatively speaking, easier to analyze.

To put the output of (1) in context, a more general input/output baseband signal representation is given by

$$Z_t = \int_{\mathbb{R}} X(t - \tau) h(t, \tau) d\tau,$$

where  $X(t)$  denotes the baseband input signal and  $h(t, \tau)$  the impulse response function. For the output in (1), the baseband input is  $X := 1$ , which corresponds to the zero bandwidth passband input signal  $X(t) = \exp(-i2\pi f_c t)$  demodulated to a 0 frequency signal, and the impulse response function is given by

$$h(t, \tau) = \frac{1}{\sqrt{M}} \sum_{m=1}^M a(\alpha_m, t) e^{-i(2\pi f_c \tau(\alpha_m, t) + \theta_m)} \delta(t - \tau(\alpha_m, t))$$

with  $\delta$  denoting the Dirac delta function. Note that although we model channels with single frequency input signal, Doppler effects deriving from the receiver moving relative to its surroundings will result in output signals having non-zero frequency bandwidth.

**1.1. The amplitude random process.** Local shadowing by reflecting objects in motion, for example cars, pedestrians and shaking leaves, causes scatterers to flip from being active to passive and vice versa. Attempting to include local shadowing in our MFC model, we define the amplitude random flip process  $a(\alpha, t)$  which flips on when it changes value from 0 to  $a^+(\alpha, t) \geq 0$  and off when the opposite change occurs. It is here assumed that the mapping  $a^+ : [0, 2\pi) \times \mathbb{R} \rightarrow \mathbb{R}_+$ , which represents the active state of a reflector, is piecewise continuous (it may for example be piecewise constant or depend on the distance from scatterer to receiver). The flip process (of a scatterer) is modeled as Poisson process with constant flip rate  $C$ :

$$P(a(\alpha_m, t) \text{ flips } k \text{ times on time step } \Delta t) = \frac{(C\Delta t)^k \exp(-C\Delta t)}{k!},$$

where flips are independent from the phase shift processes  $\{\theta_m(t)\}_{m=1}^M$  and from the scatterers' state  $\bigotimes_{m=1}^M \{0, a^+(\alpha_m, t)\}$ . The scatterers' initial state  $\{a(\alpha_m, 0)\}_{m=1}^M$  is sampled according to the i.i.d. Bernoulli distribution  $P(a(\alpha_m, 0) = a^+(\alpha, 0)) = 1/2$  and  $P(a(\alpha_m, 0) = 0) = 1/2$  which is consistent with the steady state distribution as  $t \rightarrow \infty$ . For an illustration of the effect of a positive flip rate on an output signal realization and a comparison to a real life signal measurement, see Figure 2.

The phase shift processes  $\{\theta_m(t)\}_{m=1}^M$  are at all times i.i.d. uniformly distributed in  $[0, 2\pi)$ . This is motivated from the assumption that wave lengths are very small compared to the wave paths' travel distance from transmitter to receiver. The phase of  $\theta_m(t)$  is updated by sampling  $\theta_m(t) \sim U(0, 2\pi)$  every time scatterer  $m$  flips, since we assume that a scatterer flipping at a given angle  $\alpha_m$  is equivalent to a new scatterer appearing at the given angle; a new scatterer requiring a new i.i.d. phase shift uniformly distributed in  $[0, 2\pi)$ .

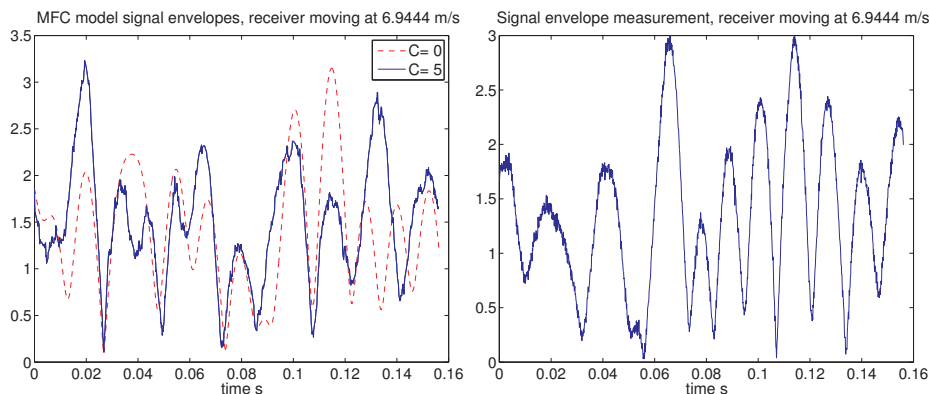


FIGURE 2. **Left plot:** Signal envelopes of two computer generated signal realizations using Algorithm 1 with the same random seed initialization. Both realizations are generated with  $f_c = 1.8775\text{GHz}$ ,  $v = 6.944\text{m/s}$ ,  $a^+ = 2$ ,  $p = (2\pi)^{-1}$  and  $\tau(\alpha, t) = -vt \cos(\alpha)/c$ . The only varying model parameter is the flip rate with  $C = 0$  (dashed line) and  $C = 5$  (full line). **Right plot:** Measured urban environment signal envelope whose carrier frequency and receiver speed are identical to the corresponding values for the realizations in the left plot.

**Remark 1.1.** *The presence of measurement noise in the data might be difficult to distinguish from flipping scatterer noise and can affect modeling parameter estimates, such as the flip rate. In our model, we assume that measurement noise is negligible relative to the noise generated by flipping scatterers.*

**1.2. Motivation for Gaussian processes model.** The MFC model with flips introduced here can be linked to the Master Equation, which is an equation often used in chemistry to describe the evolution of state space probabilities, cf. [13]. In a Master Equation setting, one typically assigns finite state spin variables on a lattice with probabilistic spin interaction dynamics and describe the time evolution of the probability to occupy each one of a discrete set of states through a differential

equation called the Master Equation. In our setting, we have the lattice  $\{\alpha_m\}_{m=1}^M$  and the possibly time dependent state space  $\bigotimes_{m=1}^M \{0, a^+(\alpha_m, t)\}$ .

Creating discrete signal realizations of (1) straightforwardly by the MFC algorithm described in Section 2 is computationally very costly. This motivated us to try to construct an algorithm which reduced the computational cost while at the same time preserved the desired signal properties by using a new result in the Master Equation setting: In [11], Katsoulakis and Szepeszy developed theory for coarse graining two state spin variables on a micro scale lattice into a macro scale Stochastic Differential Equation (SDE) representation which reduces the computational cost for such problems considerably. Using their theory, we developed a similar transition from the MFC model with flip state space on a micro scale lattice to a coarse grained SDE supposed to represent the output signal. However, comparing signal realizations of the Master Equation developed SDE to fine realizations of the MFC algorithm it was clear that the phase shift processes  $\{\theta_m(t)\}_{m=1}^M$  were not resolved in the SDE model. This observation made us believe that a Master Equation SDE is not suitable coarse graining of (1) and lead us to instead try coarse graining with Gaussian processes, which are more general stochastic processes than SDEs.

The first outcome of coarse graining with Gaussian processes is Theorem 3.4 which shows that as  $M \rightarrow \infty$ , the signal  $Z_{t,M}$  converges in distribution to a Gaussian process  $Z_t$ . Based on this theorem, we develop three algorithms, Algorithm 2, 3 and 4, using covariance and spectral properties to generate realizations of the limit Gaussian process. The developed algorithms are studied in terms of accuracy and computational cost, and a summary of this study, presented in Section 5, indicate that that the Gaussian process algorithms generate signal realizations more efficiently than the MFC algorithm, Algorithm 1, does.

**1.3. Related works and historical remarks.** In 1968, Clarke introduced the MFC model now known as Clarke's model in his seminal paper [7]. He considered the superposition of  $M$  incoming waves

$$(2) \quad \xi_{t,M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M e^{-i(2\pi f_c v \cos(\alpha_m)t/c + \theta_m)},$$

where, as before,  $f_c$  is the carrier frequency,  $\{\theta_m\}_{m=1}^M$  are i.i.d. initial phase shifts and  $\alpha_m$  is the arrival angle of the  $m^{\text{th}}$  component wave distributed according to a given azimuth density  $p(\alpha)$ . Note in particular that this is a WSS model with delay function  $\tau(\alpha_m, t) = -v \cos(\alpha_m)t/c$  which is identical to the delay function we mainly use in our MFC model. Considering the scenario with angle density  $p(\alpha) = (2\pi)^{-1}$ , Clarke noted that the auto correlation function,  $E[\xi_{t,M}\xi_{0,M}^*]$ , converges to the zeroth-order Bessel function of the first kind,  $J_0(2\pi f_c vt/c)$ , as  $M \rightarrow \infty$ . And, further, he showed that its Power Spectral Density (PSD), which describes the frequency spread around the carrier frequency and is defined by the Fourier transform of the autocorrelation function, is on the form

$$(3) \quad S(f) = \begin{cases} \frac{c}{\pi\sqrt{(vf_c)^2 - (cf)^2}} & |f| < vf_c/c \\ 0 & |f| \geq vf_c/c. \end{cases}$$

This particular PSD is often referred to as Jakes' spectrum. For plots of these results, see Figure 3.

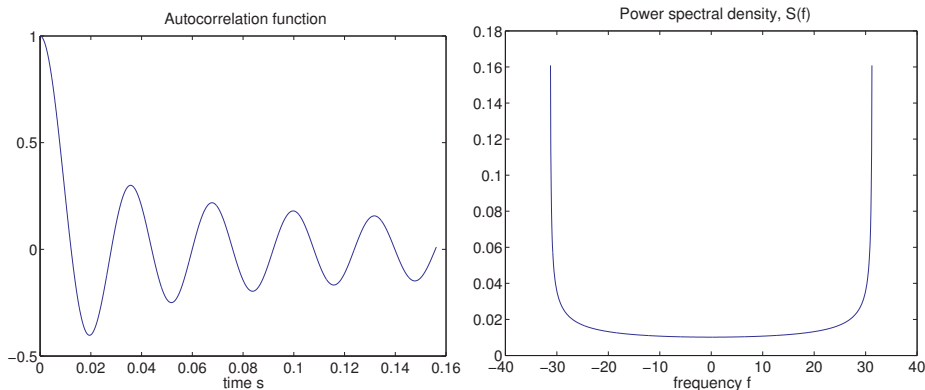


FIGURE 3. **Left plot:** The autocorrelation function for Clarke's model with azimuth density  $p(\alpha) = (2\pi)^{-1}$ ,  $v = 5\text{m/s}$  and  $f_c = 1.8775\text{GHz}$ . **Right plot:** The power spectral density of Clarke's model, often called Jakes' spectrum, with the same model parameters as for the left plot.

Among the papers linking the MFC model to SDE which motivated this work, Feng, Field and Haykin [10] studied the output signal

$$(4) \quad \xi_{t,M} = \sum_{m=1}^M a_m e^{i\phi_m(t)}.$$

Here  $a_m = O(M^{-1/2})$  are constant amplitudes and  $\phi_m$  are phases with uniform random initialization on the interval  $[0, 2\pi)$ , and, with  $dW_{t,m}$  denoting the Wiener increment and  $B$  a positive constant, they assumed the phases also satisfied the increment relation

$$(5) \quad d\phi_m(t) = \sqrt{B}dW_{t,m}.$$

Since this is a model for stationary wireless channels, the phases do not include Doppler shift terms. As  $M \rightarrow \infty$ , they derived that  $\xi_{t,M}$  converges to a complex valued Ornstein-Uhlenbeck process.

In [9], Feng and Field consider the MFC model with output

$$\xi_{t,M} = \sum_{m=1}^M a_m \exp\left(i(2\pi f_m t + \phi_t^{(m)})\right).$$

where the amplitudes  $a_m$  are i.i.d. random variables and  $f_m = f_D \cos(\alpha_m)$  are Doppler shifts with  $f_D$  denoting the maximum Doppler shift and  $\alpha_m$  the angle. The phases  $\phi_t^{(m)}$  are independent Wiener processes with uniform initial distribution in  $[0, 2\pi)$ :

$$d\phi_t^{(m)} = \sqrt{B}dW_t^{(m)}, \quad \phi_0^{(m)} \sim U[0, 2\pi),$$

where  $B$  is a constant with the dimension of frequency. With this model they obtain the auto-correlation function

$$E[\varepsilon_{t,M}\varepsilon_{0,M}^*] = \sum_{m=1}^M E[a_m^2]e^{-B|t|/2}J_0(2\pi f_D t),$$



and a corresponding PSD which they claim resemble measurements more than Clarke's model PSD. Although our modeling assumptions are quite different than Feng and Field's, we obtain similar autocorrelation and PSD results with our Gaussian Process model, cf. Section 4.

The development of flipping scatterers in our model is linked to the application of Poisson counting process as birth-death modeling of wave paths. Among the papers considering such models is Charalambous et al. [5] which study MFC models with output

$$\xi(t) = \sum_{m=1}^{M(T_s)} a_m(\tau_m) e^{i\phi_m(\tau_m, t)} x(t - \tau_m),$$

where  $x(t)$  is the input signal, and  $\tau_m$ ,  $a_m(\tau_m)$  and  $\phi_m(\tau_m, t)$  denote the propagation time delay, amplitude, and phase, respectively. The number of received paths,  $M(T_s)$ , is a Poisson counting process with  $T_s > 0$  being a fixed stopping time. Under certain conditions, they obtain explicit expressions for the autocorrelation, PSD and moment-generating functions. In [16], Zwick et al. develop a temporally dynamic indoor channel model where birth and deaths of active wave paths are used to model the varying scattering environment, and in [6], Chong et al. introduce an indoor channel model where births and deaths of wave paths are generated by a Markov transition matrix.

All papers mentioned above present channel models with relatively few parameters, but there are also reports which have a highly physical approach thereby needing many parameters. The industry standard modeling report by 3GPP [1], considers reflectors gathered in clusters and models the received signal as a superposition of wave paths hitting the reflectors. Among its input parameters are: transmitter antenna and receiver antenna orientation, line of sight angle of departure from transmitter to receiver, angle of departure for every path from transmitter, angle of arrival for every path at receiver, transmitter velocity, angle of transmitter velocity, etc.

The rest of this paper is organized as follows. In Section 2, we present a numerical algorithm for generating signal realizations of our proposed MFC model. Section 3 motivates theoretically approximating the output signal by a Gaussian process model and two algorithms for generating Gaussian process signal realizations based on the covariance matrix is developed. An error and complexity analysis of Algorithm 1, 2 and 3 is also included. In Section 4, we investigate the relation between the signal's autocorrelation and PSD for WSS signals. A method for estimating the flip rate and scatterer density from PSD measurements is described, and a PSD based algorithm for generating Gaussian process signal realizations is presented. Section 5 summarizes the complexity analysis results obtained for the five algorithms considered in this paper. Section 6 concludes the paper with a numerical examples of signal realizations for time independent and time dependent scattering environments.

## 2. THE MFC ALGORITHM

We first present a numerical algorithm for generating signal realizations by the MFC model equation (1) on a set of sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ .

**Remark 2.1.** *It is possible to extend the MFC model and algorithm to scenarios including line of sight wave components, often called specular rays. Suppose you*

**Algorithm 1** The MFC algorithm

---

**Input:** Amplitude function  $a^+$ , flip rate  $C$ , carrier frequency  $f_c$ , scatterer density  $p$ , receiver speed  $v$ , sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ .

**Output:** Signal realization  $Z_{\mathbf{t},M} = (Z_{t_1,M}, Z_{t_2,M}, \dots, Z_{t_N,M})$ .

Generate a set of i.i.d angles of arrival  $\{\alpha_k\}_{k=1}^M$  distributed according to the density  $p(\alpha)$ .

Generate a set of i.i.d. phase shifts  $\{\theta_k(0)\}_{k=1}^M$  with  $\theta_k(0) \sim U(0, 2\pi)$ .

Generate the initial state of the amplitudes  $\{a(\alpha_k, t_1)\}_{k=1}^M$  which are i.i.d. restricted to the steady state initial condition  $P(a(\alpha, t_1) = 0) = P(a(\alpha, t_1) = a^+(\alpha, t_1)) = 1/2$ .

Compute  $Z_{t_1,M}$  according to (1).

**for**  $j = 2$  to  $N$  **do**

**for**  $k = 1$  to  $M$  **do**

    Generate  $n_k \sim \text{Poisson}(C(t_j - t_{j-1}))$ .

    Flip the value of  $a(\alpha_k, t_j)$   $n_k$  times.

    If  $n_k > 0$ , update the phase shift process by generating a new random phase shift  $\theta_k(t_j) \sim U(0, 2\pi)$ .

**end for**

  Compute  $Z_{t_j,M}$  according to (1).

**end for**

---

have the input/output relation consisting of many diffuse ray contributions in  $Z_{t,M}$  and one specular ray incoming from the angle 0 with amplitude  $V$ . Then an output signal can be generated by

$$Z_{t,M} + V e^{-i(2\pi f_c \tau(0,t))},$$

where  $V e^{-i(2\pi f_c \tau(\alpha,t))}$  is a deterministic term modeling the specular ray contribution.

### 3. STOCHASTIC MODEL FOR THE SIGNAL WITH STATIC SCATTERERS

In this section we will show that the normalized signal  $Z_{t,M}$  of equation 1 converges in distribution to a complex Gaussian process as  $M \rightarrow \infty$ . Thereafter, the covariance of the limit Gaussian process is derived and used to construct an algorithm for generating signal realizations in Algorithm 2. But first, let us recall the definitions of multivariate complex normal distributions and Gaussian processes.

**Definition 3.1** (Multivariate complex normal distribution I). *Suppose  $X$  and  $Y$  are random vectors in  $\mathbb{R}^n$  such that  $(X, Y)$  is a  $2n$ -dimensional normal random vector. Then we say that  $Z = X + iY$  is complex normal distributed and its distribution is described by the mean  $\mu = E[Z]$ , the covariance matrix  $K = E[(Z - \mu)(Z - \mu)^H]$  and the pseudo-covariance matrix  $J = E[(Z - \mu)(Z - \mu)^T]$ , where  $T$  denotes the transpose operator and  $H$  the Hermitian operator. We write  $Z \sim \mathcal{N}_{\mathbb{C}}(\mu, K, J)$ .*

An alternative definition for multivariate complex normal distributions, which often is easier to work with, derives from the one-to-one relation between the characteristic function and the distribution:

**Definition 3.2** (Multivariate complex normal distribution II). *A random vector  $Z = (Z_1, Z_2, \dots, Z_n) \in \mathbb{C}^n$  is said to be complex multivariate normal distributed if*

the linear combination of its components,  $\mathbf{c}^H Z \in \mathbb{C}^1$  is complex normal distributed for all  $\mathbf{c} \in \mathbb{C}^n$ .

**Definition 3.3** (Complex Gaussian process). *A complex Gaussian process is a stochastic process  $\{Z_t\}_{t \in [0, T]}$ ,  $Z_t \in \mathbb{C}^1$ , for which any finite length sample vector  $(Z_{t_1}, Z_{t_2}, \dots, Z_{t_n})$  with  $0 \leq t_1 \leq t_2 \leq \dots \leq t_n < T$  is complex normal distributed.*

With the aid of the Central Limit Theorem (CLT), we now show that  $Z_{t, M}$  converges in distribution to a Gaussian process as  $M \rightarrow \infty$ .

**Theorem 3.4** (Model (1)'s distributional convergence to a Gaussian process). *Assume the amplitude function  $a^+(\cdot, t)$  and the delay function  $\tau(\cdot, t)$  is bounded and piecewise continuous on  $[0, 2\pi)$  for all times  $t \in [0, T)$ . Then the signal*

$$Z_{t, M} := \frac{1}{\sqrt{M}} \sum_{m=1}^M a(\alpha_m, t) e^{-i(2\pi f_c \tau(\alpha_m, t) + \theta_m(t))}$$

converges in distribution to a complex Gaussian process  $Z_t$  as  $M \rightarrow \infty$ .

*Proof.* Letting  $*$  denote the complex conjugate, definitions 3.2 and 3.3 imply that if sums of the kind

$$(6) \quad \Upsilon_{\mathbf{t}, \mathbf{c}, M} := \sum_{i=1}^n c_i^* Z_{t_i, M},$$

converge in distribution to a complex normal for all finite length sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_n) \subset [0, T)$  and complex valued vectors  $\mathbf{c} = (c_1, c_2, \dots, c_n)$ , then  $Z_{t, M}$  converges in distribution to a complex Gaussian process on  $[0, T)$ .

Writing out equation (6), we have

$$\Upsilon_{\mathbf{t}, \mathbf{c}, M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M \underbrace{\sum_{j=1}^n c_j^* a(\alpha_m, t_j) e^{-i(2\pi f_c \tau(\alpha_m, t_j) + \theta_m(t_j))}}_{=:\xi_m}.$$

Since both  $\{\theta_m(\cdot)\}_m$ ,  $\{a(\alpha_m, \cdot)\}_m$  are i.i.d. and mutually independent, the r.v.  $\{\xi_m\}_m$  are also i.i.d. with mean

$$(7) \quad \mu = E[\xi_m] = \sum_{j=1}^n c_j^* E \left[ \underbrace{E[a(\alpha_m, t_j) | \alpha_m] E[e^{-i(2\pi f_c \tau(\alpha_m, t_j) + \theta_m(t_j))} | \alpha_m]}_{=0} \right] = 0.$$

Before computing the covariance and pseudo-covariance of  $\xi_m$ , let us derive some useful properties. By the definition of the phase shift processes and amplitude process given in Subsection 1.1, we see that

$$P(\theta_m(t_j) = \theta_m(t_k)) = e^{-C|t_j - t_k|},$$

and

$$\begin{aligned} & E[a(\alpha_m, t_j) a(\alpha_m, t_k) | \alpha_m, \theta_m(t_j) = \theta_m(t_k)] \\ &= a^+(\alpha_m, t_j) a^+(\alpha_m, t_k) P(a(\alpha_m, t_j) = a^+(\alpha_m, t_k)) \\ &= \frac{a^+(\alpha_m, t_j) a^+(\alpha_m, t_k)}{2}. \end{aligned}$$

Consequently,

$$\begin{aligned}
(8) \quad & E \left[ a(\alpha_m, t_j) a(\alpha_m, t_k) e^{i(2\pi f_c(\tau(\alpha_m, t_k) - \tau(\alpha_m, t_j)) + \theta_m(t_k) - \theta_m(t_j))} \middle| \alpha_m \right] \\
&= P(\theta_m(t_j) = \theta_m(t_k)) E[a(\alpha_m, t_j) a(\alpha_m, t_k) | \alpha_m, \theta_m(t_j) = \theta_m(t_k)] e^{i2\pi f_c(\tau(\alpha_m, t_k) - \tau(\alpha_m, t_j))} \\
&= e^{-C|t_j - t_k|} \frac{a^+(\alpha_m, t_j) a^+(\alpha_m, t_k)}{2} e^{i2\pi f_c(\tau(\alpha_m, t_k) - \tau(\alpha_m, t_j))}.
\end{aligned}$$

and

$$(9) \quad E[a(\alpha_m, t_j) a(\alpha_m, t_k) e^{-i(2\pi f_c(\tau(\alpha_m, t_k) + \tau(\alpha_m, t_j)) + \theta_m(t_k) + \theta_m(t_j))} | \alpha_m] = 0.$$

The covariance of  $\xi_m$  is derived using (8)

$$\begin{aligned}
(10) \quad & K = E[|\xi_m|^2] \\
&= \sum_{j,k=1}^n c_j^* c_k E \left[ E \left[ a(\alpha_m, t_j) a(\alpha_m, t_k) e^{i(2\pi f_c(\tau(\alpha_m, t_k) - \tau(\alpha_m, t_j)) + \theta_m(t_k) - \theta_m(t_j))} \middle| \alpha_m \right] \right] \\
&= \sum_{j,k=1}^n c_j^* c_k \frac{e^{-C|t_j - t_k|}}{2} \int_0^{2\pi} a^+(\alpha, t_j) a^+(\alpha, t_k) e^{i2\pi f_c(\tau(\alpha, t_k) - \tau(\alpha, t_j))} p(\alpha) d\alpha,
\end{aligned}$$

and the pseudo-covariance from using (9),

$$\begin{aligned}
(11) \quad & J = E[\xi_m^2] \\
&= \sum_{j,k=1}^n c_j^* c_k E \left[ E \left[ a(\alpha_m, t_j) a(\alpha_m, t_k) e^{-i(2\pi f_c(\tau(\alpha_m, t_k) + \tau(\alpha_m, t_j)) + \theta_m(t_k) + \theta_m(t_j))} \middle| \alpha_m \right] \right] \\
&= 0.
\end{aligned}$$

Having shown that  $\Upsilon_{\mathbf{t}, \mathbf{c}, M} = \sum_{m=1}^M \xi_m / \sqrt{M}$  is a sum of i.i.d. complex valued random variables  $\xi_m$  with mean  $\mu = E[\xi_m] = 0$ , pseudo-covariance  $J$  and the bounded covariance  $K$  as given in equation (11) and (10), respectively, it follows by the CLT for complex valued random variables that the normalized sums  $\Upsilon_{\mathbf{t}, \mathbf{c}, M}$  converge in distribution to the complex normal  $\mathcal{N}_{\mathbb{C}}(\mu, J, K)$  as  $M \rightarrow \infty$ .  $\square$

Having proved that  $Z_{t, M}$  converges in distribution to a complex Gaussian process  $Z_t$ , our next goal is to construct an algorithm that generates realizations  $Z_{\mathbf{t}} = (Z_{t_1}, Z_{t_2}, \dots, Z_{t_N})^T$  of the process sampled at a set of times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ . To achieve that goal, we must first describe  $Z_{\mathbf{t}}$  in terms of its mean, pseudo-covariance and covariance. Consider, as in the proof of Theorem 3.4, a sum  $\Upsilon_{\mathbf{t}, \mathbf{c}} = \lim_{M \rightarrow \infty} \sum_{j=1}^N c_j^* Z_{t_j, M}$  for a complex valued vector  $\mathbf{c} = (c_1, c_2, \dots, c_N)$ . By choosing  $\mathbf{c}$  such that  $c_j = \delta_{jk}$ , with  $\delta_{jk}$  denoting the Kronecker delta, it follows from equation (7) and the proof of Theorem 3.4 that

$$\mu_k = E[Z_{t_k}] = E[\Upsilon_{\mathbf{t}, \mathbf{c}}] = 0, \quad \text{for all } k \in \{1, 2, \dots, N\}.$$

Choosing instead  $c_j = c_k = 1$  and the other elements of  $\mathbf{c}$  equal to 0, it follows from equations (11) and (10) that  $J(t_j, t_k) = 0$  and

$$(12) \quad K(t_j, t_k) = \frac{e^{-C|t_j - t_k|}}{2} \int_0^{2\pi} a^+(\alpha, t_j) a^+(\alpha, t_k) e^{i2\pi f_c(\tau(\alpha, t_k) - \tau(\alpha, t_j))} p(\alpha) d\alpha$$

for all  $j, k \in \{1, 2, \dots, N\}$ .

When the pseudo-covariance  $J = 0$ , the multivariate complex normal is said to be circular symmetric and, for simplicity, its distribution is represented by  $\mathcal{N}_{\mathbb{C}}(\mu, K)$ . An  $N$ -dimensional circular symmetric complex normal  $Z \sim \mathcal{N}_{\mathbb{C}}(\mu, K)$  has the density representation

$$P(Z = z) = \frac{e^{-(z-\mu)^H K^{-1}(z-\mu)}}{\det(K)\pi^N}.$$

For the sampled limit complex Gaussian process studied here, we thus see that  $Z_{\mathbf{t}} \sim \mathcal{N}_{\mathbb{C}}(0, K)$  with the terms of  $K$  given by (12).

**3.1. Gaussian process algorithm based on the covariance matrix.** We now present an algorithm that generates Gaussian process signal realizations on a set of  $N$  sampling times  $\mathbf{t}$  by multiplying the square root of the covariance matrix to a vector of i.i.d. standard complex normals, see (13). This is a standard way of generating Gaussian process realizations, c.f. [2], which requires having predetermined the covariance matrix  $K$  integral terms given in (12). Generally, however, these covariance integral terms are not solvable by pen and paper, so we approximate them using numerical integration. This results in an approximation of the covariance matrix  $K$ , which we denote  $\bar{K}$ . In the algorithm to be presented, Algorithm 2,  $\bar{K}$  is used to generate Gaussian process signal realizations  $\bar{Z}_{\mathbf{t}} = (\bar{Z}_{t_1}, \bar{Z}_{t_2}, \dots, \bar{Z}_{t_N})$ .

---

**Algorithm 2** Covariance matrix based Gaussian process algorithm

---

**Input:** Amplitude function  $a^+$ , flip rate  $C$ , carrier frequency  $f_c$ , scatterer density  $p$ , receiver speed  $v$ , sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ .

**Output:** Gaussian process realization  $\bar{Z}_{\mathbf{t}} = (\bar{Z}_{t_1}, \bar{Z}_{t_2}, \dots, \bar{Z}_{t_N})$ .

**for**  $i = 1$  to  $N$  **do**

**for**  $j = 1$  to  $N$  **do**

    Approximate the term  $K(t_i, t_j)$  of equation (12) with quadrature and store the value in the approximate covariance matrix  $\bar{K}(t_i, t_j)$ .

**end for**

**end for**

Singular value decompose  $\bar{K} = \bar{U} \bar{S} \bar{U}^H$ , and compute  $\bar{K}^{1/2} = \bar{U} \bar{S}^{1/2}$ .

Generate a vector of  $N$  i.i.d. standard complex normal elements;  $\hat{Z} \sim \mathcal{N}_{\mathbb{C}}(0, I_N)$ , and thereafter the sampled output signal realization

$$(13) \quad \bar{Z}_{\mathbf{t}} = \bar{K}^{1/2} \hat{Z}.$$


---

**Remark 3.5.** *As it were for the MFC algorithm, it is also possible to extend the covariance based Gaussian process algorithms to scenarios with specular ray contributions. Assume you have the input/output relation consisting of many diffuse ray contributions modeled by the Gaussian process output  $\bar{Z}_{\mathbf{t}}$  and one specular ray incoming from the angle  $\alpha$  with amplitude  $V$  by*

$$\bar{Z}_{\mathbf{t}} + V e^{-i(2\pi f_c \tau(\alpha, t))}.$$

**3.2. Error estimates and computational complexity.** Our next objective is to estimate the computational cost of generating signal realizations on a time grid  $\mathbf{t} = (t_1, t_2, \dots, t_N) \subset [0, T)$  satisfying an accuracy requirement when using either Algorithm 1 or 2. We first introduce a measure for how well, in distributional sense, an algorithm's stochastic realizations approximates the sampled limit Gaussian process  $Z_{\mathbf{t}}$ .

**Definition 3.6** (Distributional signal error measure). *Let  $K$  denote the covariance matrix of the sampled Gaussian process  $Z_{\mathbf{t}} = (Z_{t_1}, Z_{t_2}, \dots, Z_{t_N})$  expressed in equation (12) and let  $\bar{Z}_{\mathbf{t}}$  denote a generated signal realization that approximates  $Z_{\mathbf{t}}$  in distributional sense. Then, under the assumption that  $K$  is non-singular, the distributional error of  $\bar{Z}_{\mathbf{t}}$  is defined by*

$$e(\bar{Z}_{\mathbf{t}}) := \sup_{A \in \mathfrak{C}(\mathbb{C}^N)} |P(\bar{Z}_{\mathbf{t}} \in A) - P(Z_{\mathbf{t}} \in A)|,$$

where  $\mathfrak{C}(\mathbb{C}^N)$  denotes the class of convex sets in  $\mathbb{C}^N$ .

**3.3. Computational cost of Algorithm 1.** To estimate the error for signals generated by Algorithm 1, we first review a theorem on the Berry-Essen bound's dimensional dependency.

**Theorem 3.7** (Bentkus [3]). *Let  $\hat{X}_i$  be i.i.d. random vectors in  $\mathbb{R}^N$ . Assume  $\hat{X}_i$  has mean zero and identity covariance matrix. Write  $\beta = E[|\hat{X}_i|^3]$ , let*

$$S_M = \frac{1}{M} \sum_{m=1}^M \hat{X}_m,$$

and denote by  $\mathfrak{C}(\mathbb{R}^N)$  the class of all convex sets in  $\mathbb{R}^N$ . Then

$$\sup_{A \in \mathfrak{C}(\mathbb{R}^N)} \left| P(S_M \in A) - \int_A \frac{e^{-|x|^2/2}}{(2\pi)^{N/2}} dx \right| \leq \frac{400N^{1/4}E[|\hat{X}|^3]}{\sqrt{M}}.$$

**Remark 3.8.** *Considering the class of balls in  $\mathbb{R}^N$  instead of the class of convex sets, it is possible to reduce the upper bound of Theorem 3.7 by a factor  $O(N^{-1/4})$ , c.f. [3]. However, for the processes studied in this paper, it is of interest to include at least all ellipsoidally shaped sets.*

Using Theorem 3.7 we next derive an upper bound for  $e(Z_{\mathbf{t},M})$  as a function both of the number of scatterers  $M$  and the number of sampling times  $N$ .

**Corollary 3.9** (Error bound for Algorithm 1). *Let  $Z_{\mathbf{t},M} = (Z_{t_1,M}, Z_{t_2,M}, \dots, Z_{t_N,M})$  denote a sampled output signal generated by Algorithm 1 with sample times  $\mathbf{t}$  restricted to  $[0, T)$ . Assume further the covariance matrix  $K$  of the sampled limit Gaussian process  $Z_{\mathbf{t}}$ , cf. equation (12), is non-singular so that it may be represented by the singular value decomposition  $K = USU^H$  with  $S = \text{diag}(s_i)$ ,  $s_1 \geq s_2 \geq \dots \geq s_N > 0$ . Then*

$$e(Z_{\mathbf{t},M}) \leq \frac{2^{3/4} \cdot 400 N^{7/4} \|E[(a^+(\alpha, \cdot))^3]\|_{L^\infty([0,T])}}{s_N^{3/2} M^{1/2}}.$$

*Proof.* Since  $Z_{\mathbf{t}}$  is a circular symmetric complex normal, it holds for any set  $A \in \mathfrak{C}(\mathbb{C}^N)$  that

$$P(Z_{\mathbf{t},M} \in A) - P(Z_{\mathbf{t}} \in A) = P\left(\left(\frac{K}{2}\right)^{-1/2} Z_{\mathbf{t},M} \in A\right) - \int_{\left(\frac{K}{2}\right)^{-1/2} A} \frac{e^{-|z|^2/2}}{(2\pi)^N} dz,$$

where  $(K/2)^{-1/2}$  being a linear operator implies that also  $(K/2)^{-1/2}A \in \mathfrak{C}(\mathbb{C}^N)$ . The linearly transformed signal  $(K/2)^{-1/2}Z_{\mathbf{t},M}$  may be written

$$(K/2)^{-1/2}Z_{\mathbf{t},M} = \frac{1}{\sqrt{M}} \sum_{m=1}^M (K/2)^{-1/2}\check{Z}_m$$

with i.i.d. complex valued random vectors

$$\check{Z}_m := \begin{bmatrix} a(\alpha_m, t_1) \exp(-i(2\pi f_c \tau(\alpha_m, t_1) + \theta_m(t_1))) \\ a(\alpha_m, t_2) \exp(-i(2\pi f_c \tau(\alpha_m, t_2) + \theta_m(t_2))) \\ \vdots \\ a(\alpha_m, t_N) \exp(-i(2\pi f_c \tau(\alpha_m, t_N) + \theta_m(t_N))) \end{bmatrix}.$$

Following the proof of Theorem 3.4,  $\check{Z}_m$  has zero mean, zero pseudo-covariance matrix, and covariance matrix  $K$ , identical to the covariance matrix of  $Z_{\mathbf{t}}$ . This implies that the random vectors  $(K/2)^{-1/2}\check{Z}_m$  has mean zero, zero pseudo-covariance matrix, and covariance matrix equal to  $2I_N$ , with  $I_N$  denoting the  $N$ -dimensional identity matrix. Then, representing  $\mathbb{C}^N$  by  $\mathbb{R}^{2N}$ , we note that  $2N$ -dimensional real valued random vector  $(\text{Re}((K/2)^{-1/2}\check{Z}_m), \text{Im}((K/2)^{-1/2}\check{Z}_m))$  has mean 0 and identity covariance  $I_{2N}$ , and the proof is completed by bounding the norm of  $E[|(\text{Re}((K/2)^{-1/2}\check{Z}_m), \text{Im}((K/2)^{-1/2}\check{Z}_m))|^3] = E[|(K/2)^{-1/2}\check{Z}_m|^3]$ , and applying Theorem 3.7. By Hölder's inequality,

$$\begin{aligned} E[|(K/2)^{-1/2}\check{Z}_m|^3] &\leq \frac{2^{3/2}}{s_N^{3/2}} E \left[ \left( \sum_{j=1}^N (a(\alpha, t_j))^2 \right)^{3/2} \right] \\ &\leq \frac{2^{3/2}\sqrt{N}}{s_N^{3/2}} \sum_{j=1}^N E [(a(\alpha, t_j))^3] \\ &= 2^{1/2} \left( \frac{N}{s_N} \right)^{3/2} \|E[(a^+(\alpha, \cdot))^3]\|_{L^\infty((0,T))}. \end{aligned}$$

End of proof. □

It follows from the above corollary that to generate a realization fulfilling

$$(14) \quad e(Z_{\mathbf{t},M}) \leq \text{TOL},$$

requires  $M = O(\text{TOL}^{-2} N^{7/2} S_{2N}^{-3})$ , which has the computational cost  $O(MN) = O(\text{TOL}^{-2} (N^{3/2}/S_{2N})^3)$ . For statistical purposes it is often interesting to generate many realizations, so we also note that the cost generating  $L$  realizations fulfilling (14) is

$$\text{Cost}(\text{Algorithm 1}) = O \left( L \text{TOL}^{-2} \left( \frac{N^{3/2}}{S_{2N}} \right)^3 \right).$$

**3.4. Computational cost of Algorithm 2.** Algorithm 2 uses the covariance matrix  $\bar{K}$  to generate realizations  $\bar{Z}_{\mathbf{t}} = (\bar{Z}_{t_1}, \bar{Z}_{t_2}, \dots, \bar{Z}_{t_N})$ . In distribution, these realizations differ slightly from the realizations  $Z_{\mathbf{t}}$  that would be obtained if using  $Z_{\mathbf{t}}$ 's covariance matrix  $K$  to generate signals instead of  $\bar{K}$ . Here, we will bound

the error  $e(\bar{Z}_t)$  in terms of the difference between  $K$  and  $\bar{K}$ , a difference which is a consequence of approximating the integral terms  $K(t_j, t_k)$  of (12) by quadrature:

$$\bar{K}(t_j, t_k) = \frac{e^{-C|t_j - t_k|}}{2} \sum_{\ell=1}^M a^+(x_\ell, t_j) a^+(x_\ell, t_k) e^{i2\pi f_c(\tau(x_\ell, t_k) - \tau(x_\ell, t_j))} p(x_\ell) \nu_\ell$$

for the quadrature points  $0 = x_1 < x_2 < \dots < x_M = 2\pi$  and  $\nu_\ell$  being quadrature weights fulfilling  $\sum_{\ell=1}^M \nu_\ell = 2\pi$ . Using  $M$  quadrature points, the error bound

$$(15) \quad \max_{1 < j, k < N} |K_{j,k} - \bar{K}_{j,k}| \leq \epsilon = O(M^{-\gamma})$$

will be fulfilled,  $\gamma > 0$  here depending on the quadrature method used. Consequently, we may write  $\bar{K} = K + \delta K$  with  $\|\delta K\|_2 \leq N^{1/2}\epsilon$ . Having described the difference between the covariance matrices we now state a theorem that bounds  $e(\bar{Z}_t)$  for circular symmetric complex normals  $\bar{Z}_t$  approximating the sampled limit complex Gaussian process  $Z_t$ .

**Theorem 3.10.** *Assume the covariance matrix  $K$  of the sampled limit complex Gaussian process  $Z_t = (Z_{t_1}, Z_{t_2}, \dots, Z_{t_N})$ , given in (12), is non-singular so that it may be represented by the singular value decomposition  $K = USU^H$  with  $S = \text{diag}(s_i)$ ,  $s_1 \geq s_2 \geq \dots \geq s_N > 0$ . Let further  $\bar{Z}_t$  be a circular symmetric normal Gaussian with mean 0 and covariance matrix  $\bar{K}$ , and assume  $\|K - \bar{K}\|_2 \leq N^{1/2}\epsilon$ , with  $\epsilon = O(M^{-\gamma})$ ,  $\gamma > 0$  the convergence order of the quadrature method, and  $M$  chosen so large that  $10N^{3/2}\epsilon < s_N$ . Then*

$$e(\bar{Z}_t) = O\left(\frac{N^{3/2}}{s_N M^\gamma}\right).$$

*Proof.* By construction,  $\bar{K}$  is symmetric, so it has a singular value decomposition  $\bar{K} = \bar{U}\bar{S}\bar{U}^H$ , with  $\bar{s}_1 \geq \bar{s}_2 \geq \dots \bar{s}_N \in \mathbb{R}$ . According to Corollary A.1 and the assumption on  $M$  being sufficiently large,

$$\max_{1 \leq j \leq N} |s_j - \bar{s}_j| \leq N^{1/2}\epsilon < s_N.$$

This implies that  $\bar{s}_N > 0$  so that  $\bar{K}$  also is non-singular with

$$\bar{K}^{-1} = (K + \delta K)^{-1} = (I + K^{-1}\delta K)^{-1}K^{-1} = K^{-1} + \sum_{j=1}^{\infty} (-K^{-1}\delta K)^j K^{-1}.$$

For later reference we further note that by the assumption of  $M$  being sufficiently large,

$$(16) \quad \left\| K^{1/2} \sum_{j=1}^{\infty} (-K^{-1}\delta K)^j K^{-1/2} \right\|_2 \leq \frac{2N^{1/2}\epsilon}{s_N}.$$



For all Borel sets  $A \subset \mathbb{C}^N$ , we derive the following upper bound

$$\begin{aligned}
P(\bar{Z}_t \in A) - P(Z_t \in A) &= \frac{1}{\pi^N |\det(K)|} \int_A \frac{e^{-z^H \bar{K}^{-1} z}}{|\det(K^{-1} \bar{K})|} - e^{-z^H K^{-1} z} dz \\
&= \frac{1}{\pi^N |\det(K)|^{1/2}} \int_A e^{-z^H K^{-1} z} \left( \frac{e^{-z^H \sum_{j=1}^{\infty} (-K^{-1} \delta K)^j K^{-1} z}}{|\det(K^{-1} \bar{K})|} - 1 \right) dz \\
&= \frac{1}{(2\pi)^N} \int_{(K/2)^{-1/2} A} e^{-|z|^2/2} \left( \frac{e^{-z^H K^{1/2} \sum_{j=1}^{\infty} (-K^{-1} \delta K)^j K^{-1/2} z/2}}{|\det(K^{-1} \bar{K})|} - 1 \right) dz \\
&\stackrel{(16)}{\leq} \frac{1}{(2\pi)^N} \int_{(K/2)^{-1/2} A} e^{-|z|^2/2} \left( \frac{e^{(N^{1/2} \epsilon / s_N) |z|^2}}{\prod_{j=1}^N (1 - N^{1/2} \epsilon / s_j)} - 1 \right) dz \\
&\leq \frac{1}{(2\pi)^N} \int_{\mathbb{R}^{2N}} e^{-|x|^2/2} \left( \left( 1 - \frac{N^{1/2} \epsilon}{s_N} \right)^{-N} e^{(N^{1/2} \epsilon / s_N) |x|^2} - 1 \right) dx \\
&= \left( 1 - \frac{N^{1/2} \epsilon}{s_N} \right)^{-N} \left( 1 - 2 \frac{N^{1/2} \epsilon}{s_N} \right)^{-N} - 1 \\
&\leq e^{10N^{3/2} \epsilon / s_N} - 1 \\
&\leq \frac{20N^{3/2} \epsilon}{s_N}.
\end{aligned}$$

Here we used if  $\hat{X}$  is a  $2N$ -dimensional real valued multivariate normal with mean zero and identity covariance, then  $E[e^{t|\hat{X}|^2}] = (1 - 2t)^{-N}$  is the moment generating function to a chi-square distributed variable with  $2N$  degrees of freedom. Calculating along the same lines one obtains the lower bound

$$\begin{aligned}
P(\bar{Z}_t \in A) - P(Z_t \in A) &\geq \left( 1 + \frac{N^{1/2} \epsilon}{s_N} \right)^{-N} E[e^{-(N^{1/2} \epsilon / s_N) |\hat{X}|^2}] - 1 \\
&\geq -\frac{3N^{3/2} \epsilon}{s_N}.
\end{aligned}$$

Since the upper and lower bound obtained are both valid for the class of all Borel sets  $A \in \mathbb{C}^N$  and this class contains the class of convex sets  $\mathfrak{C}(\mathbb{C}^N)$ , the proof is finished.  $\square$

It follows from Theorem 3.10 that to fulfill the accuracy condition  $e(\bar{Z}_t) \leq \text{TOL}$  for realizations generated by Algorithm 2 requires

$$M = O\left(\left(\frac{N^{3/2}}{s_N \text{TOL}}\right)^{1/\gamma}\right).$$

The cost of generating  $L$  signal realizations fulfilling the above mentioned accuracy condition then becomes the sum of  $O(MN^2)$  to compute the matrix elements of  $\bar{K}$  by quadrature,  $O(N^3)$  to construct the square root of  $\bar{K}$  and  $O(LN^2)$  for creating  $L$  signal realizations by the matrix vector multiplication  $\bar{K}^{1/2} \hat{X}$ :

$$\text{Cost}(\text{Algorithm 2}) = O\left(N^3 + N^2 \left(\frac{N^{3/2}}{s_N \text{TOL}}\right)^{1/\gamma} + LN^2\right).$$

**3.5. A circulant-embedding based algorithm with error and computational cost analysis.** Computing the square root of  $\bar{K}$  in Algorithm 2 and generating the output signal by (13) were both computationally costly operations, and in this subsection we consider modeling settings where it is possible to improve the efficiency of these operations by circulant-embedding of the covariance matrix and application of finite Fourier methods. The idea presented here is adapted from the material in [2, Chapter XI].

In the modeling setting when

$$(17) \quad \tau(\alpha, t) = -v \cos(\alpha)t/c \quad \text{and} \quad a^+(\alpha, t) = a^+(\alpha),$$

the limit Gaussian process  $Z_t$  is a WSS process which simplifies the structure of the covariance:

$$(18) \quad K(t_i, t_j) = E[Z_{t_i} Z_{t_j}^*] = \frac{e^{-C|t_i - t_j|}}{2} \int_0^{2\pi} a^+(\alpha)^2 e^{i2\pi f_D \cos(\alpha)(t_i - t_j)} p(\alpha) d\alpha =: A(t_i - t_j),$$

where we recall that  $f_D = f_c v/c$  is the maximum Doppler shift and  $A(t)$  denotes the autocorrelation function. Sampling the limit Gaussian process  $Z_t$  on a uniform time grid  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ ,  $t_{i+1} - t_i = \Delta t$ , in setting (17), yields a circular symmetric complex normal  $Z_{\mathbf{t}}$  with a covariance matrix  $K$  that is Toeplitz. Using the short hand notation  $A_j = A(j\Delta t)$ , the covariance matrix has the structure

$$K = \begin{bmatrix} A_0 & A_1 & \cdot & A_{N-2} & A_{N-1} \\ A_{-1} & A_0 & \cdot & A_{N-3} & A_{N-2} \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ A_{-(N-2)} & A_{-(N-3)} & \cdot & A_0 & A_1 \\ A_{-(N-1)} & A_{-(N-2)} & \cdot & A_{-1} & A_0 \end{bmatrix}.$$

To illustrate circulant-embedding of  $K$ , let us for simplicity assume  $A_{-j} = A_j$ , and embed  $K$  as the upper left corner of a circulant of order  $2N - 2$ :

$$\mathcal{C} = \begin{bmatrix} A_0 & A_1 & \cdot & A_{N-2} & A_{N-1} & A_{N-2} & A_{N-3} & \cdot & A_2 & A_1 \\ A_1 & A_0 & \cdot & A_{N-3} & A_{N-2} & A_{N-1} & A_{N-2} & \cdot & A_3 & A_2 \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ A_{N-1} & A_{N-2} & \cdot & A_1 & A_0 & A_1 & A_2 & \cdot & A_{N-3} & A_{N-2} \\ A_{N-2} & A_{N-1} & \cdot & A_2 & A_1 & A_0 & A_1 & \cdot & A_{N-4} & A_{N-3} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ A_1 & A_2 & \cdot & A_{N-2} & A_{N-1} & A_{N-2} & A_{N-3} & \cdot & A_1 & A_0 \end{bmatrix}.$$

(In the general case when we only have  $A_{-j} = A_j^*$ , the circulant-embedding of  $K$  will be of size  $4N - 4$ .) The circulant matrix  $\mathcal{C}$  has the eigendecomposition

$$(19) \quad \mathcal{C} = F\Lambda F^H,$$

where  $F$  is the finite Fourier matrix of order  $2N - 2$  with elements  $F_{jk} = e^{i2\pi(j-1)(k-1)/(2N-2)}/\sqrt{2N-2}$ , and  $\Lambda$  is the diagonal matrix of eigenvalues

$$(20) \quad \text{diag}(\Lambda) = F\mathbf{a} \quad \text{with} \quad \mathbf{a} = (A_0, A_1, \dots, A_{N-1}, A_{N-2}, \dots, A_2, A_1)^T.$$

The representation (19) shows that provided all entries of  $\Lambda$  are non-negative, we may write  $\mathcal{C}^{1/2} = F\Lambda^{1/2}$ . The signal generated by

$$Z_{\mathbf{t}} = R F \Lambda^{1/2} \widehat{Z},$$

with  $\widehat{Z} \sim \mathcal{N}_{\mathbb{C}}(0, I_{2N-2})$  and  $R: \mathbb{C}^{2N-2} \rightarrow \mathbb{C}^N$  defined by  $Rz = (z_1, z_2, \dots, z_N)^T$  for all  $z \in \mathbb{C}^{2N-2}$ , will have the sought covariance  $K$ . This leads us to the following algorithm for generating WSS Gaussian process signal realizations.

---

**Algorithm 3** Covariance based algorithm for WSS processes

---

**Input:** Flip rate  $C$ , maximum Doppler shift  $f_D$ , scatterer density  $p$ , uniform grid of sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ .

**Output:** Gaussian process realization  $\bar{Z}_{\mathbf{t}} = (\bar{Z}_{t_1}, \bar{Z}_{t_2}, \dots, \bar{Z}_{t_N})$ .

**for**  $j = -(N-1)$  to  $N-1$  **do**

Approximate the term  $A_j$  by quadrature of equation (18) and store the approximation in  $\bar{A}_j$ .

**end for**

Determine  $\bar{\Lambda}$  by computing

$$\text{diag}(\bar{\Lambda}) = F(\bar{A}_0, \bar{A}_1, \dots, \bar{A}_{N-1}, \bar{A}_{N-2}, \dots, \bar{A}_2, \bar{A}_1)^T,$$

using the Fast Fourier Transform (FFT).

Generate a multivariate complex normal vector  $\widehat{Z} \sim \mathcal{N}_{\mathbb{C}}(0, I_{2N-2})$  and thereafter use the FFT to compute the output realization

$$(21) \quad \bar{Z}_{\mathbf{t}} = R F \bar{\Lambda}^{-1/2} \widehat{Z}.$$


---

Since a sampled realization  $\bar{Z}_{\mathbf{t}}$  generated by Algorithm 3 is multivariate circular symmetric complex Gaussian with mean zero and covariance  $\bar{K}$ , we may use Theorem 3.10 to conclude that, under the assumptions there stated,  $e(\bar{Z}_{\mathbf{t}}) = O(N^{3/2}/(s_N M^\gamma))$ . Alternatively, under the assumptions that  $\Lambda$  is positive definite,  $\|K - \bar{K}\|_2 \leq N^{1/2}\epsilon$  with  $\epsilon = O(M^\gamma)$  and  $M$  chosen sufficiently large so that  $N^{1/2} \text{tr}(\Lambda^{-1})\epsilon < 1/2$ , one may derive the second error bound

$$e(\bar{Z}_{\mathbf{t}}) = O\left(\frac{N^{1/2} \text{tr}(\Lambda^{-1})}{M^\gamma}\right)$$

from observing that for any Borel set  $A \in \mathbb{C}^N$

$$\begin{aligned} & \left| P\left(F\bar{\Lambda}^{-1/2}\widehat{Z} \in A \times \mathbb{C}^{N-2}\right) - P\left(F\Lambda^{1/2}\widehat{Z} \in A \times \mathbb{C}^{N-2}\right) \right| \\ &= \left| P\left((\Lambda^{-1}\bar{\Lambda})^{1/2}\widehat{Z} \in \Lambda^{-1/2}F^H(A \times \mathbb{C}^{N-2})\right) - P\left(\widehat{Z} \in \Lambda^{-1/2}F^H(A \times \mathbb{C}^{N-2})\right) \right| \\ &= O\left(N^{1/2}\epsilon \text{tr}(\Lambda^{-1})\right), \end{aligned}$$

where  $\text{tr}(\Lambda^{-1}) := \sum_j \Lambda_j^{-1}$ . The cost of generating  $L$  signal realizations with Algorithm 3 fulfilling the accuracy  $e(\bar{Z}_{\mathbf{t}}) \leq \text{TOL}$  becomes the sum of  $O(MN)$  to compute the covariance elements  $\{\bar{A}_j\}_{j=-N+1}^{N-1}$  of  $\bar{K}$  by quadrature,  $O(N \log(N))$  to compute  $\bar{\Lambda}$  and  $O(LN \log(N))$  for creating  $L$  signal realizations using FFT:

$$\text{Cost}(\text{Algorithm 3}) = O\left\{ N \left[ \frac{N^{1/2}}{\text{TOL}} \min\left(\frac{N}{s_N}, \text{tr}(\Lambda^{-1})\right) \right]^{1/\gamma} + LN \log(N) \right\}.$$

For a tentative comparison of the magnitude of  $N/s_N$  and  $\text{tr}(\Lambda^{-1})$ , see Section 5.

## 4. APPLICATIONS OF THE POWER SPECTRAL DENSITY

Being the Fourier dual of the covariance function, the PSD can, as the covariance, be used in algorithms to construct signal realizations and to study properties process properties. Here we derive the PSD for in the WSS process modeling setting

$$(22) \quad \tau(\alpha, t) = -v \cos(\alpha)t/c, \quad a^+(\alpha, t) = a^+(\alpha), \quad \text{and} \quad C = O(1).$$

The shape of the PSD is linked to the flip rate and the scattering density and we will describe a method for estimating the flip rate  $C$  from the PSD computed from real life signal measurements and also look into the relation between PSD and scattering density. At the end of this section, an algorithm for generating Gaussian process signal realizations using the PSD function is presented with cost and error analysis included.

We first derive an expression for the PSD, using ideas from [9] and [7]. The PSD is given by  $S_C(f) = \mathcal{F}\{A(\cdot)\}$ , where  $A(t)$  is the autocorrelation of the limit complex Gaussian process  $Z_t$ ,  $\mathcal{F}(\cdot)$  denotes the Fourier Transform, and the subscript  $C$  in  $S_C$  emphasizes that the PSD depends on the flip rate. In the setting (22), it follows from (18) that the autocorrelation is on the form

$$(23) \quad A(t) = \frac{e^{-C|t|}}{2} \int_0^{2\pi} (a^+(\alpha))^2 e^{i2\pi f_D \cos(\alpha)t} p(\alpha) d\alpha,$$

where we recall that  $f_D = vt/c$ . By the Convolution theorem for Fourier transforms,

$$\begin{aligned} S_C(f) &= \mathcal{F}\{A(\cdot)\}(f) \\ &= \mathcal{F}\left\{\frac{e^{-C|\cdot|}}{2}\right\} * \mathcal{F}\left\{\int_0^{2\pi} (a^+(\alpha))^2 e^{i2\pi f_D \cos(\alpha)\cdot} p(\alpha) d\alpha\right\}(f), \end{aligned}$$

where  $*$  denotes the convolution operator. Observe next that

$$\mathcal{F}\{e^{-C|\cdot|}\}(f) = \frac{2C}{C^2 + (2\pi f)^2}, \quad \forall C > 0,$$

and

$$\begin{aligned} &\mathcal{F}\left\{\int_0^{2\pi} (a^+(\alpha))^2 e^{i2\pi f_D \cos(\alpha)\cdot} p(\alpha) d\alpha\right\}(f) \\ &= \int_{\mathbb{R}} \int_0^{2\pi} (a^+(\alpha))^2 e^{i2\pi f_D \cos(\alpha)t} p(\alpha) d\alpha e^{-i2\pi ft} dt \quad (\alpha = -\cos^{-1}(s/f_D) + \{0, \pi\}) \\ &= \int_{\mathbb{R}} \int_{-f_D}^{f_D} \frac{1}{f_D \sqrt{1 - (s/f_D)^2}} \left( e^{i2\pi st} (a^+(-\cos^{-1}(s/f_D)))^2 p(-\cos^{-1}(s/f_D)) \right. \\ &\quad \left. + e^{-i2\pi st} (a^+(-\cos^{-1}(s/f_D) + \pi))^2 p(-\cos^{-1}(s/f_D) + \pi) \right) e^{-i2\pi ft} dt ds \\ &= \frac{1_{|f| < f_D}(f)}{f_D \sqrt{1 - (f/f_D)^2}} \left( (a^+(-\cos^{-1}(f/f_D)))^2 p(-\cos^{-1}(f/f_D)) \right. \\ &\quad \left. + (a^+(-\cos^{-1}(-f/f_D) + \pi))^2 p(-\cos^{-1}(-f/f_D) + \pi) \right). \end{aligned}$$

Convolving the last two expressions leads to the following PSD integral expression

$$(24) \quad S_C(f) = \int_0^{2\pi} \frac{C(a^+(\alpha))^2 p(\alpha)}{C^2 + (2\pi(f - f_D \cos(\alpha)))^2} d\alpha.$$

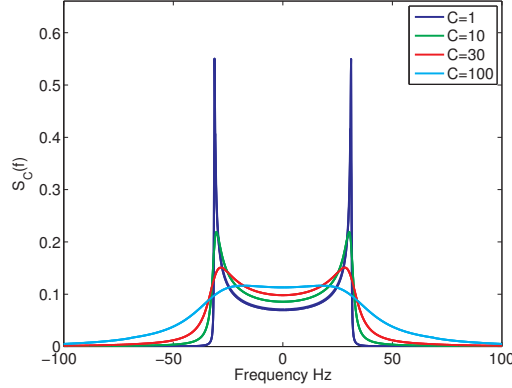


FIGURE 4. The PSD function  $S_C(f)$  is plotted for different flip rate values  $C$  with  $f_D = 50\text{Hz}$  when  $(a^+)^2 p = (2\pi)^{-1}$ .

We see that  $S_C(f)$  depends on the flip rate  $C$  and the term  $(a^+)^2 p$ , which we hereafter refer to as a scaled scatterer density. Figure 4 illustrates  $S_C(f)$ 's dependency on the flip rate  $C$  in the modeling settings  $(a^+)^2 p = (2\pi)^{-1}$ . In that setting  $S_0(f)$  is Jakes' spectrum and  $S_C(f)$  is a progressively mollified version of Jakes' spectrum the higher the value of  $C$  is. In [9] it is remarked that the PSD  $S_C(f)$  with  $C = O(1)$  is more in accordance with real life measurements than Jakes' spectrum is.

**4.1. A link to Feng and Field's model.** Our MFC model's output signal auto-correlation and PSD results are very similar to what Feng and Field obtained under quite different modeling assumptions in [9]. They considered the modified Clarke's model

$$\varepsilon_t = \sum_{n=1}^N a_n \exp\left(i(2\pi f_n t + \phi_t^{(n)})\right),$$

where the amplitudes  $a_n$  are i.i.d. random variables and  $f_n = f_D \cos(\alpha_n)$  are Doppler shifts where, as in this paper,  $f_D = f_c v/c$ . The phases  $\phi_t^{(n)}$  are independent Wiener processes with uniform initial distribution in  $[0, 2\pi)$ :

$$d\phi_t^{(n)} = \sqrt{B} dW_t^{(n)}, \quad \phi_0^{(n)} \sim U[0, 2\pi), \quad \forall n \in \{1, 2, \dots, N\},$$

where  $B$  is a constant with the dimension of frequency. For this model they obtain the autocorrelation function

$$(25) \quad E[\varepsilon_t \varepsilon_0^*] = \sum_{n=1}^N E[a_n^2] e^{-B|t|/2} \int_0^{2\pi} e^{i2\pi f_D \cos(\alpha)t} d\alpha.$$

For our model in the setting (22), we recall from (23) that the autocorrelation is given by the quite similar expression

$$E[Z_t Z_0^*] = \frac{e^{-C|t|}}{2} \int_0^{2\pi} (a^+(\alpha))^2 e^{i2\pi f_D \cos(\alpha)t} p(\alpha) d\alpha.$$

**Remark 4.1.** *It should be noted that the autocorrelation similarity is obtained although the modeling assumptions are quite different: in our MFC model the amplitude is governed by a Poisson flip process, Feng and Field's model has time invariant amplitude functions; our model updates a phase when the corresponding amplitude flips, their model has Wiener process phase evolution.*

**4.2. Model parameter estimation from PSD measurements.** Feng and Field's publication [9] presents a method for estimating the flip rate constant  $C$  from measurements which can be used when the scaled scatterer density  $(a^+)^2 p$  is known. Namely, given a measurement of the PSD,  $\widehat{S}(f)$ , estimate  $C$  by

$$(26) \quad C = \arg \min_{x>0} \cos^{-1} \left( \int_{\mathbb{R}} \sqrt{\widehat{S}(f) S_x(f)} df \right),$$

where the PSDs are scaled so that  $\|\widehat{S}\|_1 = 1$  and  $\|S_x\|_1 = 1$ . However, from (24) we see that the shape of the PSD depends both on the flip rate and the scaled scatterer density on  $(a^+)^2 p$ , see Figure 5. So to estimate the flip rate from PSD measurements,  $(a^+)^2 p$  either has to be known or it has to be estimated. It is more difficult and costly to estimate  $(a^+)^2 p$  than the flip rate, we restrict ourselves to a tentative approach to this problem. Let  $\widehat{S}$  denote a real life PSD measurement for a moving receiver, and consider relation (24) as the inhomogeneous Fredholm integral equation

$$\int_0^{2\pi} \mathcal{K}_C(f, \alpha) (a^+(\alpha))^2 p(\alpha) d\alpha = \widehat{S}(f),$$

with unknown  $(a^+(\alpha))^2 p(\alpha)$  and parametrized by the flip rate  $C$ . For a fixed flip rate, this problem might be discretized to a system of linear equations and, if expecting noise in the measurement  $\widehat{S}$ , solved as an inverse problem using Tikhonov or similar regularization techniques, c.f. [14]. For example, if we denote the discretized integral kernel  $\mathcal{K}_C \in \mathbb{R}^{m \times n}$ , then the Tikhonov regularized solution has the variational representation

$$(27) \quad (a_C^+)^2 p_C = \arg \min_{x \in \mathbb{R}^n} \|\mathcal{K}_C x - \widehat{S}\|^2 + \beta \|x\|^2$$

for a regularizing parameter  $\beta > 0$  to be chosen. If the flip rate  $C$  and  $(a^+(\alpha))^2 p(\alpha)$  have to be determined simultaneously, we suggest coupling the solution of (27) with a minimization problem on the form

$$C = \arg \min_{x>0} \|\mathcal{K}_x (a_x^+)^2 p_x - \widehat{S}\|.$$

**4.3. A PSD based Gaussian process algorithm.** We now present a PSD based algorithm for generating signal realizations of WSS complex Gaussian processes for the modeling the setting (22). The PSD (24) can be used to represent the Gaussian process spectrally by the Itô integral

$$(28) \quad Z_t = \int_{\mathbb{R}} e^{i2\pi f t} \sqrt{S_C(f)} dW_f,$$

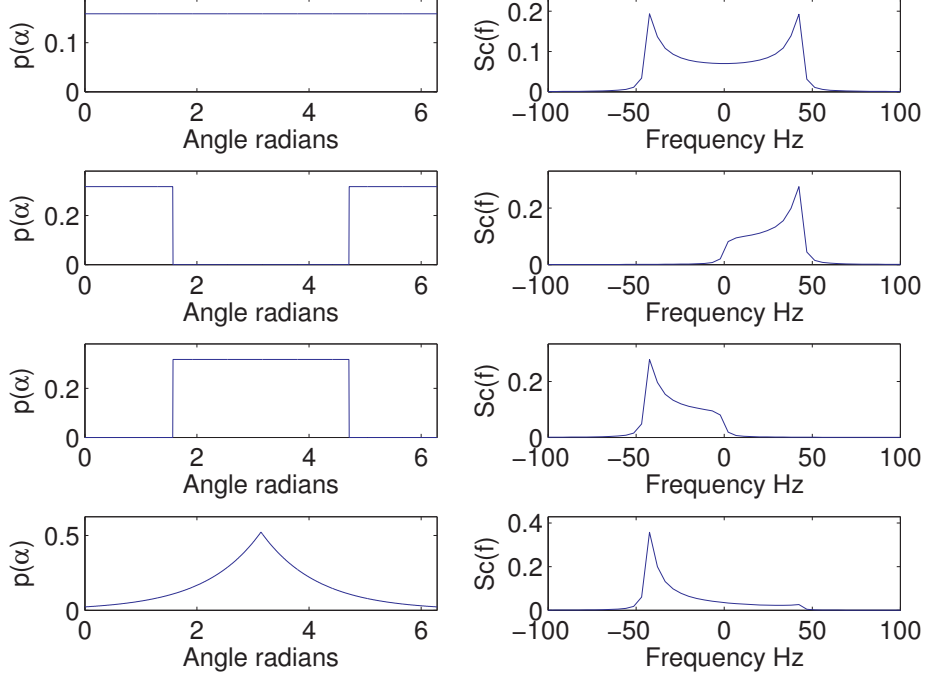


FIGURE 5. Right column plots illustrate the PSD  $S_C(f)$  obtained by equation (24) when the scaled scatterer density  $(a^+)^2 p(\alpha)$  is given by respective left column plots,  $C = 10$ , and  $f_D = 50$ . **Top row:**  $(a^+)^2 p = 1/(2\pi)$  yields the mollified Jakes' spectrum PSD. **Second row:**  $(a^+)^2 p(\alpha) = 1_{[-\pi/2, \pi/2]}(\alpha)/\pi$  yields a PSD consisting almost exclusively of positive Doppler shifts since the receiver's moves towards the active scatterers. **Third row:**  $(a^+)^2 p(\alpha) = 1_{[\pi/2, 3\pi/2]}(\alpha)/\pi$  yields a PSD consisting almost exclusively of negative Doppler shifts since the receiver moves away from the active receivers. **Last row:**  $(a^+)^2 p(\alpha) = \exp(-|\alpha - \pi|)/\text{scale}$ .

where  $W_f \in \mathbb{C}$  is a complex Wiener processes with independent infinitesimal increments  $dW_f \sim \mathcal{N}_{\mathbb{C}}(0, df)$ . Seeking to approximate the integral (28) by quadrature we first approximate the integrand by a compactly supported function by choosing a cut off frequency  $f_M \gg f_D$  and introducing the  $C^3(\mathbb{R})$  window function

$$(29) \quad \varrho(f) = \begin{cases} 0 & \text{if } f \leq -f_M \\ \varrho_I(2(f_M + f)/f_M) & \text{if } -f_M < f < -f_M/2 \\ 1 & \text{if } |f| \leq f_M/2 \\ \varrho_I(2(f_M - f)/f_M) & \text{if } f_M/2 < f < f_M \\ 0 & \text{else,} \end{cases}$$

where  $\varrho_I(f) = -20f^7 + 70f^6 - 84f^5 + 35f^4$  is the unique solution of the Birkhoff-Hermite interpolation with conditions  $\varrho_I^{(j)}(0) = 0$  for  $j = 0, 1, \dots, 4$  and  $\varrho_I(1) = 1$ ,

$\varrho_f^{(j)}(1) = 0$  for  $j = 1, \dots, 4$ . Replacing the the integrand term  $\sqrt{S_C(f)}$  in (28) with  $\sqrt{\varrho(f)S_C(f)}$  yields the integral

$$(30) \quad \int_{-f_M}^{f_M} e^{i2\pi ft} \sqrt{\varrho(f)S_C(f)} dW_f.$$

Quadrature approximating this integral with the Inverse Discrete Fourier Transform (IDFT) gives output signal realizations by

$$(31) \quad \bar{Z}_t = \sum_{j=1}^{M_1} e^{i2\pi f_j t} \sqrt{\varrho(f_j)\bar{S}_C(f_j)} \Delta W_j.$$

Here  $f_j = -f_M + (j-1)\Delta f$ ,  $\Delta f = 2f_M/(M_1-1)$ ,  $\Delta W_j \sim \mathcal{N}_{\mathbb{C}}(0, \Delta f)$  are i.i.d. complex Wiener increments and  $\bar{S}_C(f_j)$  approximates  $S_C(f_j)$  by quadrature of the integral (24),

$$(32) \quad \bar{S}_C(f_k) = \sum_{j=1}^{M_2} \frac{C(a^+(x_j))^2 p(x_j) \tilde{\nu}_j}{C^2 + (2\pi(f_k - f_D \cos(x_j)))^2}$$

on a grid  $0 = x_1 < x_2 < \dots < x_{M_2} = 2\pi$  with integration weights  $\tilde{\nu}_j \geq 0$  which satisfies  $\sum_{j=1}^{M_2} \tilde{\nu}_j = 2\pi$ .

Algorithms for generating Gaussian process signal realizations using the IDFT have been described in [12, 15], and in Algorithm 4, we present a version suited for our setting.

---

**Algorithm 4** A PSD based Gaussian process algorithm

---

**Input:** Flip rate  $C$ , maximum Doppler shift  $f_D$ , spectral cutoff  $f_M$ , sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$  and scatterer density  $p$ .

**Output:** Gaussian process realization  $\{\bar{Z}_{t_j}\}_{j=1}^N$ .

Construct a grid  $f_j = -f_M + (j-1)\Delta f$ , for  $j = 1, 2, \dots, M_1$  with  $\Delta f = 2f_M/(M_1-1)$ .

Construct a grid  $0 = \alpha_1 < \alpha_2 < \dots < \alpha_{M_2} = 2\pi$  and quadrature weights  $\tilde{\nu}_j \geq 0$  for which  $\sum_{j=1}^{M_2} \tilde{\nu}_j = 2\pi$ .

**for**  $j = 1$  to  $M_1$  **do**

    Compute  $\bar{S}_C(f_j)$  according to (32).

**end for**

Generate i.i.d. Wiener increments  $\{\Delta W_j\}_{j=1}^{M_2}$  distributed according to  $\Delta W_j \sim \mathcal{N}_{\mathbb{C}}(0, \Delta f)$ .

**for**  $k = 1$  to  $N$  **do**

    Compute  $\bar{Z}_{t_k}$  by the IDFT (31).

**end for**

---

**4.4. Computational cost of Algorithm 4.** The error analysis for this algorithm is quite similar to the error analysis of Algorithm 2 in the sense that the distributional error can be bounded in terms of the difference between the covariance matrix of realizations of Algorithm 4, denoted  $\bar{K}$ , and the covariance matrix  $K$  for the sampled limit complex Gaussian process  $Z_t$  derived in (12). Our procedure for obtaining an error bound is to first derive a bound of  $\|\bar{K} - K\|_2$  and thereafter use



Theorem 3.10 to bound  $e(\bar{\mathbf{Z}}_{\mathbf{t}})$  from above. We start with computing the elements of  $\bar{\mathbf{K}}$ .

The representation (31) implies that for two arbitrary times  $t_j, t_k$  from the sampling times  $\mathbf{t} = (t_1, t_2, \dots, t_N)$ ,

$$\begin{aligned} \bar{\mathbf{K}}(t_j, t_k) &= E \left[ \bar{\mathbf{Z}}_{t_j} \bar{\mathbf{Z}}_{t_k}^* \right] \\ &= \sum_{l,m=1}^{M_1} e^{i2\pi(f_l t_j - f_m t_k)} \sqrt{\bar{S}_C(f_l) \bar{S}_C(f_m)} \sqrt{\varrho(f_l) \varrho(f_m)} E[\Delta W_l \Delta W_m^*] \\ &= \sum_{l=1}^{M_1} e^{i2\pi f_l (t_j - t_k)} \varrho(f_l) \bar{S}_C(f_l) \Delta f, \end{aligned}$$

where we recall that  $\bar{S}_C$  is an approximation of the PSD  $S_C$ . The error of the covariance terms is split into three parts,

(33)

$$\begin{aligned} |K(t_j, t_k) - \bar{\mathbf{K}}(t_j, t_k)| &= \left| \int_{\mathbb{R}} e^{i2\pi f (t_j - t_k)} S_C(f) df - \sum_{l=1}^{M_1} e^{i2\pi f_l (t_j - t_k)} \varrho(f_l) \bar{S}_C(f_l) \Delta f \right| \\ &\leq \left| \int_{\mathbb{R}} e^{i2\pi f (t_j - t_k)} (1 - \varrho(f)) S_C(f) df \right| \\ &\quad + \left| \int_{-f_M}^{f_M} e^{i2\pi f (t_j - t_k)} \varrho(f) S_C(f) df - \sum_{l=1}^{M_1} e^{i2\pi f_l (t_j - t_k)} \varrho(f_l) S_C(f_l) \nu_l \right| \\ &\quad + \left| \sum_{l=1}^{M_1} e^{i2\pi f_l (t_j - t_k)} \varrho(f_l) (S_C(f_l) - \bar{S}_C(f_l)) \nu_l \right| \\ &= I + II + III. \end{aligned}$$

Under the assumption  $C = O(1)$ , we see from equation (24) that  $S_C(f) = O((1 + f^2)^{-1})$ . Therefore  $I = O(f_M^{-1})$ . For the second term [4, thm. 6.3] implies that since  $\varrho S_C$  up third derivative is continuous and (can be considered)  $2f_M$  periodic on  $[-f_M, f_M]$ ,  $II = O((f_M/M_1)^4)$ , where the implicit constant in the error term depends on the derivatives of  $\varrho S_C$  up to third order. (Increasing the number of derivative conditions in the Birkhoff interpolation  $\varrho_I$ , will increase the convergence order in  $II$ , but, as a trade off, the error term's implicit constant increases as well.) For the last term, the approximation  $\bar{S}_C(f)$  is obtained from the quadrature (32) using  $M_2$  integration points, so that  $|S_C(f) - \bar{S}_C(f)| = O((f^2 + 1)^{-1} M_2^{-\gamma})$ , where the convergence order  $\gamma > 0$  depends on the numerical integrator used and is the same order as that obtained in (15) for Algorithm 2. Thereby

$$(34) \quad III \leq \sum_{l=1}^{M_1} |S_C(f_l) - \bar{S}_C(f_l)| \Delta f = O \left( \int_{-f_M}^{f_M} (f^2 + 1)^{-1} df M_2^{-\gamma} \right) = O(M_2^{-\gamma}).$$

Adding the three error terms yields the covariance error bound

$$(35) \quad |K(t_j, t_k) - \bar{\mathbf{K}}(t_j, t_k)| \leq \varepsilon = O(f_M^{-1} + (f_M/M_1)^4 + M_2^{-\gamma}),$$

and by applying Theorem 3.10 we obtain the following error bound for this algorithm.

**Corollary 4.2.** *Let  $\bar{K}$  denote the covariance matrix of the stochastic process generated by Algorithm 4 in the modeling setting (22) and assume the covariance matrix  $K$  of the Gaussian process  $Z_{\mathbf{t}}$ , as given by (12), is non-singular and representable by the SVD  $K = USU^H$  with  $S = \text{diag}(s_i)$ ,  $s_1 \geq s_2 \geq \dots \geq s_N > 0$ . Then, if  $M_1$ ,  $M_2$  and  $f_M$  are chosen sufficiently large so that the error bound (35)  $\|K - \bar{K}\|_2 \leq \sqrt{N}\epsilon$  is fulfilled with  $10N^{3/2}\epsilon < s_N$ , realizations of Algorithm 4*

$$e(\bar{Z}_{\mathbf{t}}) = O\left(\frac{N^{3/2}(f_M^{-1} + (f_M/M_1)^4 + M_2^{-\gamma})}{s_N}\right).$$

*Proof.* The result follows from Theorem 3.10.  $\square$

To fulfill the accuracy condition  $e(\bar{Z}_{\mathbf{t}}) \leq \text{TOL}$ , Corollary 4.2 implies that  $f_M = O((S_N \text{TOL})^{-1}N^{3/2})$ ,  $M_1 = O(((S_N \text{TOL})^{-1}N^{3/2})^{5/4})$  and  $M_2 = O(((S_N \text{TOL})^{-1}N^{3/2})^{1/\gamma})$ . The cost of generating one signal realization for this algorithm is the sum of the cost of computing  $S_C(f_k)$  for  $k = 1, 2, \dots, M_1$  by (32), which generally amounts to  $O(M_1 M_2)$ , and the cost  $O(NM_1)$  for computing (31) to generate a realization  $\bar{Z}_{\mathbf{t}}$ . The cost of generating  $L$  signal realizations thus becomes

$$\text{Cost}(\text{Algorithm 4}) = O\left(\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4+1/\gamma} + LN\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4}\right).$$

**Remark 4.3.** *In many settings, particularly the setting when generating realizations on uniformly sampled grid points on an interval  $[0, T)$  with  $\Delta t = T/N$ , it is possible to apply the FFT techniques to speed up the quadrature computations of the discrete convolution (32) and the IDFT (31) so that the cost of generating  $L$  realizations with Algorithm 4 in terms of  $L, N$  and  $\text{TOL}$  instead amounts to*

$$\text{Cost}(\text{Alg 4, FFT}) = O\left(\left(L\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4} + \left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{1+1/\min(4, \bar{\gamma})}\right)\log\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)\right).$$

Here  $\bar{\gamma}$  is the convergence order of the discrete convolution (32) which has the relation  $\bar{\gamma} \leq \gamma$  since when computing (32) with FFT, the quadrature grid has to fulfill  $\cos(x_{j+1}) - \cos(x_j) = \Delta f / f_D$ , which generally does not optimize the accuracy of the computation.

## 5. SUMMARY OF THE COMPLEXITY ESTIMATES

Having estimated upper bounds for the computational cost of generating output realizations by four different algorithms, we now summarize the results. We do however stress that none of the error bounds for which the cost estimates are based are proven to be sharp, so the following cost comparison should not be considered conclusive. In Table 1 we present cost estimates for non-WSS modeling settings, which occurs if the amplitude function  $a^+$  is time dependent or/and if modeling with the full delay function  $\tau(\alpha, t)$ . The results of the table indicates that for settings when  $\gamma \geq 3/5$ , the covariance based Gaussian process algorithm, Algorithm 2, is asymptotically the most efficient algorithm.

In modeling settings resulting in WSS output processes,

$$(36) \quad \tau(\alpha, t) = -v \cos(\alpha)t/c, \quad \text{and} \quad a^+(\alpha, t) = a^+(\alpha),$$

we derived the computational cost presented in Table 2. The results of the table indicate that in the WSS setting, Algorithm 3 outperforms the other algorithms in terms of efficiency.

Computational cost	
Algorithm 1	$O\left(L \text{TOL}^{-2} \left(\frac{N^{3/2}}{S_N}\right)^3\right)$
Algorithm 2	$O\left(N^3 + N^2 \left(\frac{N^{3/2}}{s_N \text{TOL}}\right)^{1/\gamma} + LN^2\right)$

TABLE 1. Computational cost for non-WSS modeling settings. Recall that  $\gamma$  is the convergence order for the quadrature (15).

Computational cost	
Alg 1	$O\left(L \text{TOL}^{-2} \left(\frac{N^{3/2}}{S_N}\right)^3\right)$
Alg 2	$O\left(N^3 + N \left(\frac{N^{3/2}}{s_N \text{TOL}}\right)^{1/\gamma} + LN^2\right)$
Alg 3	$O\left(N \left[\frac{N^{1/2}}{\text{TOL}} \min\left(\frac{N}{s_N}, \text{tr}(\Lambda^{-1})\right)\right]^{1/\gamma} + LN \log(N)\right)$
Alg 4	$O\left(\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4+1/\gamma} + LN \left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4}\right)$
Alg 4, FFT	$O\left(\left(L \left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{5/4} + \left(\frac{N^{3/2}}{S_N \text{TOL}}\right)^{1+1/\min(4, \bar{\gamma})}\right) \log\left(\frac{N^{3/2}}{S_N \text{TOL}}\right)\right)$

TABLE 2. Computational cost for the WSS setting (36). Recall that  $\gamma$  is the quadrature convergence order obtained in (15) and (34) for the respective algorithms, and  $\bar{\gamma}$  with  $\bar{\gamma} \leq \gamma$  is the convergence order obtained when computing (34) with FFT.

For better understanding of the cost estimates, we would like an estimate on how  $N/S_N$  and  $\text{tr}(\Lambda^{-1})$  depend on the sample times  $\mathbf{t} = (0, t_2, \dots, t_N)$  used in the generation of signal realizations. A general answer to this question is however too demanding; we restrict ourselves to the WSS modeling setting (22).

The eigenvalue  $\Lambda_n$  is a DFT approximation of  $S_C((n-1)/T)$ , and we may, by introducing a cutoff function as in (29) and arguing similar as in the error analysis of (33), derive the error bound

$$|\Lambda_n - S_C((n-1)/t_N)| \leq O(\Delta t^4 + e^{-C|t_N|}), \quad n = 1, 2, \dots, N,$$

and

$$|\Lambda_{N+n} - S_C(-(N-1-n)/t_N)| \leq O(\Delta t^4 + e^{-C|t_N|}), \quad n = 1, 2, \dots, N-2.$$

Then it follows that for settings  $(4/C) \log(N) \leq t_N = o(N)$  and  $\Delta t = t_N/(N-1)$ , we obtain the sharp bound

$$\begin{aligned}
 \sum_{j=1}^{2N-2} \Lambda_j^{-1} &= O\left(\frac{t_N^2}{N}\right) + \sum_{j=2-N}^{N-1} S_C(j/t_N)^{-1} \\
 (37) \qquad \qquad \qquad &= O\left(\frac{t_N^2}{N} + \sum_{j=2-N}^{N-1} (j/t_N)^2\right) \\
 &= O\left(\frac{N^3}{t_N^2}\right),
 \end{aligned}$$

Figure 6 illustrates this bound for a numerical example. To relate the magnitude of

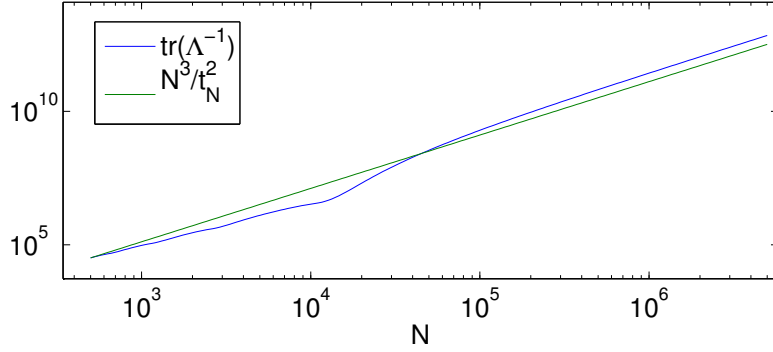


FIGURE 6. Numerical study of how  $\text{tr}(\Lambda^{-1})$  varies as a function of  $N$  in the modeling setting  $C = 10$ ,  $f_D = 50\text{Hz}$ ,  $(a^+)^2 p = \pi^{-1}$ ,  $t_N = \sqrt{N}$  and  $N$  uniform sample times on the interval  $[0, t_N)$ .

$s_N$  to  $\Lambda$ , we note that since  $K$  is circulant embedded into  $\mathcal{K}$ , Cauchy's interlacing theorem says that  $\min_j \Lambda_j \leq s_N \leq \tilde{\Lambda}_N$  with  $\tilde{\Lambda}_j$  denoting the  $j$ th largest eigenvalue of  $\Lambda$ . In settings where estimate (37) is valid, we thereby derive the upper bound  $N/s_N = O(N^3/t_N^2)$ , giving the relation  $N/s_N = O(\text{tr}(\Lambda^{-1}))$ . We end these informal estimates with the numerical study in Figure 7, showing a setting where  $N/s_N$  is orders of magnitude smaller than  $\text{tr}(\Lambda^{-1})$ .

## 6. NUMERICAL EXAMPLES

**6.1. Example 1.** The first numerical example compares realizations generated by Algorithm 1 and 2 in a WSS model setting with the scaled scatterer density  $(a^+)^2 p(\alpha) = 1/\pi$  and the delay function  $\tau(\alpha, t) = -v \cos(\alpha)t/c$ . The flip rate is determined from a real life signal measurement from a receiver moving with the speed  $v = 6.9\text{m/s}$  and carrier frequency  $f_c = 1.8775\text{GHz}$  sampled uniformly the time interval  $[-0.08, 0.08]$  using  $N = 1419$  samples, and signal realizations are thereafter generated using the same modeling parameters. From the measured signal  $\hat{Z}(t)$ , we approximate the PSD by the FFT of  $\hat{Z}(t)\hat{Z}(0)^* =: \hat{S}(f)$  and determine the flip

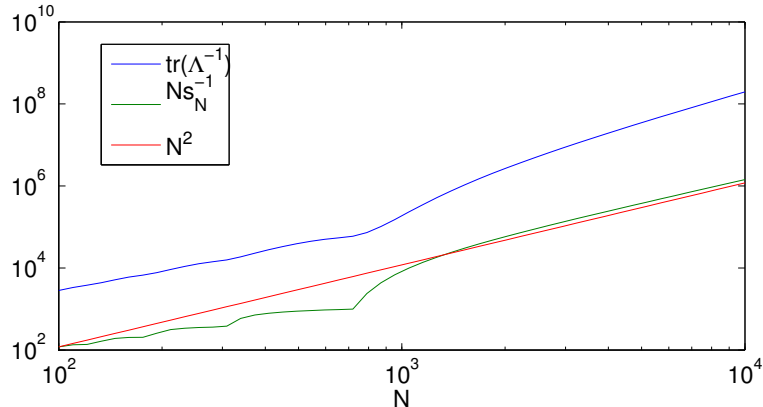


FIGURE 7. Numerical comparison of  $N/s_N$  and  $\text{tr}(\Lambda^{-1})$  in the modeling setting  $C = 15$ ,  $f_D = 50\text{Hz}$ ,  $(a^+)^2 p = \pi^{-1}$ ,  $t_N = \log(N)$  and  $N$  sample uniform sample times on the  $[0, t_N]$ . (Since  $s_N$  is computed by the SVD, it cannot be computed for as large  $N$  values as  $\text{tr}(\Lambda)$ , cf. Figure 6.)

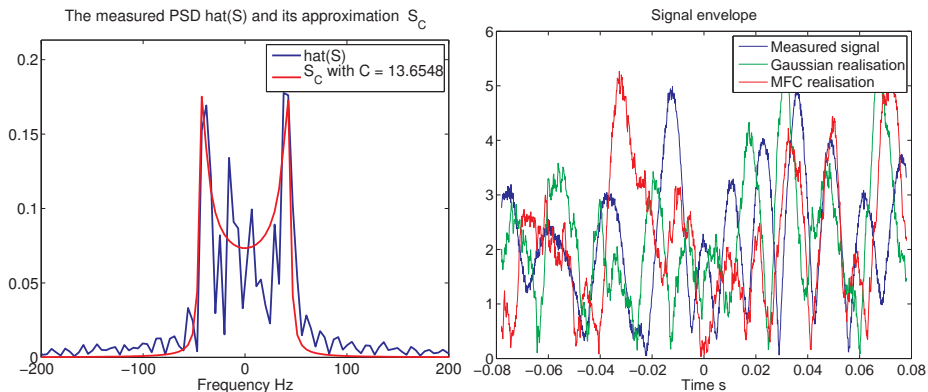


FIGURE 8. **Left plot:** The best fit of  $S_C(f)$  to the measured signal's PSD  $\hat{S}(f)$  as a function of the flip rate  $C$ , which is determined by (38). **Right plot:** Scaled signal envelopes of a measured signal, a signal realization from the MFC algorithm, and a realization from the Gaussian process algorithm.

rate by

$$(38) \quad C = \arg \min_{C > 0} \|(\hat{S} - S_{\bar{c}})\varrho\|_1,$$

which is a slight modification of Feng and Field's idea (26). See Figure 8 for an comparison of  $S_C$  and  $\hat{S}$  for the best fit flip rate  $C \approx 13.65$ .

Having determined the flip rate,  $L = 2000$  signal realizations are generated by the MFC algorithm using  $M = 2000$  scatterers and for the covariance based

Gaussian Process algorithm, one may derive using pen and paper that  $K(t_j, t_k) = J_0(2\pi f_D(t_k - t_j))$ ,  $J_0$  denoting the Bessel function of the first kind. Autocorrelation and PSD functions computed for respective algorithms are plotted in figures 8 and 9. For signal realizations of Algorithm 1, the autocorrelation is approximated by taking the sample average of

$$(39) \quad E[\bar{Z}_t \bar{Z}_0^*] \approx \sum_{j=1}^L \frac{\bar{Z}_t(\omega_j) \bar{Z}_0(\omega_j)^*}{L}$$

with  $\bar{Z}_t(\omega_j)$  denoting the  $j$ th signal realization. The approximation 39 is compared with the autocorrelation of signal realizations of Algorithm 2 (which is given by the first row of the covariance matrix  $K$ ).

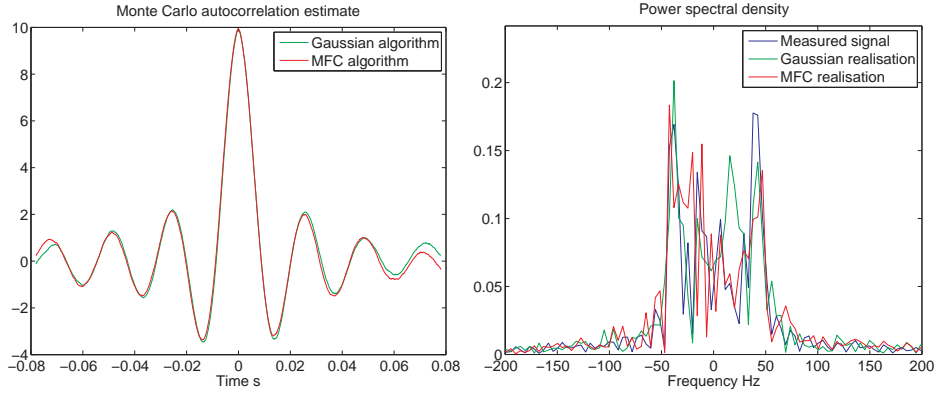


FIGURE 9. **Left plot:** Monte Carlo estimated autocorrelation for realizations generated by Algorithm 1, cf. (39), and the autocorrelation for Gaussian process signal realizations  $A(t) = J_0(2\pi f_D t)$ . **Right plot:** Scaled PSD of the measured signal and of one signal realization from each the algorithms studied.

**6.2. Example 2.** In the second example, we model the temporally varying scattering environment with a mobile receiver moving from left to right through a thin opening of a non-reflecting wall, as sketched in Figure 10. When situated on the left side of the opening, the mobile receiver receives scattered rays at its rear, and when the mobile receiver is on the right side of the opening, it receives rays at its front. As a model for this change in scattering environment we consider the the time interval  $[0, 2)$  seconds, assume that the receiver moves through the opening at  $t = 1$  and set

$$(40) \quad a^+(\alpha, t) = \begin{cases} \cos^2(\alpha) 1_{(\pi/2, 3\pi/2)}(\alpha) & \text{for } 0 \leq t < 1 \\ \cos^2(\alpha) 1_{(-\pi/2, \pi/2)}(\alpha) & \text{for } 1 \leq t \leq 2. \end{cases}$$

Other modeling parameters are set to  $p = (2\pi)^{-1}$ ,  $C = 15$  and  $f_D = 43.5\text{Hz}$ . Figure 11 contains snapshots of the time dependent PSD for a single stochastic signal realization created by Algorithm 2. It shows that when the mobile receiver is situated on the left side of the opening, the the PSD is concentrated around  $-f_D$ , and when the receiver is on the right side of the opening, the PSD is concentrated

around  $f_D$ . By the discussion of the relation between  $(a^+)^2 p$  and  $S_C(f)$  given in Section 4, the snapshotted PSDs in Figure 11 are reasonable.

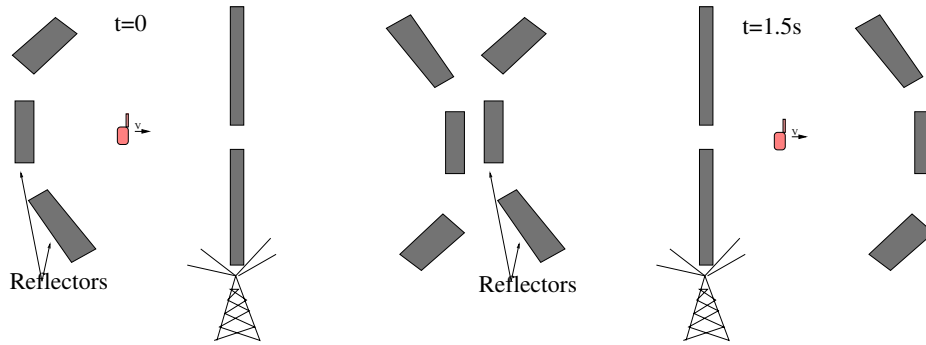


FIGURE 10. A receiver moving rightwards through a thin opening in a non-scattering wall and thereby experiencing a change in the scattering environment.

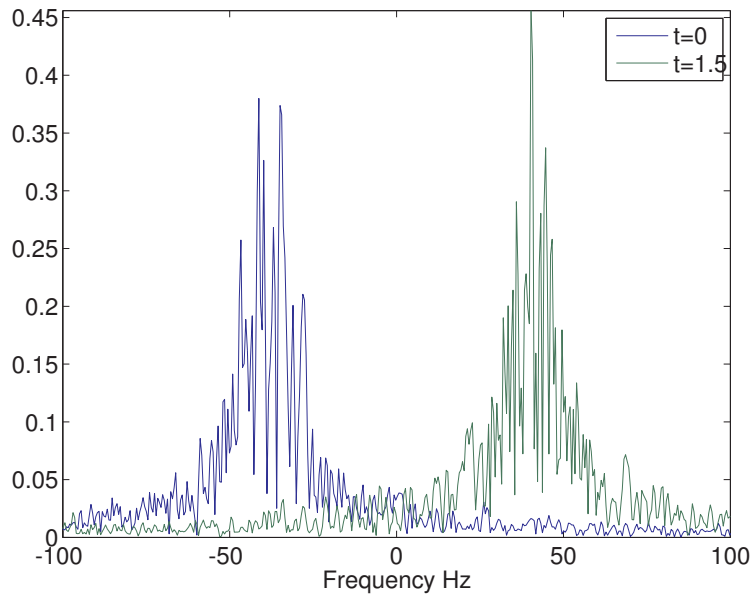


FIGURE 11. Snapshots of the time dependent PSD for a signal realization at  $t = 0.5s$  when the receiver is situated on the left side of the wall opening (blue line) and at  $t = 1.5s$  situated at to the right side of the wall opening (green line).

## APPENDIX A. THEOREMS

**Corollary A.1** ([8][p. 198]). *Let  $G$  and  $F$  be arbitrary matrices (of the same size) where  $\sigma_1 \geq \dots \geq \sigma_n$  are the singular values of  $G$  and  $\sigma'_1 \geq \dots \geq \sigma'_n$  are the singular values of  $G + F$ . Then  $|\sigma_i - \sigma'_i| \leq \|F\|_2$ .*

## APPENDIX B. APPROXIMATION OF THE DELAY FUNCTION

The delay function primarily considered in this paper,  $\tau(\alpha, t) = -vt \cos(\alpha)/c$ , is a first order approximation of the full delay function. Here we describe how this approximation is obtained.

Assuming the scattering boundary is described by  $\{(\alpha, R(\alpha)) | 0 \leq \alpha \leq 2\pi\}$  and that the receiver is moving in the direction  $(v, 0)$ , the analytical delay function is given by

$$\begin{aligned} \tau(\alpha, t) &= \frac{\sqrt{(R(\alpha) \cos(\alpha) - vt)^2 + R(\alpha)^2 \sin(\alpha)^2}}{c} \\ &= \frac{\sqrt{(vt)^2 - 2vtR(\alpha) \cos(\alpha) + R(\alpha)^2}}{c}. \end{aligned}$$

A Taylor expansion of this function with respect to  $vt$  yields

$$\tau(\alpha, t) = \frac{R(\alpha)}{c} - \frac{\cos(\alpha)}{c} vt + O\left(\frac{vt}{cR(\alpha)}\right).$$

Assuming that for the times of interest  $vt/(R(\alpha)c)$  is a small term, a good approximation of the delay function as a function of time is

$$\tau(\alpha, t) = -\frac{\cos(\alpha)}{c} vt.$$

Here the term  $R(\alpha)/c$  is removed since it is constant with respect to time and thus can be considered as part of the random phase shift term  $\theta_\alpha$  of  $\exp(-i(2\pi f_c \tau(\alpha, t) + \theta_\alpha))$ . As this first order approximation considers the Doppler effect from a given scatterer to be constant, it is only valid when the receiver is far away from the scattering boundary.

**Acknowledgement.** The authors would like to thank professor Anders Szepessy for fruitful discussions and Henrik Asplund at Ericsson Research for providing signal measurements and modeling extension ideas.

## REFERENCES

- [1] TR 25.996 3GPP. *3rd Generation Partnership Project; Technical Specification Group Radio Access Network; Spatial channel model for Multiple Input Multiple Output (MIMO) simulations (Release 6)*. 3GPP, 2003.
- [2] Søren Asmussen and Peter W. Glynn. *Stochastic simulation: algorithms and analysis*, volume 57 of *Stochastic Modelling and Applied Probability*. Springer, New York, 2007.
- [3] V. Bentkus. On the dependence of the Berry-Esseen bound on dimension. *J. Statist. Plann. Inference*, 113(2):385–402, 2003.
- [4] William L. Briggs and Van Emden Henson. *The DFT*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1995. An owner's manual for the discrete Fourier transform.
- [5] Charalambos D. Charalambous, Seddik M. Djouadi, and Christos Kourtellaris. Statistical analysis of multipath fading channels using generalizations of shot noise. *EURASIP J. Wirel. Commun. Netw.*, 2008:1–9, 2008.



- [6] Chia-Chin Chong, Chor-Min Tan, D.I. Laurenson, S. McLaughlin, M.A. Beach, and A.R. Nix. A novel wideband dynamic directional indoor channel model based on a markov process. *Wireless Communications, IEEE Transactions on*, 4(4):1539 – 1552, jul. 2005.
- [7] R. H. Clarke. A statistical theory of mobile-radio reception. *Bell Sys. Tech.*, 47:957–1000, 1968.
- [8] James W. Demmel. *Applied numerical linear algebra*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 1997.
- [9] Tao Feng and T.R. Field. Statistical analysis of mobile radio reception: an extension of clarke’s model. *Communications, IEEE Transactions on*, 56(12):2007–2012, december 2008.
- [10] Tao Feng, T.R. Field, and S. Haykin. Stochastic differential equation theory applied to wireless channels. *Communications, IEEE Transactions on*, 55(8):1478 –1483, aug. 2007.
- [11] Markos A. Katsoulakis and Anders Szepessy. Stochastic hydrodynamical limits of particle systems. *Commun. Math. Sci.*, 4(3):513–549, 2006.
- [12] J.I. Smith. A computer generated multipath fading simulation for mobile radio. *Vehicular Technology, IEEE Transactions on*, 24(3):39 – 40, aug 1975.
- [13] N. G. van Kampen. *Stochastic processes in physics and chemistry*, volume 888 of *Lecture Notes in Mathematics*. North-Holland Publishing Co., Amsterdam, 1981.
- [14] Curtis R. Vogel. *Computational methods for inverse problems*, volume 23 of *Frontiers in Applied Mathematics*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002. With a foreword by H. T. Banks.
- [15] D.J. Young and N.C. Beaulieu. The generation of correlated rayleigh random variates by inverse discrete fourier transform. *Communications, IEEE Transactions on*, 48(7):1114 – 1127, jul 2000.
- [16] T. Zwick, C. Fischer, and W. Wiesbeck. A stochastic multipath channel model including path directions for indoor environments. *Selected Areas in Communications, IEEE Journal on*, 20(6):1178 – 1192, aug. 2002.

HÅKON HOEL, DEPARTMENT OF NUMERICAL ANALYSIS AND COMPUTER SCIENCE, KTH, SE-100 44, STOCKHOLM, SWEDEN.

*E-mail address:* haakonah1@gmail.com

HENRIK NYBERG, ERICSSON RESEARCH, SE-164 80, STOCKHOLM, SWEDEN.

*E-mail address:* Henrik.L.Nyberg@ericsson.com



## Paper II

# Implementation and Analysis of an Adaptive Multilevel Monte Carlo Algorithm

Håkon Hoel\*, Erik von Schwerin†, Anders Szepessy‡ and Raúl Tempone§

April 23, 2012

## Abstract

This work generalizes a multilevel Monte Carlo (MLMC) method introduced in [7] for the approximation of expected values of functions depending on the solution to an Itô stochastic differential equation. The work [7] proposed and analyzed a forward Euler MLMC method based on a hierarchy of uniform time discretizations and control variates to reduce the computational effort required by a standard, single level, forward Euler Monte Carlo method from  $\mathcal{O}(\epsilon^{-3})$  to  $\mathcal{O}((\epsilon^{-1} \log(\epsilon^{-1}))^2)$  for a mean square error of size  $\epsilon^2$ . This work uses instead a hierarchy of adaptively refined, non uniform, time discretizations, generated by an adaptive algorithm introduced in [20, 19, 5]. Given a prescribed accuracy TOL in the weak error, this adaptive algorithm generates time discretizations based on a posteriori expansions of the weak error first developed in [24]. A theoretical analysis gives results on the stopping, the accuracy, and the complexity of the resulting adaptive MLMC algorithm. In particular, it is shown that: the adaptive refinements stop after a finite number of steps; the probability of the error being smaller than TOL is under certain assumptions controlled by a given confidence parameter, asymptotically as  $\text{TOL} \rightarrow 0$ ; the complexity is essentially the expected for MLMC methods, but with better control of the constant factors. We also show that the multilevel estimator is asymptotically normal using the Lindeberg-Feller Central Limit Theorem. These theoretical results are based on previously developed single level estimates, and results on Monte Carlo stopping from [3]. Our numerical tests include cases, one with singular drift and one with stopped diffusion, where the complexity of uniform single level method is  $\mathcal{O}(\text{TOL}^{-4})$ . In both these cases the results confirm the theory by exhibiting savings in the computational cost to achieve an accuracy of  $\mathcal{O}(\text{TOL})$ , from  $\mathcal{O}(\text{TOL}^{-3})$  for the adaptive single level algorithm to essentially  $\mathcal{O}(\text{TOL}^{-2} \log(\text{TOL}^{-1})^2)$  for the adaptive MLMC.

---

\*CSC, Royal Institute of Technology (KTH), Stockholm, Sweden

†Applied Mathematics and Computational Sciences, KAUST, Thuwal, Saudi Arabia

‡Mathematics, Royal Institute of Technology (KTH), Stockholm, Sweden

§Applied Mathematics and Computational Sciences, KAUST, Thuwal, Saudi Arabia  
(raul.tempone@kaust.edu.sa)

**Key words:** computational finance, Monte Carlo, multi-level, adaptivity, weak approximation, error control, Euler–Maruyama method, a posteriori error estimates, backward dual functions, adjoints

**AMS subject classification:** 65C30, 65Y20, 65L50, 65H35, 60H35, 60H10

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	A single level posteriori error expansion . . . . .	6
<b>2</b>	<b>Adaptive Algorithms and Multilevel Variance Reduction</b>	<b>10</b>
2.1	Path independent time stepping . . . . .	10
2.1.1	Generating the mesh hierarchy . . . . .	11
2.1.2	Multilevel simulations on a given hierarchy . . . . .	13
2.2	Fully stochastic time stepping . . . . .	14
2.3	Algorithm Listings . . . . .	17
<b>3</b>	<b>Numerical Experiments</b>	<b>20</b>
3.1	A Linear SDE . . . . .	21
3.2	Drift singularity, linear SDE . . . . .	24
3.3	Stopped diffusion . . . . .	25
<b>4</b>	<b>Theoretical results</b>	<b>32</b>
4.1	Single level results . . . . .	34
4.2	multilevel results . . . . .	41
<b>5</b>	<b>Conclusions</b>	<b>53</b>
<b>A</b>	<b>Theorems</b>	<b>54</b>

## 1 Introduction

This work develops multilevel versions of adaptive algorithms for weak approximation of Itô stochastic differential equations (SDEs)

$$dX(t) = a(t, X(t))dt + b(t, X(t))dW(t), \quad 0 < t < T, \quad (1.1)$$

where  $X(t; \omega)$  is a stochastic process in  $\mathbb{R}^d$ , with randomness generated by a  $k$ -dimensional Wiener process with independent components,  $W(t; \omega)$ ; see [15], [23]. The functions  $a(t, x) \in \mathbb{R}^d$  and  $b(t, x) \in \mathbb{R}^{d \times k}$  are given drift and diffusion fluxes.

Our goal is to, for any given sufficiently well behaved function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , approximate the expected value  $E[g(X(T))]$  by adaptive multilevel Monte Carlo (MLMC) methods. A typical example of such an expected value is to compute

option prices in mathematical finance; see [14] and [9]. Other related models based on stochastic dynamics are used for example in molecular dynamics simulations at constant temperature, for stochastic climate prediction, and for wave propagation in random media; cf. [2], [18], and [1].

The computational complexity of a Monte Carlo method is determined by the number of generated samples approximating  $g(X(T))$  and their average cost. When a standard Monte Carlo method based on a uniform time stepping scheme of weak order 1 is used to compute  $E[g(X(T))]$  to an accuracy TOL with high probability, the cost is asymptotically proportional to  $\text{TOL}^{-3}$ , provided that the functions  $a$ ,  $b$ , and  $g$  are sufficiently regular. A Monte Carlo method can not do better than a cost proportional to  $\text{TOL}^{-2}$ , since this is the cost if each sample of  $g(X(T))$  can be generated exactly at a unit cost. The goal of this work is to combine two techniques for improving the standard Monte Carlo: one technique is to use adaptive time stepping which retains the complexity  $\mathcal{O}(\text{TOL}^{-3})$  for a wider set of problems than uniform time stepping, and which can reduce the proportionality constant for other problems with widely varying scales; the other is the MLMC method, which in many cases can reduce the complexity to nearly the optimal  $\mathcal{O}(\text{TOL}^{-2})$  when based on the Euler–Maruyama time stepping scheme, and which can achieve the optimal rate using the Milstein time stepping scheme.

In the context of weak approximation of SDEs, the MLMC method based on uniform time stepping was introduced by Giles in [7]. A similar MLMC idea has been used before in the different context of parametric integration; [11, 12]. In [7], Giles developed a clever control variate type variance reduction technique for a numerical method, denoted here by  $\bar{X}$ , that approximates the solution of the SDE (1.1). The key to this variance reduction, which is an extension of a two-level control variate technique in [16], is to compute approximate solutions,  $\bar{X}_\ell$ , on hierarchies of uniform time meshes with size

$$\Delta t_\ell = C^{-\ell} \Delta t_0, \quad C \in \{2, 3, \dots\} \quad \text{and} \quad \ell \in \{0, 1, \dots, L\}, \quad (1.2)$$

thereby generating sets of realizations on different mesh levels. After computing numerical approximations on a mesh hierarchy, the expected value  $E[g(X(T))]$  is approximated by the multilevel Monte Carlo estimator

$$\begin{aligned} \mathcal{A}_{\text{Mc}}(g(\bar{X}(T)); M_0) &= \sum_{i=1}^{M_0} \frac{g(\bar{X}_0(T; \omega_{i,0}))}{M_0} \\ &+ \sum_{\ell=1}^L \sum_{i=1}^{M_\ell} \frac{g(\bar{X}_\ell(T; \omega_{i,\ell})) - g(\bar{X}_{\ell-1}(T; \omega_{i,\ell}))}{M_\ell}. \end{aligned} \quad (1.3)$$

This estimator combines  $L+1$  sample averages, based on mutually independent sample sets on the respective meshes, each with  $M_\ell$  independent samples. The number of samples  $\{M_\ell\}_{\ell=1}^L$  have a given relation to the number of samples on the coarsest mesh,  $M_0$ , which is the only free parameter in (1.3), for a fixed number of levels. To reduce the variance in the above estimator, the

realization pairs  $\bar{X}_\ell(T; \omega_{i,\ell})$  and  $\bar{X}_{\ell-1}(T; \omega_{i,\ell})$  of the summands  $g(\bar{X}_\ell(T; \omega_{i,\ell})) - g(\bar{X}_{\ell-1}(T; \omega_{i,\ell}))$  for each level  $\ell > 0$  are generated by the same Brownian path,  $W(t; \omega_{i,\ell})$ , but they are realized on different temporal grids with uniform time steps,  $\Delta t_\ell$  and  $\Delta t_{\ell-1}$ , respectively. The efficiency of this computation relies on an a priori known order of strong convergence for the numerical method employed on each level of the hierarchy.

Let  $\text{TOL} > 0$  be a desired accuracy in the approximation of  $\mathbb{E}[g(X(T))]$ . The main result of Giles' work [7] is that the computational cost needed to achieve the Mean Square Error (MSE)

$$\mathbb{E}\left[\left(\mathcal{A}_{\text{MC}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(X(T))]\right)^2\right] = \mathcal{O}(\text{TOL}^2), \quad (1.4)$$

when using the Forward Euler method to create the approximate realizations  $\bar{X}_\ell(T; \omega)$ , can be reduced from  $\mathcal{O}(\text{TOL}^{-3})$  with the standard Monte Carlo method to

$$\mathcal{O}((\text{TOL}^{-1} \log(\text{TOL}^{-1}))^2)$$

with Giles' MLMC method. Furthermore, whenever the function  $g$  is Lipschitz and for scalar Itô stochastic differential equations, the computational cost can be further reduced to  $\mathcal{O}(\text{TOL}^{-2})$  using the first order strong convergence Milstein method. In addition, the work [6] shows how to apply the Milstein method for several scalar SDE cases where the Lipschitz condition is not fulfilled and still obtain the cost  $\mathcal{O}(\text{TOL}^{-2})$ .

In this work we use the Euler–Maruyama method with non uniform time steps. Let  $0 = t_0 < t_1 < \dots < t_N = T$  denote a given time discretization, without reference to its place in the hierarchies, and  $\{0 = W(t_0; \omega), W(t_1; \omega), \dots, W(t_N; \omega)\}$  denote a generated sample of the Wiener process on that discretization. Then the Euler–Maruyama approximation to the true solution of (1.1) is given by the scheme

$$\begin{aligned} \bar{X}(t_0; \omega) &= X(0), \\ \bar{X}(t_{n+1}; \omega) &= a(\bar{X}(t_n; \omega), t_n) \Delta t_n + b(\bar{X}(t_n; \omega), t_n) \Delta W(t_n; \omega), \quad n \geq 0, \end{aligned} \quad (1.5)$$

where  $\Delta t_n = t_{n+1} - t_n$  and  $\Delta W(t_n; \omega) = W(t_{n+1}; \omega) - W(t_n; \omega)$  are the time steps and Wiener increments, respectively.

The contribution of the present paper to the MLMC method is the development and analysis of two novel algorithms with adaptive, non uniform time steps. One of the algorithms uses adaptive mesh refinements to stochastically create a path dependent mesh for each realization; the other algorithm constructs the meshes adaptively based on sample averaged error densities and then uses the same mesh hierarchy for all realizations. We refer to the previous algorithm as the *stochastic* time step algorithm and the latter as the *deterministic* time step algorithm. The construction and analysis of the adaptive algorithms is inspired by the work on single level adaptive algorithms for weak approximation of ordinary stochastic differential equations [19], and uses the adjoint weighted global error estimates first derived in [24]. The goal of these

adaptive algorithms is to choose the time steps and the number of realizations such that the event

$$|\mathcal{A}_{\mathcal{M}_\varepsilon}(g(\overline{X}(T)); M_0) - \mathbb{E}[g(X(T))]| \leq \text{TOL}, \quad (1.6)$$

holds with probability close to one. We measure computational complexity as the work needed to meet a given accuracy. On some problems where the computational complexity of methods based on uniform time steps deteriorates due to lacking regularity, these adaptive mesh refinement algorithms can regain the same rate of convergence that uniform methods have on regular problems; see [20].

It should be noted that in the setting of adaptive mesh refinement there is no given notion of mesh size, so a hierarchy of meshes can no longer be described as in the constant time step case (1.2). Instead, we generate a hierarchy of meshes by successively increasing the accuracy in our computations, introducing the time discretization error tolerance levels<sup>1</sup>

$$\text{TOL}_{\text{T},\ell} = 2^{\ell-L} \text{TOL}_{\text{T}}, \quad \text{for } \ell \in \{0, 1, \dots, L\}, \quad (1.7)$$

and (by adaptive refinements based on error indicators) determining the corresponding meshes so that for each level  $\ell \in \{0, 1, \dots, L\}$ ,

$$|\mathbb{E}[g(X(T))] - \mathbb{E}[g(\overline{X}_\ell(T))]| \lesssim \text{TOL}_{\text{T},\ell}.$$

In Section 4, we prove that also for the adaptive algorithms the computational cost for obtaining the error estimate (1.4), with probability close to one, is close to  $\mathcal{O}(\text{TOL}^{-2} \log(\text{TOL}^{-1})^2)$ . More precisely, there are two main results on efficiency and accuracy described in Section 4. Regarding accuracy with probability close to one, Theorem 2 states that the approximation errors in (1.2) are asymptotically bounded by the specified error tolerance times a problem independent factor as the tolerance parameter tends to zero. For the efficiency, Theorem 3 states that, depending on technical assumptions on the error density of the adaptive algorithm, the complexity is close to or equal to the standard complexity in the setting of uniform time steps, with explicitly given constants. A completely analogous result holds for the adaptive algorithm with deterministic steps and it is not included here for the sake of brevity; confer [19].

This work also includes three numerical examples, the most relevant ones being one with a drift singularity and one stopped diffusion. In the singularity example multilevel Monte Carlo based on adaptive time steps requires a computational work  $\mathcal{O}(\text{TOL}^{-2} \log(\text{TOL}^{-1})^2)$  and in the stopped diffusion example  $\mathcal{O}(\text{TOL}^{-2.1} \log(\text{TOL}^{-1})^2)$ . For both of these examples the observed complexity is close to the optimal, and more efficient than the single level version of the adaptive algorithm.

The rest of this paper is organized as follows: Section 1.1 introduces the notion of error density and error indicators, and recalls useful results for single

---

<sup>1</sup>For error control, the tolerance is split into a statistical error tolerance and a time discretization error tolerance;  $\text{TOL} = \text{TOL}_{\text{S}} + \text{TOL}_{\text{T}}$ , cf. Section 2.



level adaptive forward Euler Monte Carlo methods. Section 2 describes the new adaptive multilevel Monte Carlo algorithms. Section 3 presents results from numerical experiments. Finally, Section 4 proves results on accuracy, stopping, and efficiency for a simplified version of the adaptive multilevel algorithms.

### 1.1 A single level posteriori error expansion

Here we recall previous single level results that are useful for the multilevel analysis in Section 4. In particular, we recall adjoint based error expansions with computable leading order term. Assume that the process  $X$  satisfies (1.1) and its approximation,  $\bar{X}$ , is given by (1.5); then the error expansions in Theorem 1.2 and 2.2 of [24] have the form

$$\mathbb{E}[g(X(T)) - g(\bar{X}(T))] = \mathbb{E}\left[\sum_{n=1}^N \rho_n \Delta t_n^2\right] + \text{higher order terms}, \quad (1.8)$$

where  $\rho_n \Delta t_n^2$  are computable error indicators, that is they provide information for further improvement of the time mesh and  $\rho_n$  measures the density of the global error in (1.8). A typical adaptive algorithm does two things iteratively:

1. if the error indicators satisfy an accuracy condition then it stops; otherwise
2. the algorithm chooses where to refine the mesh based on the error indicators and then makes an iterative step to 1.

In addition to estimating the global error  $\mathbb{E}[g(X(T)) - g(\bar{X}(T))]$  in the sense of equation (1.8), the error indicators  $\rho_n \Delta t_n^2$  also give simple information on where to refine to reach an optimal mesh, based on the almost sure convergence of the density  $\rho_n$  as we refine the discretization, see Section 4 in [20].

In the remaining part of this section we state in Theorem 1 a single level error expansion from [24], which can be used with either stochastic or deterministic time steps.

Given an initial time discretization  $\Delta t[0](t)$  and, for the stochastic time steps algorithm, refining until<sup>2</sup>

$$|\rho(t, \omega)|(\Delta t(t))^2 < \text{constant}, \quad (1.9)$$

we construct a partition  $\Delta t(t)$  by repeated halving of intervals so that it satisfies

$$\Delta t(t) = \Delta t[0](t)/2^n \quad \text{for some natural number } n = n(t, \omega).$$

The criterion (1.9) uses an approximate error density function  $\rho$ , satisfying for  $t \in [0, T]$  and all outcomes  $\omega$  the uniform upper and lower bounds

$$\rho_{low}(\text{TOL}_T) \leq |\rho(t, \omega)| \leq \rho_{up}(\text{TOL}_T). \quad (1.10)$$

---

<sup>2</sup>The precise expressions including the constants are given in (2.7) and (2.21) below.

The positive functions  $\rho_{low}$  and  $\rho_{up}$  are chosen so that  $\rho_{up}(\text{TOL}_T) \rightarrow +\infty$  as  $\text{TOL} \downarrow 0$  while  $\rho_{low}(\text{TOL}_T) \rightarrow 0$  such that  $\text{TOL}_T/\rho_{low}(\text{TOL}_T) \rightarrow 0$ . We further make the assumption that for all  $s, t \in [0, T]$  the sensitivity of the error density to values of the Wiener process can be bounded,

$$|\partial_{W(t)}\rho(s, \omega)| \leq D\rho_{up}(\text{TOL}_T), \quad (1.11)$$

for some positive function  $D\rho_{up}$  such that  $D\rho_{up}(\text{TOL}_T) \rightarrow +\infty$  as  $\text{TOL}_T \downarrow 0$ . For each realization successive subdivisions of the steps yield the largest time steps satisfying (1.9). The corresponding stochastic increments  $\Delta W$  will have the correct distribution, with the necessary independence, if the increments  $\Delta W$  related to the new steps are generated by Brownian bridges [15]; that is the time steps are generated by conditional expected values of the Wiener process.

We begin now by stating in the next lemma the regularity conditions to be used in the analysis of the adaptive multilevel algorithms; the lemma corresponds to Lemma 2.1 in [24], while this formulation is given<sup>3</sup> in Lemma 2.1 in [21].

**Lemma 1** (Regularity). (a) *Assume that the following regularity conditions hold:*

- (1) *The functions  $a(t, x)$  and  $b(t, x)$  are continuous in  $(t, x)$  and are twice continuously differentiable with respect to  $x$ .*
- (2) *The partial derivatives of first and second order with respect to  $x$  of the functions  $a$  and  $b$  are uniformly bounded.*
- (3) *The function  $g$  is twice continuously differentiable, and together with its partial derivatives of first and second order it is uniformly bounded.*

*Then the cost to go function, defined by*

$$u(t, x) = E[g(X(T)) | X(t) = x], \quad (1.12)$$

*satisfies the Kolmogorov equation*

$$\partial_t u(t, x) + a_k \partial_k u(t, x) + d_{kn} \partial_{kn} u(t, x) = 0, \quad u(T, \cdot) = g, \quad (1.13)$$

*where we have used Einstein summation convention<sup>4</sup>, and where  $d_{kn} = \frac{1}{2} b_k^l b_n^l$ .*

(b) *Furthermore, if the following regularity conditions are satisfied:*

- (1) *The functions  $\partial_\beta a(t, \cdot)$  and  $\partial_\beta b(t, \cdot)$  are bounded uniformly in  $t$  for multi-indices  $\beta$  with  $1 \leq |\beta| \leq 8$ ;*
- (2) *The functions  $a(\cdot, x)$ ,  $b(\cdot, x)$  have continuous and uniformly bounded first order time derivatives;*

<sup>3</sup>Including a jump term omitted here.

<sup>4</sup>When an index variable appears twice in a single term this means that a summation over all possible values of the index takes place; for example  $a_k \partial_k u(t, x) = \sum_{k=1}^d a_k \partial_k u(t, x)$ , where  $d$  is the space dimension of the SDE ( $a, x \in \mathbb{R}^d$ ).

(3) The function  $g$  has spatial derivatives  $\partial_\beta g$ , with polynomial growth for  $|\beta| \leq 8$ ;

then the function  $u$  has continuous partial derivatives with respect to  $x$  up to the order 8, satisfying the following polynomial growth condition: for all  $i \in \{0, 1, 2\}$  and  $\alpha \in \mathbb{N}^d$  with  $i + |\alpha| \leq 8$  there exists  $p_{\alpha,i} \in \mathbb{N}$  and  $C_{\alpha,i} > 0$  such that

$$\max_{0 \leq t \leq T} |\partial_t^i \partial_\alpha u(t, x)| \leq C_{\alpha,i} (1 + |x|^{p_{\alpha,i}}) \quad \forall x \in \mathbb{R}^d.$$

In what follows, Lemma 2 and Theorem 1 show that although the steps adaptively generated to satisfy (1.9)–(1.11) are not adapted to the natural Wiener filtration, the method indeed converges to the correct limit, which is the same as the limit of the forward Euler method with adapted time steps. Lemma 2 and Theorem 1 correspond to Lemma 3.1 and Theorem 3.3 in [24].

**Lemma 2** (Strong Convergence). *For  $X$  the solution of (1.1) suppose that  $a$ ,  $b$ , and  $g$  satisfy the assumptions in Lemma 1, that  $\bar{X}$  is constructed by the forward Euler method, based on the stochastic time stepping algorithm defined in Section 2, with step sizes  $\Delta t_n$  satisfying (1.9)–(1.11), and that their corresponding  $\Delta W_n$  are generated by Brownian bridges. Then*

$$\sup_{0 \leq t \leq T} \mathbb{E} \left[ |X(t) - \bar{X}(t)|^2 \right] = \mathcal{O}(\Delta t_{\text{sup}}) = \mathcal{O} \left( \frac{\text{TOL}_T}{\rho_{\text{low}}(\text{TOL}_T)} \right) \rightarrow 0 \quad (1.14)$$

as  $\text{TOL}_T \downarrow 0$ , where  $\Delta t_{\text{sup}} \equiv \sup_{n,\omega} \Delta t_n(\omega)$ .

In Theorem 1 and the rest of this work, we will use Einstein summation convention with respect to functional and spatial indices, but not with respect to the temporal one (usually denoted  $t_n$ ).

**Theorem 1** (Single level stochastic time steps error expansion). *Given the assumptions in Lemma 2 and a deterministic initial value  $X(0)$ , the time discretization error in (1.8) has the following expansion, based on both the drift and diffusion fluxes and the discrete dual functions  $\varphi$ ,  $\varphi'$ , and  $\varphi''$  given in (1.17)–(1.22), with computable leading order terms:*

$$\begin{aligned} \mathbb{E}[g(X(T))] - \mathbb{E}[g(\bar{X}(T))] &= \mathbb{E} \left[ \sum_{n=0}^{N-1} \tilde{\rho}(t_n, \omega) (\Delta t_n)^2 \right] \\ &+ \mathcal{O} \left( \left( \frac{\text{TOL}_T}{\rho_{\text{low}}(\text{TOL}_T)} \right)^{1/2} \left( \frac{\rho_{\text{up}}(\text{TOL}_T)}{\rho_{\text{low}}(\text{TOL}_T)} \right)^\epsilon \right) \mathbb{E} \left[ \sum_{n=0}^{N-1} (\Delta t_n)^2 \right], \end{aligned} \quad (1.15)$$

for any  $\epsilon > 0$  and where

$$\begin{aligned} \tilde{\rho}(t_n, \omega) &\equiv \frac{1}{2} \left( (\partial_t a_k + \partial_j a_k a_j + \partial_{ij} a_k d_{ij}) \varphi_k(t_{n+1}) \right. \\ &+ (\partial_t d_{km} + \partial_j d_{km} a_j + \partial_{ij} d_{km} d_{ij} + 2\partial_j a_k d_{jm}) \varphi'_{km}(t_{n+1}) \\ &\left. + (2\partial_j d_{km} d_{jr}) \varphi''_{kmr}(t_{n+1}) \right) \end{aligned} \quad (1.16)$$

and the terms in the sum of (1.16) are evaluated at the a posteriori known points  $(t_n, \bar{X}(t_n))$ , i.e.,

$$\partial_\alpha a \equiv \partial_\alpha a(t_n, \bar{X}(t_n)), \quad \partial_\alpha b \equiv \partial_\alpha b(t_n, \bar{X}(t_n)), \quad \partial_\alpha d \equiv \partial_\alpha d(t_n, \bar{X}(t_n)).$$

Here  $\varphi \in \mathbb{R}^d$  is the solution of the discrete dual backward problem

$$\begin{aligned} \varphi_i(t_n) &= \partial_i c_j(t_n, \bar{X}(t_n)) \varphi_j(t_{n+1}), \quad t_n < T, \\ \varphi_i(T) &= \partial_i g(\bar{X}(T)), \end{aligned} \quad (1.17)$$

with

$$c_i(t_n, x) \equiv x_i + \Delta t_n a_i(t_n, x) + \Delta W_n^\ell b_i^\ell(t_n, x) \quad (1.18)$$

and its first and second variation

$$\varphi'_{ij} \equiv \partial_{x_j(t_n)} \varphi_i(t_n) \equiv \frac{\partial \varphi_i(t_n; \bar{X}(t_n) = x)}{\partial x_j}, \quad (1.19)$$

$$\varphi''_{ikm}(t_n) \equiv \partial_{x_m(t_n)} \varphi'_{ik}(t_n) \equiv \frac{\partial \varphi'_{ik}(t_n; \bar{X}(t_n) = x)}{\partial x_m}, \quad (1.20)$$

which satisfy

$$\begin{aligned} \varphi'_{ik}(t_n) &= \partial_i c_j(t_n, \bar{X}(t_n)) \partial_k c_p(t_n, \bar{X}(t_n)) \varphi'_{jp}(t_{n+1}) \\ &\quad + \partial_{ik} c_j(t_n, \bar{X}(t_n)) \varphi_j(t_{n+1}), \quad t_n < T, \\ \varphi'_{ik}(T) &= \partial_{ik} g(\bar{X}(T)), \end{aligned} \quad (1.21)$$

and

$$\begin{aligned} \varphi''_{ikm}(t_n) &= \partial_i c_j(t_n, \bar{X}(t_n)) \partial_k c_p(t_n, \bar{X}(t_n)) \partial_m c_r(t_n, \bar{X}(t_n)) \varphi''_{jpr}(t_{n+1}) \\ &\quad + \partial_{im} c_j(t_n, \bar{X}(t_n)) \partial_k c_p(t_n, \bar{X}(t_n)) \varphi'_{jp}(t_{n+1}) \\ &\quad + \partial_i c_j(t_n, \bar{X}(t_n)) \partial_{km} c_p(t_n, \bar{X}(t_n)) \varphi'_{jp}(t_{n+1}) \\ &\quad + \partial_{ik} c_j(t_n, \bar{X}(t_n)) \partial_m c_p(t_n, \bar{X}(t_n)) \varphi'_{jp}(t_{n+1}) \\ &\quad + \partial_{ikm} c_j(t_n, \bar{X}(t_n)) \varphi_j(t_{n+1}), \quad t_n < T, \\ \varphi''_{ikm}(T) &= \partial_{ikm} g(\bar{X}(T)), \end{aligned} \quad (1.22)$$

respectively.

The previous result can also be directly applied to the particular case of deterministic time steps. Observe that the constant in  $\mathcal{O}$  that appears in (1.15) may not be uniform with respect to the value  $\epsilon$ . Thus, in practice one chooses  $\epsilon = \epsilon(\text{TOL})$  to minimise the contribution of the remainder term to the error expansion (1.15).

Let us now discuss how to modify the error density  $\tilde{\rho}(t_n, \omega)$  in (1.16) to satisfy the bounds (1.10) and at the same time guarantee that  $\Delta t_{\text{sup}} \rightarrow 0$  as  $\text{TOL}_T \downarrow 0$ ; see Lemma 2.

Consider, for  $t \in [t_n, t_{n+1})$  and  $n = 1, \dots, N$ , the piecewise constant function

$$\rho(t) \equiv \text{sign}(\tilde{\rho}(t_n)) \min \left( \max(|\tilde{\rho}(t_n)|, \rho_{\text{low}}(\text{TOL}_T)), \rho_{\text{up}}(\text{TOL}_T) \right), \quad (1.23)$$

where

$$\begin{aligned}\rho_{low}(\text{TOL}_T) &= \text{TOL}_T^{\bar{\gamma}}, & 0 < \bar{\gamma} < \frac{\alpha}{\alpha+2}, & \quad 0 < \alpha < \frac{1}{2}, \\ \rho_{up}(\text{TOL}_T) &= \text{TOL}_T^{-r}, & r > 0,\end{aligned}\tag{1.24}$$

and with the standard notation for the function sign, that is  $\text{sign}(x) = 1$  for  $x \geq 0$  and  $-1$  for  $x < 0$ . The function  $\rho$  defined by (1.23) measures the density of the time discretisation error; the stochastic time stepping algorithm uses it in (2.20) and (2.21) to guide the mesh refinements, while the deterministic (path independent) time stepping algorithm uses a sample average of it for the same purpose; see (2.6) and (2.7). From now on, with a slight abuse of notation,  $\rho(t_n) = \rho_n$  denotes the modified density (1.23).

Following the error expansion in Theorem 1, the time discretization error is approximated by

$$|\mathcal{E}_T| = |\mathbb{E}[g(X(T)) - g(\bar{X}(T))]| \lesssim \mathbb{E}\left[\sum_{n=1}^N r(n)\right]\tag{1.25}$$

using the error indicator,  $r(n)$ , defined by

$$r(n) \equiv |\rho(t_n)|\Delta t_n^2\tag{1.26}$$

with the modified error density defined by (1.23). According to Corollary 4.3 and Theorem 4.5 in [19], we have the almost sure convergence of the error density to a limit density denoted by  $\hat{\rho}$ ,  $\rho \rightarrow \hat{\rho}$  as  $\text{TOL}_T \downarrow 0$ .

## 2 Adaptive Algorithms and Multilevel Variance Reduction

In this section we will describe two versions of the adaptive MLMC algorithm suitable for different problem settings. The first algorithm version which we present in Section 2.1 is designed for problems with deterministic time dependence on the drift and the diffusion. This algorithm version constructs a mesh hierarchy by adaptive refinements based on comparatively small sample sets and then performs a greater number of realizations on the constructed meshes to control the statistical error. The second algorithm version which we present in Section 2.2 is designed for problems where stochastic effects motivate the adaptive refinements. For this algorithm version the computational meshes are constructed by adaptive refinements for each individual realization of the underlying Wiener process.

### 2.1 Path independent time stepping

We recall that for a given SDE (1.1), function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$ , and tolerance  $\text{TOL} > 0$ , our goal is to construct an MLMC algorithm for which the event

$$|\mathcal{A}_{\text{MLC}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(X(T))]| \leq \text{TOL},$$

holds with probability close to one for the MLMC estimator  $\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0)$  defined by (1.3). We approach this goal by splitting the above approximation error as follows

$$\begin{aligned} & |\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(X_T)]| \\ & \leq \underbrace{|\mathbb{E}[g(\bar{X}_L(T)) - g(X(T))]|}_{=:\mathcal{E}_T} + \underbrace{|\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(\bar{X}_L(T))]|}_{=:\mathcal{E}_S}, \end{aligned}$$

and control the total error by requiring that the time discretization error fulfills  $\mathcal{E}_T \lesssim \text{TOL}_T$  and the statistical error fulfills  $\mathcal{E}_S \leq \text{TOL}_S$ , where the tolerance also have been split into a time discretization error tolerance and a statistical error tolerance,

$$\text{TOL} = \text{TOL}_T + \text{TOL}_S.$$

The computations then naturally divides into two phases. The first phase, consisting of Algorithm 1 and Algorithm 2, constructs a hierarchy of grids to control the time discretization error  $\mathcal{E}_T$ . The second phase, consisting of Algorithm 3 and Algorithm 4, computes a sufficiently large number of Euler–Maruyama realizations (1.5) on the constructed hierarchy of grids to ensure that  $\mathcal{E}_S \leq \text{TOL}_S$ , with probability close to one.

### 2.1.1 Generating the mesh hierarchy

We start with generating a hierarchy of meshes  $\{\Delta t_\ell\}_{\ell=0}^L$  for numerical approximation of the SDE (1.1). The meshes are sequentially and adaptively refined from a given initial mesh  $\Delta t_{-1}$  such that  $\Delta t_{\ell-1} \subset \Delta t_\ell$  for all  $\ell \in \{0, 1, \dots, L\}$ . On level  $\ell$  the grid is constructed with the aim that the time discretization error in the approximation of  $\mathbb{E}[g(\bar{X}_\ell(T))]$  fulfills

$$\left| \mathbb{E}[g(\bar{X}_\ell(T)) - g(X(T))] \right| < 2^{L-\ell} \text{TOL}_T =: \text{TOL}_{T,\ell}, \quad (2.1)$$

where  $\bar{X}_\ell(T)$  denotes an Euler–Maruyama approximation of the SDE (1.1) on the grid  $\Delta t_\ell$ . The number of grid levels  $L$  is chosen so that the largest tolerance

$$\text{TOL}_{T,0} = 2^L \text{TOL}_T, \quad (2.2)$$

is much larger than  $\text{TOL}_T$  and results in quite coarse meshes on the lower levels  $\Delta t_0, \Delta t_1, \dots$ . To be more precise, with a rough estimate of the magnitude of  $\mathbb{E}[g(X(T))]$  taken into account we prescribe an upper bound  $\text{TOL}_{T,\text{Max}}$  for  $\text{TOL}_{T,0}$  and determine  $L$  by the equation

$$L = \lfloor \log_2(\text{TOL}_{T,\text{Max}}/\text{TOL}_T) \rfloor. \quad (2.3)$$

The inputs in Algorithm 1 are: initial mesh  $\Delta t_{-1}$ , initial number of sample realizations  $M_{-1}$ , time discretization error tolerance  $\text{TOL}_T$ , grid levels  $L$ , initial estimate of the number of time steps on the accepted coarse mesh  $\bar{N}_0$ , and the three parameters  $C_R$ ,  $C_S$ , and  $R$  which are all used in the refinement and

stopping conditions (2.7), (2.6), and (2.10), respectively. We choose the initial estimated number of time steps  $\bar{N}_0$  as a small integer no smaller than the number of steps in  $\Delta t_{-1}$ . For smaller tolerances,  $\ell > 0$ , the initial guess is calculated according to  $\bar{N}_{\ell, \text{init}} = 2\bar{N}_{\ell-1, \text{accepted}}$ .

The grid refinement algorithm uses sample averages and sample variances of computed error indicators (1.23). With this in mind, we introduce some notation. Consider a set of  $M$  independent, identically distributed samples from the probability domain, and for such a sample set, introduce the sample average operator

$$\mathcal{A}(f; M) := \frac{1}{M} \sum_{i=1}^M f(\omega_i) \quad (2.4)$$

and, similarly, define the sample variance operator

$$\mathcal{V}(f; M) := \frac{1}{M-1} \sum_{i=1}^M \left( f(\omega) - \mathcal{A}(f; M) \right)^2. \quad (2.5)$$

When refining the grid  $\Delta t_\ell$ , the sampled error indicators  $r_\ell(n)$ , as defined in equation (1.26), are computed for all the time steps of the grid. Let  $N_\ell$  denote the number of time steps and  $\bar{N}_\ell$  be an estimate of  $N_\ell$  as described above. Then, if the stopping condition

$$\max_{1 \leq n \leq N_\ell} \mathcal{A}(r_\ell(n); M_\ell) < C_S \frac{\text{TOL}_\ell}{\bar{N}_\ell}, \quad (2.6)$$

is fulfilled, the grid is accepted and the refinements stop; otherwise the  $n$ -th time step is refined by splitting it into two equal parts if

$$\mathcal{A}(r_\ell(n); M_\ell) \geq C_R \frac{\text{TOL}_\ell}{\bar{N}_\ell}. \quad (2.7)$$

Normally, the value for  $C_R$  would be around 2, and one must take  $C_S > C_R$  following the theory developed in [20, 19].

The adaptive refinements of the computational grid are based on the sample averaged error indicators  $\mathcal{A}(r_\ell(n); M_\ell)$ . To estimate the mean error indicators  $\mathbb{E}[r_\ell(n)]$  with sufficient accuracy, we need a mechanism for deciding how many samples to use in the sample averages. With  $\mathbf{E}_{\Delta t_\ell}$  denoting the computed estimate of the time discretization error, i.e.,

$$\mathbf{E}_{\Delta t_\ell} = \sum_{n=1}^{N_\ell} \mathcal{A}(r_\ell(n); M_\ell), \quad (2.8)$$

a reasonable reliability requirement is

$$\sqrt{\text{Var}(\mathbf{E}_{\Delta t_\ell})} < R \mathbb{E}[\mathbf{E}_{\Delta t_\ell}], \quad (2.9)$$

for some suitably chosen  $0 < R < 1$ . (In our numerical examples, for instance, we use  $R = 0.2$ .) The variance of  $\mathbf{E}_{\Delta t_\ell}$  is however unknown, but the i.i.d.

distribution of the error indicators sampled motivates the approximation

$$\text{Var}(\mathbf{E}_{\Delta t_\ell}) \approx \frac{\mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)}{M_\ell} \quad \text{for } \ell = 0, 1, \dots, L.$$

We consequently approximate the reliability requirement (2.9) by

$$\sqrt{\frac{\mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)}{M_\ell}} < R \mathbf{E}_{\Delta t_\ell}, \quad \text{for } \ell = 0, 1, \dots, L, \quad (2.10)$$

where the number of sample realizations  $M_\ell$  used on level  $\ell$  in the grid construction phase is increased by repeated doubling, i.e.,  $M_{\ell, \text{new}} = 2 M_{\ell, \text{old}}$ , until inequality (2.10) is satisfied. The initial batch size at each level is set by  $M_\ell = M_{\ell-1}$ , where we for the moment let  $M_{\ell-1}$  denote the stopped number of samples at level  $\ell - 1$ , and for level  $\ell = 0$  it turns out to be sufficient to use initial batch size  $M_0 = M_{-1}$  with

$$M_{-1} = \text{const} \cdot \text{TOL}_T^{-1}. \quad (2.11)$$

The adaptive algorithm that generates the above described mesh hierarchy for the deterministic time step version of the MLMC algorithm is presented in Algorithm 1–2 in Section 2.3.

### 2.1.2 Multilevel simulations on a given hierarchy

In the second phase we will describe the algorithms which ensure that our MLMC estimate of  $\mathbf{E}[g(\bar{X}_L(T))]$  fulfills the statistical error bound

$$\mathcal{E}_S = |\mathcal{A}_{\mathcal{MC}}(g(\bar{X}(T)); M_0) - \mathbf{E}[g(\bar{X}_L(T))]| \leq \text{TOL}_S, \quad (2.12)$$

with probability close to one. We recall from (1.3) that the MLMC estimator is defined by

$$\mathcal{A}_{\mathcal{MC}}(g(\bar{X}(T)); M_0) = \mathcal{A}(g(\bar{X}_0(T)); M_0) + \sum_{\ell=1}^L \mathcal{A}(g(\bar{X}_\ell(T)) - g(\bar{X}_{\ell-1}(T)); M_\ell), \quad (2.13)$$

where the realization pairs  $\bar{X}_\ell(T; \omega_{i,\ell})$  and  $\bar{X}_{\ell-1}(T; \omega_{i,\ell})$  of the summands  $g(\bar{X}_\ell(T; \omega_{i,\ell})) - g(\bar{X}_{\ell-1}(T; \omega_{i,\ell}))$  for each level  $\ell > 0$  are generated by the Euler–Maruyama method (1.5) using the same Brownian path  $W(t; \omega_{i,\ell})$  on the respective *different* temporal meshes  $\Delta t_\ell$  and  $\Delta t_{\ell-1}$  that were computed by Algorithm 1. Furthermore, all Brownian paths  $\{W(t; \omega_{i,\ell})\}_{i,\ell}$  are independent, and the number of samples at the coarsest level is set to  $M_0 = 2^{L + \lceil C_{\mathcal{MC}} L \rceil + 1}$  for a suitable constant  $C_{\mathcal{MC}} \in (0, 1)$ , cf. Remark 3, and the number of samples on higher levels is expressed in terms of  $M_0$  by the ratio

$$M_\ell = \frac{M_0}{2^L} \left[ \frac{2^L \rho_{\text{low}}(\text{TOL}_{T,0}) \text{TOL}_{T,\ell}}{\rho_{\text{low}}(\text{TOL}_{T,\ell}) \text{TOL}_{T,0}} \right], \quad \ell = 1, \dots, L, \quad (2.14)$$



where  $\rho_{low}$  is the lower bound for the error density introduced in (1.24) and  $\lceil \cdot \rceil$  denotes rounding upwards to the nearest integer. The enforced lower bound for the sample sets  $\{M_\ell\}_{\ell=0}^L$  implies that  $M_L \rightarrow \infty$  as  $\text{TOL} \downarrow 0$ , and this motivates the approximation of

$$\frac{\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(\bar{X}_L(T))]}{\sqrt{\text{Var}(\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0))}}$$

by a normal distributed random variable; see Lemma 8 in Section 4 for a justification of this approximation for the stochastic time step setting. Relying on this approximation, the statistical error (2.12) will be controlled by bounding the MLMC estimator variance  $\sqrt{\text{Var}(\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0))} \leq C_C \text{TOL}_S$ , for a given positive confidence parameter  $C_C$ . The variance  $\text{Var}(\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0))$  is however unknown, so we introduce the following approximation

$$\text{Var}(\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0)) \approx \underbrace{\frac{\mathcal{V}(g(\bar{X}_0(T)); M_0)}{M_0} + \sum_{\ell=1}^L \frac{\mathcal{V}(g(\bar{X}_\ell(T)) - g(\bar{X}_{\ell-1}(T)); M_\ell)}{M_\ell}}_{=:\sigma^2}. \quad (2.15)$$

Our stopping criterion for the Monte Carlo simulations then becomes

$$\sigma < \frac{\text{TOL}_S}{C_C}. \quad (2.16)$$

Until this condition is fulfilled, the number of samples are iteratively doubled ( $M_0 = 2M_0$ ) and the number of samples at the levels  $\{M_\ell\}_{\ell=1}^L$  are updated according the ratio (2.14), and a new sample estimate  $\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0)$  is generated using the MLMC estimator (2.13). Having determined  $M_0$ , we lastly generate and return the output MLMC estimate  $\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0)$ .

The probability of controlling the statistical error, i.e., fulfilling the event (2.12) depends on the chosen value for the confidence parameter  $C_C$ . For example, with  $C_C = 1.65$  the event

$$|\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(\bar{X}_L(T))]| < C_C \sigma,$$

occurs with probability greater than 0.9, asymptotically as  $\text{TOL} \downarrow 0$ . See Algorithm 3–4 in Section 2.3 for more details on the MLMC algorithms approximating  $\mathbb{E}[g(\bar{X}_L(T))]$  in the path independent time step setting.

## 2.2 Fully stochastic time stepping

In this subsection we describe the MLMC algorithm for approximating  $\mathbb{E}[g(X(T))]$  in the setting with adaptive stochastic time steps.

Quite similar to the setting of path independent time steps, the error control of the MLMC estimate  $|\mathcal{A}_{\mathcal{MLC}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(X(T))]|$  is in this setting

based on constructing numerical realizations  $\bar{X}_\ell(t)$  on *stochastic* adaptively refined meshes  $\Delta t_\ell$  so that the time discretization errors

$$\left| \mathbb{E}[g(\bar{X}_\ell(T)) - g(X(T))] \right| \lesssim \text{TOL}_{\text{T},\ell}, \quad \text{for } \ell = 0, 1, \dots, L, \quad (2.17)$$

are fulfilled, and by determining the number of samples  $M_0$  to ensure that the statistical error

$$|\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0) - \mathbb{E}[g(\bar{X}_L(T))]| \leq \text{TOL}_{\text{S}}, \quad (2.18)$$

is fulfilled, with a given confidence.

The statistical error (2.18) is controlled very similarly as in the setting of path independent time steps:

1. set the initial of samples used in the MLMC estimator (2.13) as  $M_0 = 2^{L + \lceil C_{\mathcal{M}\mathcal{L}} L \rceil + 1}$  with  $C_{\mathcal{M}\mathcal{L}} \in (0, 1)$ , cf. Remark 3;
2. configure the number of samples  $M_\ell$  on higher levels in terms of  $M_0$  by the ratio (2.14);
3. generate realizations  $\{\bar{X}_\ell(T)\}$  for the MLMC estimator  $\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0)$  and compute the sample variance  $\sigma^2$  as defined in (2.15);
4. If the stopping condition (2.16) is fulfilled, generate a last output MLMC estimate  $\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0)$  and break. Otherwise, set  $M_0 = 2M_0$ , update the stochastic algorithm parameters estimating the average number of time steps on each grid level,<sup>5</sup> and return to step 2.

For the  $\ell$ -th sample average summand of  $\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0)$ , i.e.,  $\mathcal{A}(g(\bar{X}_0(T)); M_0)$  if  $\ell = 0$  and  $\mathcal{A}(g(\bar{X}_\ell(T)) - g(\bar{X}_{\ell-1}(T)); M_\ell)$  if  $\ell > 0$ , the algorithm generates  $M_\ell$  Euler–Maruyama realization pairs<sup>6</sup>,  $(\bar{X}_{\ell-1}(T), \bar{X}_\ell(T))$  according to (1.5) with the time discretization errors respectively bounded by  $\text{TOL}_{\text{T},\ell-1}$  and  $\text{TOL}_{\text{T},\ell}$  in the sense (2.17). The realization pairs are constructed by stochastic adaptive refinements of a given initial mesh  $\Delta t_{-1}$ . The realizations in a realization pair  $(\bar{X}_{\ell-1}(T), \bar{X}_\ell(T))$  are respectively generated on the adaptively refined meshes  $\Delta t_{\ell-1}$  and  $\Delta t_\ell$ . These meshes are determined by iteratively refining an initial mesh  $\Delta t_{-1}$ . First,  $\Delta t_{-1}$  is adaptively refined to a mesh  $\Delta t_0$  on which  $|\mathbb{E}[g(\bar{X}_0(T)) - g(X(T))]| \lesssim \text{TOL}_{\text{T},0}$  is fulfilled. Thereafter,  $\Delta t_0$  is adaptively refined to a mesh  $\Delta t_1$  on which  $|\mathbb{E}[g(\bar{X}_1(T)) - g(X(T))]| \lesssim \text{TOL}_{\text{T},1}$  is fulfilled. This iterative refinement procedure continues until the mesh  $\Delta t_{\ell-2}$  is adaptively refined to generate the first output mesh  $\Delta t_{\ell-1}$  and, lastly,  $\Delta t_{\ell-1}$  is adaptively refined to generate the second output mesh  $\Delta t_\ell$ .

The iterative adaptive mesh refinement procedure in Algorithm 5 ensures that a mesh  $\Delta t_\ell$  for the fine realization in a pair  $(\bar{X}_{\ell-1}(T), \bar{X}_\ell(T))$  is determined in the same way as a mesh  $\Delta t_\ell$  for the coarse realization in pair

<sup>5</sup>See Algorithm 6 for details on the parameter update.

<sup>6</sup>Observe that for the level  $\ell = 0$  only the realizations of  $\bar{X}_0(T)$  are generated.

$(\bar{X}_\ell(T), \bar{X}_{\ell+1}(T))$ , and consequently that  $\mathbb{E}[\bar{X}_\ell(T)]$  when computed from the finer realization in a pair  $(\bar{X}_{\ell-1}(T), \bar{X}_\ell(T))$  is equal to  $\mathbb{E}[\bar{X}_\ell(T)]$  when computed from the coarse realization in a pair  $(\bar{X}_\ell(T), \bar{X}_{\ell+1}(T))$ . This property is useful since it implies that the following consistency condition for our MLMC estimator is fulfilled

$$\mathbb{E}[\mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)] = \mathbb{E}[g(\bar{X}_L(T))].$$

Let us next take a closer look at the mesh refinement. Due to the stochastic nature of SDEs, each realization pair  $(\bar{X}_{\ell-1}(T), \bar{X}_\ell(T))$  may refine the initial mesh  $\Delta t_{-1}$  differently. In particular, meshes corresponding to different realizations on a given level  $\ell$  may differ. To describe the mesh refinement, taking this feature into account, we introduce some notation. Let  $N_\ell$  and  $\bar{\mathcal{N}}_\ell$  denote the number of time steps and the approximate average number of time steps for realizations at level  $\ell$ , respectively; see Algorithm 6 for details on the approximation technique and its update through the iteration. Further, denote the grid corresponding to one realization at level  $\ell$  by

$$\Delta t_\ell = [\Delta t_\ell(0), \dots, \Delta t_\ell(N_\ell - 1)], \quad (2.19)$$

and its corresponding Wiener increments by

$$\Delta W_\ell = [\Delta W_\ell(0), \dots, \Delta W_\ell(N_\ell - 1)].$$

The refinement condition is based on the error indicator  $r_\ell$ , defined in (1.26), and uses similar refinements to those in the single level method. The refinements of  $\Delta t_\ell$  are stopped when

$$\max_{1 \leq n \leq N_\ell} r_\ell(n) < C_S \frac{\text{TOL}_{\text{T},\ell}}{\bar{\mathcal{N}}_\ell}. \quad (2.20)$$

As long as inequality (2.20) is violated, the  $n^{\text{th}}$  time step of  $\Delta t_\ell$  is refined if

$$r_\ell(n) \geq C_R \frac{\text{TOL}_{\text{T},\ell}}{\bar{\mathcal{N}}_\ell}. \quad (2.21)$$

Normally, the value for  $C_R$  would be around 2, and  $C_S > C_R$  following the theory developed in [20, 19].

A detailed description of the adaptive MLMC algorithm is given in Algorithm 5 with subroutines Algorithm 6–7 in Section 2.3.

The inputs in Algorithm 5 are:  $\text{TOL}_S$ ,  $\text{TOL}_T$ , initial number of sample realizations  $M_0$ ,  $L$ ,  $\Delta t_{-1}$ , initial guesses for the mean number of time steps  $\{\bar{\mathcal{N}}_{\ell,\text{init}}\}_{\ell=0}^L$  of the accepted meshes from the adaptive refinements, and the three parameters  $C_R$ ,  $C_C$ , and  $C_S$  used in the refinement condition (2.21) and stopping conditions (2.16) and (2.20), respectively. In this algorithm the initial estimate of the mean number of time steps are chosen as  $\bar{\mathcal{N}}_{\ell,\text{init}} = c \text{TOL}_{\text{T},\ell}^{-1}$ , for  $\ell = 0, \dots, L$  and a constant  $c$  such that  $\bar{\mathcal{N}}_{0,\text{init}}$  is a small integer; in the numerical examples in Section 3 the constant was chosen so that  $\bar{\mathcal{N}}_{0,\text{init}} \approx 10$ .

## 2.3 Algorithm Listings

---

### Algorithm 1: Adaptive Generation of a Mesh Hierarchy

---

**Input** :  $\text{TOL}_T, M_{-1}, \Delta t_{-1}, L, \bar{N}_0, C_R, C_S, R$   
**Output**:  $\{\Delta t_\ell\}_{\ell=0}^L, M_L$

**for**  $\ell = 0, 1, \dots, L$  **do**  
  Set  $\text{keep\_sampling} = \text{TRUE}$ ,  $\text{keep\_refining} = \text{TRUE}$ ,  
   $\Delta t_\ell = \Delta t_{\ell-1}$ ,  $M_\ell = M_{\ell-1}$ , and  $\text{TOL}_{T,\ell} = 2^{L-\ell}\text{TOL}_T$ .  
  **while**  $\text{keep\_sampling}$  **or**  $\text{keep\_refining}$  **do**  
    Set  $\text{keep\_sampling} = \text{FALSE}$ ,  $\text{keep\_refining} = \text{FALSE}$   
    Compute  $r_\ell, \mathbf{E}_{\Delta t_\ell}$ , and  $\mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)$  by calling  
    Algorithm 2: Euler( $M_\ell, \Delta t_\ell$ )  
    **if**  $\mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)$  and  $\mathbf{E}_{\Delta t}$  violates (2.10) **then**  
      Set  $\text{keep\_sampling} = \text{TRUE}$   
      Update the number of samples by  
       $M_\ell = 2M_\ell$   
    **else**  
      **if**  $r_\ell$  violates (2.6) **then**  
        Set  $\text{keep\_refining} = \text{TRUE}$   
        Refine  $\Delta t_\ell$  by  
        **forall** intervals  $n = 1, 2, \dots, N_\ell$  **do**  
          **if**  $r_\ell(n)$  satisfies (2.7) **then**  
            divide the interval  $n$  into two equal parts  
          **end**  
        **end**  
        Set  $\bar{N}_\ell = \max\{\bar{N}_\ell, \text{new } N_\ell\}$ .  
      **end**  
    **end**  
  Set  $\bar{N}_{\ell+1} = 2\bar{N}_\ell$   
**end**

---

### Algorithm 2: Euler

---

**Input** :  $M_\ell, \Delta t_\ell$   
**Output**:  $r_\ell, \mathbf{E}_{\Delta t_\ell}, \mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)$

Compute  $M_\ell$  new realizations of  $\bar{X}_\ell$  on  $\Delta t_\ell$  by Euler–Maruyama method (1.5) and use them to compute the error indicators  $r_\ell(n)$  on  $\Delta t_\ell$  by equation (1.26),  $\mathbf{E}_{\Delta t_\ell}$  by equation (2.8), and  $\mathcal{V}\left(\sum_{n=1}^{N_\ell} r_\ell(n); M_\ell\right)$  by equation (2.5).

---

---

**Algorithm 3:** Multilevel Monte Carlo on a Mesh Hierarchy

---

**Input** :  $\text{TOL}_S, M_0, L, \{\Delta t_\ell\}_{\ell=0}^L, C_C$

**Output:**  $\mu = \mathcal{A}_{\mathcal{MC}}(g(\bar{X}(T)); M_0)$

Set  $k = 0$ .

**while**  $k < 1$  **or** (2.16) **is violated do**

Set  $\mu = 0$  and  $\sigma^2 = 0$ .

**for**  $\ell = 0, 1, \dots, L$  **do**

Set  $M_\ell$  as in (2.14)

**if**  $\ell = 0$  **then**

Call Algorithm 4: Euler( $M_0, \{\Delta t_0\}$ ).

Set  $\mu = \mu + \mathcal{A}(g(\bar{X}_0(T)); M_\ell)$

and  $\sigma^2 = \sigma^2 + \frac{\mathcal{V}(g(\bar{X}_0(T)); M_0)}{M_0}$ .

**else**

Call Algorithm 4: Euler( $M_\ell, \{\Delta t_\ell, \Delta t_{\ell-1}\}$ ).

Set  $\mu = \mu + \mathcal{A}(g(\bar{X}_\ell(T)) - g(\bar{X}_{\ell-1}(T)); M_\ell)$

and  $\sigma^2 = \sigma^2 + \frac{\mathcal{V}(g(\bar{X}_\ell(T)) - g(\bar{X}_{\ell-1}(T)); M_\ell)}{M_\ell}$ .

**end**

**end**

**if**  $\sigma$  violates (2.16) **then**

Update the number of samples by

$M_0 = 2M_0$

**end**

Increase  $k$  by 1

**end**

Generate and return the output  $\mu = \mathcal{A}_{\mathcal{MC}}(g(\bar{X}(T)); M_0)$  according to (2.13).

---

---

**Algorithm 4:** Euler

---

**Input** :  $M, \{\Delta t_\ell\}_{\ell=l_0, l_1}$

**Output:**  $\mathcal{V}(g(\bar{X}_0(T)); M), \mathcal{A}(g(\bar{X}_0(T)); M)$  if  $l_0 = l_1 = 0$  **or**

$\mathcal{V}(g(\bar{X}_{l_1}(T)) - g(\bar{X}_{l_0}(T)); M), \mathcal{A}(g(\bar{X}_{l_1}(T)) - g(\bar{X}_{l_0}(T)); M)$  if  $l_0 \neq l_1$

Simulate  $M$  new outcomes of the Wiener process  $W(t)$  on  $\Delta t_{l_1} \supseteq \Delta t_{l_0}$ .

**if**  $l_0 = l_1 = 0$  **then**

Compute the corresponding realizations of  $\bar{X}_0$  on  $\Delta t_0$  and use them to compute  $\mathcal{A}(g(\bar{X}_0(T)); M)$  and  $\mathcal{V}(g(\bar{X}_0(T)); M)$  by (2.4) and (2.5).

**else**

Compute the corresponding realizations of  $\bar{X}_{l_1}$  and  $\bar{X}_{l_0}$  on  $\Delta t_{l_1}$  and  $\Delta t_{l_0}$  and use them to compute  $\mathcal{A}(g(\bar{X}_{l_1}(T)) - g(\bar{X}_{l_0}(T)); M)$  and  $\mathcal{V}(g(\bar{X}_{l_1}(T)) - g(\bar{X}_{l_0}(T)); M)$  by (2.4) and (2.5).

**end**

---

---

**Algorithm 5:** Multilevel Monte Carlo with stochastic time stepping

---

**Input** :  $\text{TOL}_S, \text{TOL}_T, M_0, \Delta t_{-1}, L, \{\overline{\mathcal{N}}_\ell\}_{\ell=0}^L, C_R, C_S, C_C$

**Output:**  $\mu = \mathcal{A}_{\mathcal{MC}}(g(\overline{X}(T)); M_0)$

Set  $k = 0$ .

**while**  $k < 1$  **or** (2.16) is violated **do**

    Compute  $M_0$  new realizations of  $g(\overline{X}_0(T))$

    and their corresponding number of time steps,  $\{N_{0,m}\}_{m=1}^{M_0}$ ,

    by generating Wiener increments  $\{\Delta W_{-1,m}\}_{m=1}^{M_0}$  on the mesh  $\Delta t_{-1}$   
    (independently for each realization  $m$ ) and calling Algorithm 7:

**ATSSE**( $\Delta t_{-1}, \Delta W_{-1,m}, \text{TOL}_T 2^L, \overline{\mathcal{N}}_0$ ).

    Set  $\mu = \mathcal{A}(g(\overline{X}_0(T)); M_0)$  and  $\sigma^2 = \frac{\mathcal{V}(g(\overline{X}_0(T)); M_0)}{M_0}$ .

    Compute the average number of time steps  $\mathcal{A}(N_0; M_0)$ .

**for**  $\ell = 1, \dots, L$  **do**

        Set  $M_\ell$  as in (2.14)

        Compute  $M_\ell$  new realizations of  $g(\overline{X}_{\ell-1}(T))$ ,

        their corresponding number of time steps,  $\{N_{\ell-1,m}\}_{m=1}^{M_\ell}$ , and

        Wiener increments,  $\{\Delta W_{\ell-1,m}\}_{m=1}^{M_\ell}$ , by generating Wiener steps  
         $\{\Delta W_{-1,m}\}_{m=1}^{M_0}$  on the mesh  $\Delta t_{-1}$  (independently for each

        realization  $m$ ) and using the loop

**for**  $\hat{\ell} = 0, \dots, \ell - 1$  **do**

            compute  $\Delta t_{\hat{\ell},m}$  and  $\Delta W_{\hat{\ell},m}$  by calling Algorithm 7:

**ATSSE**( $\Delta t_{\hat{\ell}-1,m}, \Delta W_{\hat{\ell}-1,m}, \text{TOL}_T 2^{L-\hat{\ell}}, \overline{\mathcal{N}}_{\hat{\ell}}$ ).

**end**

        Compute the corresponding  $M_\ell$  realizations of  $g(\overline{X}_\ell(T))$  and

        their number of time steps,  $\{N_{\ell,m}\}_{m=1}^{M_\ell}$ , by calling Algorithm 7:

**ATSSE**( $\Delta t_{\ell-1,m}, \Delta W_{\ell-1,m}, \text{TOL}_T 2^{L-\ell}, \overline{\mathcal{N}}_\ell$ ).

        Set  $\mu = \mu + \mathcal{A}(g(\overline{X}_\ell(T)) - g(\overline{X}_{\ell-1}(T)); M_\ell)$  and

$\sigma^2 = \sigma^2 + \frac{\mathcal{V}(g(\overline{X}_\ell(T)) - g(\overline{X}_{\ell-1}(T)); M_\ell)}{M_\ell}$ .

        Compute the average number of time steps  $\mathcal{A}(N_{\ell-1}; M_\ell)$  and  
         $\mathcal{A}(N_\ell; M_\ell)$ .

**end**

**if**  $\sigma$  violates (2.16) **then**

        Update the number of samples by

$M_0 = 2M_0$ .

        Update the values of  $\{\overline{\mathcal{N}}_\ell\}_{\ell=0}^L$  by calling Algorithm 6:

**UMNT** ( $\{M_\ell\}_{\ell=0}^L, \{\mathcal{A}(N_\ell; M_\ell)\}_{\ell=0}^L, \{\mathcal{A}(N_{\ell-1}; M_\ell)\}_{\ell=1}^L$ ).

**end**

    Increase  $k$  by 1.

**end**

Generate and return the output  $\mu = \mathcal{A}_{\mathcal{MC}}(g(\overline{X}(T)); M_0)$  according  
to (2.13).

---

---

**Algorithm 6:** Update for the mean number of time steps, (**UMNT**)

---

**Input** :  $\{M_\ell\}_{\ell=0}^L, \{\mathcal{A}(N_\ell; M_\ell)\}_{\ell=0}^L, \{\mathcal{A}(N_{\ell-1}; M_\ell)\}_{\ell=1}^L$

**Output:**  $\{\bar{N}_\ell\}_{\ell=0}^L$

**for**  $\ell = 0, 1, \dots, L$  **do**

**if**  $\ell < L$  **then**

    Set  $\bar{N}_\ell = \frac{M_\ell \mathcal{A}(N_\ell; M_\ell) + M_{\ell+1} \mathcal{A}(N_\ell; M_{\ell+1})}{M_\ell + M_{\ell+1}}$ .

**else**

    Set  $\bar{N}_L = \mathcal{A}(N_L; M_L)$ .

**end**

**end**

---

---

**Algorithm 7:** Adaptive Time Step Stochastic Euler (**ATSSE**)

---

**Input** :  $\Delta t_{in}, \Delta W_{in}, \text{TOL}, \bar{N}_{in}$

**Output:**  $\Delta t_{out}, \Delta W_{out}, N_{out}, g_{out}$

Set  $k = 0, \Delta t_{[0]} = \Delta t_{in}, \Delta W_{[0]} = \Delta W_{in}, N_{[0]} =$  number of steps in  $\Delta t_{in}$

**while**  $k < 1$  **or**  $(r_{[k-1]}; \text{TOL}, \bar{N}_{in})$  violates (2.20) **do**

  Compute the Euler approximation  $\bar{X}_{[k]}$  and the error indicators  $r_{[k]}$  on  $\Delta t_{[k]}$  with the known Wiener increments  $\Delta W_{[k]}$ .

**if**  $(r_{[k]}; \text{TOL}, \bar{N}_{in})$  violates (2.20) **then**

    Refine the grid  $\Delta t_{[k]}$  by

**forall** intervals  $n = 1, 2, \dots, N_{[k]}$  **do**

**if**  $r_{[k]}(n)$  satisfies (2.21) **then**

        divide the interval  $n$  into two equal parts

**end**

**end**

    and store the refined grid in  $\Delta t_{[k+1]}$ .

    Compute  $\Delta W_{[k+1]}$  from  $\Delta W_{[k]}$  using Brownian bridges on  $\Delta t_{[k+1]}$ .

    Set  $N_{[k+1]} =$  number of steps in  $\Delta t_{[k+1]}$ .

**end**

  Increase  $k$  by 1.

**end**

Set  $\Delta t_{out} = \Delta t_{[k-1]}, \Delta W_{out} = \Delta W_{[k-1]}, N_{out} = N_{[k-1]}, g_{out} = g(\bar{X}_{[k-1]})$ .

---

### 3 Numerical Experiments

This section presents numerical results from implementations of the algorithms of Section 2. We have selected problems to indicate the use of the adaptive methods. Specifically, uniform time steps are suitable for problem 3.1, adaptively refined deterministic time steps are suitable for problem 3.2, and fully stochastic time steps are suitable problem 3.3. In both problems 3.2 and 3.3 the use of the multilevel adaptive algorithms is much more efficient than the use of the corresponding single level versions of the algorithms, which is in turn much more efficient than using a single level uniform time stepping method.

For those problems the complexity is close to that of uniform MLMC, since the observed order of strong convergence remains close to  $1/2$  even though the order of weak convergence is reduced using uniform time steps. The current adaptive algorithm is not optimized with respect to the strong error, but is subject to ongoing research as an extension of the present adaptive algorithm.

The computations were performed in `Matlab 7` using the built in pseudo random number generator `randn` for simulating sampling from the normal distribution. For the parameter in the refinement criteria (2.7) and (2.21) we used  $C_R = 2$  in all the cases, and for the parameter in the stopping criteria (2.6) and (2.20) we used  $C_S = 5$  in problems 3.2 and 3.3, and  $C_S = 3$  in problem 3.1 where we expect uniform refinements and all error indicators of the same size. In all examples the error tolerance was split equally,  $TOL_S = TOL_T = TOL/2$ ; the proof of Theorem 3 indicates that this choice is not optimal.

### 3.1 A Linear SDE

*Consider first the standard geometric Brownian Motion,*

$$\begin{aligned} dX(t) &= rX(t)dt + \sigma X(t) dW(t), & t \in (0, T), \\ X(0) &= 1, \end{aligned}$$

*using  $r = 1$  and  $\sigma = 0.5$  with a final time  $T = 1$  and  $g(x) = x$ .*

In this simple example adaptive time stepping is not expected to improve the time discretization error. In fact, the path independent adaptive algorithm produces a hierarchy of uniform grids, and when the fully stochastic adaptive algorithm is applied to this problem all generated meshes are uniform but different realizations of the driving Wiener process may result in different step sizes. The computational cost, measured as the total number of time steps, in all stages in the adaptive refinements, for all realizations of the Euler approximation  $\bar{X}$ , is shown in Figure 1. For both versions of the algorithm, the computational cost is consistent with, but slightly better than, the main complexity result in Theorem 3 of Section 4, from which we get an upper bound  $cost \leq C(TOL^{-1} \log(TOL_{T,0}/TOL))^2$ . The work measured this way is very similar in the two versions of the algorithm. However, the version in Section 2.1 is more efficient in this case since it only computes dual solutions in the construction of the mesh hierarchy which is of negligible cost<sup>7</sup>, while the version in Section 2.2 computes both primal and dual for every realization. The accuracy of both versions of the algorithm is shown in Figure 2.

The work we measure in Figure 1 is greater than the work analyzed in Section 4, which only counts the number of Euler steps used on the accepted meshes. The comparison made in Table 1 shows the same growth rate as  $TOL \downarrow 0$  when the fully stochastic adaptive algorithm is applied to problem 3.1.

---

<sup>7</sup>See Figure 3 for problem 3.2.



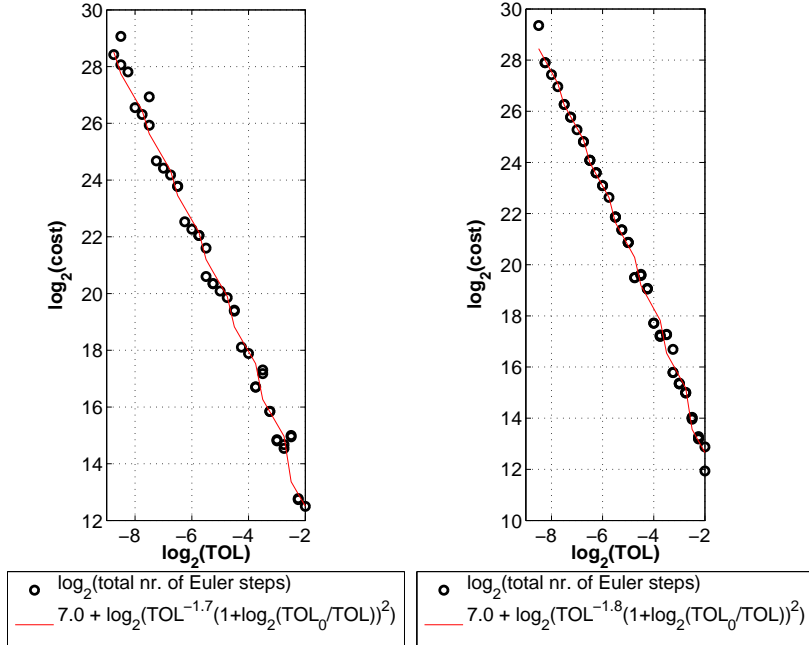


Figure 1: Experimental complexity for both versions of the algorithm applied to the geometrical Brownian motion example of Section 3.1; to the left the version of mesh creation followed by sampling on fixed meshes, in Section 2.1, and to the right the path dependent sampling version in Section 2.2. The computational cost is measured as the total number of Euler time steps taken in all refinement iterations on all levels for all realizations. The graphs show three independent realizations of the underlying Wiener processes for each prescribed tolerance. A least squares fit, in  $\log_2$ - $\log_2$ -scale, of the model  $cost = c_1 \text{TOL}^{-c_2} (1 + \log_2(\text{TOL}_{T,0}/\text{TOL}))^2$  gives  $c_2 = 1.7$  and  $c_2 = 1.8$  in the two cases respectively; this is slightly better than the prediction of Theorem 3 of Section 4

The pronounced clustering of data points in the left graph is primarily due to the fact that the tolerance changes by a constant factor  $2^{1/4}$  and for this example the adaptive algorithm generates uniform meshes by repeated halving of the same initial uniform mesh; the effect is that several consecutive tolerances result in identical mesh hierarchies, while only the number of samples changes. With the pathwise adaptive algorithm, to the right, the generated meshes are again uniform for geometric Brownian motion, but the resolution may be different for each individual outcome of the Wiener process; again the number of levels will be constant for four consecutive tolerances since it is decided by the requirement that  $\text{TOL}_{T,0} = 2^L \text{TOL}_{T,L}$  is smaller than a fixed constant.

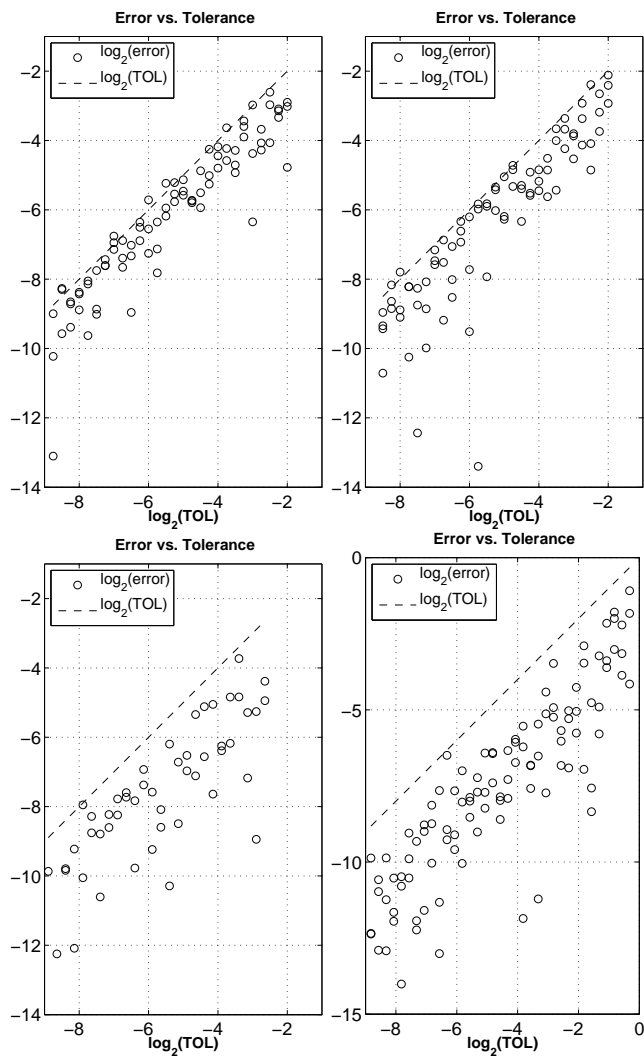


Figure 2: These accuracy tests show the error versus the prescribed tolerance when the adaptive MLMC algorithm is applied to the test examples of Section 3; to the left the version of Section 2.1 applied to the geometric Brownian motion in Section 3.1 (top) and the singularity problem in Section 3.2 (bottom), and to the right the version of Section 2.2 applied to the geometric Brownian motion in Section 3.1 (top) and the stopped diffusion problem in Section 3.3 (bottom).

Problem	Version	sampled randn's		accepted Euler steps		all Euler steps	
		$c_1$	$c_2$	$c_1$	$c_2$	$c_1$	$c_2$
GBM,	Sec. 2.1	6.9	1.7	–	–	7.0	1.7
GBM,	Sec. 2.2	6.1	1.7	6.5	1.8	7.0	1.8
Sing.,	Sec. 2.1	13.2	2.0	–	–	13.4	2.0
Barrier,	Sec. 2.2	8.2	1.9	8.8	1.9	9.5	2.1

Table 1: Complexity estimates for the three different problems: the geometric Brownian motion of Section 3.1, the deterministic singularity problem of Section 3.2, and the stopped diffusion problem of Section 3.3. The tabulated values are least square fits of the parameters  $c_1$  and  $c_2$  in the model  $\log_2(\text{cost}) = c_1 + \log_2(\text{TOL}^{-c_2}(1 + \log_2(\text{TOL}_{T,0}/\text{TOL}))^2)$  when the cost is measured in three different ways: by counting the total number of sampled random variables, by the number of accepted Euler steps  $\mathcal{A}[\text{cost}] = M_0 \mathcal{A}(N_0; M_0) + \sum_{\ell=1}^L M_\ell \{\mathcal{A}(N_\ell; M_\ell) + \mathcal{A}(N_{\ell-1}; M_\ell)\}$  which is approximated by the work estimate defined in (4.2), and by counting the total number of Euler steps performed when solving the primal problem in *all* refinement stages for all levels in the multilevel algorithms. The cost when counting Euler steps in the accepted meshes only was not recorded in the experiments using the path independent version of the algorithm, but it is by necessity bounded from below and above by the other two measures of the work.

### 3.2 Drift singularity, linear SDE

Consider for a real constant  $\alpha \in (0, T)$  the linear stochastic differential equation

$$dX(t) = \begin{cases} X(t) dW(t), & t \in [0, \alpha], \\ \frac{X(t)}{2\sqrt{t-\alpha}} dt + X(t) dW(t), & t \in (\alpha, T], \end{cases} \quad (3.1)$$

$$X(0) = 1,$$

with the unique solution

$$X(t) = \begin{cases} \exp(W(t) - t/2), & t \in [0, \alpha], \\ \exp(W(t) - t/2) \exp(\sqrt{t-\alpha}), & t \in (\alpha, T]. \end{cases}$$

The goal is to approximate the expected value  $\mathbb{E}[X(T)] = \exp(\sqrt{T-\alpha})$ . Here we choose  $T = 1$  and  $\alpha = T/3$ . To avoid evaluating arbitrarily large values of the drift in (3.1) we modify the drift to be

$$a(t, x) = \begin{cases} 0, & t \in [0, \alpha], \\ \frac{x}{2\sqrt{t-\alpha+\text{TOL}^4}}, & t \in (\alpha, T], \end{cases} \quad (3.2)$$

yielding a higher order perturbation  $\mathcal{O}(\text{TOL}^2)$  in the computed result and in the size of the optimal time steps. This regularization was applied to maintain consistency with the numerical tests in [19], but is not strictly necessary with

the upper bound,  $\rho \leq \rho_{up}(\text{TOL})$ , on the error density in (1.24). Due to the time discontinuity of the drift function and to ensure optimal convergence of the adaptive algorithms, we modify the Euler method to

$$\bar{X}_{n+1} - \bar{X}_n = a(\hat{t}, \bar{X}_n) \Delta t_n + \bar{X}_n \Delta W_n, \quad n = 0, 1, 2, \dots, \quad (3.3)$$

where we choose the stochastic evaluation time  $\hat{t} \in \{t_n, t_{n+1}\}$  so that

$$|a(\hat{t}, \bar{X}_n)| = \max(|a(t_n, \bar{X}_n)|, |a(t_{n+1}, \bar{X}_n)|).$$

Observe that the use of  $\hat{t}$  does not change the adapted nature of the Euler method.

Since the added difficulty compared to example 3.1 is a singularity in the drift at a deterministic time, the path independent adaptive algorithm described in Section 2.1 is the most suitable, and it is used in this example. The goal here is to verify that the adaptive multilevel algorithms of Section 2 give the same improvement from the single level adaptive algorithm as multilevel Monte Carlo does in the uniform case for regular problems.

The accuracy test in Figure 2 shows good agreement between observed error and prescribed tolerance. As shown in the complexity study in Table 1 and Figure 3 the computational costs grow like  $\text{TOL}^{-1.8}(1 + \log(\text{TOL}_{T,0}/\text{TOL}))^2$  which is slightly better than the predicted complexity  $\text{TOL}^{-2}(\log(\text{TOL}_{T,0}/\text{TOL}))^2$ . The cost of the mesh construction phase of the algorithm is seen to be negligible compared to the total work.

In this example the weak rate of convergence for the Euler–Maruyama method with uniform time steps is only 1/2, so the total cost for a single level uniform time stepping algorithm is proportional to  $\text{TOL}^{-4}$ . The left part of Figure 4 shows that the single level version of the adaptive algorithm improves that complexity to approximately  $\text{TOL}^{-3}$ , while the multilevel version improves the complexity by nearly one order more. With the regularization (3.2) the observed order of strong convergence of the Euler–Maruyama method with uniform time steps is still 1/2, so the complexity estimate in Theorem 1 of [7] for uniform multilevel simulations applies, and we should get the ideal complexity  $(\epsilon^{-1} \log(\epsilon^{-1}))^2$  for a mean square error of size  $\epsilon^2$ . The right part of Figure 4 shows that this is approximately true for the cost as a function of the maximal observed error over 11 independent realizations.

### 3.3 Stopped diffusion

Here we compute the solution of a more challenging problem that motivates the use of stochastic time steps that are adaptively refined for each sample path.

The additional difficulty of the problem is that we now wish to compute approximations of an expected value

$$\mathbb{E}[g(X(\tau), \tau)], \quad (3.4)$$

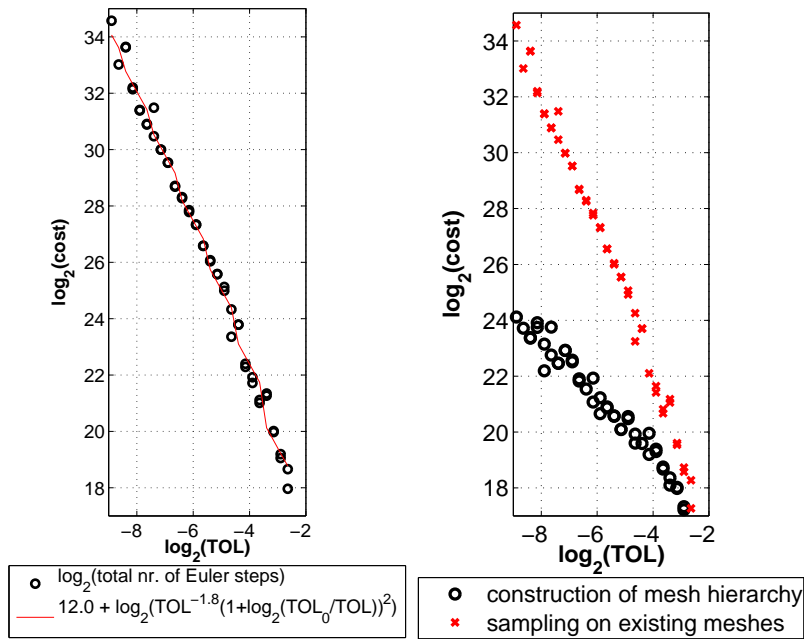


Figure 3: Experimental complexity when the algorithm in Section 2.1 is applied to the drift singularity problem in Section 3.2. To the left is shown the cost of both phases of the algorithm, and to the right the contribution from the generation of the mesh hierarchy and the subsequent sampling to reduce the statistical error; it is clear that the cost of the first phase is negligible compared to the second for small tolerances. The computational cost is measured as the total number of Euler time steps taken in all refinement iterations on all levels for all realizations. The graphs show three independent realizations of the underlying Wiener processes for each prescribed tolerance. A least squares fit, in  $\log_2 - \log_2$ -scale, of the model  $\text{cost} = c_1 \text{TOL}^{-c_2} (1 + \log_2(\text{TOL}_{T,0}/\text{TOL}))^2$  gives  $c_2 = 1.8$ .

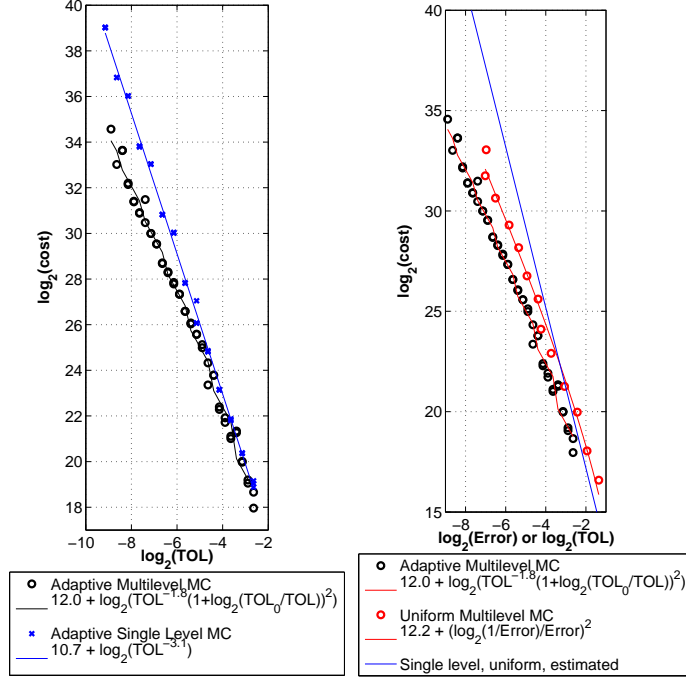


Figure 4: The computational cost of the path independent adaptive algorithm of Section 2.1, applied to the deterministic singularity problem 3.2, is compared to several alternatives. Left: the multilevel version improves the computational complexity of the single level version of the same adaptive algorithm from approximately proportional to  $\text{TOL}^{-3}$  to approximately proportional to  $\text{TOL}^{-2} \log(\text{TOL}^{-1})^2$ . Right: the cost of a standard, uniform time step, Monte Carlo method would be proportional to  $\text{TOL}^{-4}$ ; here the work was estimated from a Central Limit Theorem type confidence interval based on the time discretization errors and sample variances. The cost of the uniform MLMC method is shown as a function of the maximal error over 11 realizations. The observed cost oscillates around  $(\epsilon^{-1} \log(\epsilon^{-1}))^2$ , which is expected since the observed observed strong order of convergence is  $1/2$ .

For the adaptive algorithm the cost is estimated by the total number of Euler steps taken on all levels in all stages of the adaptive refinement process.

where  $X(t)$  solves the SDE (1.1) as before, but where the function  $g : D \times [0, T] \rightarrow \mathbb{R}$  is evaluated at the first exit time

$$\tau := \inf\{t > 0 : (X(t), t) \notin D \times (0, T)\}$$

from a given open domain  $D \times (0, T) \subset \mathbb{R}^d \times (0, T)$ . This kind of stopped (or killed) diffusion problems arise for example in mathematical finance when pricing barrier options and for boundary value problems in physics.

The main difficulty in the approximation of the stopped diffusion on the boundary  $\partial D$  is that a continuous sample path may exit the given domain  $D$  even though a discrete approximate solution does not cross the boundary of  $D$ . Due to this hitting of the boundary the order of weak convergence of the Euler–Maruyama method is reduced from 1 to 1/2, in terms of the step size of uniform meshes; see [10]. In this subsection we combine the adaptive multilevel algorithm of Section 2.2 with an error estimate derived in [5] that also takes into account the hitting error. This error estimate, and the adaptive algorithm, can be used also when  $D$  is multi dimensional even if the boundary  $\partial D$  has corners for example.

The hitting error is accounted for by an extra contribution to the error density in (1.23); this contribution can be expressed in terms of exit probabilities for individual time steps, conditioned on the computed path at the beginning and the end of the time steps, and of the change in the goal function,  $g$ , when evaluated at a possible exit point within the time step instead of the actually computed exit  $(\bar{X}(\bar{\tau}), \bar{\tau})$ . The full expression of the resulting error indicators is given in equation (50) of [5]. Since the differential  $\partial_i g(\bar{X}(T), T)$  in the discrete dual backward problem (1.17) does not exist if  $T$  is replaced by  $\bar{\tau} < T$  this initial value must be alternatively defined; this can be done using difference quotients with restarted computed trajectories as described, both for the discrete dual and for its first and second variations, in equations (20-25) of [5]. Note that for this modified error density the proof in [20] of almost sure convergence to a limit density does not apply.

In addition to the modification of the error density a lower bound is introduced on the step size to avoid excessive refinements near the barrier,

$$\Delta t_n \geq \min \left\{ \text{TOL}_{T,\ell}^{1.5}, \frac{\text{dist}_n \text{dist}_{n+1} / b(\bar{X}(t_n; \omega), t_n)^2}{-3 \log(\text{TOL}_{T,\ell})} \right\}, \quad (3.5)$$

where  $\text{dist}_j$  denotes the distance from  $\bar{X}(t_j; \omega)$  to the barrier.

*For the numerical example we consider the stopped diffusion problem*

$$dX(t) = \frac{11}{36}X(t) dt + \frac{1}{6}X(t) dW(t), \quad \text{for } t \in [0, 2] \text{ and } X(t) \in (-\infty, 2), \quad (3.6)$$

$$X(0) = 1.6.$$

*For  $g(x, t) = x^3 e^{-t}$  with  $x \in \mathbb{R}$ , this problem has the exact solution  $\mathbb{E}[g(X_\tau, \tau)] = u(X(0), 0) = X(0)^3$ , where the solution,  $u$ , of the Kolmogorov backward equation is  $u(x, t) = x^3 e^{-t}$ . We chose an example in one space dimension for simplicity, although it is only in high dimension that Monte Carlo methods are more*

efficient than deterministic finite difference or finite element methods to solve stopped diffusion problems. The comparison here between the standard Monte Carlo and the Multilevel Monte Carlo methods in the simple one dimensional example indicates that the Multilevel Monte Carlo method will also be more efficient in high dimensional stopped diffusion problems, where a Monte Carlo method is a good choice. In the case of a scalar SDE, where  $D$  is an interval on the real line, the strong order of convergence of the Euler-Maruyama scheme for barrier problems can be close to  $1/2$ . In fact, it is shown in [8] that  $\text{Var}(g(\bar{X}_\ell) - g(\bar{X}_{\ell-1})) = \mathcal{O}(\Delta t^{1-\delta})$ , for any  $\delta > 0$ , using the Euler-Maruyama method with uniform step size  $\Delta t$  on a class of options including some barrier options. In this case Theorem 3.1 of [7] tells us that, for any choice of  $\delta > 0$ , uniform MLMC simulations can be performed at a cost  $\mathcal{O}(\epsilon^{-2(1+\delta)})$  for a mean square error of order  $\epsilon^2$ .

In the remainder of this section we present results on the accuracy and cost of the adaptive multilevel algorithm of Section 2.2, applied to (3.6), with the error estimate modified for the barrier problem, and with the lower bound (3.5) on the step size. The algorithm was applied with a sequence of tolerances with three simulations for each tolerance using different initial states in the pseudo random number generator. The observed errors are scattered below the corresponding tolerances in Figure 2, showing that the algorithm achieves the prescribed accuracy.

The experimental complexity is illustrated in Figure 5 and Table 1. A least squares fit of the model

$$\text{cost} = c_1 \left( \frac{1}{\text{TOL}} \right)^{c_2} \left( 1 + \log \left( \frac{\text{TOL}_{T,0}}{\text{TOL}} \right) \right)^2 \quad (3.7)$$

in the  $\log_2$ - $\log_2$ -scale of the graph using equal weights on all data points gives  $c_2 = 1.9$  when the work is estimated by the total number of time steps on all the *accepted meshes*; this is the measure of work that is estimated by (4.2) in Section 4. When all Euler steps, in all refinement stages are included, the least squares fit gives  $c_2 = 2.1$ , which is also close to the rate predicted in Theorem 3. However, the corresponding cost using the single level adaptive algorithm with just one data point per tolerance used grows faster than  $\text{TOL}^{-3}$  in this example; see Figure 6.

In conclusion the observed convergence of the adaptive MLMC method applied to the barrier problem (3.6) is close to the predicted rate in Theorem 3. This shows an improved convergence compared to the single level version of the adaptive Monte Carlo algorithm where the cost grows approximately like  $\text{TOL}^{-3}$ , which in itself is a better order of weak convergence than the one obtained using a single level Monte Carlo method with constant time steps where the cost grows like  $\text{TOL}^{-4}$ .



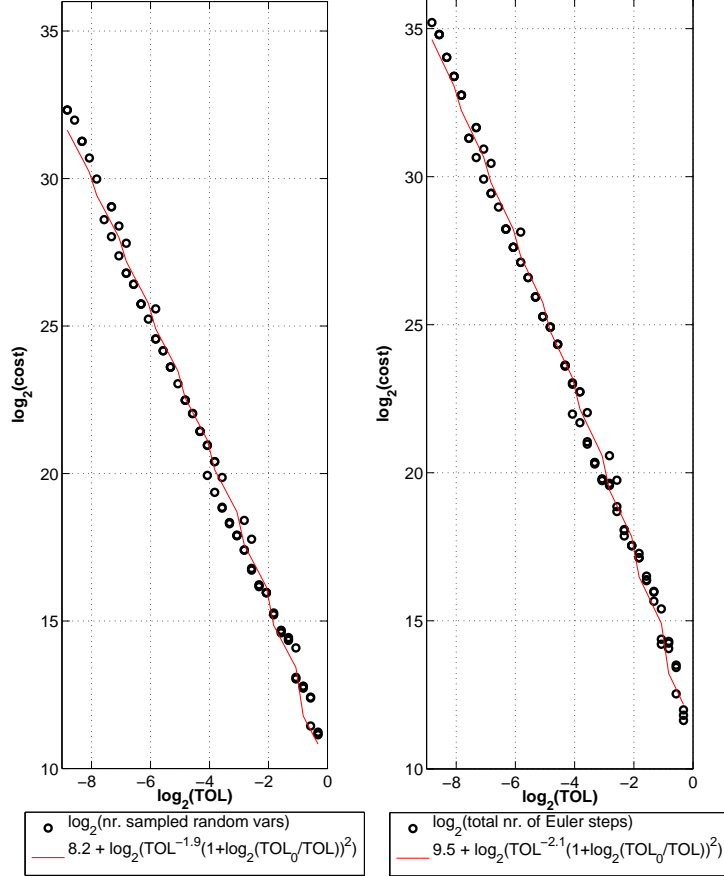


Figure 5: Experimental complexity for the barrier example in Section 3.3. The computational cost of the multilevel adaptive algorithm is shown for varying tolerances using three different initial states in the pseudo random number algorithm. To the left is shown the work estimate based on the number of Euler steps in the accepted meshes,  $\mathcal{A}[cost] = M_0 \mathcal{A}(N_0; M_0) + \sum_{\ell=1}^L M_\ell \{\mathcal{A}(N_\ell; M_\ell) + \mathcal{A}(N_{\ell-1}; M_\ell)\}$ , which is the work measure closest to (4.2) used in Section 4; to the right is shown the estimate based on all Euler steps taken in all stages in the adaptive mesh refinement process. A least squares fit, in  $\log_2$ - $\log_2$ -scale, of the model  $cost = c_1 \text{TOL}^{-c_2} (\log_2(1 + \text{TOL}_{T,0}/\text{TOL}))^2$  with equal weight on all observations results in  $c_2 = 1.9$  and  $c_2 = 2.1$  in the two cases. This is in agreement with the prediction in Theorem 3.

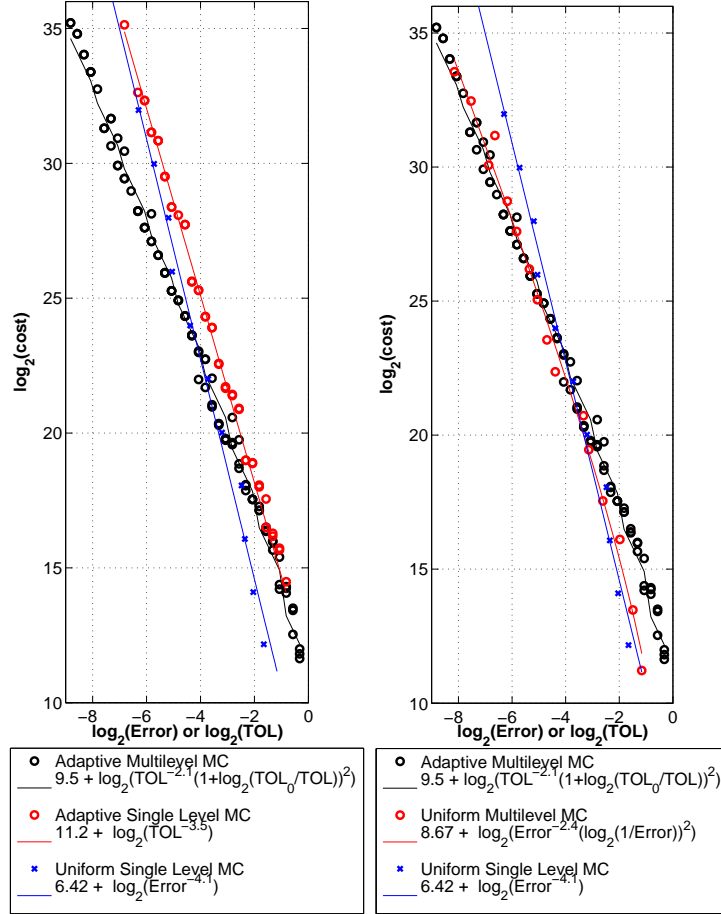


Figure 6: Left: The multilevel version of the path dependent adaptive algorithm of Section 2.2 applied to the barrier problem 3.3 improves the computational complexity of the single level version of the same adaptive algorithm; a single level method based on uniform time steps has even worse complexity with the computational cost growing like  $\text{TOL}^{-4}$ . Right: The cost of the uniform MLMC method is shown as a function of the maximal error over 16 realizations. The observed cost is close to that of adaptive multilevel Monte Carlo, which is expected since the observed strong order of convergence is  $1/2$ , but oscillates around a slightly worse fitted complexity  $\epsilon^{-2.4}(\log(\epsilon^{-1}))^2$ . The cost is estimated by the total number of Euler steps taken on all levels in all stages of the adaptive refinement process.

## 4 Theoretical results

In this section we will study the asymptotic accuracy and complexity of the MLMC algorithm in the setting of stochastic adaptive time steps introduced in Section 2.2. We recall that for a sought accuracy  $\text{TOL} > 0$ , the goal of the MLMC algorithm is to construct a Monte Carlo approximation of  $\mathbb{E}[g(X(T))]$  that with probability close to one fulfills

$$|\mathbb{E}[g(X(T))] - \mathcal{A}_{\text{MC}}(g(\bar{X}(T)); M_0)| \leq \text{TOL}.$$

Our main result on asymptotic accuracy for MLMC algorithm, proved in Subsection 4.2, is

**Theorem 2** (Multilevel accuracy). *Suppose that the modeling assumptions of Lemma 1 hold, that (4.7) holds, and that  $\text{TOL}_{\text{T}} \leq \text{TOL}_{\text{S}}$ . Then the adaptive MLMC algorithm with confidence parameter  $C_C > 0$  and stochastic time steps (2.20) and (2.21) satisfies*

$$\liminf_{\text{TOL} \downarrow 0} P(|\mathbb{E}[g(X(T))] - \mathcal{A}_{\text{MC}}(g(\bar{X}(T)); M_0)| \leq \text{TOL}) \geq \int_{-C_C}^{C_C} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx. \quad (4.1)$$

The motivation for introducing multiple levels in the MC algorithm is to reduce the computational complexity. To analyze the asymptotic complexity of the adaptive MLMC algorithm we define

$$\text{WORK}(\text{TOL}) = \sum_{\ell=0}^L \mathbb{E}[M_\ell] \mathbb{E}[N_\ell], \quad (4.2)$$

recalling that  $M_\ell$  denotes the number of realization samples  $g(\bar{X}_\ell(T; \omega))$  at level  $\ell$  required to control the statistical error, and  $N_\ell$  denotes the number of adaptive time steps required in the construction of a numerical realization  $g(\bar{X}_\ell(T; \omega))$  to control the time discretization error at level  $\ell$ . The function  $\text{WORK}(\text{TOL})$  is thus our estimate of the average number of arithmetic operations required in the generation and sampling of  $\{g(\bar{X}_\ell(T))\}_{\ell=0}^L$  to approximate  $\mathbb{E}[g(X(T))]$  for the prescribed confidence  $C_C$  and accuracy  $\text{TOL}$ . By analyzing the asymptotics of  $\mathbb{E}[M_\ell]$  and  $\mathbb{E}[N_\ell]$  separately, our main complexity theorem is as follows.

**Theorem 3** (Multilevel computational complexity). *Suppose the assumptions of Lemma 1 and (4.7) hold and suppose the lower bound for the error density is on the form  $\rho_{\text{low}}(\text{TOL}_{\text{T}}) = \text{TOL}_{\text{T}}^{\bar{\gamma}}$ , cf. (1.24). Then the work for the MLMC algorithm using stochastic time steps, as defined in (4.2), fulfills the following bounds:*

(I) *If  $\rho_{\text{low}}(\text{TOL}_{\text{T}}) = \rho_{\text{min}} \in \mathbb{R}_+$  (i.e.  $\bar{\gamma} = 0$ ) and*

$$\min_{\tau \in [0, T]} |\hat{\rho}(\tau)| \geq \rho_{\text{min}} \text{ a.s.} \quad (4.3)$$

then

$$\limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}^2}{L^2} \leq \frac{16 C_C^2 C_G}{\text{TOL}_{T,Max} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\widehat{\rho}(\tau)|} d\tau \right] \right)^2. \quad (4.4)$$

(II) Let  $\bar{\gamma}$  be constant valued and satisfying  $0 < \bar{\gamma} < \alpha/(2 + \alpha)$  with  $0 < \alpha < 1/2$ , cf. (1.24). Then

$$\begin{aligned} \limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}^{2+\bar{\gamma}}}{L} &\leq \frac{(2 + \bar{\gamma})^{2+\bar{\gamma}} 2^{2+\bar{\gamma}}}{\bar{\gamma}^{1+\bar{\gamma}} \log(2)} \cdot \frac{C_C^2 C_G C_S^{\bar{\gamma}}}{(\text{TOL}_{T,Max})^{1-\bar{\gamma}} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\widehat{\rho}(\tau)|} d\tau \right] \right)^2. \end{aligned} \quad (4.5)$$

(III) If  $\bar{\gamma} \rightarrow 0$  and  $L\bar{\gamma} \rightarrow \infty$  as  $\text{TOL} \downarrow 0$ , then

$$\limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}^2 \bar{\gamma}}{L 2^{\bar{\gamma}L}} \leq \frac{16 C_C^2 C_G}{\log(2) \text{TOL}_{T,Max} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\widehat{\rho}(\tau)|} d\tau \right] \right)^2. \quad (4.6)$$

Here, the number of levels  $L = \mathcal{O}(\log(\text{TOL}^{-1}))$ ,  $C_C$  is the confidence parameter,  $C_R$  and  $C_S$  are refinement parameters described by (2.20) and (2.21),  $C_G$  is the constant in the second moment bound (4.43), where  $\text{TOL}_{T,Max}$  is the upper bound of the time discretization tolerance at level  $\ell = 0$ , and  $\bar{\gamma}$  is the lower bound error density exponent;  $\rho_{low}(\text{TOL}_T) = \text{TOL}_T^{\bar{\gamma}}$ , cf. (1.24).

**Remark 1** (Complexity example). Case (III) of Theorem 3 implies that if the exponent of the lower error density  $\rho_{low}$  is given by  $\bar{\gamma}(\text{TOL}) = \log_2(\log_2(L))/L$ , then the following complexity bound, notably close to the standard complexity in the setting of uniform time steps, is achieved:

$$\begin{aligned} \limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}^2 \log_2(\log_2(L))}{L^2 \log_2(L)} &\leq \frac{16 C_C^2 C_G}{\log(2) \text{TOL}_{T,Max} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\widehat{\rho}(\tau)|} d\tau \right] \right)^2. \end{aligned}$$

To introduce the reader gently to our proofs of Theorem 2 and 3, we have chosen to first prove analogous results for the SLMC algorithm in Subsection 4.1. With the single level proofs fresh in mind, we move on to the more daunting task of proving Theorem 2 and 3 in Subsection 4.2. We restrict ourselves to proving Theorem 2 and 3 for the stochastic time step setting only. The stochastic time step setting is however the most general setting, so one can easily prove corresponding results for the deterministic time step setting as well.

In addition to Lemma 1, the analysis in this section will be derived relying on one more technical assumption on strong approximations.

For  $p = 2$  and  $4$ , we have that

$$\begin{aligned} \mathbb{E} \left[ |g(X(T)) - g(\bar{X}(T))|^p \right] &= \mathcal{O} \left( \frac{\text{TOL}_T}{\rho_{low}(\text{TOL}_T)} \right)^{p/2} \\ \mathbb{E} \left[ |g(\bar{X}(T))|^p \right] &= \mathcal{O}(1), \end{aligned} \quad (4.7)$$

with  $\rho_{low}(\text{TOL}_T)$  denoting the lower bound for the error density, as defined in (1.24). The work [24] gives conditions under which (4.7) is fulfilled.

#### 4.1 Single level results

The SLMC algorithm we consider in this subsection, first described and analyzed in [24], is reanalyzed here with the goal to construct proofs for the asymptotic accuracy and complexity of the SLMC algorithm that are later extended to the MLMC algorithm. In the first lemma we show that the adaptive refinement Algorithm 7 stops after a finite number of iterations. This property allows us to later bound the amount of computational work in the single level adaptive algorithm. It also has another important implication: the imposed lower bound on the error density,  $\rho_{low}(\text{TOL}_T)$  in (1.10), ensures that the maximum mesh size of the mesh generated by Algorithm 7,  $\Delta t_{sup}(\text{TOL}_T)$  introduced in Lemma 2, tends to zero as  $\text{TOL}_T$  tends to zero. This in turn implies the almost sure convergence of the error density, which is crucial in the proofs of the main results of this section. A similar result for the multilevel case is direct to obtain and will not be stated for the sake of brevity.

**Lemma 3** (Stopping). *Suppose the adaptive Algorithm 7 applies the mesh refinement strategy (2.20) and (2.21) on a set of realizations having the same uniform initial mesh of step size  $\Delta t_0$ . Further, assume that the estimated average number of time steps,  $\bar{N}_{in}$ , satisfies*

$$\bar{N}_{in} < N_{up} := \frac{T^2 \rho_{up}(\text{TOL}_T)}{C_R \text{TOL}_T}. \quad (4.8)$$

*Then, given a prescribed accuracy parameter  $\text{TOL}_T > 0$ , the adaptive refinement Algorithm 7 stops after a finite number of iterations.*

*Proof.* The main idea of the proof is to use the uniform upper bound on the error density,  $\bar{\rho} \leq \rho_{up}(\text{TOL}_T)$ , according to (1.10). Given an initial mesh size  $\Delta t_0$  with  $N_0$  time steps, Algorithm 7 satisfies both the stopping and the non-refinement conditions, (2.20) and (2.21), for the uniform mesh size

$$\tilde{\Delta}t(\text{TOL}_T) = \frac{\Delta t_0}{\max\{1, 2^k\}}, \quad \text{with } k = \left\lceil \log_2 \left( \frac{\rho_{up}(\text{TOL}_T) T \Delta t_0}{C_R \text{TOL}_T} \right) \right\rceil. \quad (4.9)$$

Furthermore, if during the refinement process one time step reaches the mesh size  $\tilde{\Delta}t(\text{TOL}_T)$  then it cannot be refined further, according to (2.21) and (4.9).

Since the number of possible refinements from the initial mesh size  $\Delta t_0$  to the mesh  $\tilde{\Delta}t(\text{TOL}_T)$  is bounded by  $N_0 \max\{1, 2^k\}$ , there is only a finite number of possible refinements. The proof is concluded by observing that the Algorithm 7 either stops or makes at least one refinement during each iteration.  $\square$

The work [20] also proves a similar stopping result, cf. Theorem 3.2 in [20], based on the assumption that the initial mesh is sufficiently refined so that the error density does not vary too much between refinement levels. Then, when the single level adaptive algorithm stops, one can prove asymptotic accuracy and efficiency estimates on the resulting approximation. In contrast, here we make essentially no assumption on the initial mesh size  $\Delta t_0$ : although the quality of the resulting approximation for the lower levels of the multilevel estimator may be poor, they have no influence in the bias of the multilevel approximation, which is only determined by the finest level,  $L$ . Since  $L \rightarrow \infty$  as  $\text{TOL} \downarrow 0$  we can still prove asymptotic accuracy and efficiency estimates. Finally, we observe that assumption (4.8) is fulfilled in all practical cases since one should start the adaptive algorithm with  $\bar{N}_{in}$  of the order of  $\text{TOL}_T^{-1}$ , which is much smaller than  $N_{up}$ .

The following proofs are inspired by the treatment by Chow and Robbins [3] on the accuracy and complexity of sequential stopping rules for sampling i.i.d. random variables.

We denote the SLMC sample average estimator of  $\mathbb{E}[g(X(T))]$  by

$$\mathcal{A}(g(\bar{X}(T)); M) = \sum_{i=1}^M \frac{g(\bar{X}(T; \omega_i))}{M},$$

where the realizations of  $\bar{X}(T)$  are generated on adaptive meshes and fulfill the weak error bound  $|\mathbb{E}[g(\bar{X}(T)) - g(X(T))]| \lesssim \text{TOL}_T$ . Here the total tolerance  $\text{TOL}$  is split into a time discretization error tolerance and a statistical error tolerance,  $\text{TOL} = C_S \text{TOL}_T + \text{TOL}_S$  (Remark 2 discusses the optimal splitting). Let  $2^{\mathbb{N}}$  denote the set  $\{2^n | n \in \mathbb{N}\}$ . For the SLMC estimator, the number of samples used in the sample average estimator to control the statistical error  $|\mathcal{A}(g(\bar{X}(T)); M) - \mathbb{E}[g(\bar{X}(T))]| \leq \text{TOL}_S$  is a stochastic process  $M : \mathbb{R}_+ \rightarrow 2^{\mathbb{N}}$  defined by

$$\begin{aligned} M(\text{TOL}_S) &:= \text{the smallest } k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil} \\ &\text{such that } \mathcal{V}(g(\bar{X}(T)); k) < \frac{k \text{TOL}_S^2}{C_C^2}, \end{aligned} \quad (4.10)$$

where the sample variance is defined by

$$\mathcal{V}(g(\bar{X}(T)); k) = \sum_{i=1}^k \frac{(g(\bar{X}(T; \omega_i)) - \mathcal{A}(g(\bar{X}(T)); k))^2}{k-1}. \quad (4.11)$$

Restricting the initial value of  $M$  to the set  $2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil}$  implies that that  $\lim_{\text{TOL} \downarrow 0} M = \infty$ . The asymptotic behavior of  $M$  as  $\text{TOL} \downarrow 0$  is crucial in our

proofs of the asymptotic accuracy and complexity. When proving the asymptotically accuracy result of Theorem 4,  $M$  should increase sufficiently fast to obtain the sought confidence. For the complexity result of Theorem 5, it is on the other hand useful to bound  $M$  from above and ensure that it does not grow too fast.

**Lemma 4.** *Suppose the assumptions of Lemma 1 hold and (4.7) holds for at least  $p = 2$ . Then*

$$\liminf_{\text{TOL} \downarrow 0} \frac{M \text{TOL}_S^2}{\text{Var}(g(\bar{X}(T))) C_C^2} = 1 \quad \text{a.s.} \quad \text{and} \quad \limsup_{\text{TOL} \downarrow 0} \frac{M \text{TOL}_S^2}{\text{Var}(g(\bar{X}(T))) C_C^2} = 2 \quad \text{a.s.} \quad (4.12)$$

*Proof.* The strong convergence (4.7) for  $p = 2$ , gives  $\lim_{\text{TOL} \downarrow 0} \text{Var}(g(\bar{X}(T))) = \text{Var}(g(X(T)))$ , which in particular means that there exists a constant  $\widetilde{\text{TOL}} > 0$  such that

$$\frac{\text{Var}(g(X(T)))}{2} < \text{Var}(g(\bar{X}(T))) < 2\text{Var}(g(X(T))), \quad \forall \text{TOL} \in (0, \widetilde{\text{TOL}}]. \quad (4.13)$$

The strong law of large numbers then implies that for all

$$\lim_{k \rightarrow \infty} \mathcal{V}(g(\bar{X}(T)); k) = \text{Var}(g(\bar{X}(T))) \quad \text{a.s.} \quad \forall \text{TOL} \in (0, \widetilde{\text{TOL}}]. \quad (4.14)$$

In order to prove (4.12), introduce the sequence of stochastic processes  $y_k : \mathbb{R}_+ \rightarrow \mathbb{R}_+$  sub-indexed by  $k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil}$  and defined by

$$y_k(\text{TOL}) = \frac{\mathcal{V}(g(\bar{X}(T)); k)}{\text{Var}(g(\bar{X}(T)))}. \quad (4.15)$$

Using  $y_k$ , definition (4.10) of  $M(\text{TOL}_S)$  is equivalent to

$$M(\text{TOL}_S) := \text{the smallest } k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil} \\ \text{such that } y_k(\text{TOL}_S) < \frac{k \text{TOL}_S^2}{\text{Var}(g(\bar{X}(T))) C_C^2}.$$

This gives rise to the bounds

$$y_M(\text{TOL}_S) < \frac{M \text{TOL}_S^2}{\text{Var}(g(\bar{X}(T))) C_C^2} \leq 2y_{M/2}(\text{TOL}_S). \quad (4.16)$$

Combining (4.14) with definition (4.10), which ensures that  $\lim_{\text{TOL} \downarrow 0} M = \infty$ , we conclude that

$$\lim_{\text{TOL} \downarrow 0} \mathcal{V}(g(\bar{X}(T)); M(\text{TOL}_S)) = \text{Var}(g(X(T))) > 0 \quad \text{a.s.}$$

which implies that also  $\lim_{\text{TOL} \downarrow 0} y_M(\text{TOL}_S) = 1$  a.s. Statement (4.12) then follows by taking limits in (4.16).  $\square$

Having obtained asymptotic bounds for  $M$ , we are ready to prove the main accuracy result for the SLMC algorithm.

**Theorem 4** (Single level accuracy). *Suppose the modeling assumptions of Lemma 1 holds, that (4.7) holds for at least  $p = 2$ , and that  $\text{TOL}_T \leq \text{TOL}_S$ . Then, the adaptive SLMC algorithm with confidence refinement parameter  $C_C > 0$ , and stochastic time steps (2.20) and (2.21), satisfies*

$$\liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T))] - \mathcal{A}(g(\bar{X}(T))); M| \leq \text{TOL}) \geq \int_{-C_C}^{C_C} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx. \quad (4.17)$$

*Proof.* For a given  $\delta > 0$ , we first bound the probability in (4.17) from below as follows:

$$\begin{aligned} & \liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T))] - \mathcal{A}(g(\bar{X}(T))); M| \leq \text{TOL}) \\ & \geq \liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T)) - g(\bar{X}(T))]| \\ & \quad + |E[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T))); M| \leq C_S \text{TOL}_T + \text{TOL}_S) \\ & \geq \liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T)) - g(\bar{X}(T))]| \leq (C_S + \delta) \text{TOL}_T) \quad (4.18) \\ & \quad \text{and } |E[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T))); M| \leq (1 - \delta) \text{TOL}_S) \\ & = \liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T)) - g(\bar{X}(T))]| \leq (C_S + \delta) \text{TOL}_T) \\ & \quad \times P(|E[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T))); M| \leq (1 - \delta) \text{TOL}_S) \end{aligned}$$

The proof is continued by analyzing the two product terms of the last line of the inequality above separately:

**The time discretization error.** From the treatment of the time discretization error for the single level case we refer to the proof of Theorem 3.4, p. 530 in [19] which shows that

$$\limsup_{\text{TOL} \downarrow 0} \frac{|E[g(X(T)) - g(\bar{X}(T))]|}{\text{TOL}_T} \leq C_S.$$

Thereby,

$$\liminf_{\text{TOL} \downarrow 0} P(|E[g(X(T)) - g(\bar{X}(T))]| \leq (C_S + \delta) \text{TOL}_T) = 1.$$

**The statistical error.** For the above introduced  $\delta > 0$ , define the family of sets

$$\Omega_\delta(\text{TOL}_S) = \left\{ k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil} \mid 1 - \delta < \frac{k \text{TOL}_S^2}{\text{Var}(g(\bar{X}(T))) C_C^2} \leq 2 + \delta \right\}. \quad (4.19)$$

By the convergence (4.12), we conclude that

$$\lim_{\text{TOL} \downarrow 0} P(M \in \Omega_\delta) = 1.$$



Recall that for the SLMC algorithm, the number of samples  $M$  is determined in the step prior to generating the output  $\mathcal{A}(g(\bar{X}(T)); M)$ , so that  $M$  is independent from  $\mathcal{A}(g(\bar{X}(T)); M)$ . Using this independence property, Fatou's Lemma, and Lindeberg-Feller's version of the Central Limit Theorem, cf. Theorem 6, yield that

$$\begin{aligned}
& \liminf_{\text{TOL} \downarrow 0} P(|\mathbb{E}[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T)); M)| \leq (1 - \delta)\text{TOL}_S) \\
&= \liminf_{\text{TOL} \downarrow 0} \sum_{k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil}} P(|\mathbb{E}[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T)); k)| \leq (1 - \delta)\text{TOL}_S) P(M = k) \\
&\geq \liminf_{\text{TOL} \downarrow 0} \sum_{k \in \Omega_\delta} P(|\mathbb{E}[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T)); k)| \leq (1 - \delta)\text{TOL}_S) P(M = k) \\
&\quad + \sum_{k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil} \setminus \Omega_\delta} \liminf_{\text{TOL} \downarrow 0} P(|\mathbb{E}[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T)); k)| \leq (1 - \delta)\text{TOL}_S) P(M = k) \\
&\geq \liminf_{\text{TOL} \downarrow 0} \sum_{k \in \Omega_\delta} P\left(\sqrt{k} \frac{|\mathbb{E}[g(\bar{X}(T))] - \mathcal{A}(g(\bar{X}(T)); k)|}{\text{Var}(g(\bar{X}(T)))} \leq (1 - \delta)^{3/2} C_C\right) P(M = k) \\
&\geq \int_{-(1-\delta)^{3/2} C_C}^{(1-\delta)^{3/2} C_C} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx.
\end{aligned} \tag{4.20}$$

The proof is finished by noting that the argument leading to inequality (4.20) is valid for all  $\delta > 0$ .  $\square$

We conclude this subsection with a complexity analysis of the SLMC algorithm. Similar to the definition of the work for the MLMC algorithm given in (4.2), we define the SLMC work by

$$\text{WORK}(\text{TOL}) = \mathbb{E}[M]\mathbb{E}[N], \tag{4.21}$$

where we recall that  $M$  denotes the number of samples of  $g(\bar{X}(T))$  required to control the statistical error and  $N$  denotes the number of adaptive time steps required in the construction of a numerical realization  $g(\bar{X}(T; \omega))$  to control the time discretization error  $|\mathbb{E}[g(\bar{X}(T)) - g(X(T))]| \leq \text{TOL}_T$ . We start by bounding  $\mathbb{E}[M]$ .

**Lemma 5.** *Suppose the assumptions of Lemma 1 and (4.7) hold. Then the number of samples used in the approximation of  $\mathbb{E}[g(X(T))]$  is bounded by*

$$\limsup_{\text{TOL} \downarrow 0} \frac{\mathbb{E}[M]\text{TOL}_S^2}{\text{Var}(g(\bar{X}(T)))C_C^2} \leq 2. \tag{4.22}$$

*Proof.* For a given  $\delta > 0$ , define the deterministic function

$$\tilde{M}(\text{TOL}_S) = \min \left\{ k \in 2^{\mathbb{N} + \lceil \log_2(\text{TOL}^{-1}) \rceil} \mid \frac{k\text{TOL}_S^2}{\text{Var}(g(\bar{X}(T)))C_C^2} > 1 + \delta \right\}.$$

Assuming TOL is sufficiently small so that (4.13) holds, the relation (4.16), the fourth moment bound (4.7) and k-Statistics bounds on the variance of the sample variance, cf. [17], yield

$$\begin{aligned}
P(M = 2\widetilde{M}) &\leq P\left(\frac{\mathcal{V}(g(\overline{X}(T)); \widetilde{M})}{\text{Var}(g(\overline{X}(T)))} > \widetilde{M} \frac{\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2}\right) \\
&\leq P\left(\frac{\mathcal{V}(g(\overline{X}(T)); \widetilde{M})}{\text{Var}(g(\overline{X}(T)))} > 1 + \delta\right) \\
&\leq P\left(|\mathcal{V}(g(\overline{X}(T)); \widetilde{M}) - \text{Var}(g(\overline{X}(T)))| > \delta \text{Var}(g(\overline{X}(T)))\right) \\
&\leq 2\mathbb{E}\left[\frac{|\mathcal{V}(g(\overline{X}(T)); \widetilde{M}) - \text{Var}(g(\overline{X}(T)))|^2}{\delta^2 \text{Var}(g(\overline{X}(T)))^2}\right] \\
&< \frac{C}{\delta^2 \widetilde{M}},
\end{aligned}$$

Furthermore, for  $\ell = 1, 2, \dots$  we get that

$$\begin{aligned}
P(M = 2^{\ell+1}\widetilde{M}) &\leq P\left(|\mathcal{V}(g(\overline{X}(T)); 2^\ell \widetilde{M}) - \text{Var}(g(\overline{X}(T)))| > 2^{\ell-1} \text{Var}(g(\overline{X}(T)))\right) \\
&\leq 2\mathbb{E}\left[\frac{|\mathcal{V}(g(\overline{X}(T)); 2^\ell \widetilde{M}) - \text{Var}(g(\overline{X}(T)))|^2}{2^{2(\ell-1)} \text{Var}(g(\overline{X}(T)))^2}\right] \\
&< \frac{C}{2^{2\ell} \widetilde{M}}.
\end{aligned}$$

Consequently,

$$\begin{aligned}
\frac{\mathbb{E}[M]\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2} &\leq \frac{[P(M \leq \widetilde{M}) + \sum_{\ell=1}^{\infty} 2^\ell P(M = 2^\ell \widetilde{M})] \widetilde{M}\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2} \\
&\leq 2(1 + \delta) \left[ P(M \leq \widetilde{M}) + P(M = 2\widetilde{M}) + \sum_{\ell=1}^{\infty} 2^{\ell+1} P(M = 2^{\ell+1}\widetilde{M}) \right] \quad (4.23) \\
&\leq 2(1 + \delta) \left[ P(M \leq \widetilde{M}) + \frac{C}{\delta^2 \widetilde{M}} + \frac{C}{\widetilde{M}} \sum_{\ell=1}^{\infty} 2^{-\ell} \right].
\end{aligned}$$

By taking limits in the above inequality, we obtain

$$\limsup_{\text{TOL} \downarrow 0} \frac{\mathbb{E}[M]\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2} \leq 2(1 + \delta),$$

and noting that this result holds for any  $\delta > 0$ , the proof is finished.  $\square$

For an asymptotic bound on  $E[N]$ , we recall Theorem 3.5 of [19]. The bound given in this theorem is derived by studying the asymptotic form of the error indicators obtained by the stopping condition (2.20). The theorem further shows that up to a multiplicative constant, the mesh refinement scheme (2.20)-(2.21) yields stochastic meshes which are optimal in mean sense. The theorem is here stated as a lemma.

**Lemma 6** (Single level asymptotic average number of time steps). *Suppose that the assumptions of Lemma 1 hold. Then the final number of adaptive steps generated by the algorithm (2.20) and (2.21) satisfies asymptotically*

$$\limsup_{\text{TOL} \downarrow 0} \text{TOL}_T E[N] \leq \frac{4}{C_R} \left( E \left[ \int_0^T \sqrt{|\hat{\rho}(t)|} dt \right] \right)^2. \quad (4.24)$$

The product of the asymptotic upper bounds for  $E[M]$  and  $E[N]$  and an optimization of the choice of  $\text{TOL}_T$  and  $\text{TOL}_S$  gives the following upper bound on the computational complexity for the SLMC algorithm.

**Theorem 5** (SLMC computational complexity). *Suppose the assumptions of Lemma 1 and (4.7) hold. Then the work for the SLMC algorithm with stochastic time steps, defined in (4.21), satisfies*

$$\limsup_{\text{TOL} \downarrow 0} \text{WORK}(\text{TOL}) \text{TOL}^3 \leq \frac{2 \cdot 3^3 \text{Var}(g(X(T))) C_C^2 C_S}{C_R} \left( E \left[ \int_0^T \sqrt{|\hat{\rho}(t)|} dt \right] \right)^2, \quad (4.25)$$

where  $C_C$  is the confidence parameter and  $C_R$  and  $C_S$  are refinement parameters described by (2.20) and (2.21).

*Proof.* Lemma 5 and 6 straightforwardly yield the upper bound

$$\limsup_{\text{TOL} \downarrow 0} \text{WORK}(\text{TOL}) \text{TOL}_S^2 \text{TOL}_T \leq \frac{2^3 \text{Var}(g(X(T))) C_C^2}{C_R} \left( E \left[ \int_0^T \sqrt{|\hat{\rho}(t)|} dt \right] \right)^2.$$

So  $\text{WORK}(\text{TOL}) = \mathcal{O}(\text{TOL}_S^{-2} \text{TOL}_T^{-1})$ . Minimizing  $\text{TOL}_S^{-2} \text{TOL}_T^{-1}$  subject to the restriction  $C_S \text{TOL}_T + \text{TOL}_S = \text{TOL}$ , yields

$$\text{TOL}_T = \frac{\text{TOL}}{3C_S} \text{ and } \text{TOL}_S = \frac{2\text{TOL}}{3}.$$

These values for  $\text{TOL}_T$  and  $\text{TOL}_S$  lead to the upper bound (4.25).  $\square$

**Remark 2.** *The optimal choices of  $\text{TOL}_T$  and  $\text{TOL}_S$  for minimizing  $\text{WORK}(\text{TOL})$  are derived in the proof of Theorem 5 to be*

$$\text{TOL}_T = \frac{\text{TOL}}{3C_S} \text{ and } \text{TOL}_S = \frac{2\text{TOL}}{3}.$$

## 4.2 multilevel results

We recall from the description of the MLMC algorithm in Section 2.2 that given an accuracy  $\text{TOL} = C_S \text{TOL}_T + \text{TOL}_S$ , the MLMC algorithm generates realizations  $g(\bar{X}_\ell(T))$  on adaptive meshes fulfilling the weak error bounds  $|\mathbb{E}[g(\bar{X}_\ell(T)) - g(X(T))]| \lesssim \text{TOL}_{T,\ell}$  on the levels  $\ell = 0, 1, \dots, L$ . The time discretization tolerance levels are given by  $\text{TOL}_{T,\ell} = 2^\ell \text{TOL}_T$ , and the number of levels is set by  $L = \lfloor \log_2(\text{TOL}_{T,\text{Max}}/\text{TOL}_T) \rfloor$ , where  $\text{TOL}_{T,\text{Max}}$  is a predetermined max time discretization tolerance value, cf. (2.3). The MLMC sample average estimator of  $\mathbb{E}[g(X(T))]$  is denoted by

$$\mathcal{A}_{\mathcal{MC}}(g(\bar{X}(T)); M_0) = \sum_{i=1}^{M_0} \frac{g(\bar{X}_0(T; \omega_{0,i}))}{M_0} + \sum_{\ell=1}^L \sum_{i=1}^{M_\ell} \frac{\Delta_\ell g(\bar{X}(T; \omega_{\ell,i}))}{M_\ell},$$

where  $M_0 \in 2^{L+\lceil C_{\mathcal{MC}}L \rceil} 2^{\mathbb{N}}$  denotes the number of samples on the coarsest level with the constant  $C_{\mathcal{MC}} \in (0, 1)$ , and the number of samples on higher levels is expressed in terms of  $M_0$  by the ratio

$$M_\ell = \frac{M_0}{2^L} \left[ 2^L \frac{\rho_{\text{low}}(\text{TOL}_{T,0}) \text{TOL}_{T,\ell}}{\rho_{\text{low}}(\text{TOL}_{T,\ell}) \text{TOL}_{T,0}} \right] = \frac{M_0}{2^L} \left[ 2^{L+(\bar{\gamma}-1)\ell} \right] \quad \ell = 1, 2, \dots, L. \quad (4.26)$$

The number of samples at the coarsest level is a stochastic process  $M_0 : \mathbb{R}_+ \rightarrow 2^{\mathbb{N}+L+\lceil C_{\mathcal{MC}}L \rceil}$  defined by

$$\begin{aligned} M_0(\text{TOL}_S) &= \text{the smallest } k_0 \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{MC}}L \rceil} \text{ such that} \\ \mathcal{V}_{\mathcal{MC}}(g(\bar{X}(T)); k_0) &< \frac{k_0 \text{TOL}_S^2}{C_C^2}, \end{aligned} \quad (4.27)$$

where

$$\begin{aligned} \mathcal{V}_{\mathcal{MC}}(g(\bar{X}(T)); k_0) &= \sum_{i=1}^{k_0} \frac{(g(\bar{X}_0(T; \omega_{0,i})) - \mathcal{A}(g(\bar{X}_0(T; \omega_{0,\cdot})); k_0))^2}{k_0 - 1} \\ &+ \sum_{\ell=1}^L \frac{k_0}{k_\ell} \sum_{i=1}^{k_\ell} \frac{(\Delta_\ell g(\bar{X}(T; \omega_{\ell,i})) - \mathcal{A}(\Delta_\ell g(\bar{X}(T; \omega_{\ell,\cdot})); k_\ell))^2}{k_\ell - 1} \\ &= \mathcal{V}(g(\bar{X}_0(T; \omega_{0,\cdot})); k_0) + 2^L \sum_{\ell=1}^L \frac{\mathcal{V}(\Delta_\ell g(\bar{X}_0(T; \omega_{\ell,\cdot})); k_\ell)}{[2^{L+\ell(\bar{\gamma}-1)}]} \end{aligned} \quad (4.28)$$

and, analogous to the definition of  $M_\ell$ ,

$$k_\ell := \frac{k_0}{2^L} \left[ 2^{L+(\bar{\gamma}-1)\ell} \right], \quad \ell = 1, 2, \dots, L. \quad (4.29)$$

**Remark 3.** In the analysis of the SLMC Algorithm, the requirement  $M_0 \in 2^{\mathbb{N}+\lceil \log(1/\text{TOL}) \rceil}$  ensured that the number of samples used in the MC estimate fulfilled  $\liminf_{\text{TOL} \downarrow 0} M = \infty$ . For the MLMC Algorithm, we analogously ensure that  $\liminf_{\text{TOL} \downarrow 0} M_L = \infty$  by requiring that  $M_0 \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{MC}}L \rceil}$  for any positive constant  $C_{\mathcal{MC}}$ .

The stochastic process  $M_0$  is defined in a similar way as the stochastic process  $M$  was defined for SLMC algorithm, cf. (4.10). For the SLMC algorithm, asymptotic accuracy and complexity results were easily obtained by applying the asymptotic bounds of  $M$ , cf. Lemma 4. Applying the same strategy for the MLMC algorithm, we will derive asymptotic bounds for  $M_0$  and use these bounds to prove the accuracy and complexity results of Theorem 2 and 3.

**Lemma 7** (Asymptotic bounds for  $M_0$ ). *Let*

$$\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) := \text{Var}(g(\bar{X}_0(T))) + 2^L \sum_{\ell=1}^L \frac{\text{Var}(\Delta_\ell g(\bar{X}(T)))}{\lceil 2^{L+\ell(\bar{\gamma}-1)} \rceil} \quad (4.30)$$

and suppose the assumptions of Lemma 1 and (4.7) hold. Further, assume that  $\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) > 0$  for all sufficiently small  $\text{TOL} > 0$ . Then  $M_0(\text{TOL}_S)$  defined according to (4.27) fulfills

$$\begin{aligned} \liminf_{\text{TOL} \downarrow 0} \frac{M_0 \text{TOL}_S^2}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) C_C^2} &= 1 \quad \text{in probability, and} \\ \limsup_{\text{TOL} \downarrow 0} \frac{M_0 \text{TOL}_S^2}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) C_C^2} &= 2 \quad \text{in probability.} \end{aligned} \quad (4.31)$$

*Proof.* The definition of  $M_0$  given in (4.27) implies that the following inequalities hold:

$$\frac{\mathcal{V}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)); M_0)}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))} \leq \frac{M_0 \text{TOL}_S^2}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) C_C^2} \leq 2 \frac{\mathcal{V}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)); M_0/2)}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))}.$$

So to conclude the proof, we will show that

$$\lim_{\text{TOL} \downarrow 0} \frac{\mathcal{V}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)); M_0)}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))} = 1, \quad \text{in probability.} \quad (4.32)$$

Define the deterministic function  $\tilde{k}_0(\text{TOL}_T) = 2^{L(\text{TOL}_T) + \lceil C_{\mathcal{M}\mathcal{C}} L(\text{TOL}_T) \rceil + 1}$  and let  $\{\tilde{k}_\ell\}_{\ell=1}^L$  be the corresponding level functions defined according to (4.29).

Then, for a given  $\epsilon > 0$ , let us consider

$$\begin{aligned}
& P \left( \left| \frac{\mathcal{V}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)); \tilde{k}_0)}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))} - 1 \right| > \epsilon \right) \\
&= P \left( \left| \mathcal{V}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)); \tilde{k}_0) - \text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T))) \right| > \text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon \right) \\
&\leq P \left( \left| \mathcal{V}(g(\bar{X}_0(T)); \tilde{k}_0) - \text{Var}(g(\bar{X}_0(T))) \right| + \sum_{\ell=1}^L 2^L \left[ 2^{L+\ell(\bar{\gamma}-1)} \right]^{-1} \right. \\
&\quad \left. \times \left| \mathcal{V}(\Delta_\ell g(\bar{X}(T)); \tilde{k}_\ell) - \text{Var}(\Delta_\ell g(\bar{X}(T))) \right| > \text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon \right) \\
&\leq P \left( \left| \mathcal{V}(g(\bar{X}_0(T)); \tilde{k}_0) - \text{Var}(g(\bar{X}_0(T))) \right| > \frac{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon}{L+1} \right) \\
&\quad + \sum_{\ell=1}^L P \left( 2^{(1-\bar{\gamma})\ell} \left| \mathcal{V}(\Delta_\ell g(\bar{X}(T)); \tilde{k}_\ell) - \text{Var}(\Delta_\ell g(\bar{X}(T))) \right| > \frac{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon}{L+1} \right)
\end{aligned}$$

From the fourth moment bound (4.7), Chebycheff's inequality and k-Statistics bounds on the variance of the sample variance, cf. [17], we get that

$$\begin{aligned}
& P \left( \left| \mathcal{V}(g(\bar{X}_0(T)); \tilde{k}_0) - \text{Var}(g(\bar{X}_0(T))) \right| > \frac{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon}{L+1} \right) \\
&\leq \frac{C(L+1)^2}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))^2 \epsilon^2 \tilde{k}_0}.
\end{aligned}$$

The equality  $2^{(1-\bar{\gamma})\ell} = \frac{\rho_{low}(\text{TOL}_{\mathcal{T},\ell})\text{TOL}_{\mathcal{T},0}}{\rho_{low}(\text{TOL}_{\mathcal{T},0})\text{TOL}_{\mathcal{T},\ell}}$  combined with (4.7) further yields that

$$\begin{aligned}
& P \left( 2^{(1-\bar{\gamma})\ell} \left| \mathcal{V}(\Delta_\ell g(\bar{X}(T)); \tilde{k}_\ell) - \text{Var}(\Delta_\ell g(\bar{X}(T))) \right| > \frac{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))\epsilon}{L+1} \right) \\
&\leq \frac{C(L+1)^2}{\text{Var}_{\mathcal{M}\mathcal{C}}(g(\bar{X}(T)))^2 \epsilon^2 \tilde{k}_\ell}.
\end{aligned}$$

Since  $\tilde{k}_0 = 2^{L+\lceil C_{\mathcal{M}\mathcal{C}}L \rceil + 1}$ , the definition of  $\tilde{k}_\ell$  in (4.29) implies that  $\tilde{k}_\ell \geq 2^{L+\lceil C_{\mathcal{M}\mathcal{C}}L \rceil + 1 + (\bar{\gamma}-1)\ell}$  for  $\ell = 1, 2, \dots, L$ , with  $\bar{\gamma} \geq 0$  denoting the lower error

density exponent in  $\rho_{low}(\text{TOL}_T) = \text{TOL}_T^{\tilde{\gamma}}$ , cf. (1.24). Consequently,

$$\begin{aligned} P \left( \left| \frac{\mathcal{V}_{\mathcal{M}C}(g(\bar{X}(T)); \tilde{k}_0)}{\text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))} - 1 \right| > \epsilon \right) &\leq \frac{C(L+1)^2}{\text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))^2 \epsilon^2 \tilde{k}_0} \sum_{\ell=0}^L \frac{\tilde{k}_0}{\tilde{k}_\ell} \\ &\leq \frac{C(L+1)^2}{\text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))^2 \epsilon^2 \tilde{k}_0} \sum_{\ell=0}^L 2^{(1-\tilde{\gamma})\ell} \\ &< \frac{C(L+1)^2}{2^{\lceil C_{\mathcal{M}C}L \rceil + \tilde{\gamma}L} \text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))^2 \epsilon^2} \end{aligned}$$

which implies that for any  $\epsilon > 0$ ,

$$\begin{aligned} \lim_{\text{TOL} \downarrow 0} P \left( \left| \frac{\mathcal{V}_{\mathcal{M}C}(g(\bar{X}(T)); \tilde{k}_0)}{\text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))} - 1 \right| > \epsilon \right) \\ < \lim_{\text{TOL} \downarrow 0} \frac{C(L+1)^2}{2^{\lceil C_{\mathcal{M}C}L \rceil + \tilde{\gamma}L} \text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))^2 \epsilon^2} = 0. \end{aligned}$$

Since  $M_0 \geq \tilde{k}_0$  by definition, we conclude that also (4.32) holds, i.e.

$$\lim_{\text{TOL} \downarrow 0} P \left( \left| \frac{\mathcal{V}_{\mathcal{M}C}(g(\bar{X}(T)); M_0)}{\text{Var}_{\mathcal{M}C}(g(\bar{X}(T)))} - 1 \right| > \epsilon \right) = 0,$$

for any  $\epsilon > 0$ . □

With the asymptotic bounds on  $M_0$  we are ready to prove the main asymptotic accuracy result, Theorem 2. For the convenience of the reader, we first recall its formulation.

**Theorem 2** (Multilevel accuracy). *Suppose that the modeling assumptions of Lemma 1 hold, that (4.7) holds, and that  $\text{TOL}_T \leq \text{TOL}_S$ . Then the adaptive MLMC algorithm with confidence parameter  $C_C > 0$  and stochastic time steps (2.20) and (2.21) satisfies*

$$\liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T))] - \mathcal{A}_{\mathcal{M}C}(g(\bar{X}(T)); M_0)| \leq \text{TOL} \right) \geq \int_{-C_C}^{C_C} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx. \quad (4.33)$$

*Proof.* This proof is quite similar to the proof of Theorem 4 for the asymptotic accuracy in the single level setting, but for the sake of the differing details, a full proof is included in this setting also. For a given  $\delta > 0$ , we start by bounding the left-hand side of (4.33) by a product of the statistical error and the time

discretization error.

$$\begin{aligned}
& \liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T))] - \mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)| \leq \text{TOL} \right) \\
& \geq \liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T)) - g(\bar{X}_L(T))]| \right. \\
& \quad \left. + |\mathbb{E}[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)| \leq C_S \text{TOL}_T + \text{TOL}_S \right) \\
& \geq \liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T)) - g(\bar{X}_L(T))]| \leq (C_S + \delta) \text{TOL}_T \right. \\
& \quad \left. \text{and } |\mathbb{E}[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)| \leq (1 - \delta) \text{TOL}_S \right) \\
& = \liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T)) - g(\bar{X}_L(T))]| \leq (C_S + \delta) \text{TOL}_T \right) \\
& \quad \times P \left( |\mathbb{E}[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)| \leq (1 - \delta) \text{TOL}_S \right)
\end{aligned}$$

**The time discretization error.** The proof of Theorem 3.4, p. 530 in [19] shows that the following bound is fulfilled

$$\lim_{\text{TOL} \downarrow 0} \frac{|\mathbb{E}[g(X(T)) - g(\bar{X}_L(T))]|}{\text{TOL}_{T,L}} \leq C_S.$$

By construction  $\text{TOL}_{T,L} = \text{TOL}_T$ , and this implies by the above that

$$\liminf_{\text{TOL} \downarrow 0} P \left( |\mathbb{E}[g(X(T)) - g(\bar{X}_L(T))]| \leq (1 + \delta) C_S \text{TOL}_T \right) = 1.$$

**The statistical error.** From the above introduced  $\delta > 0$ , define the family of sets

$$\Omega_\delta(\text{TOL}_S) = \left\{ k \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{M}c}L \rceil} \left| 1 - \delta < \frac{k \text{TOL}_S^2}{\text{Var}_{\mathcal{M}c}(g(\bar{X}(T))) C_C^2} \leq 2 + \delta \right. \right\}, \quad (4.34)$$

indexed by  $\text{TOL}_S > 0$ . Lemma 7 then implies that  $\lim_{\text{TOL} \downarrow 0} P(M_0 \in \Omega_\delta) = 1$ . Recall further that for the MLMC algorithm, the number of samples  $M_0$  is determined in the step prior to generating the output  $\mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)$ , so that  $M_0$  is independent from  $\mathcal{A}_{\mathcal{M}c}(g(\bar{X}(T)); M_0)$ . Using this independence property



and Fatou's Lemma, the statistical error is bounded from below as follows:

$$\begin{aligned}
& \liminf_{\text{TOL}\downarrow 0} P(|E[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); M_0)| \leq (1-\delta)\text{TOL}_S) \\
&= \liminf_{\text{TOL}\downarrow 0} \sum_{k_0 \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{M}\mathcal{L}}L \rceil}} P(|E[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)| \leq (1-\delta)\text{TOL}_S) P(M_0 = k_0) \\
&\geq \liminf_{\text{TOL}\downarrow 0} \sum_{k_0 \in \Omega_\delta} P(|E[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)| \leq (1-\delta)\text{TOL}_S) P(M_0 = k_0) \\
&\quad + \sum_{k_0 \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{M}\mathcal{L}}L \rceil} \setminus \Omega_\delta} \liminf_{\text{TOL}\downarrow 0} P(|E[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)| \leq (1-\delta)\text{TOL}_S) P(M_0 = k_0) \\
&\geq \liminf_{\text{TOL}\downarrow 0} \sum_{k_0 \in \Omega_\delta} P\left(\sqrt{k_0} \frac{|E[g(\bar{X}_L(T))] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)|}{\sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} \leq (1-\delta)^{3/2} C_C\right) P(M_0 = k_0) \\
&\geq \int_{-(1-\delta)^{3/2} C_C}^{(1-\delta)^{3/2} C_C} \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx.
\end{aligned} \tag{4.35}$$

The last inequality above follows from the application of Lindeberg-Feller's Central Limit Theorem (CLT) which is justified by Lemma 8 and the observation that  $E[\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)] = E[g(\bar{X}_L(T))]$ . The reasoning leading to inequality (4.35) is valid for any  $\delta > 0$ , so the proof of (4.33) is finished.  $\square$

Next we derive the weak convergence CLT result for the multilevel estimator  $\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)$  which is needed in Theorem 2.

**Lemma 8** (A CLT result). *Suppose the assumptions of Lemma 1 and (4.7) hold and, for a given  $\delta > 0$ , let*

$$k_0(\text{TOL}_S) := \min \left\{ k \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{M}\mathcal{L}}L \rceil} \left| \frac{k \text{TOL}_S^2}{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T))) C_C^2} > 1 - \delta \right. \right\},$$

in correspondence with the set defined in (4.34). Then for any  $z \in \mathbb{R}_+$ , we have that

$$\lim_{\text{TOL}\downarrow 0} P\left(\sqrt{k_0} \frac{|E[\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)|}{\sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} \leq z\right) = \int_{-z}^z \frac{e^{-x^2/2}}{\sqrt{2\pi}} dx. \tag{4.36}$$

*Proof.* This Lemma will be proved by verifying that the assumptions of the Lindeberg-Feller CLT are fulfilled, cf. Theorem 6. Let us write

$$\sqrt{k_0} \frac{E[\mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)] - \mathcal{A}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)); k_0)}{\sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} = \sum_{i=1}^K Y_{K,i}$$

where  $K := \sum_{\ell=0}^L k_\ell$  and the elements of  $Y_{K,i}$  are independent and defined by

$$Y_{K,i} := \begin{cases} \frac{\mathbb{E}[g(\bar{X}_0(T))] - g(\bar{X}_0(T; \omega_i))}{\sqrt{k_0} \sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} & \text{for } i = 1, 2, \dots, k_0, \\ \frac{\sqrt{\frac{k_0}{k_1}} (\mathbb{E}[\Delta_1 g(\bar{X}(T))] - \Delta_1 g(\bar{X}(T; \omega_i)))}{\sqrt{k_1} \sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} & \text{for } i = k_0 + 1, \dots, k_0 + k_1 \\ \vdots & \vdots \\ \frac{\sqrt{\frac{k_0}{k_L}} (\mathbb{E}[\Delta_L g(\bar{X}(T))] - \Delta_L g(\bar{X}(T; \omega_i)))}{\sqrt{k_L} \sqrt{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}} & \text{for } i = k_{L-1} + 1, \dots, K. \end{cases}$$

Then it follows that

$$\sum_{i=1}^K \mathbb{E}[Y_{K,i}^2] = \frac{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))} = 1, \quad \forall \text{TOL} > 0,$$

so condition (a) of Theorem 6 is fulfilled. To verify that condition (b) of Theorem 6 is fulfilled, one must show that for any  $\epsilon > 0$ ,

$$\limsup_{\text{TOL} \rightarrow 0} \sum_{i=1}^K \mathbb{E}[Y_{K,i}^2 1_{|Y_{K,i}| > \epsilon}] = 0.$$

The definition of  $k_\ell$ , cf. (4.29), combined with the moment bound (4.7) implies that there exists a  $C > 0$  such that

$$\mathbb{E} \left[ \left( \frac{k_0}{k_\ell} \right)^2 |\Delta_\ell g(\bar{X}(T)) - \mathbb{E}[\Delta_\ell g(\bar{X}(T))]|^4 \right] \leq C, \quad \forall \ell \in \{1, 2, \dots, L\}.$$

Using Chebycheff's inequality and the fact that  $k_L \geq 2^{\lceil C_{\mathcal{M}\mathcal{L}} L \rceil + \bar{\gamma} L + 1}$ , cf. (4.29), we derive that

$$\begin{aligned} \sum_{i=1}^K \mathbb{E}[Y_{K,i}^2 1_{|Y_{K,i}| > \epsilon}] &\leq \sum_{i=1}^K \epsilon^{-2} \mathbb{E}[Y_{K,i}^4] \\ &= \frac{1}{\epsilon^2 \text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))^2} \left\{ \frac{1}{k_0} \mathbb{E} \left[ |g(\bar{X}_0(T)) - \mathbb{E}[g(\bar{X}_0(T))]|^4 \right] \right. \\ &\quad \left. + \sum_{\ell=1}^L \frac{1}{k_\ell} \mathbb{E} \left[ \left( \frac{k_0}{k_\ell} \right)^2 |\Delta_\ell g(\bar{X}(T)) - \mathbb{E}[\Delta_\ell g(\bar{X}(T))]|^4 \right] \right\} \\ &\leq \frac{C}{\epsilon^2 \text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))^2} \sum_{\ell=0}^L k_\ell^{-1} \\ &\leq \frac{C L}{k_L \epsilon^2 \text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))^2} \rightarrow 0, \quad \text{as } \text{TOL} \downarrow 0. \end{aligned}$$

This verifies that condition (b) is fulfilled.  $\square$

We conclude the analysis of the MLMC algorithm by estimating the work required to fulfill the accuracy estimate (4.1). We recall that  $\text{WORK}(\text{TOL})$ , defined in (4.2) by

$$\text{WORK}(\text{TOL}) = \sum_{\ell=0}^L \mathbb{E}[M_\ell] \mathbb{E}[N_\ell],$$

is an estimate of the average number of operations required in the generation of  $\mathcal{A}_{\mathcal{M}_\ell}(g(\bar{X}(T)); M_0)$  to approximate  $\mathbb{E}[g(X(T))]$  with the prescribed confidence  $C_C$  and accuracy  $\text{TOL}$ . First, let us derive an asymptotic bound for  $\mathbb{E}[M_0]$ .

**Lemma 9.** *Suppose the assumptions of Lemma 1 and (4.7) hold. Then the number of samples  $M_0$  used at the base level of the MLMC algorithm approximation of  $\mathbb{E}[g(X(T))]$  satisfies*

$$\limsup_{\text{TOL} \downarrow 0} \frac{\mathbb{E}[M_0] \text{TOL}_S^2}{\text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) C_C^2} \leq 2. \quad (4.37)$$

*Proof.* For given  $\delta > 0$ , define the deterministic function

$$\widetilde{M}_0(\text{TOL}) = \min \left\{ k \in 2^{\mathbb{N}+L+\lceil C_{\mathcal{M}_\ell} L \rceil} \mid \frac{k_0 \text{TOL}^2}{\text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) C_C^2} > 1 + \delta \right\}$$

By the relation (4.16), the moment bound assumption (4.7), Hölder's inequality, and k-Statistics bounds on the variance of the sample variance, cf. [17], we derive that

$$\begin{aligned} P(M_0 = 2\widetilde{M}_0) &\leq P \left( \frac{\mathcal{V}_{\mathcal{M}_\ell}(g(\bar{X}(T)); \widetilde{M}_0)}{\text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T)))} > \widetilde{M}_0 \frac{\text{TOL}^2}{\text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) C_C^2} \right) \\ &\leq P \left( \frac{\mathcal{V}_{\mathcal{M}_\ell}(g(\bar{X}(T)); \widetilde{M}_0)}{\text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T)))} > 1 + \delta \right) \\ &\leq P \left( \mathcal{V}_{\mathcal{M}_\ell}(g(\bar{X}(T)); \widetilde{M}_0) - \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) > \delta \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) \right) \\ &\leq \mathbb{E} \left[ \frac{\left| \mathcal{V}_{\mathcal{M}_\ell}(g(\bar{X}(T)); \widetilde{M}_0) - \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T))) \right|^2}{\delta^2 \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T)))^2} \right] \\ &\leq \frac{\text{Var} \left( \mathcal{V}(g(\bar{X}_0(T)); \widetilde{M}_0) \right) + \sum_{\ell=1}^L \text{Var} \left( \mathcal{V}(\Delta_\ell g(\bar{X}(T)); \widetilde{M}_\ell) \right)}{\delta^2 \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T)))^2} \\ &\leq \frac{CL}{\delta^2 \text{Var}_{\mathcal{M}_\ell}(g(\bar{X}(T)))^2 \widetilde{M}_L}, \end{aligned}$$

and for  $\ell = 1, 2, \dots$  that

$$\begin{aligned}
& P(M_0 = 2^{\ell+1}\widetilde{M}_0) \\
& \leq P\left(\mathcal{V}_{\mathcal{M}_C}(g(\overline{X}(T)); 2^\ell\widetilde{M}_0) - \text{Var}_{\mathcal{M}_C}(g(\overline{X}(T))) > 2^{\ell-1}\text{Var}_{\mathcal{M}_C}(g(\overline{X}(T)))\right) \\
& \leq \mathbb{E}\left[\frac{\left|\mathcal{V}_{\mathcal{M}_C}(g(\overline{X}(T)); 2^\ell\widetilde{M}_0) - \text{Var}_{\mathcal{M}_C}(g(\overline{X}(T)))\right|^2}{2^{2(\ell-1)}\text{Var}_{\mathcal{M}_C}(g(\overline{X}(T)))^2}\right] \\
& < \frac{CL}{2^{3\ell}\text{Var}_{\mathcal{M}_C}(g(\overline{X}(T)))^2\widetilde{M}_L}.
\end{aligned}$$

Consequently,

$$\begin{aligned}
\frac{\mathbb{E}[M_0]\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2} & \leq \left[P(M_0 \leq \widetilde{M}_0) + \sum_{\ell=1}^{\infty} 2^\ell P(M_0 = 2^\ell\widetilde{M}_0)\right] \frac{\widetilde{M}_0\text{TOL}_S^2}{\text{Var}(g(\overline{X}(T)))C_C^2} \\
& \leq 2(1+\delta) \left[P(M_0 \leq \widetilde{M}_0) + P(M_0 = 2\widetilde{M}_0) + \sum_{\ell=1}^{\infty} 2^{\ell+1} P(M_0 = 2^{\ell+1}\widetilde{M}_0)\right] \\
& \leq 2(1+\delta) \left[P(M_0 \leq \widetilde{M}_0) + \frac{CL}{\delta^2\widetilde{M}_L} + \frac{CL}{\widetilde{M}_L} \sum_{\ell=1}^{\infty} 2^{-2\ell}\right].
\end{aligned}$$

Taking limits in the above inequality leads to

$$\limsup_{\text{TOL} \downarrow 0} \frac{\mathbb{E}[M_0]\text{TOL}_S^2}{\text{Var}_{\mathcal{M}_C}(g(\overline{X}(T)))C_C^2} \leq 2(1+\delta).$$

Finally, observe that since the obtained inequality holds true for any  $\delta > 0$ , the proof is finished.  $\square$

An asymptotic bound on  $\mathbb{E}[N_\ell]$  may be deduced from the single level result of Lemma 6. For the convenience of the reader we present the result of Lemma 6 in a way that is fitting for the multilevel setting.

**Lemma 10** (Multilevel asymptotical average number of time steps). *Suppose that the assumptions of Lemma 1 hold. Then the final number of adaptive steps generated by the MLMC algorithm with stochastic time steps (2.20) and (2.21) and  $\text{TOL}_{T,\ell} = 2^{-\ell}\text{TOL}_{T,0}$  satisfies*

$$\limsup_{\ell \uparrow \infty} \text{TOL}_{T,\ell} \mathbb{E}[N_\ell] \leq \frac{4}{C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(t)|} dt \right] \right)^2. \quad (4.38)$$

With bounds for  $\mathbb{E}[M_0]$  and  $\mathbb{E}[N_\ell]$  in hand, we are ready to prove the main complexity theorem for the MLMC algorithm, Theorem 3. The proof is divided into three cases.

*Proof of case (I).* Lemma 10 implies that for any given  $\delta > 0$ , there exists an  $\hat{L}(\delta)$  not depending on TOL such that

$$\text{TOL}_{\text{T},\ell} \mathbb{E}[N_\ell] \leq (1 + \delta) \frac{4}{C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2, \quad \forall \ell \geq \hat{L}. \quad (4.39)$$

Furthermore, for  $\rho_{low}(\text{TOL}_{\text{T}}) = \rho_{\min}$ ,  $M_\ell$  as defined in (4.26) fulfills

$$\mathbb{E}[M_\ell] = 2^{-\ell} \mathbb{E}[M_0], \quad \forall \ell \in \{0, 1, \dots, L\}. \quad (4.40)$$

By inequality (4.39), equation (4.40), and the monotonic relation  $N_\ell \leq N_{\ell+1}$ , we derive that

$$\begin{aligned} \sum_{\ell=0}^L \mathbb{E}[M_\ell] \mathbb{E}[N_\ell] &\leq \mathbb{E}[N_{\hat{L}}] \text{TOL}_{\text{T},\hat{L}} \sum_{\ell=0}^{\hat{L}} \frac{\mathbb{E}[M_\ell]}{\text{TOL}_{\text{T},\hat{L}}} + \sum_{\ell=\hat{L}+1}^L \frac{\mathbb{E}[M_\ell]}{\text{TOL}_{\text{T},\ell}} \mathbb{E}[N_\ell] \text{TOL}_{\text{T},\ell} \\ &\leq \frac{4(1+\delta)\mathbb{E}[M_0]}{C_R \text{TOL}_{\text{T},0}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2 \left( 2^{\hat{L}} \sum_{\ell=0}^{\hat{L}-1} 2^{-\ell} + \sum_{\ell=\hat{L}}^L 1 \right) \\ &\leq \frac{4(1+\delta)\mathbb{E}[M_0]}{C_R \text{TOL}_{\text{T},0}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2 (2^{\hat{L}+1} + (L - \hat{L})). \end{aligned} \quad (4.41)$$

Recalling the definition  $L = \lfloor \log_2(\text{TOL}_{\text{T},\text{Max}}/\text{TOL}_{\text{T}}) \rfloor$  and that  $\hat{L}$  is fixed, it follows that

$$\lim_{\text{TOL} \downarrow 0} \frac{2^{\hat{L}+1} + (L - \hat{L})}{L} = 1.$$

Using (4.41) combined with Lemma 9 and recalling that  $\text{TOL}_{\text{T},0} > \text{TOL}_{\text{T},\text{Max}}/2$ , we obtain the bound

$$\limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}_{\text{S}}^2}{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T))) L} \leq \frac{16(1+\delta)\mathbb{E}[M_0]}{C_R \text{TOL}_{\text{T},\text{Max}}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2. \quad (4.42)$$

We observe that  $\text{WORK}(\text{TOL}) = \mathcal{O}(\text{TOL}_{\text{S}}^{-2} L \text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T))))$ , and to obtain a bound on more explicit form, we note that by the assumption (4.7) on  $L^p$  convergence, there exists a  $C_G > 0$  such that

$$\limsup_{\ell \uparrow \infty} \frac{\rho_{low}(\text{TOL}_{\text{T},\ell})}{\text{TOL}_{\text{T},\ell}} \mathbb{E} \left[ |\Delta_\ell g(\bar{X}(T))|^2 \right] \leq C_G. \quad (4.43)$$

See for instance Remark 4 for a discussion on how to estimate  $C_G$ . Inequality (4.43) further implies that

$$\limsup_{\text{TOL} \downarrow 0} \frac{\text{Var}_{\mathcal{M}\mathcal{L}}(g(\bar{X}(T)))}{L} \leq C_G. \quad (4.44)$$

To approximately minimize the complexity we introduce the splitting choice

$$\begin{aligned}\text{TOL}_S &= \frac{\log(\text{TOL}^{-1})}{\log(\text{TOL}^{-1}) + \log(\log(\text{TOL}^{-1}))} \text{TOL}, \quad \text{and} \\ \text{TOL}_T &= \frac{\log(\log(\text{TOL}^{-1}))}{(\log(\text{TOL}^{-1}) + \log(\log(\text{TOL}^{-1})))C_S} \text{TOL},\end{aligned}$$

which we see fulfills the restrictions  $C_S \text{TOL}_T + \text{TOL}_S = \text{TOL}$  and  $\text{TOL}_T \leq \text{TOL}_S$ . Combining (4.43) with the above splitting choice in inequality (4.42), and noting that this bounding procedure is valid for any  $\delta > 0$  leads to (4.4).  $\square$

*Proof of case (II).* Recall that  $M_\ell$  as defined in (4.26) fulfills

$$\mathbb{E}[M_\ell] \leq (2^{\ell(\bar{\gamma}-1)} + 2^{-L})\mathbb{E}[M_0], \quad \forall \ell \in \{0, 1, \dots, L\}.$$

By this property, inequality (4.39) and recalling that by construction  $\text{TOL}_{T,0} > \text{TOL}_{T,\text{Max}}/2$ , we derive the following

$$\begin{aligned}\sum_{\ell=0}^L \mathbb{E}[M_\ell] \mathbb{E}[N_\ell] &\leq \mathbb{E}[N_{\hat{L}}] \text{TOL}_{T,\hat{L}} \sum_{\ell=0}^{\hat{L}} \frac{\mathbb{E}[M_\ell]}{\text{TOL}_{T,\hat{L}}} + \sum_{\ell=\hat{L}+1}^L \frac{\mathbb{E}[M_\ell]}{\text{TOL}_{T,\ell}} \mathbb{E}[N_\ell] \text{TOL}_{T,\ell} \\ &\leq \frac{(1+\delta)4\mathbb{E}[M_0]}{C_R \text{TOL}_{T,0}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2 \\ &\quad \times \left( 2^{\hat{L}} \sum_{\ell=0}^{\hat{L}-1} (2^{\ell(\bar{\gamma}-1)} + 2^{-L}) + \sum_{\ell=\hat{L}}^L (2^{\ell\bar{\gamma}} + 2^{-L+\ell}) \right) \\ &\leq \frac{(1+\delta)8\mathbb{E}[M_0]}{C_R \text{TOL}_{T,\text{Max}}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2 \left( \frac{2^{\hat{L}}}{1-2^{\bar{\gamma}-1}} + \hat{L}2^{\hat{L}-L} + \frac{2^{(L+1)\bar{\gamma}}}{\log(2^{\bar{\gamma}})} + 2 \right).\end{aligned}\tag{4.45}$$

By noting that  $\hat{L}(\delta)$  is fixed, we have

$$\lim_{\text{TOL} \downarrow 0} \frac{\frac{2^{\hat{L}}}{1-2^{\bar{\gamma}-1}} + \hat{L}2^{\hat{L}-L} + \frac{2^{(L+1)\bar{\gamma}}}{\log(2^{\bar{\gamma}})} + 2}{2^{\bar{\gamma}L}} = \frac{2^{\bar{\gamma}}}{\bar{\gamma} \log(2)},$$

and by Lemma 9 and (4.44), we derive the bound

$$\begin{aligned}\limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}_S^2}{2^{\bar{\gamma}L} L} \\ \leq (1+\delta) \frac{2^{4+\bar{\gamma}} C_C^2 C_G}{\bar{\gamma} \log(2) \text{TOL}_{T,\text{Max}} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2.\end{aligned}$$

Since  $2^{-L} \geq \text{TOL}_T / \text{TOL}_{T,\text{Max}}$  by construction, we further obtain

$$\begin{aligned} \limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}_S^2 \text{TOL}_T^{\bar{\gamma}}}{L} \\ \leq (1 + \delta) \frac{2^{4+\bar{\gamma}} C_C^2 C_G}{\bar{\gamma} \log(2) (\text{TOL}_{T,\text{Max}})^{1-\bar{\gamma}} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2. \end{aligned} \quad (4.46)$$

To approximately minimize the complexity, we introduce the splitting choice

$$\text{TOL}_S = \frac{2}{2 + \bar{\gamma}} \text{TOL}, \quad \text{and}, \quad \text{TOL}_T = \frac{\bar{\gamma}}{(2 + \bar{\gamma}) C_S} \text{TOL},$$

which fulfills the constraints  $C_S \text{TOL}_T + \text{TOL}_S = \text{TOL}$  and  $\text{TOL}_T \leq \text{TOL}_S$ . Using the splitting choice in inequality (4.46) leads to (4.4) when noting that the bounding procedure is valid for any  $\delta > 0$ .  $\square$

*Proof of case (III).* First, we note that the conditions  $\bar{\gamma} \rightarrow 0$  and  $L\bar{\gamma} \rightarrow \infty$  as  $\text{TOL} \downarrow 0$  yields a consistent lower error density, since it leads to

$$\rho_{low}(\text{TOL}_T) = \text{TOL}_T^{\bar{\gamma}} = \mathcal{O}(2^{-L\bar{\gamma}}),$$

which implies that  $\rho_{low}(\text{TOL}_T) \rightarrow 0$  as  $\text{TOL} \downarrow 0$ .

For proving case (III), recall inequality (4.45) from the proof of case (II) which is valid in the present setting as well, namely

$$\begin{aligned} \sum_{\ell=0}^L \mathbb{E}[M_\ell] \mathbb{E}[N_\ell] &\leq \frac{(1 + \delta) 8 \mathbb{E}[M_0]}{C_R \text{TOL}_{T,\text{Max}}} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2 \\ &\times \left( \frac{2^{\hat{L}}}{1 - 2^{\bar{\gamma}-1}} + \hat{L} 2^{\hat{L}-L} + \frac{2^{(L+1)\bar{\gamma}}}{\log(2^{\bar{\gamma}})} + 2 \right). \end{aligned}$$

The asymptotic conditions on  $\bar{\gamma}$  imply that

$$\lim_{\text{TOL} \downarrow 0} \frac{\bar{\gamma}}{2^{\bar{\gamma}L}} \left( \frac{2^{\hat{L}}}{1 - 2^{\bar{\gamma}-1}} + \hat{L} 2^{\hat{L}-L} + \frac{2^{(L+1)\bar{\gamma}}}{\log(2^{\bar{\gamma}})} + 2 \right) = \frac{1}{\log(2)}.$$

This property, Lemma 9 and (4.44), yield

$$\begin{aligned} \limsup_{\text{TOL} \downarrow 0} \frac{\text{WORK}(\text{TOL}) \text{TOL}_S^2 \bar{\gamma}}{L 2^{\bar{\gamma}L}} \\ \leq (1 + \delta) \frac{16 C_C^2 C_G}{\log(2) \text{TOL}_{T,\text{Max}} C_R} \left( \mathbb{E} \left[ \int_0^T \sqrt{|\hat{\rho}(\tau)|} d\tau \right] \right)^2. \end{aligned} \quad (4.47)$$

To approximately minimize the complexity, we introduce the splitting choice

$$\text{TOL}_S = \frac{2}{2 + \bar{\gamma}(\text{TOL})} \text{TOL} \quad \text{and} \quad \text{TOL}_T = \frac{\bar{\gamma}(\text{TOL})}{(2 + \bar{\gamma}(\text{TOL})) C_S} \text{TOL}.$$

Inserting the splitting into (4.47) and noting that this argument is valid for all  $\delta > 0$  leads to (4.6).  $\square$

**Remark 4** (Particular estimate for the constant  $C_G$ ). *It is possible to estimate the constant  $C_G$ . For instance, in the particular case when the exact error density is bounded away from zero, i.e. there exist a constant  $\rho_{min}$  such that  $\hat{\rho} > \rho_{min} > 0$  a.s. and, considering the equation*

$$\begin{aligned} dX(t) &= b(X(t))dW(t), \quad t > 0 \\ X(0) &= X_0, \end{aligned}$$

we have

$$C_G \leq C_S \mathbb{E} \left[ \left\| \frac{(b'b)^2(X(t))(\varphi)^2(t)}{\hat{\rho}(t)} \right\|_{L^\infty([0,T])} \right].$$

Here  $\varphi(t) = g'(X(T)) \frac{X'(T)}{X'(t)}$  and the first variation  $X'(s)$  solves, for  $s > 0$ , the linear equation

$$dX'(s) = b'(X(s))X'(s)dW(s),$$

with initial condition  $X'(0) = 1$ . The constant  $C_S$  is the parameter in the stopping condition (2.20).

**Remark 5** (Jump Diffusions). *It is possible to extend these results of adaptive multilevel weak approximation for diffusions to the case of jump diffusions with time dependent jump measure analyzed in [22].*

## 5 Conclusions

In this paper we have presented and analyzed an adaptive multilevel Monte Carlo algorithm, where the multilevel simulations are performed on adaptively generated mesh hierarchies based on computable a posteriori weak error estimates. The theoretical analysis of the adaptive algorithm showed that the algorithm stops after a finite number of steps, and proceeded to show accuracy and efficiency results under natural assumptions in Theorems 2 and 3. In particular, Theorem 2 states that the probability of the weak error being bounded by the specified tolerance TOL is asymptotically bounded by any desired probability through the confidence parameter. Theorem 3 states computational complexity results where the involved constants are explicitly given in terms of algorithm parameters and problem properties. It shows that the  $L^{1/2}$ -quasi norm of the error density appears as a multiplicative constant in the complexity bounds, instead of the larger  $L^1$ -norm of the same error density that would appear using a uniform time step MLMC algorithm; the difference between these two factors can be arbitrarily large even in problems with smooth coefficients where they are both finite. Disregarding the constants the result shows that, depending on assumptions on the limit error density and the lower bound on the computed error density used by the adaptive algorithm, the complexity can be either the



same as or nearly the same as the complexity uniform MLMC has in cases where the order of strong convergence of the Euler-Maruyama method is  $1/2$ .

Numerical results for scalar SDEs confirmed the theoretical analysis. For the two problems with reduced weak convergence order a simple single level Monte Carlo method has complexity  $\mathcal{O}(\text{TOL}^{-4})$  while the adaptive MLMC method has the improved complexity  $\mathcal{O}(\text{TOL}^{-2} \log_2(\text{TOL}_0/\text{TOL})^2)$ . The use of advanced Monte Carlo methods such as the adaptive MLMC algorithm presented in this paper is most attractive for SDEs in higher dimension, where the corresponding PDE-based computational techniques are not competitive. Numerical experiments using adaptive MLMC is ongoing work and will be presented in a later report; this also makes for an interesting comparison to uniform MLMC method for barrier problems, since it is not clear that the order of strong convergence of the Euler-Maruyama method will be  $(1 - \delta)/2$ , for any positive  $\delta$ , in that case. The fact that computational complexity of uniform multilevel Monte Carlo, disregarding constants, depends on the strong convergence indicates that adaptive mesh refinements based on strong error estimates can also be used to improve the computational efficiency; such methods are also subjects of ongoing research.

In this paper the adaptive algorithms were presented with global error control in the quantity of interest, starting from a given coarse mesh. Alternatively we can use local error estimates to guide adaptive time step control in the computation of the forward problem. This approach can be used on its own when global error control is deemed unnecessary or too computationally expensive, but it can also be used together with the global error control in situations with stiff SDEs where any given initial mesh can be too coarse depending on the realization. This is particularly relevant for MLMC simulations where stability issues in the computations on the coarsest level can destroy the results of the whole multilevel simulation, as was pointed out by Hutzenthaler, Jentzen, and Kloeden in [13].

## A Theorems

**Theorem 6** (Lindeberg-Feller Theorem [4][p. 114]). *For each  $n$ , let  $X_{n,m}$ ,  $1 \leq m \leq n$ , be independent random variables with  $E[X_{n,m}] = 0$ . Suppose*

(a)

$$\sum_{m=1}^n E[X_{n,m}^2] \rightarrow \sigma^2 > 0, \text{ and}$$

(b) for all  $\delta > 0$ ,

$$\lim_{n \rightarrow \infty} \sum_{m=1}^n E[X_{n,m}^2 1_{|X_{n,m}| > \delta}] = 0.$$

*Then the Central Limit Theorem holds, i.e., the random variable*

$$S_n := \sum_{m=1}^n X_{n,m} \rightarrow \sigma \Xi, \text{ as } n \rightarrow \infty,$$

where  $\Xi$  is a standard normal distributed random variable.

#### Acknowledgments.

This work was partially supported by the Dahlquist fellowship at the Royal Institute of Technology in Stockholm, Sweden, the department of Scientific Computing in Florida State University, the University of Austin Subcontract (Project Number 024550, Center for Predictive Computational Science), the VR project "Effektiva numeriska metoder för stokastiska differentialekvationer med tillämpningar", and King Abdullah University of Science and Technology (KAUST). The authors would like to thank Mike Giles, and also two anonymous reviewers, for useful comments.

## References

- [1] F. Kh. Abdullaev, J. C. Bronski, and G. Papanicolaou. Soliton perturbations and the random Kepler problem. *Phys. D*, 135(3-4):369–386, 2000.
- [2] Eric Cancès, Frédéric Legoll, and Gabriel Stoltz. Theoretical and numerical comparison of some sampling methods for molecular dynamics. *M2AN Math. Model. Numer. Anal.*, 41(2):351–389, 2007.
- [3] Y. S. Chow and Herbert Robbins. On the asymptotic theory of fixed-width sequential confidence intervals for the mean. *The Annals of Mathematical Statistics*, 36(2):pp. 457–462, 1965.
- [4] Richard Durrett. *Probability: theory and examples*. Duxbury Press, Belmont, CA, second edition, 1996.
- [5] Anna Dzougoutov, Kyoung-Sook Moon, Erik von Schwerin, Anders Szepessy, and Raúl Tempone. Adaptive Monte Carlo algorithms for stopped diffusion. In *Multiscale methods in science and engineering*, volume 44 of *Lect. Notes Comput. Sci. Eng.*, pages 59–88. Springer, Berlin, 2005.
- [6] Michael B. Giles. Improved multilevel monte carlo convergence using the milstein scheme. *Monte Carlo and Quasi-Monte Carlo Methods 2006*, pages 343–358.
- [7] Michael B. Giles. Multilevel Monte Carlo path simulation. *Oper. Res.*, 56(3):607–617, 2008.
- [8] Michael B. Giles, Desmond J. Higham, and Xuerong Mao. Analysing multilevel Monte Carlo for options with non-globally Lipschitz payoff. *Finance and Stochastics*, 13:403–413, 2009.
- [9] Paul Glasserman. *Monte Carlo methods in financial engineering*, volume 53 of *Applications of Mathematics (New York)*. Springer-Verlag, New York, 2004. Stochastic Modelling and Applied Probability.

- [10] Emmanuel Gobet. Weak approximation of killed diffusion using Euler schemes. *Stochastic Process. Appl.*, 87(2):167–197, 2000.
- [11] S. Heinrich. Monte Carlo complexity of global solution of integral equations. *J. Complexity*, 14(2):151–175, 1998.
- [12] S. Heinrich and E. Sindambiwe. Monte Carlo complexity of parametric integration. *J. Complexity*, 15(3):317–341, 1999. Dagstuhl Seminar on Algorithms and Complexity for Continuous Problems (1998).
- [13] Martin Hutzenthaler, Arnulf Jentzen, and Peter E. Kloeden. Divergence of the multilevel monte carlo method. *arXiv/1105.0226v1*, 2011.
- [14] E. Jouini, J. Cvitanić, and Marek Musiela, editors. *Option pricing, interest rates and risk management*. Handbooks in Mathematical Finance. Cambridge University Press, Cambridge, 2001.
- [15] Ioannis Karatzas and Steven E. Shreve. *Brownian motion and stochastic calculus*, volume 113 of *Graduate Texts in Mathematics*. Springer-Verlag, New York, second edition, 1991.
- [16] Ahmed Kebaier. Statistical Romberg extrapolation: A new variance reduction method and applications to option pricing. *Ann. Appl. Probab.*, 15(4):2681–2705, 2005.
- [17] E. S. Keeping. *Introduction to statistical inference*. D. Van Nostrand Co., Inc., Princeton, N.J.-Toronto-London-New York, 1962.
- [18] Andrew J. Majda, Ilya Timofeyev, and Eric Vanden Eijnden. A mathematical framework for stochastic climate models. *Comm. Pure Appl. Math.*, 54(8):891–974, 2001.
- [19] Kyoung-Sook Moon, Anders Szepessy, Raúl Tempone, and Georgios E. Zouraris. Convergence rates for adaptive weak approximation of stochastic differential equations. *Stoch. Anal. Appl.*, 23(3):511–558, 2005.
- [20] Kyoung-Sook Moon, Erik von Schwerin, Anders Szepessy, and Raúl Tempone. An adaptive algorithm for ordinary, stochastic and partial differential equations. In *Recent advances in adaptive computation*, volume 383 of *Contemp. Math.*, pages 325–343. Amer. Math. Soc., Providence, RI, 2005.
- [21] E. Mordecki, A. Szepessy, R. Tempone, and G. E. Zouraris. Adaptive weak approximation of diffusions with jumps. *SIAM J. Numer. Anal.*, 46(4):1732–1768, 2008.
- [22] E. Mordecki, A. Szepessy, R. Tempone, and G. E. Zouraris. Adaptive weak approximation of diffusions with jumps. *SIAM J. Numer. Anal.*, 46(4):1732–1768, 2008.
- [23] Bernt Øksendal. *Stochastic differential equations*. Universitext. Springer-Verlag, Berlin, fifth edition, 1998. An introduction with applications.

- [24] Anders Szepessy, Raúl Tempone, and Georgios E. Zouraris. Adaptive weak approximation of stochastic differential equations. *Comm. Pure Appl. Math.*, 54(10):1169–1214, 2001.

## Paper III

# ON NON-ASYMPTOTIC OPTIMAL STOPPING CRITERIA IN MONTE CARLO SIMULATIONS

CHRISTIAN BAYER, HÅKON HOEL, ERIK VON SCHWERIN, AND RAÚL TEMPONE

ABSTRACT. We consider the setting of estimating the mean of a random variable by a sequential stopping rule Monte Carlo (MC) method. The performance of a typical second moment based sequential stopping rule MC method is shown to be unreliable in such settings both by numerical examples and through analysis. By analysis and approximations, we construct a higher moment based stopping rule which is shown in numerical examples to perform more reliably and only slightly less efficiently than the second moment based stopping rule.

## 1. INTRODUCTION

Given i.i.d. random variables  $X_1, X_2, \dots$  the typical way of approximating their expected value  $\mu = E[X]$  using  $M$  samples is the sample average

$$\bar{X}_M := \sum_{i=1}^M \frac{X_i}{M}.$$

We consider the objective of choosing  $M$  sufficiently large so that the error probability satisfies

$$P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta, \quad (1)$$

for some fixed small constants  $\text{TOL} > 0$  and  $\delta > 0$ . Clearly,  $P(|\bar{X}_M - \mu| > \text{TOL})$  decreases as  $M$  increases, but at the same time the cost of computing  $\bar{X}_M$  increases. From an application and cost point of view it is therefore of interest to derive theory or algorithms that will give upper bounds on  $M$  satisfying (1) that are not far too large. When a-priori information about the distribution of  $X$  is available, for example if  $X$  is a bounded r.v. with an explicitly given bound, it is possible to derive good theoretical upper bounds for  $M$  using Hoeffding type inequalities, cf. Hoeffding [5]. In the general case when no or little information of the distribution is given, little theory is however known, and the typical way of estimating  $E[X]$  using a sufficiently large number of samples  $M$  is through a sequential stopping rule. Below we give the general structure of the class of sequential stopping rules we have in mind.

- (I) Generate a batch of  $M$  i.i.d. samples  $X_1, X_2, \dots, X_M$ .
- (II) Infer distributive properties of  $\bar{X}_M$  from the generated batch of samples through higher order sample moments, e.g., the sample mean and the sample variance.
- (III) Based on the sample moments, estimate the error probability. When, based on the estimated probability, (1) is violated, increase the number of samples  $M$  and return to step (I).  
Else, break and accept  $M$ .

---

2010 *Mathematics Subject Classification.* Primary 65C05; Secondary 62L12, 62L15.

*Key words and phrases.* Monte Carlo methods, optimal stopping, sequential stopping rules, non-asymptotic.

---

**Algorithm 1** Sample Variance Based Stopping Rule

---

**Input:** Number of samples  $M_0$ , accuracy TOL, confidence  $\delta$ , the cumulative distribution function of the standard normal distributed r.v.  $\Phi(x)$ .

**Output:**  $\bar{X}_{M_{\bar{n}}}$ .

Set  $k = 0$ , generate  $M_k$  samples  $\{X_i\}_{i=1}^{M_k}$  and compute the sample variance

$$\bar{\sigma}_{M_k}^2 := \frac{1}{M_k - 1} \sum_{i=1}^{M_k} (X_i - \bar{X}_{M_k})^2. \quad (2)$$

**while**  $2\left(1 - \Phi(\sqrt{M_k} \text{TOL} / \bar{\sigma}_{M_k})\right) > \delta$  **do**

Set  $k = k + 1$  and  $M_k = 2M_{k-1}$ .

Generate a batch of  $M_k$  i.i.d. samples  $\{X_i\}_{i=1}^{M_k}$ .

Compute the sample variance  $\bar{\sigma}_{M_k}^2$  as given in (2).

**end while**

Set  $M_{\bar{n}} = M_k$ , generate samples  $\{X_i\}_{i=1}^{M_{\bar{n}}}$  and compute the output sample mean  $\bar{X}_{M_{\bar{n}}}$ . (See Section 2 for a motivation of the choice of the stopping criterion in the while loop above.)

---

Certainly the most natural and important representative of this class of algorithms is given in Algorithm 1. The algorithm estimates the error probability by appealing to the Central Limit Theorem (CLT). Consequently, it only relies on the sample variance in addition to the sample mean. In particular, the algorithm only requires mild additional assumptions on  $X$ , namely square integrability.

In the literature, various second moment based sequential stopping rules have been introduced to estimate the steady-state mean of stochastic processes, see for example Law and Kelton [7, 8] for comparisons of the performance of different stopping rules and Bratley, Bennet, and Fox [1] for an overview. Second moment based sequential stopping rules generally tend to perform well in the asymptotic regime when  $\text{TOL} \rightarrow 0$ . In fact, Chow and Robbins [2] proved that under very loose restrictions, second moment based sequential stopping rules such as Algorithm 1 are asymptotically consistent, meaning that for a fixed  $\delta$ ,

$$\lim_{\text{TOL} \rightarrow 0} \text{P}(|\bar{X}_M - \mu| > \text{TOL}) = \delta,$$

and in Glynn and Whitt [4] the consistency property is proven to hold for such stopping rules applied to more general stochastic processes. The performance for second moment based stopping rules in the non-asymptotic regime – when TOL and  $\delta$  are fixed values – is however not as well understood. This is unsatisfactory, as in applications this is precisely the interesting regime, in particular since very often we have  $\text{TOL} \gg \delta$ . While consistency is clearly a re-assuring property in any case, in many situations one is in dire need of quantitative estimates of the error probability in the non-asymptotic regime, for instance when one tries to optimize the computational cost needed to meet a certain accuracy target using an adaptive algorithm. We could not find such a quantitative, non-asymptotic analysis of algorithms like Algorithm 1 in the literature.

In this note we demonstrate by numerical examples that second moment based stopping rules can fail convincingly in the non-asymptotic regime, especially when the underlying distribution  $X$  is heavy-tailed, see Section 2. We proceed by giving an error analysis of Algorithm 1 specifically in the non-asymptotic regime. We note a-priori that there are two obvious approximation errors in the underlying assumptions of Algorithm 1:

(I) The algorithm appeals to the CLT to approximate the tail probabilities for  $\overline{X}_M$  even though  $M$  is finite.

(II) In doing so, it uses the sample variance  $\overline{\sigma}_M^2$  instead of the true variance  $\sigma^2$ .

To get a hold on the error probability (1) despite the fact that the distribution of the sample mean  $\overline{X}_M$  is unknown, we again appeal to the central limit theorem, but we adjust the estimate by adding a Berry-Esseen type term, which extends the validity of the estimate to the non-asymptotic case, thereby dealing with the first approximation error. As the error probability (1) is a tail probability for the distribution of the sample mean and the Berry-Esseen theorem itself is rather aimed at being sharp at the center of the distribution, we appeal to non-uniform versions of the Berry-Esseen theorem, see Theorem 1.1 and Corollary 1.2 below. However, both intuition and numerical tests suggest that the approximation of the tail probabilities by the non-uniform Berry-Esseen theorem is far too pessimistic at least when the second approximation error is small, i.e., when the computed sample variance is actually close to the true variance. In this case, we adjust the normal distribution by a less stringent extra term, which is obtained from an Edgeworth expansion of the distribution function of the sample mean  $\overline{X}_M$ , c.f. Feller [3].<sup>1</sup>

Having identified possible origins of failure of Algorithm 1, we propose an improvement of Algorithm 1. However, this variant requires third and fourth sample moments, see Section 4. Finally, in Section 5, we test the new algorithm numerically. We find that the new stopping Algorithm 2 indeed satisfies the desired confidence level  $\delta$  on the error probability (1) even when  $\delta \ll \text{TOL}$ .

As already discussed above, we need to approximate the unknown distribution of a sample mean in a general, non-asymptotic regime. The uniform and non-uniform Berry-Esseen theorems provide quantitative bounds for the difference between the true distribution of the sample mean and its asymptotic limit, namely the normal distribution. The following classical theorem can be found, for instance, in Petrov [10].

**Theorem 1.1** (Uniform and Non-Uniform Berry-Esseen). *Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with  $E[X] = 0$ ,  $\sigma^2 = \text{Var}(X)$  and  $\beta = \frac{E[|X|^3]}{\sigma^3} < \infty$ . Then, for a positive constant  $C_0$ , the following uniform bound*

$$\left| P\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq C_0 \frac{\beta}{\sqrt{n}}$$

*holds. For another positive constant  $C_1$ , the following non-uniform bound holds*

$$\left| P\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq C_1 \frac{\beta}{\sqrt{n}(1+|x|^3)}.$$

Up to our knowledge, the best upper bounds presently known for the Berry-Esseen constants are  $C_0 < 0.4785$ , cf. Tyurin [11], and  $C_1 < C_0 + 8(1+e^1) < 30.2338$ ,

---

<sup>1</sup>Note that here we are introducing a gap in the analysis: the estimate based on the non-uniform Berry-Esseen theorem is reliable in the sense that it always leads to an upper bound of the error probability (1). For the Edgeworth expansion, however, there might be situations when the true error probability is underestimated, and, consequently, the accuracy target might still be missed. Numerical evidence, however, suggests that the estimate obtained from solely relying on the non-uniform Berry-Esseen theorem is usually by orders of magnitude too pessimistic. Apart from intrinsic reasons, one reason might be that the constants known in the non-uniform Berry-Esseen theorems might be far from being optimal. In the end, we think that the above compromise between Berry-Esseen type estimations and estimations based on the Edgeworth expansion might be a good compromise which retains the goal of reliably meeting the accuracy target – except maybe for very extreme situations – while keeping a certain level of efficiency. We note, however, that it is also possible to construct even more conservative stopping rules which are only based on the Berry-Esseen theorem.



cf. Michel [9]. For the purpose of this paper, it will be useful to combine the uniform and non-uniform Berry-Esseen bound as follows.

**Corollary 1.2** (Berry-Esseen). *Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with  $E[X] = 0$ ,  $\sigma^2 = \text{Var}(X)$  and  $\beta = E[|X|^3] / \sigma^3 < \infty$ . Then*

$$\left| P\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) - \Phi(x) \right| \leq C_{\text{BE}}(x) \frac{\beta}{\sqrt{n}}$$

where the bound function  $C_{\text{BE}} : \mathbb{R} \rightarrow [0, C_0]$  is defined by

$$C_{\text{BE}}(x) := \min\left(C_0, \frac{C_1}{(1 + |x|)^3}\right).$$

In the asymptotic regime, the distribution of  $P\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right)$  can be expressed by so called Edgeworth expansions. Here we present the one-term Edgeworth expansion.

**Theorem 1.3** (Edgeworth expansion, cf. Feller [3]). *Suppose  $X_1, X_2, \dots$  are i.i.d. r.v. with a distribution which is not a lattice distribution and  $E[X] = 0$ ,  $\sigma^2 = \text{Var}(X)$  and  $E[X^3] < \infty$ . Then*

$$P\left(\frac{1}{\sigma\sqrt{n}} \sum_{i=1}^n X_i < x\right) = \Phi(x) + \frac{(x^2 - 1)e^{-x^2/2} E[X^3]}{6\sqrt{2\pi n} \sigma^3} + o(n^{-1/2}),$$

uniformly for  $x \in \mathbb{R}$ .

## 2. STOPPING RULE FAILURES

Suppose we seek to estimate  $\mu = E[X]$  using Monte Carlo simulation and we actually *do know* the variance  $\sigma^2 = \text{Var}(X)$ . As before, our objective is to achieve  $P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta$ , for some fixed, small constants  $\text{TOL}, \delta > 0$ . The CLT motivates the stopping rule

$$M = \frac{C_{\text{CLT}}^2 \sigma^2}{\text{TOL}^2}, \quad C_{\text{CLT}} := \Phi^{-1}\left(\frac{2 - \delta}{2}\right), \quad (3)$$

which would exactly fulfill our objective (1) in the asymptotic regime  $M \rightarrow \infty$ . Of course, this conflicts with our choice (3) for  $M$ , since we treat  $\delta$  and  $\text{TOL}$  as finite constants. However, we can still estimate the probability in (1) using Corollary 1.2 and obtain

$$\begin{aligned} P(|\bar{X}_M - \mu| > \text{TOL}) &= P\left(\sqrt{M} \frac{|\bar{X}_M - \mu|}{\sigma} > \frac{\sqrt{M} \text{TOL}}{\sigma}\right) \\ &\leq 2 \left(1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right)\right) + 2C_{\text{BE}}\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) \frac{\beta}{\sqrt{M}} \quad (4) \\ &= 2(1 - \Phi(C_{\text{CLT}})) + 2C_{\text{BE}}(C_{\text{CLT}}) \frac{\beta}{\sqrt{M}} \text{TOL} \\ &= \delta + 2 \frac{C_{\text{BE}}(C_{\text{CLT}}) \beta}{\sigma C_{\text{CLT}}} \text{TOL} \end{aligned}$$

This means that in the worst case, the actual error probability could be  $\delta + \mathcal{O}(\text{TOL})$  instead of  $\delta^2$ . For instance, in situations where the statistical confidence in the result is more stringent than the accuracy so that  $\delta \ll \text{TOL}$ , the asymptotically motivated choice of  $M$  in (3) could, granted the bound (4) is sharp, fail to deliver

<sup>2</sup>Note that  $C_{\text{CLT}}$  as defined in (3) grows only very slowly as  $\delta$  decreases, since we have  $C_{\text{CLT}} < \sqrt{2 \log(\delta^{-1})}$ . Thus, the factor in front of  $\text{TOL}$  in the error probability can almost be neglected.

the expected level of confidence. For most r.v. however, the bound (4) is far too conservative, and one might ask whether it is reasonable to fear overshooting the error probability in the fashion we have described. The following numerical example shows the existence of r.v. for which the stopping rule (3) fails in the non-asymptotic regime

**Example 2.1.** The heavy-tailed Pareto-distribution has the probability distribution function

$$f(x) = \begin{cases} \alpha x_m^\alpha x^{-(\alpha+1)} & \text{if } x \geq x_m \\ 0 & \text{else,} \end{cases} \quad (5)$$

where  $\alpha, x_m \in \mathbb{R}_+$  are respectively the shape and the scale parameter. The moments of  $E[X^n]$  for the Pareto r.v. only exists for  $n < \alpha$  and, supposing  $\alpha > 2$ , its mean and variance are given by

$$\mu = \frac{\alpha x_m}{\alpha - 1} \text{ and } \sigma^2 = \frac{x_m^2 \alpha}{(\alpha - 1)^2 (\alpha - 2)}.$$

It is further easy to derive that for a Pareto r.v. with  $\alpha = 3 + \gamma$  and  $0 < \gamma < 1$ ,

$$\beta = \frac{E[|X - \mu|^3]}{\sigma^3} = \mathcal{O}(\gamma^{-1}).$$

This implies that there exists r.v. for which the second summand of the bound (4) can become arbitrary large. So for such r.v. the stopping rule (3) might fail. Let us investigate by numerical approximations. Considering the distribution with  $\alpha = 3.1$  (and  $x_m = 1$ ), yields a heavy-tailed r.v. with known mean, variance and third moment. For a set of accuracies  $\text{TOL} \in [0.05, 0.2]$  and confidences  $\delta = \text{TOL}^\ell$ ,  $\ell = 0.5, 1, 1.5$ , and 2 we have numerically approximated  $P(|\bar{X}_M - \mu| > \text{TOL}) \leq \delta$  using, in accordance with (3), the stopping rule

$$M = \left\lceil \frac{C_{\text{CLT}}^2 \sigma^2}{\text{TOL}^2} \right\rceil$$

The results, illustrated in Figure 1, show that when  $\delta \ll \text{TOL}$ , the sought confidence is far from met.

Example 2.1 shows that for some r.v. the confidence goal of (1) will not be met by using the stopping rule (3), at least in settings with  $\delta \ll \text{TOL}$ . Supposing we do not know the variance prior to sampling, yet another type of stopping rule failure is given in Example 2.2; it considers how the MC estimate of Algorithm 1 depends on the initial number of samples  $M_0$ .

**Example 2.2** (Premature Stopping). In this example we will sample the mean of various r.v. using Algorithm 1 and investigate how the MC estimate  $\bar{X}_M$  depends on the initial number of samples  $M_0$ . Let  $M(M_0)$  denote the number of r.v. used in the stopped estimate as a function of the initial number of samples  $M_0$ . Our MC estimate goal then becomes to achieve

$$P(|\bar{X}_{M(M_0)} - \mu| > \text{TOL}) \leq \delta. \quad (6)$$

To investigate whether this goal is fulfilled we plot  $P(|\bar{X}_{M(M_0)} - \mu| > \text{TOL})$  as a function of  $M_0$  for four different r.v. in Figure 2. Figure 2 indicates that the more heavy-tailed or skewed the distribution is, the higher  $M_0$  is needed to ensure that the goal (6) is met.

The demonstrated stopping rule failures motivated us to study and develop ways of constructing more reliable stopping rules. In Section 3, we first analyze the stopping rule of Algorithm 1, and derive an approximate upper bound for the failure probability expressed in terms of  $M$ ,  $\text{TOL}$  and  $\delta$ . In Section 4, we develop

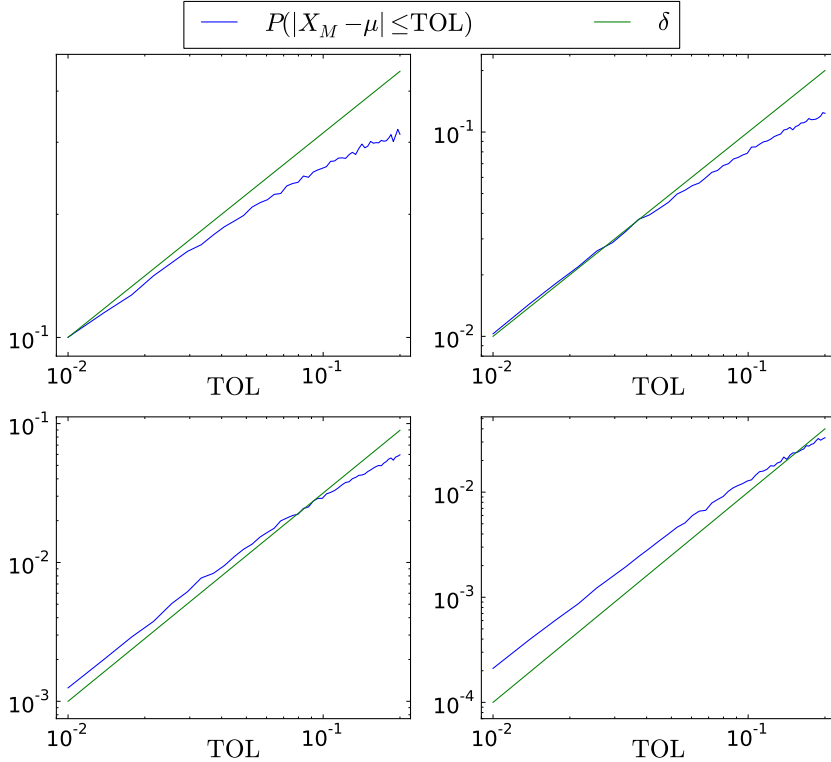


FIGURE 1. MC estimate using the stopping rule (3) for i.i.d. Pareto r.v. with parameters  $\alpha = 3.1$  and  $x_m = 1$ . The obtained failure probability  $P(|\bar{X}_M - \mu| > \text{TOL})$  (blue lines) is plotted in comparison to the sought confidence parameter  $\delta(\text{TOL}) = \text{TOL}^\ell$  (green lines), for  $\ell = 0.5$  (upper left),  $\ell = 1$  (upper right),  $\ell = 1.5$  (lower left), and  $\ell = 2$  (lower right). We observe that the smaller  $\delta$  is relative to  $\text{TOL}$ , the more apparent does the failure of the stopping rule become.

a more reliable stopping rule algorithm, which in addition to second moment of the r.v. in question, also depends on third and fourth order moments. The paper is concluded with numerical examples comparing the reliability and computational cost of Algorithm 1 with the stopping rule developed in Section 4.

### 3. ERROR ANALYSIS FOR ALGORITHM 1

Examples 2.1 and 2.2 illustrate that for some r.v. the stopping rule in Algorithm 1 does not meet the accuracy-confidence constraint (1). To construct a more reliable stopping rule, penalty terms have to be added to the stopping criterion in Algorithm 1. Some care should be taken to make the penalty terms of right size: if too large penalties are added, the new stopping rule will be reliable but very inefficient, while if too small penalty terms are added, the algorithm will of course be efficient but unreliable.

In this section, we first derive an approximate upper bound for the failure probability

$$P\left(|\bar{X}_M - \mu| > \text{TOL} \mid M\right) \quad (7)$$

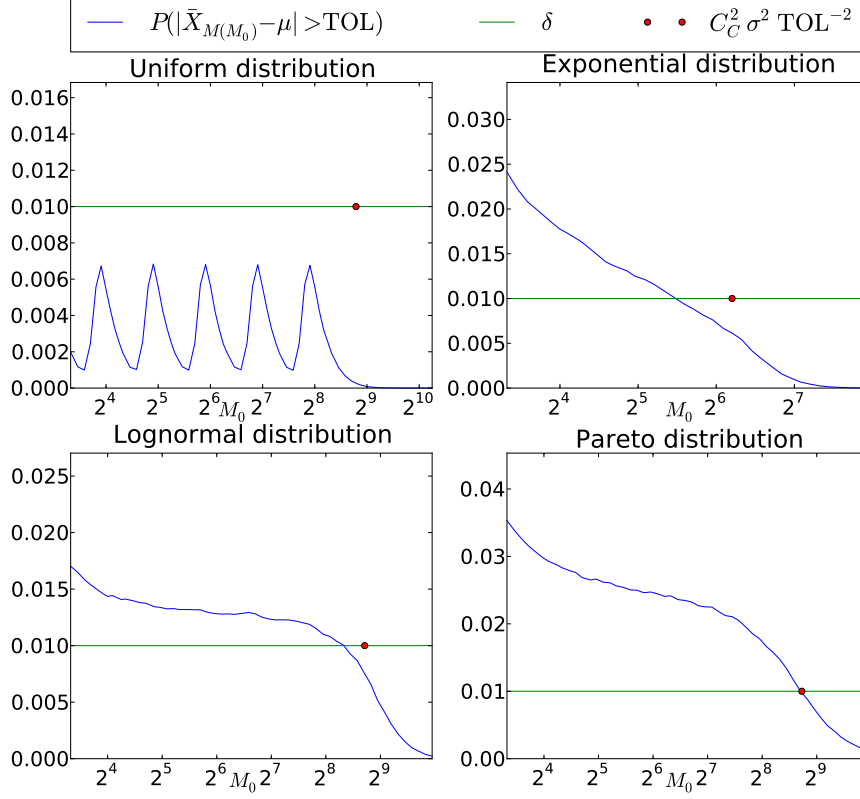


FIGURE 2. Plots of  $P(|\bar{X}_{M(M_0)} - \mu| > \text{TOL})$  as a function of  $M_0$  when using Algorithm 1 to sample  $\bar{X}_{M(M_0)}$ . The accuracy and confidence is set to  $\text{TOL} = 0.1$  and  $\delta = \text{TOL}^2$ , respectively. **Upper left:** The Uniform distribution  $X \sim U(-1, 1)$  with  $\mu = 0$  and  $\sigma^2 = 2/3$  (light-tailed). **Upper right:** The Exponential distribution with  $\mu = 1/3$  and  $\sigma^2 = 1/9$  (not heavy-tailed). **Lower left:** The Lognormal distribution  $X \sim \log(\mathcal{N}(\mu_{\mathcal{N}}, \sigma_{\mathcal{N}}^2))$  with  $\mu_{\mathcal{N}} = -1$  and  $\sigma_{\mathcal{N}}^2 = 1$ . This gives  $\mu = \exp(\mu_{\mathcal{N}} + \sigma_{\mathcal{N}}^2/2)$  and  $\sigma^2 = (\exp(\sigma_{\mathcal{N}}^2) - 1) \exp(2\mu_{\mathcal{N}} + \sigma_{\mathcal{N}}^2)$  (quite heavy-tailed). **Lower right:** The Pareto distribution with  $x_m = 1$  and  $\alpha = 3.1$ , cf. (5). (heavy-tailed).

corresponding to the stopping rule of Algorithm 1 conditional on the (random) final number of samples  $M$ . Clearly, the bound for (7) will also be a r.v. Using the bound for (7), we thereafter construct reasonable penalty terms to be added to the stopping criterion of our new stopping rule.

Let  $\bar{\sigma}_M$  denote the sample variance generated from the stopped sample batch, i.e., the samples used to generate the output MC estimate  $\bar{X}_M$ . Then our first step towards an upper bound for (7) is partitioning the probability (7) into two parts

as follows

$$\begin{aligned} & \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M\right) \\ &= \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2\}\right) \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) \\ & \quad + \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| < \sigma^2/2\}\right) \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| < \sigma^2/2 \mid M\right). \end{aligned} \quad (8)$$

The event  $|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2$  implies the estimate of the real variance is substantially wrong, and then it is likely that we use far too few samples  $M$  to ensure that our MC estimate is reliable. A relatively high penalty term should therefore be added to the stopping criterion to avoid the event  $|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2$  from occurring. To derive such a penalty term, we will first bound the factors of the product

$$\mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2\}\right) \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) \quad (9)$$

separately. For the first factor of this product,

$$\mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2\}\right),$$

we recall that in Algorithm 1 the samples used in the output estimate  $\bar{X}_M$  and for  $\bar{\sigma}_M$  are independent of the samples used to determine  $M$ . Keeping this in mind, we derive the following approximate upper bound

$$\begin{aligned} & \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2\}\right) \\ &= \mathbb{P}\left(|\bar{X}_n - \mu| > \text{TOL} \mid \{|\bar{\sigma}_n^2 - \sigma^2| \geq \sigma^2/2\}\right) \Big|_{n=M} \\ &\lesssim 2 \left( 1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) + C_{\text{BE}}\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) \frac{\beta}{\sqrt{M}} \right). \end{aligned} \quad (10)$$

Here the Berry-Esseen bound of Corollary 1.2 was used to derive the approximate bound of the last line.

For the second factor of the product (9), we obtain the following equality

$$\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) = \mathbb{P}\left(|\bar{\sigma}_n^2 - \sigma^2| \geq \sigma^2/2\right) \Big|_{n=M}.$$

Furthermore, using Chebycheff's inequality and k-Statistics to bound the variance of the sample variance, cf. Keeping [6], we derive that

$$\begin{aligned} & \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) = \mathbb{P}\left(|\bar{\sigma}_n^2 - \sigma^2| \geq \sigma^2/2\right) \Big|_{n=M} \\ &\leq 4 \mathbb{E}\left[\frac{|\sigma^2 - \bar{\sigma}_n^2|^2}{\sigma^4}\right] \Big|_{n=M} \leq 4 \frac{\sigma^4 \left(\frac{2}{M-1} + \frac{\kappa}{M}\right)}{\sigma^4} = 4 \left(\frac{2}{M-1} + \frac{\kappa}{M}\right). \end{aligned}$$

Here  $\kappa$  denotes the *kurtosis*, i.e.

$$\kappa = \frac{\mathbb{E}[|X - \mu|^4]}{\sigma^4} - 3.$$

We conclude that

$$\mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \geq \sigma^2/2 \mid M\right) \leq \min\left\{1, 4 \left(\frac{2}{M-1} + \frac{\kappa}{M}\right)\right\}. \quad (11)$$

Next, we want to bound the first factor of the second term of (8),

$$\mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| \leq \sigma^2/2\}\right).$$

The event  $|\bar{\sigma}_M^2 - \sigma^2| \leq \sigma^2/2$  indicates that the variance is not substantially wrong-estimated and thereby that it is quite likely that sufficiently many samples are used in our MC estimate. Example 2.1 however illustrated that even in settings with reasonably well-estimated  $M$ , failing to meet the confidence is still possible. A relatively weak penalty should thus be added to the stopping criterion to avoid failure in this setting. Applying the Edgeworth expansion with truncated  $o(n^{-1/2})$  as a weak penalty, cf. Theorem 1.3, we derive the approximate bound

$$\begin{aligned}
 & \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M, \{|\bar{\sigma}_M^2 - \sigma^2| > \sigma^2/2\}\right) \\
 &= \mathbb{P}\left(|\bar{X}_n - \mu| > \text{TOL} \mid \{|\bar{\sigma}_n^2 - \sigma^2| > \sigma^2/2\}\right) \Big|_{n=M} \\
 &= \mathbb{P}\left(\sqrt{n} \frac{|\bar{X}_n - \mu|}{\sigma} > \frac{\sqrt{n} \text{TOL}}{\sigma} \mid \{|\bar{\sigma}_n^2 - \sigma^2| > \sigma^2/2\}\right) \Big|_{n=M} \\
 &\lesssim 2 \left( 1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) + \frac{\left|\frac{M \text{TOL}^2}{\sigma^2} - 1\right| \exp\left(-\frac{M \text{TOL}^2}{\sigma^2}\right) |\mathbb{E}[(X - \mu)^3]|}{6\sqrt{2\pi}M\sigma^3} \right). \tag{12}
 \end{aligned}$$

Combining (10), (11) and (12), and noting that for all  $x \in \mathbb{R}_+$  and  $n \in \mathbb{N}$ ,

$$\frac{|x^2 - 1| e^{-x^2/2} |\mathbb{E}[(X - \mu)^3]|}{6\sqrt{2\pi}n\sigma^3} \leq C_{\text{BE}}(x) \frac{\beta}{\sqrt{n}},$$

we obtain the following approximate bound for failing to meet the accuracy of Algorithm 1 conditioned on the stopped number of samples  $M$ :

$$\begin{aligned}
 & \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M\right) \\
 &\lesssim 2 \left\{ 1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) + C_{\text{BE}}\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) \frac{\beta}{\sqrt{M}} \right\} \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| \leq \sigma^2/2 \mid M\right) \\
 &+ 2 \left\{ 1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) + \frac{\left|\frac{M \text{TOL}^2}{\sigma^2} - 1\right| |\mathbb{E}[(X - \mu)^3]|}{\exp\left(\frac{M \text{TOL}^2}{2\sigma^2}\right) \times 6\sqrt{2\pi}M\sigma^3} \right\} \mathbb{P}\left(|\bar{\sigma}_M^2 - \sigma^2| > \sigma^2/2 \mid M\right) \\
 &\lesssim 2 \left( 1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) \right) + 2C_{\text{BE}}\left(\frac{\sqrt{M} \text{TOL}}{\sigma}\right) \frac{\beta}{\sqrt{M}} \min\left\{1, 4\left(\frac{2}{M-1} + \frac{\kappa}{M}\right)\right\} \\
 &+ \frac{\left|\frac{M \text{TOL}^2}{\sigma^2} - 1\right| |\mathbb{E}[(X - \mu)^3]|}{\exp\left(\frac{M \text{TOL}^2}{2\sigma^2}\right) \times 3\sqrt{2\pi}M\sigma^3} \max\left\{1 - 4\left(\frac{2}{M-1} + \frac{\kappa}{M}\right), 0\right\}. \tag{13}
 \end{aligned}$$

#### 4. A HIGHER MOMENTS BASED STOPPING RULE

From the approximate stochastic error bound (13) we will in this section construct a new, more reliable stopping rule with a stopping criterion based on second, third, and fourth moments of the r.v. that is sampled. The sampled moments our new algorithm will depend on are (here represented in biased form)

$$\begin{aligned}
 \bar{\sigma}_M &:= \sqrt{\frac{\sum_{i=1}^M (X_i - \bar{X}_M)^2}{M}}, & \bar{\beta}_M &:= \sum_{i=1}^M \frac{|X_i - \bar{X}_M|^3}{M\bar{\sigma}_M^3}, \\
 \hat{\beta}_M &:= \sum_{i=1}^M \frac{(X_i - \bar{X}_M)^3}{M\bar{\sigma}_M^3}, & \text{and } \bar{\kappa}_M &:= \sum_{i=1}^M \frac{(X_i - \bar{X}_M)^4}{M\bar{\sigma}_M^4} - 3.
 \end{aligned} \tag{14}$$

Replacing moments with sample moments in (13), we obtain a computable approximate stochastic error bound

$$\begin{aligned} & \mathbb{P}\left(|\bar{X}_M - \mu| > \text{TOL} \mid M\right) \\ & \lesssim 2 \left(1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\bar{\sigma}_M}\right)\right) + 2C_{\text{BE}} \left(\frac{\sqrt{M} \text{TOL}}{\bar{\sigma}_M}\right) \frac{\bar{\beta}_M}{\sqrt{M}} \min\left\{1, 4\left(\frac{2}{M-1} + \frac{\bar{\kappa}_M}{M}\right)\right\} \\ & \quad + \frac{\left|\frac{M \text{TOL}^2}{\bar{\sigma}_M^2} - 1\right| |\hat{\beta}_M|}{\exp\left(\frac{M \text{TOL}^2}{2\bar{\sigma}_M^2}\right) \times 3\sqrt{2\pi M} \bar{\sigma}_M^3} \max\left\{1 - 4\left(\frac{2}{M-1} + \frac{\bar{\kappa}_M}{M}\right), 0\right\}. \end{aligned} \quad (15)$$

The resulting approximate stochastic error bound will be implemented as the following stopping criterion in Algorithm 2:

$$\begin{aligned} & 2 \left(1 - \Phi\left(\frac{\sqrt{M} \text{TOL}}{\bar{\sigma}_M}\right)\right) + 2C_{\text{BE}} \left(\frac{\sqrt{M} \text{TOL}}{\bar{\sigma}_M}\right) \frac{\bar{\beta}_M}{\sqrt{M}} \min\left\{1, 4\left(\frac{2}{M-1} + \frac{\bar{\kappa}_M}{M}\right)\right\} \\ & \quad + \frac{\left|\frac{M \text{TOL}^2}{\bar{\sigma}_M^2} - 1\right| |\hat{\beta}_M|}{\exp\left(\frac{M \text{TOL}^2}{2\bar{\sigma}_M^2}\right) \times 3\sqrt{2\pi M} \bar{\sigma}_M^3} \max\left\{1 - 4\left(\frac{2}{M-1} + \frac{\bar{\kappa}_M}{M}\right), 0\right\} < \delta \end{aligned} \quad (16)$$

We now present the new stopping rule algorithm.

---

**Algorithm 2** Higher Moments Based Stopping Rule

---

**Input:** Accuracy TOL, confidence  $\delta$ , and initial number of samples  $M_0$ .

**Output:**  $\bar{X}_M$ .

Set  $n = 0$ , generate i.i.d. samples  $\{X_i\}_{i=1}^{M_n}$  and compute the sample moments  $\bar{\sigma}_{M_n}$ ,  $\bar{\beta}_{M_n}$ ,  $\hat{\beta}_{M_n}$  and  $\bar{\kappa}_{M_n}$  according to (14).

**while** Inequality (16) is not fulfilled. **do**

    Set  $n = n + 1$  and  $M_n = 2M_{n-1}$ .

    Generate  $M_n$  i.i.d. samples  $\{X_i\}_{i=1}^{M_n}$  and compute the sample moments  $\bar{\sigma}_{M_n}$ ,  $\bar{\beta}_{M_n}$ , and  $\bar{\kappa}_{M_n}$ .

**end while**

Set  $M = M_n$ , generate i.i.d. samples  $\{X_i\}_{i=1}^M$  and return the sample mean  $\bar{X}_M$ .

---

## 5. NUMERICAL EXPERIMENTS

In the numerical experiments we will estimate the mean of four differently distributed r.v. by using both the the sample variance based stopping rule in Algorithm 1 and the new higher moments based stopping rule in Algorithm 2. We compare the reliability and complexity of the algorithms, with the complexity measured in terms of average number of r.v. realizations needed to generate the given MC estimate for a given accuracy-confidence pair TOL and  $\delta$ . The distributions considered here are the light-tailed Uniform distribution, the Exponential distribution, the heavier-tailed Lognormal distribution, and the heavy-tailed Pareto distribution. In all these experiments we have set the algorithm parameter initial number of samples to  $M_0 = 30$ . From Figures 3, 4, 5, and 6, which illustrate the results of the numerical experiments, we observe that for the heavy-tailed distributions Algorithm 2 performs reliably and succeeds in meeting the accuracy-confidence constraint while

Algorithm 1 does not. For the light tailed distributions considered, both Algorithms meet the accuracy-confidence constraint. Regarding the complexity of the algorithms, we see that Algorithm 2 is only slightly more costly than Algorithm 1, and, as expected, the complexities of the algorithms seem to become more similar as TOL decreases.

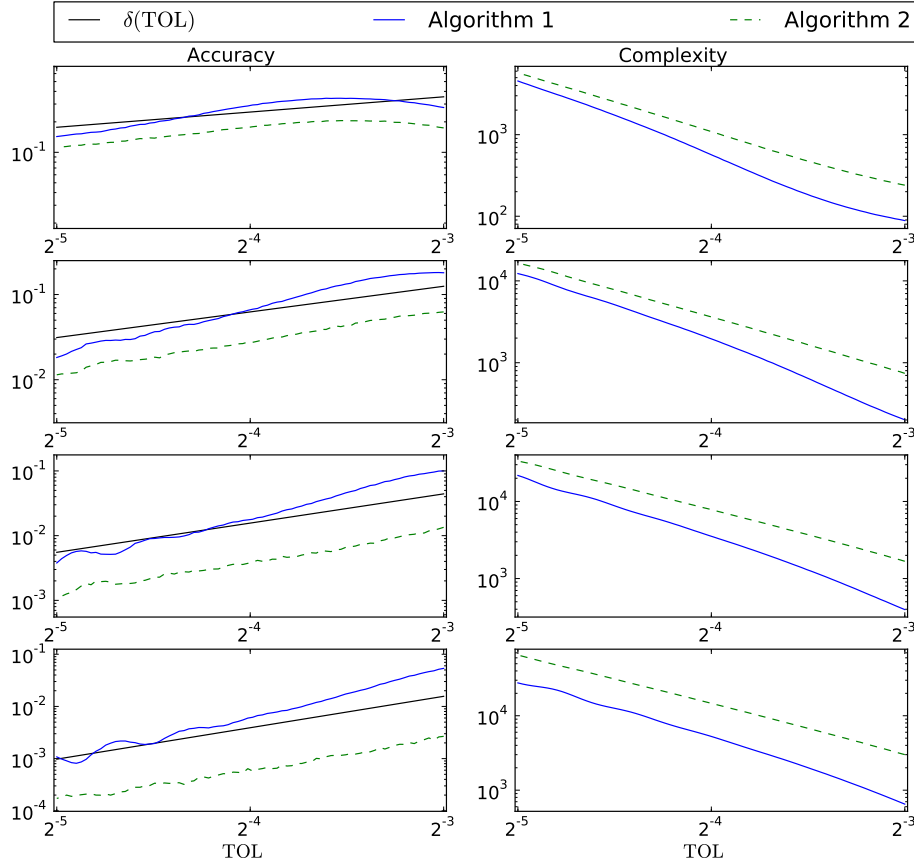


FIGURE 3. (**Pareto Distribution**) Numerical comparison of the accuracy and complexity of reaching the goal  $P(|\bar{X}_M - \mu| > \text{TOL}) < \delta$  with Algorithm 1 and 2 when sampling Pareto distributed r.v.s with parameters  $\alpha = 3.1$  and  $x_m = 1$ , cf. (5). Row plots from top to bottom is the output for the respective confidences  $\delta(\text{TOL}) = \text{TOL}^{1/2}$ ,  $\delta(\text{TOL}) = \text{TOL}$ ,  $\delta(\text{TOL}) = \text{TOL}^{3/2}$ , and  $\delta(\text{TOL}) = \text{TOL}^2$ .



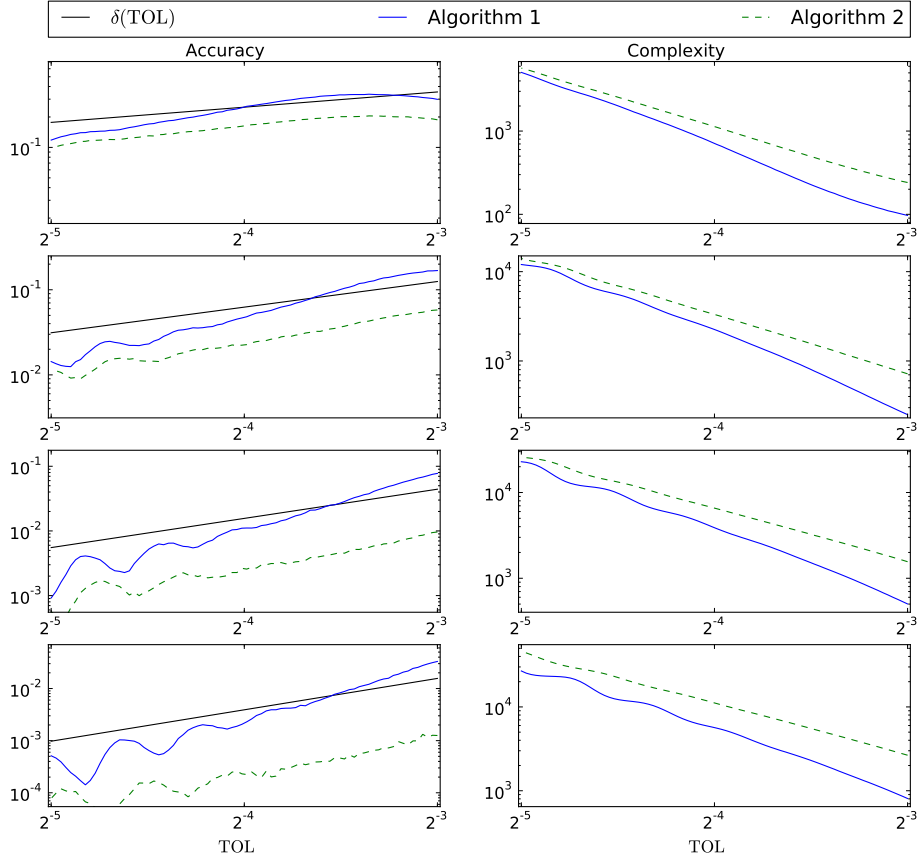


FIGURE 4. (**Lognormal Distribution**) Numerical comparison of the accuracy and complexity of reaching the goal  $\mathbb{P}(|\bar{X}_M - \mu| > \text{TOL}) < \delta$  with Algorithm 1 and 2 when sampling Lognormal distributed r.v.  $X \sim \log(\mathcal{N}(\mu_{\mathcal{N}}, \sigma_{\mathcal{N}}^2))$  with  $\mu_{\mathcal{N}} = -1$  and  $\sigma_{\mathcal{N}}^2 = 1$ . Row plots from top to bottom is the output for the respective confidences  $\delta(\text{TOL}) = \text{TOL}^{1/2}$ ,  $\delta(\text{TOL}) = \text{TOL}$ ,  $\delta(\text{TOL}) = \text{TOL}^{3/2}$ , and  $\delta(\text{TOL}) = \text{TOL}^2$ .

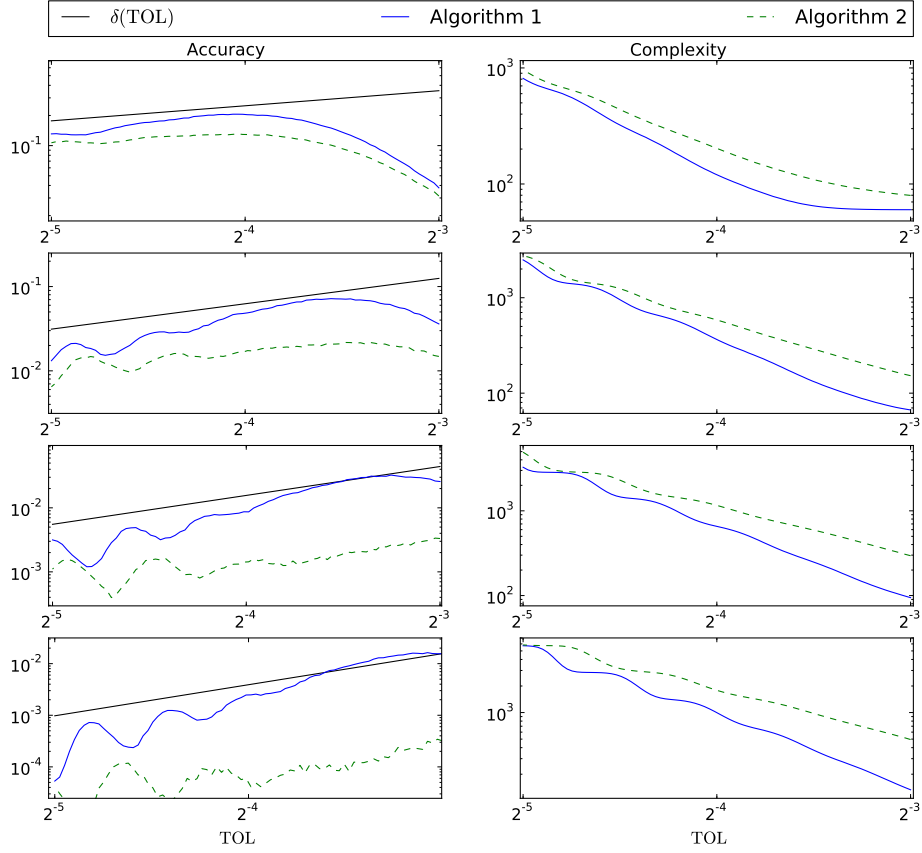


FIGURE 5. (**Exponential Distribution**) Numerical comparison of the accuracy and complexity of reaching the goal  $\mathbb{P}(|\bar{X}_M - \mu| > \text{TOL}) < \delta$  with Algorithm 1 and 2 when sampling exponentially distributed r.v. with  $\mu = 1/3$ . Row plots from top to bottom is the output for the respective confidences  $\delta(\text{TOL}) = \text{TOL}^{1/2}$ ,  $\delta(\text{TOL}) = \text{TOL}$ ,  $\delta(\text{TOL}) = \text{TOL}^{3/2}$ , and  $\delta(\text{TOL}) = \text{TOL}^2$ .

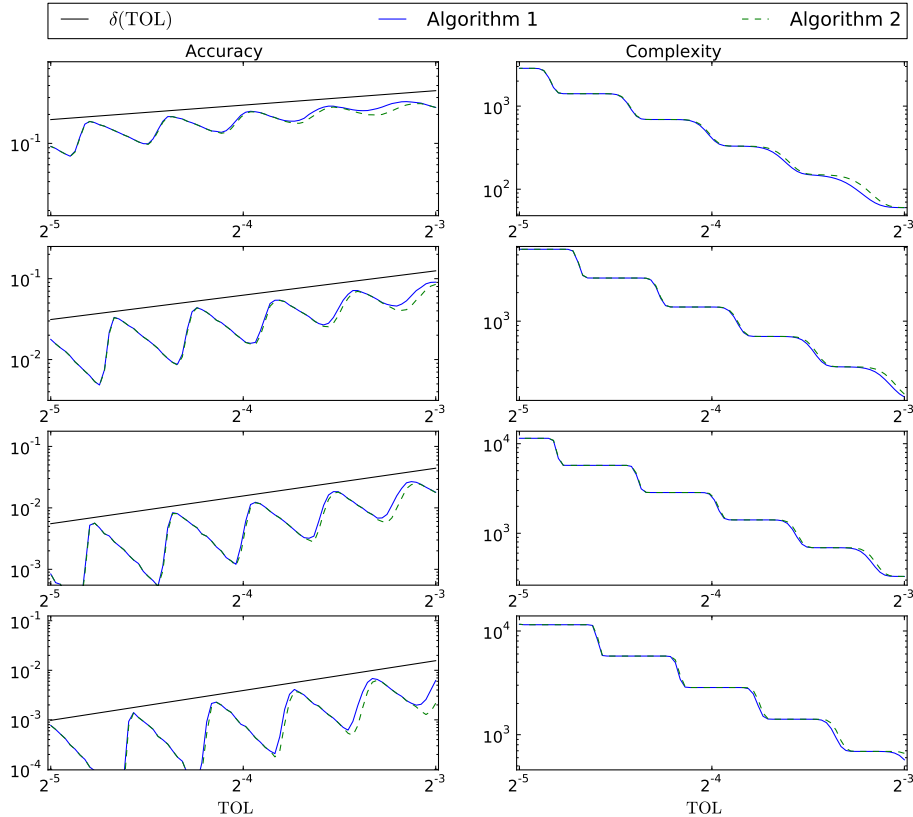


FIGURE 6. (**Uniform Distribution**) Numerical comparison of the accuracy and complexity of reaching the goal  $P(|\bar{X}_M - \mu| > \text{TOL}) < \delta$  with Algorithm 1 and 2 when sampling uniformly distributed r.v.  $X \sim U(-1, 1)$ . Row plots from top to bottom is the output for the respective confidences  $\delta(\text{TOL}) = \text{TOL}^{1/2}$ ,  $\delta(\text{TOL}) = \text{TOL}$ ,  $\delta(\text{TOL}) = \text{TOL}^{3/2}$ , and  $\delta(\text{TOL}) = \text{TOL}^2$ .

## 6. CONCLUSION

We have shown that second moment based sequential stopping rules such as Algorithm 1 run the risk of using too few samples in MC estimates, especially when sampling heavy-tailed r.v. in settings with very stringent confidence requirements, i.e.,  $\delta \ll \text{TOL}$ . Algorithm 2, a higher moment based stopping rule algorithm is proposed in this work, and, according to the numerical examples of Section 5, our new stopping rule performs much more reliable than Algorithm 1 while only slightly increasing the computational cost. In short, we believe that our new stopping rule presented in Algorithm 2 is well worth considering in settings with heavy tailed r.v. and/or  $\delta \ll \text{TOL}$ .

Note that our analysis of the original Algorithm 1 critically depends on three main ingredients:

- (I) a general, non-asymptotic estimate of the tail probabilities for the sample mean  $\bar{X}_M$ , for which we used either the non-uniform Berry-Esseen theorem given in Corollary 1.2 or the Edgeworth expansion given in Theorem 1.3,
- (II) a choice between the more conservative Berry-Esseen bound and the approximate Edgeworth bound made depending on whether the sample variance of

the samples used to generate the output MC estimate is close to, or far from the true variance,

- (III) an estimate of the conditional distribution function of the sample variance given the output  $M$  of the stopping algorithm given in (11).

There is clearly room for improvement in all these steps. First of all, the second ingredient above is dangerous as we do not know how to estimate the correlation between  $\bar{X}_M$  and the events  $|\bar{\sigma}_M^2 - \sigma^2| > \sigma^2/2$  and  $|\bar{\sigma}_M^2 - \sigma^2| \leq \sigma^2/2$ . This is problematic, as these approximations can potentially have the wrong sign, i.e., it is possible that the right-hand sides of (10) and (12) are smaller than their respective left-hand sides even though we actually seek upper bounds. It is however our hope that these approximation errors are compensated by the overly pessimistic non-uniform Berry-Esseen estimate and by using Chebycheff's inequality to bound the conditional distribution function of the sample variance. Even though the numerical evidence obtained in Section 5 seems to confirm that the compensations work well, we would prefer an analysis in which each estimation step can be controlled, at least in the sense that we indeed obtain an upper bound for the error probability.

To a lesser extent, it is not clear that the truncation of the  $o(n^{-1/2})$  of the Edgeworth expansion will lead to an upper bound for the error probability, either. In this case, the approximation error is however of higher order, so a stronger case can be made on why the effect will finally be negligible. In fact, when we used truncated Edgeworth expansion also for the estimation of (10) – instead of the non-uniform Berry-Esseen theorem – then the corresponding stopping rule turned out to be not much more reliable than Algorithm 1, indicating that there is a delicate balance between reliability in meeting the accuracy target (1) and maintaining an acceptable efficiency.

**Acknowledgments.** This work was partially supported by the Center for Predictive Computational Science), the VR project "Effektiva numeriska metoder för stokastiska differentialekvationer med tillämpningar", and King Abdullah University of Science and Technology (KAUST).

#### REFERENCES

- [1] Paul Bratley, Bennett L. Fox, and Linus E. Schrage. *A guide to simulation (2nd ed.)*. Springer-Verlag New York, Inc., New York, NY, USA, 1987.
- [2] Y. S. Chow and Herbert Robbins. On the asymptotic theory of fixed-width sequential confidence intervals for the mean. *The Annals of Mathematical Statistics*, 36(2):pp. 457–462, 1965.
- [3] W. Feller. *An introduction to probability theory and its applications. Vol II. 2nd ed.* Wiley Series in Probability and Mathematical Statistics. New York etc.: John Wiley and Sons, Inc. XXIV, 669 p. , 1971.
- [4] Peter W. Glynn and Ward Whitt. The asymptotic validity of sequential stopping rules for stochastic simulations. *The Annals of Applied Probability*, 2(1):pp. 180–198, 1992.
- [5] Wassily Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58(301):pp. 13–30, 1963.
- [6] E. S. Keeping. *Introduction to statistical inference*. D. Van Nostrand Co., Inc., Princeton, N.J.-Toronto-London-New York, 1962.
- [7] Averill M. Law and W. David Kelton. Confidence intervals for steady-state simulations, ii: A survey of sequential procedures. *Management Science*, 28(5):pp. 550–562, 1982.
- [8] Averill M. Law and W. David Kelton. Confidence intervals for steady-state simulations: I. a survey of fixed sample size procedures. *Operations Research*, 32(6):pp. 1221–1239, 1984.
- [9] R. Michel. On the constant in the nonuniform version of the berry-essen theorem. *Probability Theory and Related Fields*, 55:109–117, 1981. 10.1007/BF01013464.
- [10] Valentin V. Petrov. *Limit theorems of probability theory*, volume 4 of *Oxford Studies in Probability*. The Clarendon Press Oxford University Press, New York, 1995. Sequences of independent random variables, Oxford Science Publications.
- [11] Il'ya S Tyurin. Refinement of the upper bounds of the constants in lyapunov's theorem. *Russian Mathematical Surveys*, 65(3):586, 2010.

DEPARTMENT OF MATHEMATICS, UNIVERSITY OF VIENNA, NORDBERGSTRAE 15, 1090 WIEN,  
AUSTRIA

*E-mail address:* `christian.bayer@univie.ac.at`

DEPARTMENT OF NUMERICAL ANALYSIS, KUNGL. TEKNISKA HÖGSKOLAN, 100 44 STOCKHOLM,  
SWEDEN

*E-mail address:* `hhoel@kth.se`

DIVISION OF MATHEMATICS, KING ABDULLAH UNIVERSITY OF SCIENCE AND TECHNOLOGY,  
THUWAL 23955-6900, KINGDOM OF SAUDI ARABIA

*E-mail address:* `erik.vonschwerin@kaust.edu.sa`

DIVISION OF MATHEMATICS, KING ABDULLAH UNIVERSITY OF SCIENCE AND TECHNOLOGY,  
THUWAL 23955-6900, KINGDOM OF SAUDI ARABIA

*E-mail address:* `raul.tempone@kaust.edu.sa`



## Paper IV

# HOW ACCURATE IS MOLECULAR DYNAMICS?

CHRISTIAN BAYER, HÅKON HOEL, ASHRAFUL KADIR, PETR PLECHÁČ, MATTIAS SANDBERG,  
ANDERS SZEPESSY, AND RAUL TEMPONE

ABSTRACT. Born-Oppenheimer dynamics is shown to provide an accurate approximation of time-independent Schrödinger observables for a molecular system with an electron spectral gap, in the limit of large ratio of nuclei and electron masses, without assuming that the nuclei are localized to vanishing domains. The derivation, based on a Hamiltonian system interpretation of the Schrödinger equation and stability of the corresponding hitting time Hamilton-Jacobi equation for non ergodic dynamics, bypasses the usual separation of nuclei and electron wave functions, includes caustic states and gives a different perspective on the Born-Oppenheimer approximation, Schrödinger Hamiltonian systems and numerical simulation in molecular dynamics modeling at constant energy.

## CONTENTS

1. Motivation for error estimates of molecular dynamics	2
2. The Schrödinger and molecular dynamics models	3
3. A time-independent Schrödinger WKB-solution	7
3.1. Exact Schrödinger dynamics	7
3.2. Born-Oppenheimer dynamics	12
3.3. Equations for the density	12
3.4. Construction of the solution operator	13
4. Fourier integral WKB states including caustics	14
4.1. A preparatory example with the simplest caustic	14
4.2. A general Fourier integral ansatz	16
5. A global construction coupling caustics with single WKB-modes	22
5.1. Lagrangian manifolds	23
6. Computation of observables	25
7. Molecular dynamics approximation of Schrödinger observables	26
7.1. The Born-Oppenheimer approximation error	26
7.2. Why do symplectic numerical simulations of molecular dynamics work?	29
8. Analysis of the molecular dynamics approximation	29
8.1. Continuation of the construction of the solution operator	30
8.2. Stability from perturbed Hamiltonians	31
8.3. The Born-Oppenheimer approximation	34
9. Numerical examples	37
9.1. Example 1: A single WKB state	37
9.2. Example 2: A caustic state	38
10. The stationary phase expansion	43
Acknowledgment	45
References	45

---

2000 *Mathematics Subject Classification.* Primary: 81Q20; Secondary: 82C10.

*Key words and phrases.* Born-Oppenheimer approximation, WKB expansion, caustics, Fourier integral operators, Schrödinger operators.

The research of P.P. and A.S. was partially supported by the National Science Foundation under the grant NSF-DMS-0813893 and Swedish Research Council grant 621-2010-5647, respectively.



## 1. MOTIVATION FOR ERROR ESTIMATES OF MOLECULAR DYNAMICS

Molecular dynamics is a computational method to study molecular systems in materials science, chemistry and molecular biology. The simulations are used, for example, in designing and understanding new materials or for determining biochemical reactions in drug design, [14]. The wide popularity of molecular dynamics simulations relies on the fact that in many cases it agrees very well with experiments. Indeed when we have experimental data it is easy to verify correctness of the method by comparing with experiments at certain parameter regimes. However, if we want the simulation to predict something that has no comparing experiment, we need a mathematical estimate of the accuracy of the computation. In the case of molecular systems with few particles such studies are made by directly solving the Schrödinger equation. A fundamental and still open question in classical molecular dynamics simulations is how to verify the accuracy computationally, i.e., when the solution of the Schrödinger equation is not a computational alternative.

The aim of this paper is to derive qualitative error estimates for molecular dynamics and present new mathematical methods which could be used also for a more demanding quantitative accuracy estimation, without solving the Schrödinger solution. Having molecular dynamics error estimates opens, for instance, the possibility of systematically evaluating which density functionals or empirical force fields are good approximations and under what conditions the approximation properties hold. Computations with such error estimates could also give improved understanding when quantum effects are important and when they are not, in particular in cases when the Schrödinger equation is too computationally complex to solve.

*The first step to check the accuracy* of a molecular dynamics simulation is to know what to compare with. Here we compare with the value of any *observable*  $g(X)$ , of nuclei positions  $X$ , for the *time-independent Schrödinger* eigenvalue equation  $\mathcal{H}\Phi = E\Phi$ , so that the approximation error we study is

$$(1.1) \quad \int_{\mathbb{R}^{3(N+n)}} g(X)\Phi(x, X)^*\Phi(x, X) dx dX - \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(X_t) dt,$$

for a molecular dynamics path  $X_t$ , with total energy equal to the Schrödinger eigenvalue  $E$ . The observable can be, for instance, the local potential energy, used in [37] to determine phase-field partial differential equations from molecular dynamics simulations, see Figure 1. The time-independent Schrödinger equation has a remarkable property of accurately predicting experiments in combination with no unknown data, thereby forming the foundation of computational chemistry. However, the drawback is the high dimensional solution space for nuclei-electron systems with several particles, restricting numerical solution to small molecules. In this paper we study the *time-independent* setting of the Schrödinger equation as the reference. The proposed approach has the advantage of avoiding the difficulty of finding the initial data for the time-dependent Schrödinger equation.

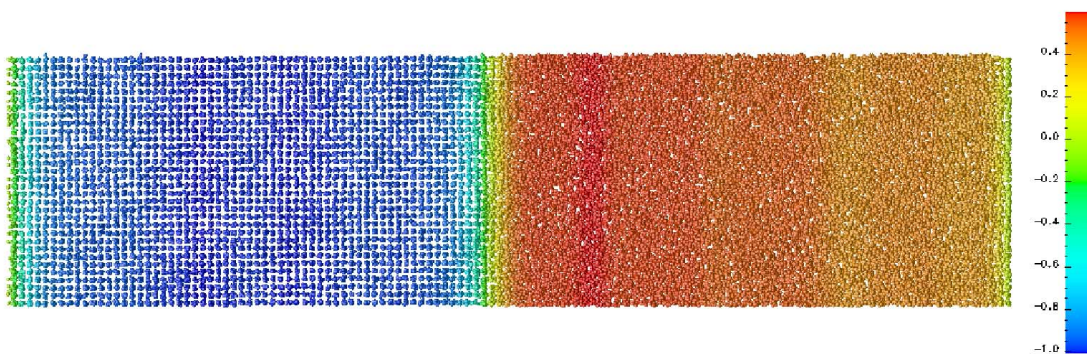


FIGURE 1. A Lennard-Jones molecular dynamics simulation of a phase transition with periodic boundary conditions, from [37]. The left part is solid and the right is liquid. The color measures the local potential energy.

*The second step to check the accuracy* is to derive error estimates. We have three types of error: time discretization error, sampling error and modeling error. The time discretization error comes from approximating

the differential equation for molecular dynamics with a numerical method, based on replacing time derivatives with difference quotients for a positive step size  $\Delta t$ . The sampling error is due to truncating the infinite  $T$  and using a finite value of  $T$  in the integral in (1.1). The modeling error (also called coarse-graining error) originates from eliminating the electrons in the Schrödinger nuclei-electron system and replacing the nuclei dynamics with their classical paths; this approximation error was first analyzed by Born and Oppenheimer in their seminal paper [2].

The time discretization and truncation error components are in some sense simple to handle by comparing simulations with different choice of  $\Delta t$  and  $T$ , although it can, of course, be difficult to know that the behavior does not change with even smaller  $\Delta t$  and larger  $T$ . The modeling error is more difficult to check since a direct approach would require to solve the Schrödinger equation. Currently the Schrödinger partial differential equation can only be solved with few particles, therefore it is not an option to solve the Schrödinger equation in general. The reason to use molecular dynamics is precisely in avoiding solution of the Schrödinger equation. Consequently the modeling error requires mathematical error analysis. In the literature there seems to be no error analysis that is precise, simple and constructive enough so that a molecular dynamics simulation can use it to assess the modeling error. Our alternative error analysis presented here is developed with the aim to give a different point of view that could help to construct algorithms that estimate the modeling error in molecular dynamics computations. Our analysis differs from previous ones by using

- the time-independent Schrödinger equation as the reference model to compare molecular dynamics with,
- an amplitude function in a WKB-Ansatz that depends only on position coordinates  $(x, X)$  (and not on momentum coordinates  $(p, P)$ ) for caustic states,
- actual solutions of the Schrödinger equation locally and asymptotic solutions globally (and not only asymptotic solutions),
- the theory of Hamilton-Jacobi partial differential equations to derive estimates for the corresponding Hamiltonian systems, i.e., the molecular dynamics systems.

Understanding both the exact Schrödinger model and the molecular dynamics model through Hamiltonian systems allows us to obtain bounds on local and some global problems for the difference of the solutions by well-established comparison results for the solutions of Hamilton-Jacobi equations, by regarding the Schrödinger Hamiltonian and the molecular dynamics Hamiltonians as perturbations of each others. The Hamilton-Jacobi theory applied to Hamiltonian systems is inspired by the error analysis of symplectic methods for optimal control problems for partial differential equations, [31]. The result is that the modeling error can be estimated based on the difference of the Hamiltonians, for the molecular dynamics system and the Schrödinger system, along the same solution path, see Theorem 7.1 and Section 8.2. The stability analysis limits the study to non ergodic dynamics.

The next section introduces the Schrödinger and molecular dynamics models. Sections 3 and 4 establish local analysis relating the Schrödinger problem to a Hamiltonian system for non caustic and caustic states, respectively. Sections 5 and 6 extend the local picture to a global construction, in the case of non ergodic dynamics. Sections 7 and 8 formulate approximation and stability results in the Hamilton-Jacobi setting. Section 9 presents numerical results.

## 2. THE SCHRÖDINGER AND MOLECULAR DYNAMICS MODELS

In deriving the approximation of the solutions to the full Schrödinger equation the heavy particles are often treated within classical mechanics, i.e., by defining the evolution of their positions and momenta by equations of motions of classical mechanics. Therefore we denote  $X_t : [0, \infty) \rightarrow \mathbb{R}^{3N}$  and  $P_t : [0, \infty) \rightarrow \mathbb{R}^{3N}$  time-dependent functions of positions and momenta with time derivatives denoted by

$$\dot{X}_t = \frac{dX_t}{dt}, \quad \ddot{X}_t = \frac{d^2X_t}{dt^2}.$$

We denote the Euclidean scalar product on  $\mathbb{R}^{3N}$  by

$$X \cdot Y = \sum_{i=1}^{3N} X^i Y^i.$$

Furthermore, we use the notation  $\nabla_X \psi(x, X) = (\nabla_{X^1} \psi(x, X), \dots, \nabla_{X^N} \psi(x, X))$ , and as customary  $\nabla_{X^i} \psi = (\partial_{X^i_1} \psi, \partial_{X^i_2} \psi, \partial_{X^i_3} \psi)$ .

On the other hand, the light particles are treated within the quantum mechanical description and the following complex valued bilinear map  $\langle \cdot, \cdot \rangle : L^2(\mathbb{R}^{3n} \times \mathbb{R}^{3N}) \times L^2(\mathbb{R}^{3n} \times \mathbb{R}^{3N}) \rightarrow L^2(\mathbb{R}^{3N})$  will be used in the subsequent calculations

$$(2.1) \quad \langle \phi, \psi \rangle = \int_{\mathbb{R}^{3n}} \phi(x, X)^* \psi(x, X) dx.$$

The notation  $\psi(x, X) = \mathcal{O}(M^{-\alpha})$  is also used for complex valued functions, meaning that  $|\psi(x, X)| = \mathcal{O}(M^{-\alpha})$  holds uniformly in  $x$  and  $X$ .

The *time-independent Schrödinger equation*

$$(2.2) \quad \mathcal{H}(x, X)\Phi(x, X) = E\Phi(x, X)$$

models many-body (nuclei-electron) quantum systems and is obtained from minimization of the energy in the solution space of wave functions, see [33, 32, 1, 35, 7]. It is an eigenvalue problem for the energy  $E \in \mathbb{R}$  of the system in the solution space, described by wave functions,  $\Phi : \mathbb{R}^{3n} \times \mathbb{R}^{3N} \rightarrow \mathbb{C}$ , depending on electron coordinates  $x = (x^1, \dots, x^n) \in \mathbb{R}^{3n}$ , nuclei coordinates  $X = (X^1, \dots, X^N) \in \mathbb{R}^{3N}$ , and the Hamiltonian operator  $\mathcal{H}(x, X)$

$$(2.3) \quad \mathcal{H}(x, X) = \mathcal{V}(x, X) - \frac{1}{2}M^{-1} \sum_{n=1}^N \Delta_{X^n}.$$

We assume that a quantum state of the system is fully described by the wave function  $\Phi : \mathbb{R}^{3n} \times \mathbb{R}^{3N} \rightarrow \mathbb{C}$  which is an element of the Hilbert space of wave functions with the standard complex valued scalar product

$$\langle\langle \Phi, \Psi \rangle\rangle = \int_{\mathbb{R}^{3n} \times \mathbb{R}^{3N}} \Phi(x, X)^* \Psi(x, X) dx dX,$$

and the operator  $\mathcal{H}$  is self-adjoint in this Hilbert space. The Hilbert space is then a subset of  $L^2(\mathbb{R}^{3n} \times \mathbb{R}^{3N})$  with symmetry conditions based on the Pauli exclusion principle for electrons, see [7, 22].

In computational chemistry the operator  $\mathcal{V}$ , the electron Hamiltonian, is independent of  $M$  and it is precisely determined by the sum of the kinetic energy of electrons and the Coulomb interaction between nuclei and electrons. We assume that the electron operator  $\mathcal{V}(\cdot, X)$  is self-adjoint in the subspace with the inner product  $\langle \cdot, \cdot \rangle$  of functions in (2.1) with fixed  $X$  coordinate and acts as a multiplication on functions that depend only on  $X$ . An essential feature of the partial differential equation (2.2) is the high computational complexity of finding the solution in an antisymmetric/symmetric subset of the Sobolev space  $H^1(\mathbb{R}^{3n} \times \mathbb{R}^{3N})$ . The mass of the nuclei, which are much greater than one (electron mass), are the diagonal elements in the diagonal matrix  $M$ .

In contrast to the Schrödinger equation, a *molecular dynamics* model of  $N$  nuclei  $X : [0, T] \rightarrow \mathbb{R}^{3N}$ , with a given potential  $V_p : \mathbb{R}^{3N} \rightarrow \mathbb{R}$ , can be computationally studied for large  $N$  by solving the ordinary differential equations

$$(2.4) \quad \ddot{X}_t = -\nabla_X V_p(X_t),$$

in the slow time scale, where the nuclei move  $\mathcal{O}(1)$  in unit time. This computational and conceptual simplification motivates the study to determine the potential and its implied accuracy compared with the the Schrödinger equation, as started already in the 1920's with the Born-Oppenheimer approximation [2]. The purpose of our work is to contribute to the current understanding of such derivations by showing convergence rates under new assumptions. The precise aim in this paper is to estimate the error

$$(2.5) \quad \frac{\int_{\mathbb{R}^{3N+3n}} g(X) \Phi(x, X)^* \Phi(x, X) dx dX}{\int_{\mathbb{R}^{3N+3n}} \Phi(x, X)^* \Phi(x, X) dx dX} - \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(X_t) dt$$

for a position dependent observable  $g(X)$  of the time-independent Schrödinger equation (2.2) approximated by the corresponding molecular dynamics observable  $\lim_{T \rightarrow \infty} T^{-1} \int_0^T g(X_t) dt$ , which is computationally cheaper to evaluate with several nuclei. The Schrödinger eigenvalue problem may typically have multiple eigenvalues and the aim is to find an eigenfunction  $\Phi$  and a molecular dynamics system that can be compared.

There may be eigenfunctions that we cannot approximate and the stability analysis we use limits the study to non ergodic dynamics.

The main step to relate the Schrödinger wave function and the molecular dynamics solution is the so-called zero-order Born-Oppenheimer approximation, where  $X_t$  solves the classical *ab initio* molecular dynamics (2.4) with the potential  $V_p : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  determined as an eigenvalue of the electron Hamiltonian  $\mathcal{V}(\cdot, X)$  for a given nuclei position  $X$ . That is  $V_p(X) = \lambda_0(X)$  and

$$\mathcal{V}(\cdot, X)\Psi_{\text{BO}}(\cdot, X) = \lambda_0(X)\Psi_{\text{BO}}(\cdot, X),$$

for an electron eigenfunction  $\Psi_{\text{BO}}(\cdot, X) \in L^2(\mathbb{R}^{3n})$ , for instance, the ground state. The Born-Oppenheimer expansion [2] is an approximation of the solution to the time-independent Schrödinger equation which is shown in [15, 19] to solve the time-independent Schrödinger equation approximately. This expansion, analyzed by the methods of multiple scales, pseudo-differential operators and spectral analysis in [15, 19, 13], can be used to study the approximation error (2.5). However, in the literature, e.g., [24], it is easier to find precise statements on the error for the setting of the time-dependent Schrödinger equation, since the stability issue is more subtle in the eigenvalue setting.

Instead of an asymptotic expansion we use a different method based on a Hamiltonian dynamics formulation of the *time-independent* Schrödinger eigenfunction and the stability of the corresponding perturbed Hamilton-Jacobi equations viewed as a hitting problem. This approach makes it possible to reduce the error propagation on the infinite time interval to finite time excursions from a certain co-dimension one hitting set. A motivation for our method is that it forms a sub-step in trying to estimate the approximation error using only information available in molecular dynamics simulations.

The related problem of approximating observables to the time-dependent Schrödinger equation by the Born-Oppenheimer expansions is well studied, theoretically in [4, 28] and computationally in [20] using the Egorov theorem. The Egorov theorem shows that finite time observables of the time-dependent Schrödinger equation are approximated with  $\mathcal{O}(M^{-1})$  accuracy by the zero-order Born-Oppenheimer dynamics with an electron eigenvalue gap. In the special case of a position observable and no electrons (i.e.,  $\mathcal{V} = V(X)$  in (2.3)), the Egorov theorem states that

$$(2.6) \quad \left| \int_{\mathbb{R}^{3N}} g(X)\Phi(X, t)^*\Phi(X, t) dX - \int_{\mathbb{R}^{3N}} g(X_t)\Phi(X_0, 0)^*\Phi(X_0, 0) dX_0 \right| \leq C_t M^{-1},$$

where  $\Phi(X, t)$  is a solution to the time-dependent Schrödinger equation

$$i\partial_t\Phi(\cdot, t) = \mathcal{H}\Phi(\cdot, t)$$

with the Hamiltonian (2.3) and the path  $X_t$  is the nuclei coordinates for the dynamics with the Hamiltonian  $\frac{1}{2}|\dot{X}|^2 + V(X)$ . If the initial wave function  $\Phi(X, 0)$  is the eigenfunction in (2.2) the first term in (2.6) reduces to the first term in (2.5) and the second term can also become the same in an ergodic limit. However, since we do not know that the parameter  $C_t$  (bounding an integral over  $(0, t)$ ) is bounded for all time we cannot directly conclude an estimate for (2.5) from (2.6).

In our perspective studying the time-independent instead of the time-dependent Schrödinger equation has the important differences that

- the infinite time study of the Born-Oppenheimer dynamics can be reduced to a finite time hitting problem,
- the computational and theoretical problem of specifying initial data for the Schrödinger equation is avoided, and
- computationally cheap evaluation of the position observable  $g(X)$  is possible using the time average  $\lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T g(X_t) dt$  along the solution path  $X_t$ .

In this paper we derive the Born-Oppenheimer approximation from the time-independent Schrödinger equation (2.2) and we establish convergence rates for molecular dynamics approximations to time-independent Schrödinger observables under simple assumptions including the so-called *caustic* points, where the Jacobian determinant  $\det J(X_t) \equiv \det(\partial X_t / \partial X_0)$  of the Eulerian-Lagrangian transformation of  $X$ -paths vanish. As mentioned previously, the main new analytical idea is an interpretation of the time-independent Schrödinger equation (2.2) as a Hamiltonian system and the subsequent analysis of the approximations by comparing Hamiltonians. This analysis employs the theory of Hamilton-Jacobi partial differential equations. The

problematic infinite-time evolution of perturbations in the dynamics is solved for non ergodic dynamics by viewing it as a finite-time hitting problem for the Hamilton-Jacobi equation, with a particular hitting set. In contrast to the traditional rigorous and formal asymptotic expansions we analyze the transport equation as a time-dependent Schrödinger equation.

The main inspiration for this paper are works [27, 6, 5] and the semi-classical WKB analysis in [25]: the works [27, 6, 5] derive the time-dependent Schrödinger dynamics of an  $x$ -system,  $i\dot{\Psi} = \mathcal{H}_1\Psi$ , from the time-independent Schrödinger equation (with the Hamiltonian  $\mathcal{H}_1(x) + \epsilon\mathcal{H}(x, X)$ ) by a classical limit for the environment variable  $X$ , as the coupling parameter  $\epsilon$  vanishes and the mass  $M$  tends to infinity; in particular [27, 6, 5] show that the time derivative enters through the coupling of  $\Psi$  with the classical velocity. Here we refine the use of characteristics to study classical *ab initio* molecular dynamics where the coupling does not vanish, and we establish error estimates for Born-Oppenheimer approximations of Schrödinger observables. The small scale, introduced by the perturbation

$$-(2M)^{-1} \sum_k \Delta_{X^k}$$

of the potential  $\mathcal{V}$ , is identified in a modified WKB eikonal equation and analyzed through the corresponding transport equation as a time-dependent Schrödinger equation along the eikonal characteristics. This modified WKB formulation reduces to the standard semi-classical approximation, see [25], in the case of the potential function  $\mathcal{V} = V(X) \in \mathbb{R}$ , depending only on nuclei coordinates, but becomes different in the case of operator-valued potentials studied here. The global analysis of WKB functions was initiated by Maslov in the 1960', [25], and lead to the subjects Geometry of Quantization and Quantum Ergodicity, relating global classical paths to eigenfunctions of the Schrödinger equation, see [10] and [38]. The analysis presented in this paper is based on a Hamiltonian system interpretation of the time-independent Schrödinger equation. Stability of the corresponding Hamilton-Jacobi equation, bypasses the usual separation of nuclei and electron wave functions in the time-dependent self-consistent field equations, [3, 23, 36].

A unique property of the time-independent Schrödinger equation we use is the interpretation that the dynamics  $X_t \in \mathbb{R}^{3N}$  can return to a co-dimension one surface  $I$  which then can reduce the dynamics to a hitting time problem with finite-time excursions from  $I$ . We assume that the (Lagrangian) manifold, generated by the visited points  $(X_t, \dot{X}_t) \in \mathbb{R}^{6N}$  in phase space is smooth, which excludes ergodic dynamics. Another advantage of viewing the molecular dynamics as an approximation of the eigenvalue problem is that stochastic perturbations of the electron ground state can be interpreted as a Gibbs distribution of degenerate nuclei-electron eigenstates of the Schrödinger eigenvalue problem (2.2), see [34]. The time-independent eigenvalue setting also avoids the issue on “wave function collapse” to an eigenstate, present in the time-dependent Schrödinger equation.

Theorem 7.1 demonstrates that observables from the zero-order Born-Oppenheimer dynamics approximate observables for the Schrödinger eigenvalue problem with the error of order  $\mathcal{O}(M^{-1+\delta})$ , for any  $\delta > 0$ , assuming that the electron eigenvalues satisfy a spectral gap condition and that the Lagrangian manifold is smooth. The result is based on the Hamiltonian (2.3) with any potential  $\mathcal{V}$  that is smooth in  $X$ , e.g., a regularized version of the Coulomb potential. The derivation does not assume that the nuclei are supported on small domains; in contrast derivations based on the time-dependent self-consistent field equations require nuclei to be supported on small domains. The reason that the small support is not needed here comes from the combination of the characteristics and sampling from an equilibrium density. In other words, the nuclei paths behave classically although they may not be supported on small domains. Section 6 shows that caustics couple the WKB modes, as is well-known from geometric optics, see [18, 25], and generate non-orthogonal WKB modes that are coupled in the Schrödinger density. On the other hand, with a spectral gap and without caustics the Schrödinger density is asymptotically decoupled into a simple sum of individual WKB densities. Section 4 constructs a WKB-Fourier integral Schrödinger solution for caustic states. Section 7.2 relates the approximation results to the accuracy of symplectic numerical methods for molecular dynamics.

We believe that these ideas can be further developed to better understanding of molecular dynamics simulations. Our study does not directly apply to ergodic dynamics, since then the Lagrangian manifold becomes dense in a set of dimension  $6N - 1$  in phase-space, which violates our assumption of a smooth Lagrangian manifold of dimension  $3N$ . It would also be desirable to have more precise conditions on the data (i.e. molecular dynamics initial data and potential  $\mathcal{V}$ ) instead of our implicit assumption on hitting

times, smooth Lagrangian manifold and convergence of the Born-Oppenheimer power series approximation in Lemma 8.2.

### 3. A TIME-INDEPENDENT SCHRÖDINGER WKB-SOLUTION

**3.1. Exact Schrödinger dynamics.** For the sake of simplicity we assume that all nuclei have the same mass. If this is not the case, we can introduce new coordinates  $M_1^{1/2} \tilde{X}^k = M_k^{1/2} X^k$ , which transform the Hamiltonian to the form we want  $\mathcal{V}(x, M_1^{1/2} M^{-1/2} \tilde{X}) - (2M_1)^{-1} \sum_{k=1}^N \Delta_{\tilde{X}^k}$ . The singular perturbation  $-(2M)^{-1} \sum_k \Delta_{X^k}$  of the potential  $\mathcal{V}$  introduces an additional small scale  $M^{-1/2}$  of high frequency oscillations, as shown by a WKB-expansion, see [29, 17, 16, 26]. We shall construct solutions to (2.2) in such a WKB-form

$$(3.1) \quad \Phi(x, X) = \phi(x, X) e^{iM^{1/2}\theta(X)},$$

where the amplitude function  $\phi : \mathbb{R}^{3n} \times \mathbb{R}^{3N} \rightarrow \mathbb{C}$  is complex valued, the phase  $\theta : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  is real valued, and the factor  $M^{1/2}$  is introduced in order to have well-defined limits of  $\phi$  and  $\theta$  as  $M \rightarrow \infty$ . Note that it is trivially always possible to find functions  $\phi$  and  $\theta$  satisfying (3.1), even in the sense of a true equality. Of course, the ansatz only makes sense if  $\phi$  and  $\theta$  do not have strong oscillations for large  $M$ . The standard WKB-construction, [25, 16], is based on a series expansion in powers of  $M^{1/2}$  which solves the Schrödinger equation with arbitrary high accuracy. Instead of an asymptotic solution, we introduce an actual solution based on a time-dependent Schrödinger transport equation. This transport equation reduces to the formulation in [25] for the case of a potential function  $\mathcal{V} = V(X) \in \mathbb{R}$ , depending only on nuclei coordinates  $X \in \mathbb{R}^{3N}$ , and modifies it for the case of a self-adjoint potential operator  $\mathcal{V}(\cdot, X)$  on the electron space  $L^2(\mathbb{R}^{3n})$  which is the primary focus of our work here. In Sections 6 and 4 we use a linear combination of WKB-eigenfunctions, but first we study the simplest case of a single WKB-eigenfunction as motivated by the following subsection.

**3.1.1. Molecular dynamics from a piecewise constant electron operator on a simplex mesh.** The purpose of this section is to convey a first formal understanding of the relation between ab initio molecular dynamics  $\ddot{X}_t = -\nabla_X \lambda_0(X_t)$  and the Schrödinger eigenvalue problem (2.2) and motivate the WKB ansatz (3.1). In subsequent sections we will describe precise analysis of error estimates for the WKB-method. The idea behind this first study is to approximate the electron operator  $\mathcal{V}$  by a finite dimensional matrix  $\mathcal{V}^h$ , which is piecewise constant on a simplex mesh in the variable  $X$ , with the mesh size  $h$ . Furthermore, we introduce the change of variables

$$\Phi = \sum_{j=0}^J \varphi_j \Psi_j =: \Psi \varphi$$

based on the piecewise constant electron eigenvalues and eigenvectors  $\mathcal{V}^h \Psi_j = \lambda_j^h \Psi_j$ ,  $\langle \Psi_j, \Psi_j \rangle = 1$ ,  $j = 0, \dots, J$ , normalized and ordered with respect to increasing eigenvalues. Then the Schrödinger equation (2.2) becomes

$$-\frac{1}{2M} \Delta_X (\Psi \varphi) + \mathcal{V}^h \Psi \varphi = E \Psi \varphi,$$

with the notation  $\Delta_X = \sum_j \Delta_{X_j}$ , so that on each simplex

$$-\frac{1}{2M} \Delta_X \varphi_j + \lambda_j^h \varphi_j = E \varphi_j,$$

which by separation of variables, for each  $j = 0, 1, 2, \dots, J$ , implies

$$(3.2) \quad \varphi_j = \sum_{P^j} a(P^j) e^{iM^{1/2} P^j \cdot X}$$

for any  $P^j \in \mathbb{C}^{3N}$  that satisfies the eikonal equation

$$\frac{1}{2} P^j \cdot P^j + \lambda_j^h = E,$$

for any  $a(P^j) \in \mathbb{C}$ , if all components of  $P^j$  are non zero. If  $P_k^j = 0$  we have  $a(P^j) = \prod_{\{k: P_k^j=0\}} (A_k X_k + B_k)$  for any  $A_k \in \mathbb{C}, B_k \in \mathbb{C}$ , since  $e^{\pm iM^{1/2} P_k^j X_k} = 1$  in this case. The solution  $\Phi$ , to (2.2), and its normal

derivative are continuous at the interfaces of the simplices. On the intersection of the faces the normal derivative is not defined but this set is of measure zero and thus negligible as seen from the  $H^1(\mathbb{R}^{3N})$  solution concept of (2.2).

We investigate a simpler, one-dimensional case,  $X \in \mathbb{R}$ , first. Then the solution  $\varphi$  simplifies to

$$\varphi_j = a_j e^{iM^{1/2}P^j \cdot X} + b_j e^{-iM^{1/2}P^j \cdot X}$$

for  $a_j, b_j, P^j \in \mathbb{C}$  and  $(P^j)^2/2 + \lambda_j = E$ . The continuity conditions

$$(3.3) \quad \begin{aligned} \lim_{X \rightarrow X_0^+} \Phi(X) &= \lim_{X \rightarrow X_0^-} \Phi(X) \\ \lim_{X \rightarrow X_0^+} \partial_X \Phi(X) &= \lim_{X \rightarrow X_0^-} \partial_X \Phi(X) \end{aligned}$$

hold for any  $X_0 \in \mathbb{R}$ , in particular, at the interval boundary where for  $X_0 = 0$

$$(3.4) \quad \begin{aligned} \lim_{X \rightarrow X_0^\pm} \Phi(X) &= \sum_j (a_{j\pm} \Psi_{j\pm} + b_{j\pm} \Psi_{j\pm}) \\ \lim_{X \rightarrow X_0^\pm} \partial_X \Phi(X) &= iM^{1/2} \sum_j (a_{j\pm} P_\pm^j \Psi_{j\pm} - b_{j\pm} P_\pm^j \Psi_{j\pm}). \end{aligned}$$

It is clear that given  $a_-$  and  $b_-$  we can determine  $a_+$  and  $b_+$  so that (3.3) holds. In order to prepare for the multi-dimensional case it is convenient to consider each incoming wave  $a_-$  and  $b_+$  separately: the incoming  $a_-$  wave is split into a refracted  $a_+$  and reflected  $b_-$  wave

$$(3.5) \quad \sum_j a_{j-} \Psi_{j-} P_-^j = \sum_j (a_{j+} \Psi_{j+} P_+^j + b_{j-} \Psi_{j-} P_-^j)$$

and similarly the incoming  $b_+$  wave is split into a refracted  $b_-$  wave and a reflected  $a_+$  wave, see Figure 2. The jump conditions at the different interfaces are coupled by the oscillatory functions  $e^{\pm iM^{1/2}P^j \cdot X}$ . The global construction of  $\varphi$  and  $\Psi$  in one dimension follows by marching in the positive  $X$ -direction to successive intervals, creating in each interval both a  $e^{iM^{1/2}P^j \cdot X} \Psi_j$  and a  $e^{-iM^{1/2}P^j \cdot X} \Psi_j$  wave.

In general each interface condition (3.4) also couples all eigenvectors  $\Psi_j$ . However, we shall see that if  $M$  is large,  $\mathcal{V}$  smooth and there is a spectral gap  $\lambda_1 - \lambda_0 > c > 0$  then, in the limit of the simplex size  $h$  tending to zero, there is an asymptotically uncoupled WKB-solution  $\Phi(x, X) = \phi(x, X) e^{iM^{1/2}\theta(X)}$ , where  $\theta : \mathbb{R}^{3N} \rightarrow \mathbb{R}$ ,  $\phi : \mathbb{R}^{3n} \times \mathbb{R}^{3N} \rightarrow \mathbb{C}$ . Under these assumptions the Born-Oppenheimer approximation in Lemma 8.2 shows that  $\phi$  is asymptotically parallel, in  $L^2(dx)$ , to the electron eigenfunction  $\Psi_0$  as  $M \rightarrow \infty$ . The gradient  $\nabla_X \theta(X) = P^0$  is obtained from the differential  $\theta(X) = \theta(X_0) + \nabla_X \theta(X_0) \cdot (X - X_0) + o(|X - X_0|)$ .

In the case of electron eigenvalue crossing, i.e.,  $\lambda_1(X) = \lambda_0(X)$  for some  $X$ , or so called avoided crossings (meaning that the eigenvalue gap  $c \ll 1$  is small and dependent on  $M$ ), a refraction will, in general, include all components  $a_j e^{iM^{1/2}P^j \cdot X} \Psi_j$ ,  $j = 1, \dots, J$  and consequently the Born-Oppenheimer approximation fails.

The construction of a solution to the Schrödinger equation with a piecewise constant potential is more involved in the multi-dimensional case for two reasons: each reflection at an interface generates, in general, an additional path in a new direction, so that many paths are needed. Furthermore, the construction of a solution to the eikonal equation is more complicated since the jump condition (3.4) implies that the tangential component  $P_t^j$  of  $P^j$  must be continuous across a simplex face and only the normal component  $P_n^j = P^j - P_t^j$  may have a jump. In multi-dimensional cases it is still possible to construct a solution of the form (3.2) by following the characteristic paths  $\dot{X}_t = P^j(X_t)$  and using the jump conditions (3.4): when the path  $X_t$  hits a simplex face, the tangential part  $P_t^j$  of  $P^j$  is continuous and the normal component  $P_n^j$  of  $P^j$  may jump. At a simplex face the new value of the  $P_n^j$  is determined by  $(P_n^j \cdot P_n^j + P_t^j \cdot P_t^j)/2 + \lambda_j^h = E$ . Analogously to the one dimensional case we treat the pair  $e^{iM^{1/2}(P_t^j + P_n^j) \cdot X}$  and  $e^{iM^{1/2}(P_t^j - P_n^j) \cdot X}$  together. However, each collision with  $e^{iM^{1/2}(P_t^j + P_n^j) \cdot X}$  on an interface now creates a reflected wave in another direction, in particular,  $e^{iM^{1/2}(P_t^j - P_n^j) \cdot X} \Psi_j$ , and we get many paths to follow. Therefore each mode  $e^{iM^{1/2}P^j \cdot X}$  follows its characteristic  $X_t$ , where  $\dot{X}_t = P^j$ , through the simplex to the adjacent simplicial faces, which the characteristic pass through when they leave the simplex, and at these outflow faces a reflected mode is created and a refracted mode continues into the adjacent simplices, see Figure 2. In this way we can

formally construct a solution of the form  $\sum_{P^j} a(P^j) e^{iM^{1/2}P^j \cdot X} \Psi_j$  to the Schrödinger equation (2.2), with possibly several different characteristic paths in each simplex.

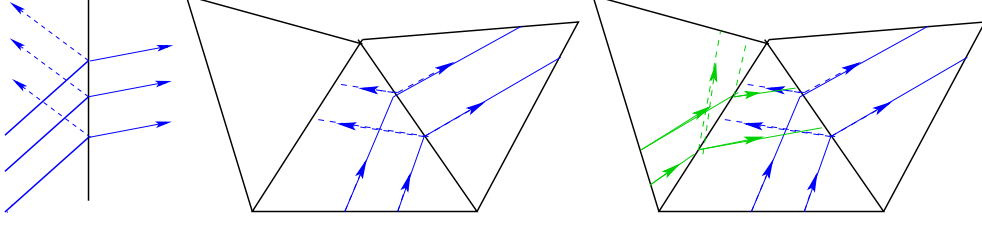


FIGURE 2. The value of  $P^j$  is constructed by following the characteristic paths  $X_t$  (the blue and green curves), based on  $\dot{X}_t = P^j$ , with a reflection-refraction at each simplex face (left) following the path through simplices (middle) and each simplex may have several  $P^j$  (right).

In conclusion, the piecewise constant electron operator shows that the solution to the Schrödinger equation (2.2) is composed of a linear combination of highly oscillatory function modes  $a_j e^{iM^{1/2}P^j \cdot X} \Psi_j$  based on the electron eigenvectors  $\Psi_j$  and eigenvalues  $\lambda_j$ , where  $P^j$  satisfies the eikonal equation  $P^j \cdot P^j / 2 + \lambda_j(X) = E$ . These modes can be followed by characteristics  $\dot{X} = P^j$  from simplex to simplex. In this paper we show that observables based on the related WKB Schrödinger solutions can be approximated by molecular dynamics time averages, when there is a spectral gap around  $\lambda_0$ .

3.1.2. *A first WKB-solution.* The WKB-solution satisfies the Schrödinger equation (2.2) provided that

$$(3.6) \quad \begin{aligned} 0 &= (\mathcal{H} - E)\phi e^{iM^{1/2}\theta(X)} \\ &= \left( \left( \frac{1}{2}|\nabla_X \theta|^2 + \mathcal{V} - E \right) \phi - \frac{1}{2M} \Delta_X \phi - \frac{i}{M^{1/2}} (\nabla_X \phi \cdot \nabla_X \theta + \frac{1}{2} \phi \Delta_X \theta) \right) e^{iM^{1/2}\theta(X)}. \end{aligned}$$

We shall see that only eigensolutions  $\Phi$  that correspond to dynamics without caustics correspond to such a single WKB-mode, as for instance when the eigenvalue  $E$  is inside an electron eigenvalue gap. Solutions in the presence of caustics use a Fourier integral of such WKB-modes, and we treat this case in detail in Section 4. To understand the behavior of  $\theta$ , we multiply (3.6) by  $\phi^* e^{-iM^{1/2}\theta(X)}$  and integrate over  $\mathbb{R}^{3n}$ . Similarly we take the complex conjugate of (3.6), and multiply by  $\phi e^{iM^{1/2}\theta(X)}$  and integrate over  $\mathbb{R}^{3n}$ . By adding these two expressions we obtain

$$(3.7) \quad \begin{aligned} 0 &= 2 \left( \frac{1}{2} |\nabla_X \theta|^2 - E \right) \langle \phi, \phi \rangle + \underbrace{\langle \phi, \mathcal{V} \phi \rangle + \langle \nabla_X \phi, \phi \rangle}_{=2\langle \phi, \mathcal{V} \phi \rangle} - \frac{1}{2M} (\langle \phi, \Delta_X \phi \rangle + \langle \Delta_X \phi, \phi \rangle) \\ &\quad - \frac{i}{M^{1/2}} \left( \underbrace{\langle \phi, \nabla_X \phi \cdot \nabla_X \theta \rangle - \langle \nabla_X \phi \cdot \nabla_X \theta, \phi \rangle}_{=2i \operatorname{Im} \langle \phi, \nabla_X \phi \cdot \nabla_X \theta \rangle} + \frac{i}{2M^{1/2}} \underbrace{(\langle \phi, \phi \rangle - \langle \phi, \phi \rangle)}_{=0} \Delta_X \theta \right). \end{aligned}$$

The purpose of the phase function  $\theta$  is to generate an accurate approximation in the limit as  $M \rightarrow \infty$ . A possible and natural definition of  $\theta$  would be the formal limit of (3.7) as  $M \rightarrow \infty$ , which is the *Hamilton-Jacobi equation*, also called the *eikonal equation*

$$(3.8) \quad \frac{1}{2} |\nabla_X \theta|^2 = E - V_0,$$

where the function  $V_0 : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  is

$$(3.9) \quad V_0 := \frac{\langle \phi, \mathcal{V} \phi \rangle}{\langle \phi, \phi \rangle}.$$

The solution to the Hamilton-Jacobi eikonal equation can be constructed from the associated Hamiltonian system

$$(3.10) \quad \begin{aligned} \dot{X}_t &= P_t \\ \dot{P}_t &= -\nabla_X V_0(X_t) \end{aligned}$$



through the characteristics path  $(X_t, P_t)$  satisfying  $\nabla_X \theta(X_t) =: P_t$ . The amplitude function  $\phi$  can be determined by requiring the ansatz (3.6) to be a solution, which gives

$$\begin{aligned} 0 &= (\mathcal{H} - E)\phi e^{iM^{1/2}\theta(X)} \\ &= \underbrace{\left( \frac{1}{2}|\nabla_X \theta|^2 + V_0 - E \right)}_{=0} \phi \\ &\quad - \frac{1}{2M}\Delta_X \phi + (\mathcal{V} - V_0)\phi - \frac{i}{M^{1/2}}(\nabla_X \phi \cdot \nabla_X \theta + \frac{1}{2}\phi \Delta_X \theta) e^{iM^{1/2}\theta(X)}, \end{aligned}$$

so that by using (3.8) we have

$$-\frac{1}{2M}\Delta_X \phi + (\mathcal{V} - V_0)\phi - \frac{i}{M^{1/2}}(\nabla_X \phi \cdot \nabla_X \theta + \frac{1}{2}\phi \Delta_X \theta) = 0.$$

The usual method for determining  $\phi$  from this so-called *transport equation* uses an asymptotic expansion  $\phi \simeq \sum_{k=0}^K M^{-k/2} \phi_k$ , see [15, 19] and the beginning of Section 8. An alternative is to write it as a Schrödinger equation, similar to work in [25]: we apply the characteristics in (3.10) to write

$$\frac{d}{dt}\phi(X_t) = \nabla_X \phi \cdot \dot{X}_t = \nabla_X \phi \cdot \nabla_X \theta,$$

and define the weight function  $G$  by

$$(3.11) \quad \frac{d}{dt} \log G_t = \frac{1}{2} \Delta_X \theta(X_t),$$

and the variable  $\psi_t := \phi(X_t)G_t$ . We use the notation  $\phi(X)$  instead of the more precise  $\phi(\cdot, X)$ , so that e.g.  $\psi_t = \psi_t(x) = \phi(x, X_t)G_t$ . Then the transport equation becomes a Schrödinger equation

$$(3.12) \quad iM^{-1/2}\dot{\psi}_t = (\mathcal{V} - V_0)\psi_t - \frac{G_t}{2M}\Delta_X \left( \frac{\psi_t}{G_t} \right).$$

In conclusion, equations (3.8)-(3.12) determine the WKB-ansatz (3.1) to be a local solution to the Schrödinger equation (2.2).

**Theorem 3.1.** *Assume the Hamilton-Jacobi equation, with the corresponding Hamiltonian,*

$$H_S(X, P) := \frac{1}{2}|P|^2 + \underbrace{\frac{\langle \psi(X), \mathcal{V}(X)\psi(X) \rangle}{\langle \psi(X), \psi(X) \rangle}}_{=: V_0(X)} - E = 0,$$

*based on the primal variable  $X$  and the dual variable  $P = P(X) = \nabla_X \theta(X)$ , has a smooth solution  $\theta(X)$  in a domain  $\mathcal{U} \subseteq \mathbb{R}^{3N}$ , then  $\theta$  generates a solution to the time-independent Schrödinger equation  $(\mathcal{H} - E)\Phi = 0$  in  $\mathcal{U}$ , in the sense that*

$$\Phi(X_t, x) = \hat{G}^{-1}(X_t)\hat{\psi}(x, X_t)e^{iM^{1/2}\theta(X_t)},$$

*solves the equation (2.2) in  $\mathcal{U}$ , where  $\hat{\psi}(X_t) := \psi_t$  satisfies the transport equation (3.12) and*

$$\hat{G}(X_t) = G_t,$$

$$\frac{d}{dt} \log G_t = \frac{1}{2} \Delta_X \theta(X_t),$$

$$(X_t, P_t) \text{ solves the Hamiltonian system (3.10) corresponding to } H_S.$$

The theorem tells us that if there is a  $C^2$  solution to the Hamilton-Jacobi equation, then we have a family of characteristic paths with the desired property. It is well known that Hamilton-Jacobi equations in general do not have global  $C^2$  solutions, due to  $X$ -paths that collide, as seen by (5.4) generating blow up in  $\partial_{XX}\theta(X)$ . However if the domain is small enough, the data on the boundary is compatible (in the sense that  $H_S(X, \nabla_X \theta(X)) = 0$  in the boundary) and noncharacteristic (in the sense that the normal derivative  $\partial_n \theta(X) \neq 0$  on the boundary) and  $V_0$  is smooth, then the converse property holds: that the characteristics generate a local solution to the Hamilton-Jacobi equation, see Ref. [12]. In Section 5 we describe Maslov's method to find a global asymptotic solution by patching together local solutions; an important ingredient is how to set up data for the Hamiltonian system, which is previewed in the next section.

3.1.3. *Data for the Hamiltonian system.* For the energy  $E$  chosen larger than the potential energy, that is such that  $E \geq V_0$ , Theorem 3.1 yields a solution  $\theta : \mathcal{U} \rightarrow \mathbb{R}$  to the eikonal equation (3.8) locally in a neighborhood  $\mathcal{U} \subseteq \mathbb{R}^{3N}$ , for regular compatible data  $(X_0, P_0)$  given on a  $3N - 1$  dimensional "inflow"-domain  $X_0 \in I \subset \bar{\mathcal{U}}$ . For a Schrödinger eigenvalue problem, the domain  $I$  and the data  $(X_0, P_0)|_I$  are not given (except that the total energy is  $E$ ). In contrast for a scattering problem, the domain  $I$  has given data. If paths leaving from  $I$  return to  $I$ , there is an additional global compatibility of data on  $I$ : assume  $X_0 \in I$  and  $X_t \in I$ , then the values  $P_t$  are determined from  $P_0$ ; continuing the path to subsequent hitting points  $X_{t_j} \in I$ ,  $j = 1, 2, \dots$  determines  $P_{t_j}$  from  $P_0$ . The characteristic path  $(X_t, P_t)$ ,  $t > 0$ , generates a  $3N$  dimensional Lagrangian manifold in the  $6N$  dimensional phase space  $(X, P)$ , which is smooth under our assumptions. This Lagrangian manifold is in general only locally of the form  $(X, P(X))$ , but in the case of no caustics it is globally of this form and then there is a phase function  $X \mapsto \theta(X)$  such that  $P(X) = \nabla_X \theta(X)$  globally. Section 5.1 reviews background material on Lagrangian manifolds used in this paper.

In Section 4 we study phase space manifolds with caustics and Section 5 presents a global construction of a Lagrangian manifold in some cases. We will use a variant of Maslov's construction [25] to obtain a global asymptotic solution to the Schrödinger equation (2.2) from local WKB-solutions and we apply a Poincare map to determine the initial Lagrangian manifold, as described in Section 5: the first step is to define a codimension one hitting plane in the phase space  $\mathbb{R}^{6N}$ ; the problem is reduced to find an initial Lagrangian manifold of dimension  $3N - 1$  in the hitting plane by following the characteristic paths extending  $\theta_j$  and  $\phi_j$  locally; around a caustic the solution is a Fourier integral of WKB solutions, described in Section 4 and the stationary phase method yields boundary conditions for the phase and amplitude functions from the Fourier integral solution; on the hitting plane the solution has to coincide with the initial data, giving a fixed point problem for the initial Lagrangian manifold.

3.1.4. *Liouville's formula.* In this section we verify Liouville's formula

$$(3.13) \quad \frac{G_0^2}{G_t^2} = e^{-\int_0^t \text{Tr}(\nabla_X P(X_t)) dt} = \left| \det \frac{\partial(X_0)}{\partial(X_t)} \right|,$$

given in [25]. The characteristic  $\dot{X}_t = P(X_t)$  implies  $\frac{d}{dt} J(X_t) = \nabla_X P J(X)$ , where  $J(X)_{ij} = \partial X_t^i / \partial X_0^j$  denotes the first variation with respect to perturbations of the initial data. The logarithmic derivative then satisfies  $d/dt(\log J(X))_{ij} = \partial_{X^j} P^i(X_t) = \partial_{X^i X^j} \theta(X)$  which implies that  $\log J(X_t)$  is symmetric and shows that (3.13) holds

$$\text{div} P = \text{Tr} \nabla_X P = \frac{d}{dt} \text{Tr} \log J(X) = \frac{d}{dt} \log \det J(X).$$

The last step uses that  $J(X)$  can be diagonalized by an orthogonal transformation and that the trace is invariant under orthogonal transformations.

3.1.5. *The density and the first variation.* Note that the nuclei density, using  $\hat{G}$ , can be written

$$(3.14) \quad \rho := \frac{\langle \phi, \phi \rangle}{\int_{\mathbb{R}^{3N}} \langle \phi, \phi \rangle dX} = \frac{\langle \hat{\psi}, \hat{\psi} \rangle \hat{G}^{-2}}{\int_{\mathbb{R}^{3N}} \langle \hat{\psi}, \hat{\psi} \rangle \hat{G}^{-2} dX},$$

and since each time  $t$  determines a unique point  $(X_t, P_t) = (X_t, \nabla_X \theta(X_t))$  in the phase space the functions  $\hat{G}$  and  $\hat{\psi}$  are well defined.

The integrating factor  $G$  and its derivative  $\partial_{X^i} G$  can be determined from  $(P, \partial_{X^i} P, \partial_{X^i X^j} P)$  along the characteristics by the following characteristic equations obtained from (3.8) by differentiation with respect

to  $X$

$$\begin{aligned}
\frac{d}{dt} \partial_{X^r} P^k &= \left[ \sum_j P^j \partial_{X^j X^r} P^k = \sum_j P^j \partial_{X^r X^k} P^j \right] \\
&= - \sum_j \partial_{X^r} P^j \partial_{X^k} P^j - \partial_{X^r X^k} V_0, \\
(3.15) \quad \frac{d}{dt} \partial_{X^r X^q} P^k &= \left[ \sum_j P^j \partial_{X^j X^r X^q} P^k + \sum_j P^j \partial_{X^r X^k X^q} P^j \right] \\
&= - \sum_j \partial_{X^r} P^j \partial_{X^k X^q} P^j - \sum_j \partial_{X^r X^q} P^j \partial_{X^k} P^j - \partial_{X^r X^k X^q} V_0,
\end{aligned}$$

and similarly  $\partial_{X^i X^j} G$  can be determined from  $(P, \partial_{X^i} P, \partial_{X^i X^j} P, \partial_{X^i X^j X^k} P)$ .

**3.2. Born-Oppenheimer dynamics.** The Born-Oppenheimer approximation leads to the standard formulation of *ab initio* molecular dynamics, in the micro-canonical ensemble with the constant number of particles, volume and energy, for the nuclei positions  $X = X_{\text{BO}}$ ,

$$\begin{aligned}
(3.16) \quad \dot{X}_t &= P_t, \\
\dot{P}_t &= -\nabla_X \lambda_0(X_t),
\end{aligned}$$

by using that the electrons are in the eigenstate  $\psi = \Psi_{\text{BO}}$  with eigenvalue  $\lambda_0$  to  $\mathcal{V}$ , in  $L^2(dx)$  for fixed  $X$ , i.e.,  $\mathcal{V}(X)\Psi_{\text{BO}} = \lambda_0(X)\Psi_{\text{BO}}$ . The corresponding Hamiltonian is  $H_{\text{BO}}(X, P) := |P|^2/2 + \lambda_0(X)$  with the eikonal equation

$$(3.17) \quad \frac{1}{2} |\nabla_X \theta_{\text{BO}}(X)|^2 + \lambda_0(X) = E.$$

**3.3. Equations for the density.** We note that

$$\phi = \hat{G}^{-1} \hat{\psi} = \left( \frac{\rho}{\langle \hat{\psi}, \hat{\psi} \rangle / \int \langle \hat{\psi}, \hat{\psi} \rangle \hat{G}^{-2} dX} \right)^{1/2} \hat{\psi},$$

shows that  $G$  and  $\psi$  determine the density

$$(3.18) \quad \rho_{\text{S}} = \rho = \frac{\langle \hat{\psi}, \hat{\psi} \rangle |\hat{G}|^{-2}}{\int \langle \hat{\psi}, \hat{\psi} \rangle |\hat{G}|^{-2} dX},$$

defined in (3.14). Using the Born-Oppenheimer approximation in Lemma 8.2 we have  $\langle \hat{\psi}, \hat{\psi} \rangle = 1 + \mathcal{O}(M^{-1})$  in the case of a spectral gap. Therefore the weight function  $|\hat{G}|^{-2}$  approximates the density and we know from Theorem 3.1 that  $|\hat{G}|^{-2}$  is determined by the phase function  $\theta$ .

The Born-Oppenheimer dynamics generates an approximate solution  $\Psi_{\text{BO}} \hat{G}_{\text{BO}}^{-1} e^{iM^{1/2}\theta_{\text{BO}}}$  which yields the density

$$(3.19) \quad \rho_{\text{BO}} = |\hat{G}_{\text{BO}}|^{-2},$$

where

$$\frac{d}{dt} \log |\hat{G}_{\text{BO}}|^{-2} = -\Delta_X \theta_{\text{BO}}(X).$$

This representation can also be obtained from the conservation of mass

$$(3.20) \quad 0 = \text{div}(\rho_{\text{BO}} \nabla_X \theta_{\text{BO}})$$

implying

$$(3.21) \quad \frac{d}{dt} \rho_{\text{BO}}(X_t) = \nabla_X \rho_{\text{BO}}(X_t) \cdot \dot{X}_t = -\rho_{\text{BO}}(X_t) \text{div} \nabla_X \theta_{\text{BO}},$$

with the solution

$$(3.22) \quad \rho_{\text{BO}}(X_t) = \frac{C}{|\hat{G}_{\text{BO}}(X_t)|^2},$$

where  $C$  is a positive constant for each characteristic. Note that the derivation of this classical density does not need a corresponding WKB equation but uses only the conservation of mass that holds for classical paths satisfying a Hamiltonian system. The classical density corresponds precisely to the Eulerian-Lagrangian change of coordinates  $|G_t|^2/|G_0|^2 = \det(\partial X_t/\partial X_0)$  in (3.13).

**3.4. Construction of the solution operator.** The WKB Ansatz (3.1) is meaningful when  $\phi$  does not include the full small scale. In Lemma 8.2 we present conditions for  $\psi$  to be smooth.

To construct the solution operator it is convenient to include a non interacting particle in the system, i.e., a particle without charge, and assume that this particle moves with a constant, high speed  $dX_1^1/dt = P_1^1 \gg 1$  (or equivalently with the unit speed and a large mass). Such a non interacting particle does not affect the other particles. The additional new coordinate  $X_1^1$  is helpful in order to simply relate the time-coordinate  $t$  and  $X_1^1$ . We add the corresponding kinetic energy  $(P_1^1)^2/2$  to  $E$  in order not to change the original problem (2.2) and write the equation (3.12) in the fast time scale  $\tau = M^{1/2}t$

$$i \frac{d}{d\tau} \psi = (\mathcal{V} - V_0) \psi - \frac{1}{2M} G \sum_j \Delta_{X^j} (G^{-1} \psi).$$

Furthermore, we change to the coordinates

$$(\tau, X_*) := (\tau, X_2^1, X_3^1, X^2, \dots, X^N) \in [0, \infty) \times I, \text{ instead of } (X^1, X^2, \dots, X^N) \in \mathbb{R}^{3N},$$

where  $X^j = (X_1^j, X_2^j, X_3^j) \in \mathbb{R}^3$ . Hence we obtain

$$(3.23) \quad i \dot{\psi} + \frac{1}{2(P_1^1)^2} \ddot{\psi} = (\mathcal{V} - V_0) \psi - \frac{1}{2M} G \sum_j \Delta_{X_*^j} (G^{-1} \psi) =: \tilde{\mathcal{V}} \psi,$$

using the notation  $\dot{w} = dw/d\tau$  in this section. In Section 8.1 we show that the left hand side can be reduced to  $i \dot{\psi}$  as  $P_1^1 \rightarrow \infty$ , by choosing special initial data. Note also that  $G$  is independent of the first component in  $X^1$ . We see that the operator

$$\bar{\mathcal{V}} := G^{-1} \tilde{\mathcal{V}} G = \underbrace{G^{-1}(\mathcal{V} - V_0)G}_{= \mathcal{V} - V_0} - \frac{1}{2M} \sum_j \Delta_{X_*^j}$$

is symmetric on  $L^2(\mathbb{R}^{3n+3N-1})$ . Assume now the data  $(X_0, P_0, Z_0)$  for  $X_0 \in \mathbb{R}^{3N-1}$  is  $(L\mathbb{Z})^{3N-1}$ -periodic, then also  $(X_\tau, P_\tau, Z_\tau)$  is  $(L\mathbb{Z})^{3N-1}$ -periodic, for  $Z_t = \theta(X_t)$  and  $P_t = \nabla_X \theta(X_t)$ . To simplify the notation for such periodic functions, define the periodic circle

$$\mathbb{T} := \mathbb{R}/(L\mathbb{Z}).$$

We seek a solution  $\Phi$  of (2.2) which is  $(L\mathbb{Z})^{3(n+N)-1}$ -periodic in the  $(x, X_*)$ -variable. The Schrödinger operator  $\bar{\mathcal{V}}(\cdot, X_\tau)$  has, for each  $\tau$ , real eigenvalues  $\{\lambda_m(\tau)\}$  with a complete set of eigenvectors  $\{\zeta^m(x, X_*, \tau)\}$  orthogonal in the space functions, a subset of  $L^2(\mathbb{T}^{3n+3N-1})$ , see [1]. The proof uses that the operator  $\bar{\mathcal{V}}_\tau + \gamma I$  generates a compact solution operator in the Hilbert space functions in  $L^2(\mathbb{T}^{3n+3N-1})$ , for the constant  $\gamma \in (0, \infty)$  chosen sufficiently large. The discrete spectrum and the compactness comes from Fredholm theory for compact operators and the fact that the bilinear form  $\int_{\mathbb{T}^{3(n+N)-1}} v \bar{\mathcal{V}}_\tau w + \gamma v w \, dx \, dX_*$  is continuous and coercive on  $H^1(\mathbb{T}^{3(n+N)-1})$ , see [12]. We see that  $\tilde{\mathcal{V}}$  has the same eigenvalues  $\{\lambda_m(\tau)\}$  and the eigenvectors  $\{G_\tau \zeta^m(\tau)\}$ , orthogonal in the weighted  $L^2$ -scalar product

$$\int_{\mathbb{T}^{3N-1}} \langle v, w \rangle \hat{G}^{-2} \, dX_*.$$

The construction and analysis of the solution operator continues in Section 8.1 based on the spectrum.

**Remark 3.2** (Boundary conditions). The Schrödinger problem (2.2) makes sense not only in the periodic setting but also with alternative boundary conditions, e.g. from interaction with an external environment in scattering problems.

#### 4. FOURIER INTEGRAL WKB STATES INCLUDING CAUSTICS

**4.1. A preparatory example with the simplest caustic.** As an example of a caustic, we study first the simplest example of a fold caustic based on the Airy function  $A : \mathbb{R} \rightarrow \mathbb{R}$  which solves

$$(4.1) \quad -\partial_{xx}A(x) + xA(x) = 0.$$

The scaled Airy function

$$u(x) = C A(M^{1/3}x)$$

solves the Schrödinger equation

$$(4.2) \quad -\frac{1}{M}\partial_{xx}u(x) + xu(x) = 0,$$

for any constant  $C$ . In our context an important property of the Airy function is the fact that it is the inverse Fourier transform of the function

$$\hat{A}(p) = \sqrt{\frac{2}{\pi}} e^{ip^3/3},$$

i.e.,

$$(4.3) \quad A(x) = \frac{1}{\pi} \int_{\mathbb{R}} e^{i(xp+p^3/3)} dp.$$

In the next section, we will consider a general Schrödinger equation and determine a WKB Fourier integral corresponding to (4.3) for the Airy function; as an introduction to the general case we show how to derive (4.3): by taking the Fourier transform of the ordinary differential equation (4.1)

$$(4.4) \quad 0 = \int_{\mathbb{R}} (-\partial_{xx} + x) A(x) e^{-ixp} dx = (p^2 + i\partial_p)\hat{A}(p),$$

we obtain an ordinary differential equation for the Fourier transform  $\hat{A}(p)$  with the solution  $\hat{A}(p) = C e^{ip^3}$ , for any constant  $C$ . Then, by differentiation, it is clear that the scaled Airy function  $u$  solves (4.2). Furthermore, the stationary phase method, cf. Section 10, shows that to the leading order  $u$  is approximated by

$$u(x) \simeq C \left(-xM^{1/3}\right)^{-1/4} \cos\left(M^{1/2}(-x)^{3/2} - \pi/4\right), \quad \text{for } x < 0,$$

and  $u(x) \simeq 0$  to any order (i.e.,  $\mathcal{O}(M^{-K})$  for any positive  $K$ ) when  $x > 0$ . The behaviour of the Airy function is illustrated in Figure 3.

**4.1.1. Molecular dynamics for the Airy function.** The eikonal equation corresponding to (4.2) is

$$p^2 + x = 0$$

with solutions for  $x \leq 0$ , which leads to the phase

$$(4.5) \quad p = \theta'(x) = \pm(-x)^{1/2}, \quad \text{and} \quad \theta(x) = \mp \frac{2}{3}(-x)^{3/2}.$$

We compute the Legendre transform

$$\theta^*(p) = xp - \theta(x)$$

where by (4.5) and  $-x = p^2$  we obtain

$$\theta^*(p) = -p^2p + \frac{2}{3}p^3 = -\frac{p^3}{3}.$$

We note that this solution is also obtained from the eikonal equation

$$p^2 + \partial_p\theta^*(p) = 0,$$

which is solved by

$$\theta^*(p) = -p^3/3.$$

Thus we recover the relation for the Legendre transform  $-xp + \theta^*(p) = -\theta(x)$ .

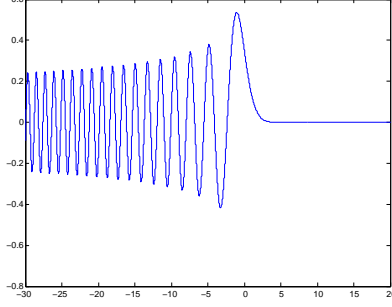


FIGURE 3. The Airy function.

4.1.2. *Observables for the Airy function.* The primary object of our analysis is an observable (a functional depending on  $u$ ) rather than the solution  $u(x)$  itself. Thus we first compute the observable evaluated on the solution obtained from the Airy function. In the following calculation we denote by  $C$  a generic constant not necessarily the same at each occurrence,

$$\begin{aligned}
\int_{\mathbb{R}} g(x) |u(x)|^2 dx &= C \int_{\mathbb{R}} g(x) \int_{\mathbb{R}} \int_{\mathbb{R}} e^{-iM^{1/2}(xp+p^3/3)} e^{iM^{1/2}(xq+q^3/3)} dq dp dx \\
&= C \int_{\mathbb{R}} \int_{\mathbb{R}} \hat{g}\left(M^{1/2}(p-q)\right) e^{iM^{1/2}(q^3/3-p^3/3)} dq dp \\
&= C \int_{\mathbb{R}} \int_{\mathbb{R}} \hat{g}\left(M^{1/2}(p-q)\right) e^{iM^{1/2}((q-p)^3/12+(q-p)(p+q)^2/4)} dq dp \\
(4.6) \quad &= C \int_{\mathbb{R}} \int_{\mathbb{R}} \underbrace{\hat{g}(-M^{1/2}\bar{q})}_{=t} e^{iM^{1/2}(\bar{q}^3/12+\bar{q}\bar{p}^2/4)} \underbrace{d\bar{q}d\bar{p}}_{\bar{q}=q-p, \bar{p}=p+q} \\
&= C \int_{\mathbb{R}} \int_{\mathbb{R}} \hat{g}(t) e^{-i(t^3/(12M)+t\bar{p}^2/4)} dt d\bar{p} \\
&= C \int_{\mathbb{R}} g * A_M\left(\underbrace{-\bar{p}^2}_{=\partial_p \theta^*(\bar{p})=x}\right) d\bar{p} \\
&= C \int_{-\infty}^0 g * A_M(x) |\partial_x p(x)| dx,
\end{aligned}$$

where

$$(4.7) \quad A_M(x) := \left(\frac{M}{4}\right)^{1/3} A\left(\left(\frac{M}{4}\right)^{1/3} x\right) \text{ is the Fourier transform of } e^{-it^3/(12M)}.$$

**Lemma 4.1.** *The scaled Airy function  $A_M$  is an approximate identity in the following sense*

$$(4.8) \quad \|g * A_M - g\|_{L^2(\mathbb{R})} \leq \frac{1}{12M} \|\partial_x^3 g\|_{L^2(\mathbb{R})}.$$

*Proof.* Plancherel's Theorem implies

$$\begin{aligned} M\|g * A_M - g\|_{L^2} &= M\|\hat{g}\hat{A}_M - \hat{g}\|_{L^2} = \|\hat{g}(e^{ip^3/(12M)} - 1)M\|_{L^2} \\ &\leq \frac{1}{12}\| |p|^3 \hat{g} \|_{L^2} = \frac{1}{12}\|\partial_x^3 g\|_{L^2}. \end{aligned}$$

The inequality follows from  $|e^{iy} - 1| \leq |y|$  which holds for all  $y \in \mathbb{R}$ .  $\square$

The classical molecular dynamics approximation corresponding to the Schrödinger equation (4.2) is the Hamiltonian system

$$\dot{X} = p, \quad \dot{p} = -\frac{1}{2}$$

with a solution  $X_t = -t^2/4$  and the corresponding approximation of the observable

$$\frac{1}{T} \int_0^T g(X_t) dt = \frac{1}{T} \int_0^T g(X_t) \frac{dX_t}{\dot{X}_t} = \frac{1}{T} \int_{-T^2/4}^0 g(x) \frac{dx}{|p(x)|}.$$

In this specific case the phase satisfies  $|p(x)| = |x|^{1/2}$  and  $|\partial_x p| = |x|^{-1/2}/2$ , and hence the non-normalized density  $|p|^{-1}$  is in this case equal to  $2|\partial_x p|$ . Equation (4.6) and Lemma 4.1 imply

$$\left| \int_{\mathbb{R}} g|u|^2 dx - \int_{\mathbb{R}} g\partial_x p(x) dx \right| = \mathcal{O}(M^{-1})$$

and consequently for two different observables  $g_1$  and  $g_2$  we have that Schrödinger observables are approximated by the classical observables with the error  $\mathcal{O}(M^{-1})$

$$(4.9) \quad \frac{\int_{\mathbb{R}} g_1 |u|^2 dx}{\int_{\mathbb{R}} g_2 |u|^2 dx} - \frac{\int_{\mathbb{R}} g_1 |\partial_{xx}\theta| dx}{\int_{\mathbb{R}} g_2 |\partial_{xx}\theta| dx} = \mathcal{O}(M^{-1}),$$

using  $\partial_x p(x) = \partial_{xx}\theta(x)$ . The reason we compare two different observables with a compact support is that  $\int_{\mathbb{R}} u^2(x) dx = \infty$  in the case of the Airy function.

We note that in (4.6) we used

$$\frac{1}{3}(q^3 - p^3) = \theta^*(p) - \theta^*(q) = (p - q)\partial_p \theta^* \left( \frac{1}{2}(p + q) \right) + \frac{1}{3}\partial^3 \theta^* \left( \frac{1}{2}(p + q) \right) \left( \frac{1}{2}(p - q) \right)^3$$

which in the next section is generalized to other caustics. For the Airy function there holds

$$\frac{1}{3}\partial^3 \theta^* \left( \frac{1}{2}(p + q) \right) = -\frac{2}{3}.$$

**4.2. A general Fourier integral ansatz.** In order to treat a more general case with a caustic of the dimension  $d$  we use the Fourier integral ansatz

$$(4.10) \quad \Phi(X, x) = \int_{\mathbb{R}^d} \phi(X, x) e^{-iM^{1/2}\Theta(\check{X}, \hat{X}, \check{P})} d\check{P}$$

and we write

$$\begin{aligned} X &= (\hat{X}, \check{X}), \quad P = (\hat{P}, \check{P}) \\ \check{X} \cdot \check{P} &= \sum_{j=1}^d \check{X}^j \check{P}^j, \quad \hat{X} \cdot \hat{P} = \sum_{j=d+1}^N \hat{X}^j \hat{P}^j \\ \Theta(\check{X}, \hat{X}, \check{P}) &= \check{X} \cdot \check{P} - \theta^*(\hat{X}, \check{P}), \end{aligned}$$

based on the Legendre transform

$$\theta^*(\hat{X}, \check{P}) = \min_{\check{X}} \left( \check{X} \cdot \check{P} - \theta(\hat{X}, \check{X}) \right).$$

If the function  $\theta^*(\hat{X}, \check{P})$  is not defined for all  $\check{P} \in \mathbb{R}^d$ , but only for  $\check{P} \in \mathcal{U} \subset \mathbb{R}^d$  we replace the integral over  $\mathbb{R}^d$  by integration over  $\mathcal{U}$  using a smooth cut-off function  $\chi(\check{P})$ . The cut-off function is zero outside  $\mathcal{U}$  and equal to one in a large part of the interior of  $\mathcal{U}$ , see Section 4.2.3. The ansatz (4.10) is inspired by Maslov's work [25], although it is not the same since our amplitude function  $\phi$  depends on  $(\hat{X}, \check{X}, x)$  but not on  $\check{P}$ . We emphasize that our modification consisting in having an amplitude function that is not dependent on  $\check{P}$  is essential in the construction of the solution and for determining the accuracy of observables based on this solution.

4.2.1. *Making the ansatz for a Schrödinger solution.* In this section we construct a solution to the Schrödinger equation from the ansatz (4.10). The constructed solution will be an *actual* solution and not only an asymptotic solution as in [25]. We consider first the case when the integration is over  $\mathbb{R}^d$  and then conclude in the end that the cut-off function  $\chi(\check{P})$  can be included in all integrals without changing the property of the Fourier integral ansatz being a solution in the  $\check{X}$ -domain where  $\check{X} = \nabla_{\check{P}}\theta^*(\hat{X}, \check{P})$  for some  $\check{P}$  satisfying  $\chi(\check{P}) = 1$ .

The requirement to be a solution means that there should hold

$$(4.11) \quad \begin{aligned} 0 &= (\mathcal{H} - E)\Phi \\ &= \int_{\mathbb{R}^d} \left( \frac{1}{2} |\nabla_{\hat{X}} \theta^*(\hat{X}, \check{P})|^2 + \frac{1}{2} |\check{P}|^2 + V_0(X) - E \right) \phi(X, x) e^{-iM^{1/2}\Theta(\check{X}, \hat{X}, \check{P})} d\check{P} \\ &\quad - \int_{\mathbb{R}^d} \left( iM^{-1/2} (\nabla_{\hat{X}} \phi \cdot \nabla_{\hat{X}} \theta^* - \nabla_{\hat{X}} \phi \cdot \check{P} + \frac{1}{2} \phi \Delta_{\hat{X}} \theta^*) - (\mathcal{V} - V_0)\phi + \frac{1}{2M} \Delta_X \phi \right) e^{-iM^{1/2}\Theta(\check{X}, \hat{X}, \check{P})} d\check{P}. \end{aligned}$$

Comparing this expression to the previously discussed case of a single WKB-mode we see that the zero order term is now  $\Delta_{\hat{X}} \theta^*$  instead of  $\Delta_X \theta$  and that we have  $-\nabla_{\hat{X}} \phi \cdot \check{P}$  instead of  $\nabla_{\hat{X}} \phi \cdot \nabla_{\hat{X}} \theta$ . However, the main difference is that the first integral is not zero (only the leading order term of its stationary phase expansion is zero, cf. (10.1)). Therefore, the first integral contributes to the second integral. The goal is now to determine a function  $F(\hat{X}, \check{X}, \check{P})$  satisfying

$$(4.12) \quad \begin{aligned} &\int_{\mathbb{R}^d} \left( \frac{1}{2} |\nabla_{\hat{X}} \theta^*|^2 + \frac{1}{2} |\check{P}|^2 + V_0(X) - E \right) e^{-iM^{1/2}\Theta(\check{X}, \hat{X}, \check{P})} d\check{P} \\ &= iM^{-1/2} \int_{\mathbb{R}^d} F(\hat{X}, \check{X}, \check{P}) e^{-iM^{1/2}\Theta(\check{X}, \hat{X}, \check{P})} d\check{P}, \end{aligned}$$

and verify that it is bounded.

**Lemma 4.2.** *There holds  $F = F_0 + F_1$  where*

$$\begin{aligned} F_0 &= \frac{1}{2} \sum_{i,j} \partial_{\hat{X}^i \hat{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P})) \partial_{\check{P}^j \check{P}^i} \theta^*(\check{P}), \\ F_1 &= iM^{-1/2} \int_0^1 \int_0^1 \int_{\mathbb{R}^d} \sum_{i,j,k} t(1-t) \partial_{\check{P}^k} \left[ \partial_{\hat{X}^i \hat{X}^j \hat{X}^k} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + s t \delta \theta^*(\check{P})) \partial_{\check{P}^j \check{P}^i} \nabla_{\check{P}} \theta^*(\check{P}) \right] dt ds. \end{aligned}$$

*Proof.* The function  $\theta^*(\hat{X}, \check{P})$  is defined as a solution to the Hamilton-Jacobi (eikonal) equation

$$(4.13) \quad \frac{1}{2} |\nabla_{\hat{X}} \theta^*(\hat{X}, \check{P})|^2 + \frac{1}{2} |\check{P}|^2 + V_0(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P})) - E = 0$$

for all  $(\hat{X}, \check{P})$ . Consequently, the integral on the left hand side of (4.12) is

$$\int_{\mathbb{R}^d} \left( V_0(\hat{X}, \check{X}) - V_0(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P})) \right) e^{-iM^{1/2}(\check{X} \cdot \check{P} - \theta^*(\hat{X}, \check{P}))} d\check{P}.$$

Let  $\check{P}_0(\check{X})$  be any solution to the stationary phase equation  $\check{X} = \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}_0)$  and introduce the notation

$$\Theta'(\check{X}, \hat{X}, \check{P}) := \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}_0) \cdot \check{P} - \theta^*(\hat{X}, \check{P}).$$



Then by writing a difference as  $V(y_1) - V(y_2) = \int_0^1 \partial_y V(y_2 + t(y_1 - y_2)) dt \cdot (y_1 - y_2)$ , identifying a derivative  $\partial_{\check{P}_i}$  and integrating by parts the integral can be written

$$\begin{aligned}
& \int_{\mathbb{R}^d} \left( V_0(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}_0)) - V_0(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P})) \right) e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} \\
&= \int_0^1 \int_{\mathbb{R}^d} \sum_i \partial_{\check{X}^i} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + t[\nabla_{\check{P}} \theta^*(\check{P}_0) - \nabla_{\check{P}} \theta^*(\check{P})]) \times \\
&\quad \times (\partial_{\check{P}_i} \theta^*(\check{P}_0) - \partial_{\check{P}_i} \theta^*(\check{P})) e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt \\
&= -iM^{-1/2} \int_0^1 \int_{\mathbb{R}^d} \sum_i \partial_{\check{X}^i} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + t[\nabla_{\check{P}} \theta^*(\check{P}_0) - \nabla_{\check{P}} \theta^*(\check{P})]) \partial_{\check{P}_i} e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt \\
&= iM^{-1/2} \int_0^1 \int_{\mathbb{R}^d} \sum_i \partial_{\check{P}_i} \partial_{\check{X}^i} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + t[\nabla_{\check{P}} \theta^*(\check{P}_0) - \nabla_{\check{P}} \theta^*(\check{P})]) e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt.
\end{aligned}$$

Therefore the leading order term in  $F =: F_0 + F_1$  is

$$\begin{aligned}
F_0 &:= \int_0^1 \sum_{i,j} (1-t) \partial_{\check{X}^i \check{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P}) dt \\
&= \frac{1}{2} \sum_{i,j} \partial_{\check{X}^i \check{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P}).
\end{aligned}$$

Denoting  $\delta\theta^*(\check{P}) = \nabla_{\check{P}} \theta^*(\check{P}_0) - \nabla_{\check{P}} \theta^*(\check{P})$  the remainder becomes

$$\begin{aligned}
& -iM^{-1/2} \int_0^1 \int_{\mathbb{R}^d} \sum_{i,j} [\partial_{\check{X}^i \check{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P})) - \partial_{\check{X}^i \check{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + t\delta\theta^*(\check{P}))] \\
&\quad \times (1-t) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P}) e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt \\
&= iM^{-1/2} \int_0^1 \int_0^1 \int_{\mathbb{R}^d} \sum_{i,j,k} t(1-t) \partial_{\check{X}^i \check{X}^j \check{X}^k} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + st\delta\theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P}) \\
&\quad \times (\partial_{\check{P}_k} \theta^*(\check{P}_0) - \partial_{\check{P}_k} \theta^*(\check{P})) e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt ds \\
&= -\frac{1}{M} \int_0^1 \int_0^1 \int_{\mathbb{R}^d} \sum_{i,j,k} t(1-t) \partial_{\check{P}^k} [\partial_{\check{X}^i \check{X}^j \check{X}^k} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + st\delta\theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P})] \\
&\quad \times e^{-iM^{1/2} \Theta'(\check{X}, \hat{X}, \check{P})} d\check{P} dt ds,
\end{aligned}$$

hence the function  $F_1$  is purely imaginary and small

$$F_1 = iM^{-1/2} \int_0^1 \int_0^1 \int_{\mathbb{R}^d} \sum_{i,j,k} t(1-t) \partial_{\check{P}^k} \left[ \partial_{\check{X}^i \check{X}^j \check{X}^k} V_0(\nabla_{\check{P}} \theta^*(\check{P}) + st\delta\theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \nabla_{\check{P}} \theta^*(\check{P}) \right] dt ds,$$

and

$$(4.14) \quad 2\text{Re } F = \sum_{i,j} \partial_{\check{X}^i \check{X}^j} V_0(\nabla_{\check{P}} \theta^*(\check{P})) \partial_{\check{P}_j \check{P}_i} \theta^*(\check{P}).$$

□

The eikonal equation (4.13) and the requirement that  $(\mathcal{H} - E)\Phi = 0$  in (4.11) then imply that

$$\begin{aligned}
(4.15) \quad 0 &= \int_{\mathbb{R}^d} \left[ iM^{-1/2} \left( \nabla_{\check{X}} \phi \cdot \nabla_{\check{X}} \theta^* - \nabla_{\check{X}} \phi \cdot \check{P} + \frac{1}{2} \phi (\Delta_{\check{X}} \theta^* - 2F(X, \check{P})) \right) \right. \\
&\quad \left. - (\mathcal{V} - V_0) \phi + \frac{1}{2M} \Delta_X \phi \right] e^{-iM^{1/2} \Theta(\check{X}, \hat{X}, \check{P})} d\check{P}.
\end{aligned}$$

The Hamilton-Jacobi eikonal equation (4.13), in the primal variable  $(\hat{X}, \check{P})$  with the corresponding dual variable  $(\hat{P}, \check{X})$ , can be solved by the characteristics

$$(4.16) \quad \begin{aligned} \dot{\hat{X}} &= \hat{P} \\ \dot{\hat{P}} &= -\nabla_{\hat{X}} V_0(\hat{X}, \check{X}) \\ \dot{\check{X}} &= -\check{P} \\ \dot{\check{P}} &= \nabla_{\check{X}} V_0(\hat{X}, \check{X}), \end{aligned}$$

using the definition

$$\begin{aligned} \nabla_{\hat{X}} \theta^*(\hat{X}, \check{P}) &= \hat{P} \\ \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}) &= \check{X}. \end{aligned}$$

The characteristics give

$$\frac{d}{dt} \phi = \nabla_{\hat{X}} \phi \cdot \nabla_{\hat{X}} \theta^* - \nabla_{\check{X}} \phi \cdot \check{P},$$

so that the Schrödinger transport equation becomes, as in (3.12),

$$(4.17) \quad iM^{-1/2} \left( \dot{\phi} + \phi \frac{\dot{G}}{G} \right) = (\mathcal{V} - V_0) \phi - \frac{1}{2M} \Delta_X \phi$$

and for  $\psi = G\phi$

$$(4.18) \quad iM^{-1/2} \dot{\psi} = (\mathcal{V} - V_0) \psi - \frac{G}{2M} \Delta_X \frac{\psi}{G}$$

with the complex valued weight function  $G$  defined by

$$(4.19) \quad \frac{d}{dt} \log G_t = \frac{1}{2} \Delta_{\hat{X}} \theta^*(\hat{X}_t, \check{P}_t) - F(\hat{X}_t, \check{P}_t).$$

This transport equation is of the same form as the transport equation for a single WKB-mode, with a modification of the weight function  $G$ .

Differentiation of the second equation in the Hamiltonian system (4.16) implies that the first variation  $\partial \check{P}_t / \partial \check{X}_0$  satisfies

$$\frac{d}{dt} \left( \frac{\partial \check{P}_t^i}{\partial \check{X}_0} \right) = \sum_{j,k} \partial_{\hat{X}^i \check{X}^j} V_0(\hat{X}, \check{X}_t) \partial_{\check{P}^j \check{P}^k} \theta^*(\check{P}) \frac{\partial \check{P}_t^k}{\partial \check{X}_0},$$

which by the Liouville formula (3.13) and the equality

$$2\text{Re } F = \sum_{i,j} \partial_{\hat{X}^i \check{X}^j} V_0 \partial_{\check{P}^j \check{P}^i} \theta^* = \text{Tr} \left( \sum_j \partial_{\hat{X}^i \check{X}^j} V_0 \partial_{\check{P}^j \check{P}^k} \theta^* \right)$$

in (4.14) yields the relation,

$$(4.20) \quad e^{-2 \int_0^t \text{Re } F dt'} = C \left| \det \frac{\partial \check{P}_t}{\partial \check{X}_0} \right|,$$

for the constant  $C := \left| \det \frac{\partial \check{X}_0}{\partial \check{P}_0} \right|$ . We use relation (4.20) to study the density in the next section.

**Remark 4.3.** The conclusion in this section holds also when all integrals over  $d\check{P}$  in  $\mathbb{R}^d$  are replaced by integrals with the measure  $\chi(\check{P}) d\check{P}$ . Then there holds  $2\text{Re } F = \sum_{i,j} \partial_{\hat{X}^i \check{X}^j} \mathcal{V} \partial_{\check{P}^i} (\chi \partial_{\check{P}^j} \theta^*)$ . We use that the observable  $g$  is zero when the cut-off function  $\chi_j$  is not one, see Section 4.2.3. In Section 5 we show how to construct a global solution by connecting the Fourier integral solutions, valid in a neighborhood where  $\det \partial(X)/\partial(P)$  vanishes (and  $\chi(\check{P}) = 1$ ), to a sum of WKB-modes, valid in neighborhoods where  $\det \partial(P)/\partial(X)$  vanishes (and  $\chi(\check{P}) < 1$ ).

4.2.2. *The Schrödinger density for caustics.* In this section we study the density generated by the solution

$$\Phi(X, x) = \int_{\mathbb{R}^d} \phi(X, x) e^{-iM^{1/2}(\tilde{X} \cdot \tilde{P} - \theta^*(\tilde{X}, \tilde{P}))} d\tilde{P}.$$

The analysis of the density generalizes the calculations for the Airy function in Section 4.1.2. We have, using the notation  $\hat{g}$  for the Fourier transform of  $\tilde{g}$  with respect to the  $\tilde{X}$  variable, and by introducing the notation  $\check{R} = \frac{1}{2}(\check{P} + \check{Q})$  and  $\check{S} = \check{P} - \check{Q}$

$$\begin{aligned} \int g(X) |\Phi(x, X)|^2 dx dX &= \int \underbrace{g(X) \langle \phi, \phi \rangle}_{=: \tilde{g}(X)} e^{iM^{1/2}(\tilde{X} \cdot \tilde{P} - \theta^*(\tilde{X}, \tilde{P}))} e^{-iM^{1/2}(\tilde{X} \cdot \check{Q} - \theta^*(\tilde{X}, \check{Q}))} d\check{P} d\check{Q} dX \\ &= \int \hat{g}(\hat{X}, M^{1/2}\check{S}) e^{iM^{1/2}(\theta^*(\hat{X}, \check{Q}) - \theta^*(\hat{X}, \check{P}))} d\check{P} d\check{Q} d\hat{X} \\ &= \int \hat{g}(\hat{X}, M^{1/2}\check{S}) e^{iM^{1/2} \frac{1}{6} (\check{S} \cdot \nabla_{\check{P}})^3 \theta^*(\hat{X}, \check{R} + \gamma\check{S}/2)} \times \\ (4.21) \quad &\quad \times e^{iM^{1/2} \check{S} \cdot \nabla_{\check{P}} \theta^*(\hat{X}, \check{R})} d\check{S} d\check{R} d\hat{X} \\ &= \left(\frac{1}{2\pi}\right)^{d/2} M^{-1/2} \int \tilde{g} * A_M(\hat{X}, \underbrace{\nabla_{\check{P}} \theta^*(\hat{X}, \check{R})}_{=\check{X}}) d\check{R} d\hat{X} \\ &= \left(\frac{1}{2\pi}\right)^{d/2} M^{-1/2} \int \tilde{g} * A_M(\hat{X}, \check{X}) \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| d\check{X}. \end{aligned}$$

In the convolution  $\tilde{g} * A_M$ , the function  $A_M$ , analogous to (4.7), is the Fourier transform of

$$e^{i\frac{1}{M}(\omega \cdot \nabla_{\check{P}})^3 \theta^*(\hat{X}, \check{P})} \Big|_{\check{P} = \check{R} + \gamma\omega}$$

with respect to  $\omega \in \mathbb{R}^d$  and the integration in  $\check{X}$  is with respect to the range of  $\nabla_{\check{P}} \theta^*(\hat{X}, \cdot)$ . As a next step we evaluate the Fourier transform and its derivatives at zero and obtain

$$\begin{aligned} \int_{\mathbb{R}^d} A_M(\check{X}) d\check{X} &= 1, \quad \int_{\mathbb{R}^d} \check{X}^i A_M(\check{X}) d\check{X} = 0, \\ \int_{\mathbb{R}^d} \check{X}^i \check{X}^j A_M(\check{X}) d\check{X} &= 0, \quad M \int_{\mathbb{R}^d} \check{X}^i \check{X}^j \check{X}^k A_M(\check{X}) d\check{X} = \mathcal{O}(1). \end{aligned}$$

Here we use that both differentiation with respect to  $(\omega \cdot \nabla_{\check{P}})^3$  and  $\theta^*(\hat{X}, \check{R} + \gamma\omega)$  yield factors of  $\omega$  which vanish. The vanishing moments of  $A_M$  imply that

$$(4.22) \quad \|\tilde{g} * A_M - \tilde{g}\|_{L^2(d\check{X})} = \mathcal{O}(M^{-1})$$

as in (4.8), so that up to  $\mathcal{O}(M^{-1})$  error the convolution with  $A_M$  can be neglected.

4.2.3. *Integration over a compact set in  $\check{P}$ .* In the case when the integration is over  $\mathcal{U} \subset \mathbb{R}^d$  instead of  $\mathbb{R}^d$ , we use a smooth cut-off function  $\chi(\check{P})$ , which is zero outside  $\mathcal{U}$  and restrict our analysis to the case when the smooth observable mapping  $\check{P} \mapsto g(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}))$  is compactly supported in the domain where  $\chi$  is one. In this way  $g(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{P}))$  is zero when  $\nabla_{\check{P}} \chi(\check{P})$  is non zero. The integrand is thus equal to

$$(g(X) \langle \phi, \phi \rangle) \chi(\check{P}) \chi(\check{Q})$$

and we use the convergent Taylor expansion

$$\chi(\underbrace{\check{R} + M^{-1/2}\omega}_{\check{P}}) \chi(\underbrace{\check{R} - M^{-1/2}\omega}_{\check{Q}}) = \sum_{k=0}^{\infty} \frac{|\omega|^{2k}}{M^k} a_k(\check{R}).$$

Then the observable becomes

$$(2\pi)^{-d/2} M^{-1/2} \sum_{k=0}^{\infty} \int (a_k(M^{-1} \Delta_{\check{X}})^k \tilde{g}) * A_M(\hat{X}, \nabla_{\check{P}} \theta^*(\hat{X}, \check{R})) d\check{R} d\hat{X}.$$

As in (4.22) we can remove the convolution with  $A_M$  by introducing an error  $\mathcal{O}(M^{-1})$  and since for  $k > 0$  we have  $a_k(\check{R})g(\hat{X}, \nabla_{\check{P}}\theta^*(\hat{X}, \check{R})) = 0$  and  $a_0 = 1$ , we obtain the same observable as before

$$\begin{aligned} & \sum_{k=0}^{\infty} \int (a_k (M^{-1} \Delta_{\check{X}})^k \tilde{g}) * A_M(\hat{X}, \nabla_{\check{P}}\theta^*(\hat{X}, \check{R})) d\check{R} d\hat{X} \\ &= \sum_{k=0}^{\infty} \int (a_k (M^{-1} \Delta_{\check{X}})^k \tilde{g}) \left( \hat{X}, \nabla_{\check{P}}\theta^*(\hat{X}, \check{R}) \right) d\check{R} d\hat{X} + \mathcal{O}(M^{-1}) \\ &= \int \tilde{g} \left( \hat{X}, \nabla_{\check{P}}\theta^*(\hat{X}, \check{R}) \right) d\check{R} d\hat{X} + \mathcal{O}(M^{-1}) \\ &= \int \tilde{g}(\hat{X}, \check{X}) \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| dX + \mathcal{O}(M^{-1}). \end{aligned}$$

4.2.4. *Comparing the Schrödinger and molecular dynamics densities.* We compare the Schrödinger density to the molecular dynamics density generated by the continuity equation

$$0 = \operatorname{div}(\rho \nabla \theta) = \nabla \rho \cdot \nabla \theta + \rho \operatorname{div}(\nabla \theta) = \dot{\rho} + \rho \operatorname{div}(\nabla \theta)$$

which yields the density

$$e^{-\int \operatorname{div}(\nabla \theta) dt}.$$

We have  $P = \nabla \theta$ , so that  $\frac{\partial(P)}{\partial(X)} = \partial_{X X} \theta$ . The Liouville formula (3.13) implies the molecular dynamics density

$$(4.23) \quad \rho_{\text{BO}} = e^{-\int_0^t \operatorname{div}(\nabla \theta) dt'} = \det \frac{\partial X_{0, \text{BO}}}{\partial X_{t, \text{BO}}}.$$

The observable for the Schrödinger equation has, by (4.21), the density

$$(g\langle \phi, \phi \rangle) * A_M \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right|.$$

We want to compare it with the molecular dynamics density  $\rho_{\text{BO}}$ . The convolution with  $A_M$  gives an error term of the order  $\mathcal{O}(M^{-1})$ , as in (4.8). The Schrödinger transport equation (4.17) and the definition of the weight  $G$  in (4.19), show that the amplitude function satisfies, by (4.18) and (4.19) and the Born-Oppenheimer approximation Lemma 8.2,

$$\langle \phi, \phi \rangle = |G|^2 \langle \psi, \psi \rangle = e^{\int 2\operatorname{Re} F - \Delta_{\check{X}} \theta^* dt} + \mathcal{O}(M^{-1}),$$

so that by (4.20)

$$\begin{aligned} (4.24) \quad & (g\langle \phi, \phi \rangle) * A_M \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| = (g\langle \phi, \phi \rangle) \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| + \mathcal{O}(M^{-1}) \\ &= g e^{\int 2\operatorname{Re} F - \Delta_{\check{X}} \theta^* dt} \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| + \mathcal{O}(M^{-1}) \\ &= g \left| \det \frac{\partial(\check{X}_0)}{\partial(\check{P})} \right| \left| \det \frac{\partial(\hat{X}_0)}{\partial(\hat{X})} \right| \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right| + \mathcal{O}(M^{-1}), \\ &= g \left| \det \frac{\partial(\check{X}_0)}{\partial(\check{X})} \right| \left| \det \frac{\partial(\hat{X}_0)}{\partial(\hat{X})} \right| + \mathcal{O}(M^{-1}), \\ &= g \left| \det \frac{\partial(X_0)}{\partial(X)} \right| + \mathcal{O}(M^{-1}). \end{aligned}$$

When we restrict the domain to  $\mathcal{U}$  with the cut-off function  $\chi$  as in Remark 4.3 we use the fact that  $g(\hat{X}, \nabla_{\check{P}}\theta^*(\hat{X}, \check{P}))$  is zero when  $\nabla_{\check{P}}\chi(\check{P})$  is non zero and obtain the same. The representations (4.24) and (4.23) show that the density generated in the caustic case with a Fourier integral also takes the same form, to the leading order, as the molecular dynamics density and the remaining discrepancy is only due to  $\theta^* = \theta_{\text{S}}^*$  and  $\theta^* = \theta_{\text{BO}}^*$  being different. This difference is, as in the single mode WKB expansion, of size  $\mathcal{O}(M^{-1})$  which is estimated by the difference in Hamiltonians of the Schrödinger and molecular dynamics eikonal

equations. The estimate of the difference of the phase functions uses the Hamilton-Jacobi equation (4.13) for  $\theta_S^*(\hat{X}, \check{P})$  and a similar Hamilton-Jacobi equation for  $\theta_{\text{BO}}^*(\hat{X}, \check{P})$  with  $V_0 = \lambda_{\text{BO}} + \mathcal{O}(M^{-1})$  replaced by  $\lambda_{\text{BO}}$ . The difference in the weight functions  $\log(|G(\hat{X}, \check{P})|^{-2})$  is estimated by the Hamilton-Jacobi equation

$$\left( \nabla_{\hat{X}} \theta_S^*(\hat{X}, \check{P}) \cdot \nabla_{\hat{X}} - \nabla_{\hat{X}} V_0(\hat{X}, \check{X}) \cdot \nabla_{\check{P}} \right) \log |G_S(\hat{X}, \check{P})|^{-2} - \Delta_{\hat{X}} \theta_S^*(X, \check{P}) + \text{Re } F(X, \check{P}) = 0,$$

where  $\text{Re } F$  is given in (4.14), and by the similar Hamilton-Jacobi equation with  $V_0 = \lambda_{\text{BO}} + \mathcal{O}(M^{-1})$  replaced by  $\lambda_{\text{BO}}$  and  $\theta_S^*$  by  $\theta_{\text{BO}}^*$ .

## 5. A GLOBAL CONSTRUCTION COUPLING CAUSTICS WITH SINGLE WKB-MODES

We use the paths generated by the Hamiltonian to construct a global asymptotic solution. More precisely, we will construct a  $3N$ -dimensional Lagrangian manifold from a Hamiltonian; for this we need some initial data, i.e. a  $3N - 1$  dimensional Lagrangian manifold, which we construct from a Poincare map in four steps. The basic properties of a Lagrangian manifold we use are reviewed in Section 5.1 and the basic assumption we make is that the Lagrangian manifold is smooth, which excludes ergodic dynamics where the Lagrangian manifold is dense in set of dimension  $6N - 1$ .

**Step 1. Define a hitting plane.** Consider a codimension one hitting plane in the  $X$ -coordinate space, e.g.  $X_{1_1} = 0$ .

**Step 2. Consider a WKB Schrödinger solution and its initial Lagrangian manifold data.** We seek initial data on the hitting plane  $X_{1_1} = 0$  in the form of a (smooth)  $3N - 1$  dimensional Lagrangian manifold  $L$ , satisfying the two constraints  $H(X, P) = E$  &  $X_{1_1} = 0$  for  $(X, P) \in L$ . For instance any smooth function  $G : \mathbb{R}^{3N} \rightarrow \mathbb{R}$  generates a local subset of a Lagrangian manifold

$$L \supseteq \{(X, P(X)) \mid \text{for all } X \text{ such that } H(X, P(X)) = E, \& X_{1_1} = 0, P(X) = \nabla_X G\}$$

and one can permute the role of  $X$  and  $P$  to obtain other parts of the set. Section 3.1.1 shows, in the case of piecewise constant potentials, that the solution to the Schrödinger equation is of the form  $\Phi(X) = \sum_j e^{iM^{1/2}P^j \cdot X} \phi_j(X)$ , where  $H(X, P^j) = E$ . Assume that a solution to the Schrödinger equation has the asymptotic WKB-form in the hitting plane

$$(5.1) \quad \Phi(X) = \sum_{\nu=1}^K \phi_\nu(X, x) e^{iM^{1/2}\theta_\nu(X)} + \mathcal{O}(M^{-n}), \text{ for } X_{1_1} = 0 \text{ and for all } n \in \mathbb{N},$$

where  $\phi_\nu$  and  $\theta_\nu$  are smooth functions, based on a finite sum of WKB-solutions (3.1) (or the caustic ansatz (4.10)).

**Step 3. Use paths to extend the initial Lagrangian manifold.** Use the characteristics paths in Theorem 3.1 for the WKB-functions  $\phi_\nu(X_t)$  and  $\theta_\nu(X_t)$  (or (4.13) and (4.16) for the caustic case  $\theta_\nu(X) = \check{X} \cdot \check{P} - \theta^*(\hat{X}, \check{P})$ ) to locally extend the guessed initial  $3N - 1$  dimensional Lagrangian manifold, in the hitting plane  $X_{1_1} = 0$ , to dimension  $3N$  outside  $X_{1_1} = 0$ , by starting a path from each point on the initial Lagrangian manifold; change coordinates in a neighborhood of a caustic and apply the stationary phase method in Step 4 to continue the solutions until the first hitting time  $\tau$ , for all possible initial  $\phi_\nu$  and  $\theta_\nu$ . Here  $\tau$  is the first time the path  $(X_t, P_t)$  leaves the set  $X_{1_1} < 0$  if it initially went into the set  $X_{1_1} > 0$  (i.e. if  $\lim_{t \rightarrow 0^+} \text{sign}(X_{1_1}(t)) = 1$ ), and similarly the first time it leaves  $X_{1_1} > 0$  if  $\lim_{t \rightarrow 0^+} \text{sign}(X_{1_1}(t)) = -1$ . For all initial  $\phi_\nu$  and  $\theta_\nu$  this yields an asymptotic Schrödinger solution

$$(5.2) \quad \Phi(X) \simeq \sum_{\nu} \phi_\nu(X, x) e^{iM^{1/2}\theta_\nu(X)}$$

for all  $X_{1_1} \neq 0$ , since by construction the WKB integral is an asymptotic solution

$$(\mathcal{H} - E) \sum_{\nu} \phi_\nu(X, x) e^{iM^{1/2}\theta_\nu(X)} = \mathcal{O}(M^{-n})$$

in the domain  $X_{1_1} \neq 0$  and we assume there exists a stable global solution  $\Phi$  (including  $X_{1_1} = 0$ ) to the Schrödinger equation (2.2), for  $E$  chosen to be an eigenvalue with a distance bigger than  $\mathcal{O}(M^{-n})$  to other eigenvalues. The particular  $\phi_\nu$  and  $\theta_\nu$  that gives the asymptotic eigensolution satisfies

$$(5.3) \quad \lim_{X_{1_1} \rightarrow 0^+} \Phi(X) = \lim_{X_{1_1} \rightarrow 0^-} \Phi(X).$$

We have used the WKB-method to characterize the global solution. In particular the initial Lagrangian manifold, given by  $\{\theta_\nu \mid \forall \nu\}$ , is by the characteristics mapped to the same manifold at the hitting surface  $X_{1_1} = 0$ , since otherwise (5.3) and  $\lim_{X_{1_1} \rightarrow 0^+} \theta_\nu(X) = \lim_{X_{1_1} \rightarrow 0^-} \theta_\nu(X)$  cannot hold. The obtained fixed point initial data  $\Phi$  is smooth, since it is a solution of the Schrödinger eigenvalue problem. We assume that also the amplitudes  $\phi_\nu$  and the phases  $\theta_\nu$  are smooth functions and that there are only finitely many of them. The construction of WKB solutions makes it possible to compare their Lagrangian manifolds by their Hamiltonians. We interpret the molecular dynamics as the formal limit  $M \rightarrow \infty$  of the WKB system.

**Step 4. Apply Maslov's matching around caustics.** We see that the weight function  $G$ , in (3.13), based on a single WKB-mode (3.1) blows up at caustics, where  $\det(\partial(\tilde{X})/\partial(\tilde{P})) = 0$ , and that the weight function  $G$  in (4.17) for the Fourier integral (4.10) blows up at points where  $\det(\partial(\tilde{P})/\partial(\tilde{X}))$  vanishes. Therefore, in neighborhoods around caustic points we need to use the representation  $\theta^*(\tilde{X}, \tilde{P})$  of the phase based on the Fourier integrals, while around points where  $\det(\partial(\tilde{P})/\partial(\tilde{X}))$  vanishes we apply the representation  $\theta(\tilde{X}, \tilde{X})$  based on the Legendre transform, as pointed out by Maslov in [25] and described in the simplifying setting of the harmonic oscillator in [11].

One way to make a global construction of a WKB solution, in the spirit of [25], is to use the characteristics and a partition of the phase-space as follows, also explained constructively by the numerical algorithm 9.2 in the next section. Start with a Fourier integral representation in a neighborhood  $\mathcal{U}$  of a caustic point, which gives a representation of the Schrödinger solution  $\Phi$  in  $\mathcal{U}$ . Then we use the stationary phase expansion, see Section 10, to find an asymptotic approximation  $\tilde{\Phi}$  (accurate to any order  $n \in \mathbb{N}$ ) at the boundary points  $\tilde{X}$  of  $\mathcal{U}$  as a sum of single WKB-modes with phase functions  $\theta_j$

$$\int_{\mathbb{R}^d} \chi(\tilde{P}) e^{-iM^{1/2}(\tilde{X} \cdot \tilde{P} - \theta^*(\tilde{X}, \tilde{P}))} d\tilde{P} = \sum_j e^{-iM^{1/2}\theta_j(X)} \phi_j(X) + \mathcal{O}(M^{-n})$$

where each phase function  $\theta_j(X) := \tilde{X} \cdot \tilde{P}_{X,j} - \theta^*(\tilde{X}, \tilde{P}_{X,j})$  corresponds to a branch of the boundary and the index  $j$  corresponds to different solutions  $\tilde{P}_{X,j}$  of the stationary phase equation  $\tilde{X} = \nabla_{\tilde{P}} \theta^*(\tilde{X}, \tilde{P}_X)$ . The single WKB-modes  $\phi(x, X) e^{iM^{1/2}\theta(X)}$  are then constructed along the characteristics to be Schrödinger solutions in a domain around the point where  $\det(\partial(\tilde{P})/\partial(\tilde{X}))$  vanishes, following the construction in Theorem 3.1 using the initial data of  $\tilde{\Phi}$  at  $\partial\mathcal{U}$ . We note that the tiny error of size  $\mathcal{O}(M^{-n})$  that we make in the initial data for  $\phi$  also yields a tiny perturbation error in  $\phi$  of size  $\mathcal{O}(M^{-n})$  along the path, due to the assumption of the  $\mathcal{O}(1)$  bounded hitting times. A small error we make in the expansion therefore leads to a negligible error in the Schrödinger solution and the corresponding density.

When a characteristic leaves the domain and enters another region around a caustic we again use the stationary phase method at the boundary to give initial data for  $(X, P, \phi, G)$ . When the characteristic finally returns to the first boundary  $\partial\mathcal{U}$ , there is a compatibility condition to have a global solution, by having the incoming final phase equal to the initial phase function in  $\mathcal{C}^1$ . We can think of this as trying to find a co-dimension one surface  $I$  in  $\mathbb{R}^{3N}$  where the incoming and outgoing phases are equal. First to have one point where they agree is possible if we restrict the possible solutions to a discrete set of energies  $E$ , i.e., the eigenvalues, and therefore the compatibility condition is called a quantization condition. Then, having one point where the difference of the two phase function is zero, we can combine this with the assumption that the Lagrangian manifold generated by the characteristics path  $(X_t, P_t)$  is continuous: the two phases have the same gradient on  $I$ , since  $(X, P) = (X, \nabla_X \theta(X)) = \left( (\tilde{X}, \nabla_{\tilde{P}} \theta^*(\tilde{X}, \tilde{P})), (\nabla_{\tilde{X}} \theta^*(\tilde{X}, \tilde{P}), \tilde{P}) \right)$  so the phases are  $\mathcal{C}^1$ . If the Lagrangian manifold would be simple connected, the compatibility condition is always satisfied by the construction of the Lagrangian manifold; here we assume that one parameter is enough to describe the non simple connectedness. In this way we define the  $(X, P, \phi, G)$  globally, for the eigenvalue energies  $E$ . To evaluate observables we use a partition of unity to restrict the observable to a domain with a single representation, either a Fourier integral representation for a caustic or a single WKB-mode when  $\det(\partial(\tilde{P})/\partial(\tilde{X})) = 0$ .

**5.1. Lagrangian manifolds.** This section presents some basic properties of Lagrangian manifolds. Given the  $3N - 1$  dimensional Lagrangian manifold  $L$  on the hitting surface  $X_{1_1} = 0$ , the solution paths

$$\{(X_t, P_t) \in \mathbb{R}^{6N} \mid 0 \leq t < \infty, \forall (X_0, P_0) \in L\}$$

of the Hamiltonian system

$$\begin{aligned}\dot{X}_t &= \nabla_P H(P_t, X_t) \\ \dot{P}_t &= -\nabla_X H(P_t, X_t)\end{aligned}$$

with a smooth and bounded Hamiltonian  $H(P, X)$  generate a  $3N$ -dimensional manifold called Lagrangian manifold. The fact that the Lagrangian manifold has dimension  $3N$  implies that it can locally be described by  $(X, P)$  with  $P$  as a function of  $X$  or with  $X$  as a function of  $P$  and in general  $3N$  coordinates are functions of the other  $3N$  coordinates.

The Lagrangian manifold generated by the tube of trajectories is defined by the phase function  $\theta(X)$  that plays the role of a generating function of the Lagrangian manifold. Thus we seek a function  $\theta : U \subset \mathbb{R}^{3N} \rightarrow \mathbb{R}$  such that  $P_t = \nabla_X \theta(X_t)$ . We show that there exists a potential function  $\theta$  by determining an equation that preserves the symmetry for the matrix  $Q_t$ , defined as  $Q^{ij}(X) := \partial_{X^j} P^i(X)$  and  $Q_t^{ij} := Q^{ij}(X_t)$ . The relations  $P_t^i = P^i(X_t)$  and  $Q_t^{ij} := \partial_{X^j} P^i(X_t)$  imply

$$\dot{P}_t^i = \frac{d}{dt} P^i(X_t) = \sum_j \dot{X}_t^j \partial_{X^j} P_t^i = \sum_j \dot{X}_t^j Q_t^{ij} = \sum_j \partial_{P^j} H(P(X_t), X_t) Q_t^{ij},$$

so that

$$\begin{aligned}\partial_{X^k} \dot{P}_t^i &= \partial_{X^k} \left( \sum_j \partial_{P^j} H(P(X_t), X_t) Q_t^{ij} \right) \\ &= \underbrace{\sum_j \dot{X}_t^j \partial_{X^k} Q_t^{ij}}_{=\sum_j \dot{X}_t^j \partial_{X^k X^j} P_t^i = \sum_j \dot{X}_t^j \partial_{X^j X^k} P_t^i = \dot{Q}_t^{ik}} + \sum_j \partial_{P^j P^l} H(P(X_t), X_t) \underbrace{\partial_{X^k} P^l}_{=Q^{lk}} Q_t^{ij} \\ &= \sum_j \dot{X}_t^j \partial_{X^k X^j} P_t^i = \sum_j \dot{X}_t^j \partial_{X^j X^k} P_t^i = \dot{Q}_t^{ik} \\ &+ \sum_j \partial_{P^j X^k} H(P(X_t), X_t) Q_t^{ij}\end{aligned}$$

and

$$\partial_{X^k} \dot{P}_t^i = -\partial_{X^k} \left( \partial_{X^i} H(P(X_t), X_t) \right) = -\partial_{X^i X^k} H(P(X_t), X_t) - \sum_j \partial_{X^i P^j} H(P(X_t), X_t) \underbrace{\partial_{X^k} P^j}_{=Q^{jk}}$$

together with the symmetry of  $Q_t$  show that

$$(5.4) \quad \begin{aligned}\dot{Q}_t^{ik} &= -\partial_{X^i X^k} H(P_t, X_t) - \sum_{j,l} \partial_{P^j P^l} H(P_t, X_t) Q_t^{kl} Q_t^{ij} \\ &- \sum_j \partial_{P^j X^k} H(P_t, X_t) Q_t^{ij} - \sum_j \partial_{P^j X^i} H(P_t, X_t) Q_t^{kj}.\end{aligned}$$

Since the Hamiltonian is assumed to be smooth it follows that the right hand side in (5.4) is symmetric and thus the matrix  $Q_t$  remains symmetric since it is initially symmetric. Hence there exists a potential function  $\theta(X)$  such that  $P(X) = \nabla_X \theta(X)$  in simple connected domains where  $Q$  is smooth. The function  $Q$  may become unbounded due to the term  $\partial_{P^j P^l} H Q^{kl} Q^{ij}$ , even though  $H$  has bounded third derivatives. Points  $X_t$  at which  $|\text{Tr}(Q_t)| = \infty$  satisfy, by Liouville's theorem (see Section 3.1.4),  $\left| \det \frac{\partial X_0}{\partial X_t} \right| = \infty$  and such points are called *caustic* points.

The same construction of a potential works for the local chart expressed as  $X = X(P)$  instead of  $P = P(X)$ . In fact any new variable  $\hat{X}$  (not including both  $X^i$  and  $P^i$  for any  $i$ ), based on  $3N$  of the  $6N$  variables  $(X, P)$ , and the remaining variables  $3N$  variables,  $\hat{P}$ , represent the same Hamiltonian system with the Hamiltonian  $\hat{H}(\hat{P}, \hat{X}) := H(P, X)$ . The Lagrangian manifold is defined by  $\hat{P} = \nabla_{\hat{X}} \hat{\theta}(\hat{X})$  in the local chart of  $\hat{P}$ -coordinates with the generating (potential) function  $\hat{\theta}(\hat{X})$  defined in domains excluding caustics, i.e., where  $\det \left| \frac{\partial \hat{X}_0}{\partial \hat{X}_t} \right| < \infty$ . Maslov, [25], realized that a Lagrangian manifold can be partitioned, by changing coordinates in the neighborhood of a caustic, into domains where  $\hat{P} = \nabla_{\hat{X}} \hat{\theta}(\hat{X})$  is smooth. He used the generating (potential) functions  $\hat{\theta}$  to construct asymptotic WKB solutions of Fourier integral type. A sketch of this general situation is depicted in Figure 4. In previous sections we have described construction of

solutions in a simpler case without caustics, i.e.,  $P_t = \nabla_X \theta(X_t)$  holds everywhere. In this section we described the global construction of WKB solutions in the general case when caustics are present.

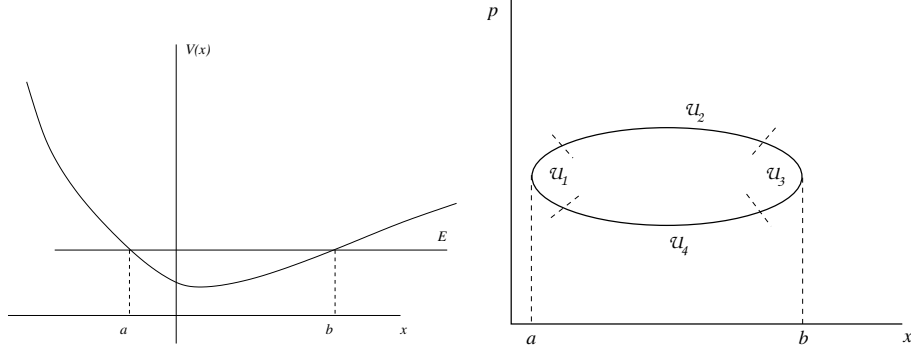


FIGURE 4. The left figure depicts a graph of the molecular dynamics potential  $\lambda(X)$  in the case which exhibits caustics at  $X = a$  and  $X = b$  for a given energy  $E$ . The right figure shows a general case of the Lagrangian manifold with two caustic points  $X = a$  and  $b$  and its covering with charts  $\mathcal{U}_i$ . In the charts  $\mathcal{U}_i$ ,  $i = 2, 4$  the manifold is defined by  $P = \nabla_X \theta_i(X)$  and the solution to Schrödinger equation is constructed by simple WKB modes. The caustics belong to the charts  $\mathcal{U}_i$ ,  $i = 1, 3$  and in this case the manifold is defined by  $X = \nabla_X \theta_i(P)$  and the solutions are given by the Fourier integrals.

## 6. COMPUTATION OF OBSERVABLES

Suppose the goal is to compute a real-valued *observable*

$$\int_{\mathbb{T}^{3N}} \langle \Phi, A\Phi \rangle dX$$

for a given bounded linear multiplication operator  $A = A(X)$  on  $L^2(\mathbb{T}^{3N})$  and a solution  $\Phi = \sum_k \phi_k e^{iM^{1/2}\theta_k}$  of (2.2). We have

$$(6.1) \quad \begin{aligned} \int_{\mathbb{T}^{3N}} \langle \Phi, A\Phi \rangle dX &= \sum_{k,l} \int_{\mathbb{T}^{3N}} \langle A\phi_k e^{iM^{1/2}\theta_k(X)}, \phi_l e^{iM^{1/2}\theta_l(X)} \rangle dX \\ &= \sum_{k,l} \int_{\mathbb{T}^{3N}} A e^{iM^{1/2}(\theta_l(X) - \theta_k(X))} \langle \phi_k, \phi_l \rangle dX. \end{aligned}$$

The integrand is oscillatory for  $k \neq l$ , hence critical points (or near critical points) of the phase difference give the main contribution. The stationary phase method, see [10, 25] and Section 10, shows that these integrals are small, bounded by  $\mathcal{O}(M^{-3N/4})$ , in the case when the phase difference has non degenerate critical points, or no critical point, and the functions  $A\langle \phi_k, \phi_l \rangle$  and  $\theta_l$  are sufficiently smooth. A critical point  $X_c \in \mathbb{R}^{3N}$  satisfies  $\nabla_X \theta_l(X_c) - \nabla_X \theta_k(X_c) = 0$ , which means that the two different paths, generated by  $\theta_l$  and  $\theta_k$ , passing through  $X = X_c$  also have the same momentum  $P$  at this point. That the critical point is degenerate means that the Hessian matrix  $\partial_{X^i X^j}(\theta_k - \theta_l)(X_c)$  is singular (or asymptotically singular for  $M \rightarrow \infty$  as for avoided crossings when the electron eigenvalues have a vanishing spectral gap depending on  $M$ ). Therefore caustics, crossing or avoided crossing electron eigenvalues may generate coupling between the WKB terms. On the other hand, without such coupling the density of a linear combination of WKB terms separates asymptotically to a sum of densities of the individual WKB terms

$$(6.2) \quad \int_{\mathbb{T}^{3N}} \langle \Phi, A\Phi \rangle dX = \int_{\mathbb{T}^{3N}} A \underbrace{\sum_{k=1}^{\bar{k}} \langle \phi_k, \phi_k \rangle}_{=\rho_k} dX + \mathcal{O}(M^{-1}),$$

=  $\rho$



in the case of multiple WKB-functions,  $\bar{k} > 1$ , and

$$\int_{\mathbb{T}^{3N}} \langle \Phi, A\Phi \rangle dX = \int_{\mathbb{T}^{3N}} A \langle \phi_1, \phi_1 \rangle dX$$

for a single WKB-function and we have seen in (4.24) that the Fourier integral WKB-solutions around caustics yields the same density as a single WKB mode.

In the presence of a caustic, the WKB terms can be asymptotically non orthogonal, since their coefficients and phases typically are not smooth enough to allow the integration by parts to gain powers of  $M^{-1/2}$ . Non-orthogonal WKB functions tell how the caustic couples the WKB modes.

Regarding the inflow density  $\rho_k|_I$  there are two situations: either the characteristics return often to the inflow domain or not. If they do not return we have a scattering problem and it is reasonable to define the inflow-density  $\rho_k|_I$  as an initial condition. If characteristics return, the dynamics can be used to estimate the return-density  $\rho_k|_I$  as follows: Assume that the following limits exist

$$(6.3) \quad \lim_{T \rightarrow \infty} \frac{1}{T} \int_0^T A(X_t) dt = \int_{\mathbb{T}^{3N}} A(X) \rho_k(X) dX$$

which bypasses the need to find  $\rho_k|_I$  and the quadrature in the number of characteristics. When there are multiple amplitudes  $\phi_k$  in (6.2) we can, away from caustics, obtain the sum of the densities  $\sum_k \rho_k$  from the time average split into the time the dynamics spends in each phase  $\theta_k$ . A way to think about the time average (6.3) is to sample the return points  $X_t \in I$  and from these samples construct an *empirical* return-density, converging to  $\rho_k|_I$  as the number of return iterations tends to infinity. We shall use this perspective to view the eikonal equation (3.8) as a hitting problem on  $I$ , with hitting times  $\tau$  (i.e., return times). The stronger property having  $\rho_k$  constant as a function of  $(X_0, P_0)$  for  $H(X_0, P_0) = E$  is called *ergodicity*, see [30]. We allow the density  $\rho_k|_I$  to depend on the initial position  $X_0$  and momentum  $P_0$  and then our observables are conditional expected values. An example of a hitting surface is the co-dimension one surface where the first component  $X_{1_1}$  in  $X_1 = (X_{1_1}, X_{1_2}, X_{1_3})$  is equal to its initial value  $X_{1_1}(0)$ . The dynamics does not always have such a hitting surface: for instance if all particles are close initially and then are scattered away from each other, as in an explosion, no co-dimension one hitting surface exists.

## 7. MOLECULAR DYNAMICS APPROXIMATION OF SCHRÖDINGER OBSERVABLES

A numerical computation of an approximation to  $\sum_k \int_{\mathbb{T}^{3N}} \langle \phi_k, A\phi_k \rangle dX$  has the main ingredients:

- (1) to approximate the exact characteristics by molecular dynamics characteristics (3.10),
- (2) to discretize the molecular dynamics equations, and
- (3a) if  $\rho|_I$  is an inflow-density, to introduce quadrature in the number of characteristics, or
- (3b) if  $\rho|_I$  is a return-density, to replace the ensemble average by a time average using the property (6.3).

This section presents a derivation of the approximation error in the step (1) in the case of a return density and comments on the time-discretization of step (2) treated in Section 7.2. The third and fourth discretization steps, which are not described here, are studied, for instance, in [8, 7, 21].

**7.1. The Born-Oppenheimer approximation error.** This section states our main result of molecular dynamics approximating Schrödinger observables based on hitting time and spectral gap assumptions. We formulate it using the assumption of the Born-Oppenheimer property

$$(7.1) \quad \|\psi_t - \Psi_{\text{BO}}(X_t)\|_{L^2(dx)} = \mathcal{O}(M^{-1/2}), \quad \text{uniformly in } t.$$

This assumption is then proved in Lemma 8.2 based on a setting with a spectral gap.

*The spectral gap condition.* The electron eigenvalues  $\{\lambda_k\}$  satisfy, for some positive  $c$ , the spectral gap condition

$$(7.2) \quad \inf_{k \neq 0, Y \in D} |\lambda_k(Y) - \lambda_0(Y)| > c,$$

where  $D := \{X_S(t) | t \geq 0\} \cup \{X_{\text{BO}}(t) | t \geq 0\}$  is the set of all nuclei positions obtained from the Schrödinger characteristics  $X = X_S$  in Theorem 3.1 and from the Born-Oppenheimer dynamics  $X = X_{\text{BO}}$  in (3.16), for all considered initial data.

*The hitting time assumption.* In the case of quasiperiodic dynamics, the points  $(X_{t_k}, P_{t_k})$  in a codimension one hitting plane, e.g.  $X_1^1 \equiv X_{1_1} = 0$ , accumulate on a set of dimension  $3N - 1$ , see Figure 5, while for ergodic dynamics the hitting points are distributed in the phase-space set  $\{(X, P) \in \mathbb{R}^{6N} \mid H(X, P) = E \ \& \ X_{1_1} = 0\}$  of dimension  $6N - 2$  with positive density, see Figure 6.

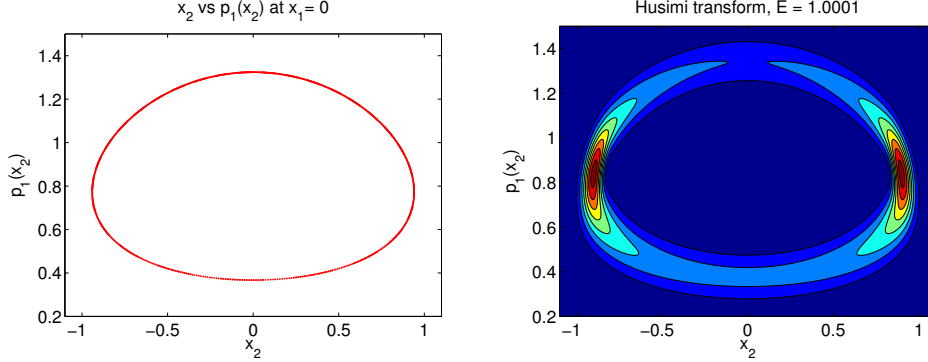


FIGURE 5. Left plot: The hitting points  $(X_2, P_1)$  in the plane  $X_1 = 0$  that accumulate on a curve of dimension one, for the Hamiltonian  $H = |P|^2/2 + |X_1|^2/2 + |X_2|^2/2^{1/2} + 0.3 \sin(X_1 X_2) = 1.0001$ . Right plot: The Husimi transform,  $M^{1/2} \left| \int_{\mathbb{R}} \Phi(0 - y, X_2) e^{-iM^{1/2}yP_1} e^{-M^{1/2}|y|^2/2} dy \right|^2$ , of the Schrödinger solution,  $\Phi(X)$ , as a function of  $X_2$  and  $P_1$  in the plane  $X_1 = 0$ , with the potential  $V(X) = |X_1|^2/2 + |X_2|^2/2^{1/2} + 0.3 \sin(X_1 X_2)$  and the eigenvalue  $E = 1.0001$  for mass  $M = 8000$ . The Husimi transform measures the density in  $X_2$  and  $P_1$  and equals the absolute value squared of the FBI transform, which here is integrated in the  $P_2$  direction.

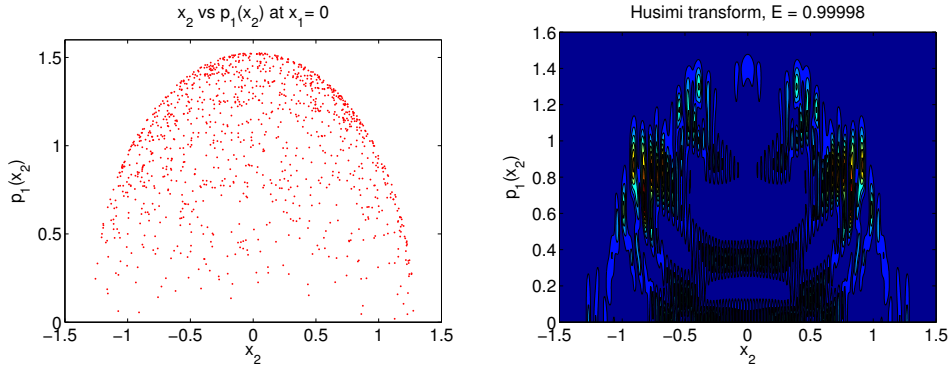


FIGURE 6. Left: The hitting points  $(X_2, P_1)$  in the plane  $X_1 = 0$  that become dense set of dimension two, for the Hamiltonian  $H = |P|^2/2 + |X_1|^2/2 + |X_2|^2/2^{1/2} + 2 \sin(X_1 X_2) = 1$ . Right: The Husimi transform of the Schrödinger solution in the plane  $X_1 = 0$  for the same potential and eigenvalue  $E = 0.99998$ .

In the quasiperiodic setting it is therefore reasonable to assume that the hitting time (i.e. the return time) is finite, since only a uniformly bounded set of phases in the hitting plane is visited.

**Theorem 7.1.** *Assume that there is an asymptotic solution to the Schrödinger equation (2.2) that can be written as a finite sum of WKB-modes (5.1) (as for caustics (4.10)) and that all phase functions  $\theta_S$  and  $\theta_{BO}$  are smooth solutions to the Schrödinger eikonal equation (written as (3.8) or (4.13)) and the molecular dynamics eikonal equation (3.17), respectively, satisfying the Born-Oppenheimer property (7.1) and the limit (6.3), with uniformly bounded hitting times  $\tau$ , in (8.12), (8.16), and (8.20), then the zero-order*

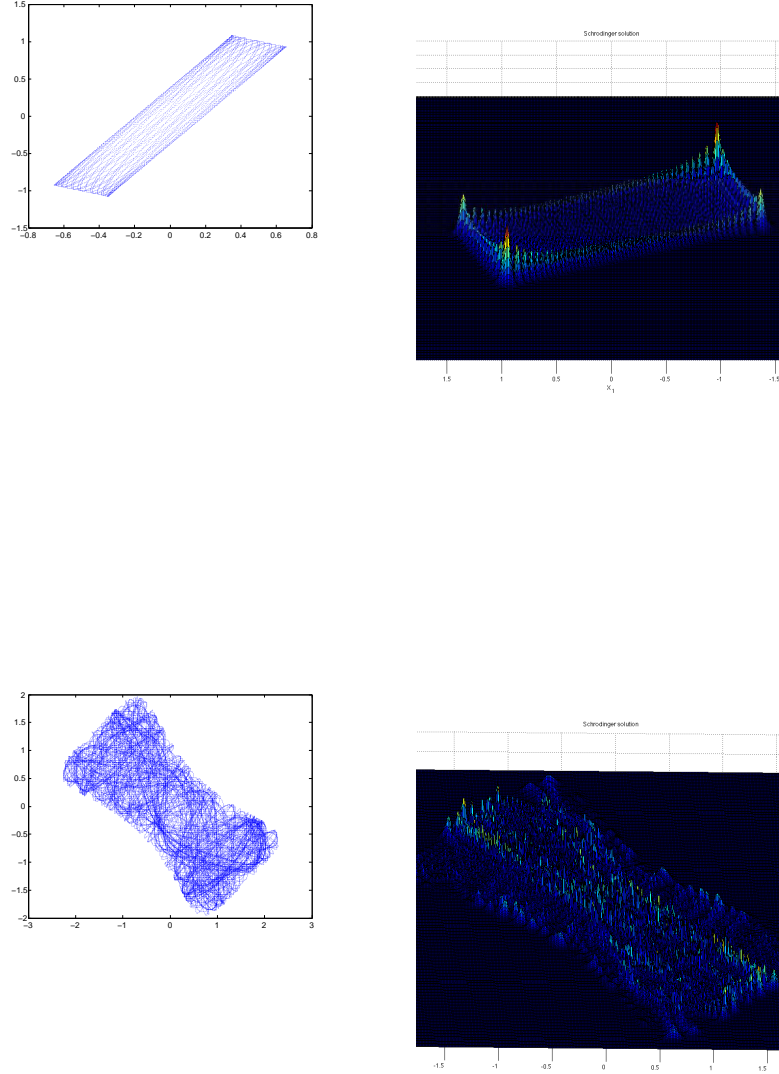


FIGURE 7. Left: the path  $(X_1(t), X_2(t))$  for  $\beta = 0.3$  and  $\beta = 2$ . Right: the Schrödinger density  $|\Phi(X)|^2$  for  $\beta = 0.3$  and  $\beta = 2$ .

*Born-Oppenheimer dynamics (3.16), with initial data generated by the construction in Steps 1-4 in Section 5, approximates time-independent Schrödinger observables, with error bounded by  $\mathcal{O}(M^{-1+\delta})$ , i.e.*

$$(7.3) \quad \int_{\mathbb{T}^{3N}} g(X) \rho_{\text{BO}}(X) dX = \int_{\mathbb{T}^{3N}} g(X) \rho_{\text{S}}(X) dX + \mathcal{O}(M^{-1+\delta}), \quad \text{for any } \delta > 0 \text{ and } g \in \mathcal{C}^3(\mathbb{R}^{3N}).$$

The proof is given in Section 8.

**7.2. Why do symplectic numerical simulations of molecular dynamics work?** The derivation of the approximation error for the Born-Oppenheimer dynamics, in Theorem 7.1, also allows to study perturbed systems. For instance, the perturbed Born-Oppenheimer dynamics

$$\begin{aligned}\dot{X}_t &= P_t + \nabla_P H^\epsilon(X_t, P_t) \\ \dot{P}_t &= -\nabla_X \lambda_0(X_t) - \nabla_X H^\epsilon(X_t, P_t),\end{aligned}$$

generated from a perturbed Hamiltonian  $H_{\text{BO}}(X, P) + H^\epsilon(X, P) = E$ , with the perturbation satisfying

$$(7.4) \quad \|H^\epsilon\|_{L^\infty} \leq \epsilon \quad \text{for some } \epsilon \in (0, \infty)$$

yields through (8.14) and (8.21) an additional error term  $\mathcal{O}(\epsilon)$  to the approximation of observables in (7.3). So called symplectic numerical methods are precisely those that can be written as perturbed Hamiltonian systems, see [31], and consequently we have a method to precisely analyze their numerical error by combining an explicit construction of  $H^\epsilon$  with the stability condition (7.4) to obtain  $\mathcal{O}((M^{-1} + \epsilon)^{1-\delta})$  accurate approximations, provided the corresponding phase function has bounded second difference quotients. The popular Störmer-Verlet method is symplectic and the positions  $X$  coincides with those of the symplectic Euler method, for which  $H^\epsilon$  is explicitly constructed in [31] with  $\epsilon$  proportional to the time step. The construction in [31] is not using the modified equation and formal asymptotics, instead a piecewise linear extension of the solution generates  $H^\epsilon$ .

## 8. ANALYSIS OF THE MOLECULAR DYNAMICS APPROXIMATION

Before we proceed with the analysis of the approximation error we motivate our results by a significantly simpler case of a system *without electrons*. We use the densities (3.18) and (3.19) and we show heuristically how the characteristics can be used to estimate the difference  $\rho_S - \rho_{\text{BO}}$ , leading to  $\mathcal{O}(M^{-1})$  accurate Born-Oppenheimer approximations of Schrödinger observables

$$\int g(X) \underbrace{\rho_S(X)}_{\langle \Phi, \Phi \rangle} dX = \int g(X) \rho_{\text{BO}}(X) dX + \mathcal{O}(M^{-1}).$$

In the special case of no electrons, the dynamics of  $X$  does not depend on  $\psi$  and therefore  $X_{\text{BO}} = X_S = X$  and consequently  $G_{\text{BO}} = G_S$ . The difference  $\psi_S - \psi_{\text{BO}}$  can be understood from iterative approximations of (3.12)

$$(8.1) \quad \frac{i}{M^{1/2}} \dot{\psi}_{k+1} - (\mathcal{V} - V_0) \psi_{k+1} = \frac{1}{2M} G \Delta_X (G^{-1} \psi_k)$$

with  $\psi_0 = 0$ . Then  $\psi_{\text{BO}} = \psi_1$  is the Born-Oppenheimer approximation and formally we have the iterations approaching the full Schrödinger solution as  $k \rightarrow \infty$ .

In the special case of no electrons, there holds  $\mathcal{V} = V_0$ , thus the transport equation  $i\dot{\psi}_1 = 0$  has constant solutions. We let  $\psi_1 = 1$  and then  $\psi_2 - \psi_1$  is imaginary with its absolute value bounded by  $\mathcal{O}(M^{-1/2})$ . We write the iterations of  $\psi_k$  by integrating (8.1) as the linear mapping

$$\psi_{k+1} = 1 + iM^{-1/2} \hat{\mathcal{S}}(\psi_k) = \sum_{l=0}^k i^l M^{-l/2} \hat{\mathcal{S}}^l(\psi_1),$$

which formally shows that

$$|\psi_S|^2 = |\psi_1|^2 + 2\text{Re} \langle \psi_S - \psi_1, \psi_1 \rangle + |\psi_S - \psi_1|^2 = 1 + \mathcal{O}(M^{-1}).$$

Consequently this special Born-Oppenheimer density satisfies

$$(8.2) \quad \rho_{\text{BO}} = \underbrace{G_S^{-2} \langle \psi_S, \psi_S \rangle}_{=\rho_S} + \mathcal{O}(M^{-1}),$$

since  $G_{\text{BO}} = G_S$  and  $X$  do not depend on  $\psi$ .

In the general case with electrons and a spectral gap, we show in Lemma 8.2 that there is a solution  $\psi_S$  satisfying

$$(8.3) \quad \|\psi_S - \Psi_{\text{BO}}\|_{L^2(dx)} = \mathcal{O}(M^{-1/2}),$$

for the electron eigenfunction  $\Psi_{\text{BO}}$ , satisfying

$$\mathcal{V}(\cdot, X)\Psi_{\text{BO}}(\cdot, X) = \lambda_0(X)\Psi_{\text{BO}}(\cdot, X)$$

and the eigenvalue  $\lambda_0(X) \in \mathbb{R}$  with a (fixed) nuclei position  $X$ . Then the state  $\psi_1$  equal to a constant, in the case of no electrons, corresponds to the electron eigenfunction  $\Psi_{\text{BO}}$  in the case with electrons present. In the general case the  $X$  dynamics for the Schrödinger and the Born-Oppenheimer dynamics are not the same, but we will show that (8.3) implies that the Hamiltonians  $H_S$  and  $H_{\text{BO}}$  are  $\mathcal{O}(M^{-1})$  close. Using stability of Hamilton-Jacobi equations, the phase functions  $\theta_S$  and  $\theta_{\text{BO}}$  are then also close in the maximum norm, which, combined with an assumption of smooth phase functions, show that  $|G_S - G_{\text{BO}}| = \mathcal{O}(M^{-1+\delta})$  for any  $\delta > 0$ . Lemma 8.2 also shows that  $|\langle \psi_S, \psi_S \rangle - 1| = \mathcal{O}(M^{-1})$  and consequently the density bound  $|\rho_S - \rho_{\text{BO}}| = \mathcal{O}(M^{-1+\delta})$  holds. To obtain the estimate (8.3) the important new property, compared to no electrons, is to use oscillatory cancellation in directions orthogonal to  $\Psi_{\text{BO}}$ .

**8.1. Continuation of the construction of the solution operator.** This section continues the construction of the solution operator started in Section 3.4. Assume for a moment that  $\tilde{\mathcal{V}}$  is independent of  $\tau$ . Then the solution to (3.23) can be written as a linear combination of the two exponentials

$$ae^{i\tau\mathcal{A}_+} + be^{i\tau\mathcal{A}_-}$$

where the two characteristic roots are the operators

$$\mathcal{A}_{\pm} = (P_1^1)^2 \left( -1 \pm (1 - 2(P_1^1)^{-2}\tilde{\mathcal{V}})^{1/2} \right).$$

We see that  $e^{i\tau\mathcal{A}_-}$  is a highly oscillatory solution on the fast  $\tau$ -scale with

$$\lim_{P_1^1 \rightarrow \infty} \frac{1}{(P_1^1)^2} \mathcal{A}_- = -2\text{Id},$$

while

$$(8.4) \quad \lim_{P_1^1 \rightarrow \infty} \mathcal{A}_+ = -\tilde{\mathcal{V}},$$

in distribution sense. Therefore we chose initial data

$$(8.5) \quad i\dot{\psi}|_{\tau=0} = -\mathcal{A}_+\psi|_{\tau=0}$$

to have  $b = 0$ , which eliminates the fast scale, and the limit  $P_1^1 \rightarrow \infty$  determines the solution by the Schrödinger equation

$$i\dot{\psi} = \tilde{\mathcal{V}}\psi.$$

The next section presents an analogous construction for the slowly, in  $\tau$ , varying operator  $\tilde{\mathcal{V}}$ .

**8.1.1. Spectral decomposition.** Write (3.23) as the first order system

$$\begin{aligned} i\dot{\psi} &= \pi \\ i\dot{\pi} &= -2(P_1^1)^2(\tilde{\mathcal{V}}\psi - \pi), \end{aligned}$$

which for  $\bar{\psi} := (\psi, \pi)$  takes the form

$$\dot{\bar{\psi}} = i\mathcal{B}\bar{\psi}, \quad \mathcal{B} := \begin{pmatrix} 0 & -1 \\ 2(P_1^1)^2\tilde{\mathcal{V}} & -2(P_1^1)^2 \end{pmatrix},$$

where the eigenvalues  $\Lambda_\pm$ , right eigenvectors  $\mathcal{Q}_\pm$  and left eigenvectors  $\mathcal{Q}_\pm^{-1}$  of the real “matrix” operator  $\mathcal{B}$  are

$$\begin{aligned}\Lambda_\pm &:= (P_1^1)^2 \left( -\text{Id} \pm \left( \text{Id} - 2(P_1^1)^{-2} \tilde{\mathcal{V}} \right)^{1/2} \right), \\ \mathcal{Q}_+ &:= \begin{pmatrix} \text{Id} \\ -\Lambda_+ \end{pmatrix}, \quad \mathcal{Q}_- := \begin{pmatrix} -\Lambda_-^{-1} \\ \text{Id} \end{pmatrix}, \\ \mathcal{Q}_+^{-1} &:= (\text{Id} - \Lambda_+ \Lambda_-^{-1})^{-1} \begin{pmatrix} \text{Id} \\ \Lambda_-^{-1} \end{pmatrix}, \quad \mathcal{Q}_-^{-1} := (\text{Id} - \Lambda_+ (\Lambda_-)^{-1})^{-1} \begin{pmatrix} \Lambda_+ \\ \text{Id} \end{pmatrix}.\end{aligned}$$

We see that  $\lim_{P_1^1 \rightarrow \infty} \Lambda_+ = -\tilde{\mathcal{V}}$  and  $\lim_{P_1^1 \rightarrow \infty} (P_1^1)^{-2} \Lambda_- = -2\text{Id}$ . The important property here is that the left eigenvector limit  $\lim_{P_1^1 \rightarrow \infty} \mathcal{Q}_+^{-1} = (\text{Id}, 0)$  is constant, independent of  $\tau$ , which implies that the  $\mathcal{Q}_+$  component  $\mathcal{Q}_+^{-1} \bar{\psi} = \psi$  decouples. We obtain in the limit  $P_1^1 \rightarrow \infty$  the time-dependent Schrödinger equation

$$\begin{aligned}i\dot{\psi}(\tau) &= i \frac{d}{d\tau} (\mathcal{Q}_+^{-1} \bar{\psi}_\tau) = i \mathcal{Q}_+^{-1} \frac{d}{d\tau} \bar{\psi}_\tau = -\mathcal{Q}_+^{-1} \mathcal{B}_\tau \bar{\psi}_\tau \\ &= -\Lambda_+(\tau) \mathcal{Q}_+^{-1} \bar{\psi}_\tau = -\Lambda_+(\tau) \psi(\tau) = \tilde{\mathcal{V}}_\tau \psi(\tau),\end{aligned}$$

where the operator  $\tilde{\mathcal{V}}_\tau$  depends on  $\tau$  and  $(x, X_0)$ , and we define the solution operator  $\mathcal{S}$

$$(8.6) \quad \psi(\tau) = \mathcal{S}_{\tau,0} \psi(0).$$

As in (8.5) we can view this as choosing special initial data for  $\psi(0)$ . From now on we only consider such data.

The operator  $\tilde{\mathcal{V}}$  can be symmetrized

$$(8.7) \quad \bar{\mathcal{V}}_\tau := G_\tau^{-1} \tilde{\mathcal{V}}_\tau G_\tau = (\mathcal{V} - V_0)_\tau - \frac{1}{2M} \sum_j \Delta_{X_j^i},$$

with real eigenvalues  $\{\check{\lambda}_m\}$  and orthonormal eigenvectors  $\{\zeta^m\}$  in  $L^2(dx dX_*)$ , satisfying

$$\bar{\mathcal{V}}_\tau \zeta^m(\tau) = \check{\lambda}_m(\tau) \zeta^m(\tau).$$

Therefore  $\tilde{\mathcal{V}}_\tau$  has the same eigenvalues and the eigenvectors  $\bar{\zeta}^m := G_\tau \zeta^m$ , which establishes the spectral representation

$$(8.8) \quad \tilde{\mathcal{V}}_\tau \psi(\cdot, \tau, \cdot) = \sum_m \check{\lambda}_m(\tau) \int_{\mathbb{T}^{3N-1}} \langle \psi(\cdot, \tau, \cdot), \bar{\zeta}^m \rangle G_\tau^{-2} dX_* \bar{\zeta}^m(\tau).$$

We note that the weight  $G^{-2}$  on the co-dimension one surface  $\mathbb{T}^{3N-1}$  appears precisely because the operator  $\tilde{\mathcal{V}}$  is symmetrized by  $G^{-2}$  and the weight  $G^{-2}$  corresponds to the Eulerian-Lagrangian change of coordinates (3.13)

$$(8.9) \quad \int_{\mathbb{T}^{3N-1}} \langle \psi, \bar{\zeta}^m \rangle G_\tau^{-2} dX_* = \int_{\mathbb{T}^{3N-1}} \langle \psi, \bar{\zeta}^m \rangle dX_0.$$

The existence of the orthonormal set of eigenvectors and real eigenvalues makes the operator  $\tilde{\mathcal{V}}$  self-adjoint in the Lagrangian coordinates and hence the solution operator  $\mathcal{S}$  becomes unitary in the Lagrangian coordinates.

**8.2. Stability from perturbed Hamiltonians.** In this section we derive error estimates of the weight functions  $G$  when the corresponding Hamiltonian system is perturbed. To derive the stability estimate we consider the Hamilton-Jacobi equation

$$(8.10) \quad H(\nabla_X \theta(X), X) = 0$$

in an optimal control perspective with the corresponding Hamiltonian system

$$\begin{aligned}\dot{X}_t &= \nabla_P H(P_t, X_t) \\ \dot{P}_t &= -\nabla_X H(P_t, X_t).\end{aligned}$$

We define the “value” function

$$\theta(X_0) = \theta(X_t) - \int_0^t h(P_s, X_s) ds,$$

where the “cost” function defined by

$$h(P, X) := P \cdot \nabla_P H(P, X) - H(P, X)$$

satisfies the Pontryagin principle (related to the Legendre transform)

$$(8.11) \quad H(P, X) = \sup_Q (P \cdot \nabla_Q H(Q, X) - h(Q, X)).$$

Let  $\theta|_I$  be defined by the hitting problem

$$\theta(X_0) = \theta(X_\tau) - \int_0^\tau h(P_s, X_s) ds$$

using the hitting time  $\tau$  on the return surface  $I$

$$(8.12) \quad \tau := \inf\{t \mid X_0 \in I, X_t \in I \& t > 0\}.$$

For a perturbed Hamiltonian  $\tilde{H}$  and its dynamics  $(\tilde{X}_t, \tilde{P}_t)$  we define analogously the value function  $\tilde{\theta}$  and the cost function  $\tilde{h}$ .

8.2.1. *Estimates of the phase functions.* We can think of the difference  $\theta - \tilde{\theta}$  as composed by a perturbation of the boundary data (on the return surface  $I$ ) and perturbations of the Hamiltonians. The difference of the value functions due to the perturbed Hamiltonian satisfies the stability estimate

$$(8.13) \quad \begin{aligned} \theta(X_0) - \tilde{\theta}(X_0) &\geq \theta(\tilde{X}_{\tilde{\tau}}) - \tilde{\theta}(\tilde{X}_{\tilde{\tau}}) + \int_0^{\tilde{\tau}} (H - \tilde{H}) \left( \nabla_X \theta(\tilde{X}_t), \tilde{X}_t \right) dt \\ \theta(X_0) - \tilde{\theta}(X_0) &\leq \theta(X_\tau) - \tilde{\theta}(X_\tau) + \int_0^\tau (H - \tilde{H}) \left( \nabla_X \tilde{\theta}(X_t), X_t \right) dt \end{aligned}$$

with a difference of the Hamiltonians evaluated along the same solution path. This result follows by differentiating the value function along a path and using the Hamilton-Jacobi equations, see Remark 8.1 and [9].

The global construction of a Hamilton-Jacobi equation in Section 5 shows that the Hamiltonian  $H(X, P)$  remains the same while the parametrization  $(X, P) = (\hat{X}, \hat{X}, \hat{P}, \hat{P})$  of the Lagrangian manifold may vary in different domains - e.g.  $P = P(X) = \nabla_X \theta(X)$  or  $(\hat{P}, \hat{X}) = (\nabla_{\hat{X}} \theta^*(\hat{X}, \hat{P}), \nabla_{\hat{P}} \theta^*(\hat{X}, \hat{P}))$ - due to the effect of caustics. At the interface to such domains the phase function satisfies a boundary condition, which determines the phase function with an asymptotically vanishing error of size  $\mathcal{O}(M^{-n})$  for any natural number  $n$ . This small error in the initial data of the phase function is negligible compared to other sources of error.

We assume that

$$(8.14) \quad \sup_{(P, X) \in L(H) \cup L(\tilde{H})} |(H - \tilde{H})(P, X)| = \mathcal{O}(M^{-1}),$$

which is verified in (8.15) for Schrödinger and Born-Oppenheimer Hamiltonians; here  $L(H)$  is the Lagrangian manifold for the Hamiltonian  $H$ . We consider the case when the return time is uniformly bounded, related to quasiperiodic dynamics.

The Hamiltonians we use are

$$\begin{aligned} H_S &= \frac{|P|^2}{2} + \frac{\langle \psi(X), \mathcal{V}(X) \psi(X) \rangle}{\langle \psi(X), \psi(X) \rangle} - E, \\ H_{BO} &= \frac{|P|^2}{2} + \lambda_0(X) - E, \end{aligned}$$

based on the cost functions

$$\begin{aligned} h_S &= E + \frac{|P|^2}{2} - \frac{\langle \psi(X), \mathcal{V}(X) \psi(X) \rangle}{\langle \psi(X), \psi(X) \rangle}, \\ h_{BO} &= E + \frac{|P|^2}{2} - \lambda_0(X). \end{aligned}$$

For the Born-Oppenheimer case the electron wave function is the eigenstate  $\Psi_{BO}$ . The Born-Oppenheimer approximation (7.1), proved in Lemma 8.2, implies that

$$(8.15) \quad \|H_S - H_{BO}\|_{L^\infty} = \mathcal{O}(M^{-1}),$$

which verifies (8.14).

*Uniformly bounded hitting times.* We choose the  $3N - 1$  dimensional hitting set as

$$(8.16) \quad I := \{X \in \mathbb{T}^{3N} \mid \theta(X) = \tilde{\theta}(X)\}$$

on which the two phases coincide. If  $\theta(X)$  is a solution also  $\theta(X) + C$  is a solution to the Hamilton-Jacobi equation (8.10), for any constant  $C$ ; therefore we can choose  $\theta$  such that  $\theta(X_0) = \tilde{\theta}(X_0)$  for any  $X_0$ . Now assume that  $I$  forms a codimension one set in  $\mathbb{T}^{3N}$  and that the maximal hitting time  $\tau$  for characteristics starting on  $I$  is bounded; the fact that  $I$  is a codimension one set holds, for instance, locally if  $|\nabla_X(\theta - \tilde{\theta})|$  is nonzero. In fact, it is sufficient to assume that there exists a function  $\gamma : \mathbb{T}^{3N} \rightarrow \mathbb{R}$ , satisfying  $\gamma = \mathcal{O}(M^{-1})$ , such that the set  $I := \{X \in \mathbb{T}^{3N} \mid \theta(X) - \tilde{\theta}(X) = \gamma(X)\}$  is a codimension one set with bounded hitting times. We also define the  $3N - 1$  dimensional hitting sets in phase space

$$\begin{aligned} I_\theta &:= \{(X, P) \mid X \in I, P = \nabla\theta(X)\}, \\ I_{\tilde{\theta}} &:= \{(X, P) \mid X \in I, P = \nabla\tilde{\theta}(X)\}. \end{aligned}$$

Then the representation (8.13), for any time  $t$  replacing  $\tau$  and  $\tilde{\tau}$ , together with the stability of the Hamiltonians (8.14) and the initial data  $(\theta - \tilde{\theta})|_I = 0$  obtained from (8.16) imply that

$$(8.17) \quad \|\theta - \tilde{\theta}\|_{L^\infty} = \mathcal{O}(M^{-1}),$$

provided the maximal hitting time  $\tau$  is bounded: we assume the hitting time is uniformly bounded in the case when the Lagrangian manifold  $L$  for the Schrödinger Hamiltonian system in Section 5 has a smooth limit, as  $M$  tends to infinity, and the hitting points on a codimension one plane in phase space accumulate on a  $3N - 1$  dimensional subset of phase space; this case is related to quasiperiodic molecular dynamics systems.

**8.2.2. Estimates of the densities.** To estimate the density, we will use the characteristic paths. When the value functions  $\theta$  and  $\tilde{\theta}$  are smoothly differentiable in  $X$ , with derivatives bounded uniformly in  $M$ , the stability estimate (8.13) implies that also the difference of the second derivatives has the bound

$$(8.18) \quad \|\Delta_X\theta - \Delta_X\tilde{\theta}\|_{L^\infty} = \mathcal{O}(M^{-1+\delta}), \quad \text{for any } \delta > 0.$$

Our goal is to analyze the density function  $\rho = |G|^{-2}\langle\psi, \psi\rangle$  with  $G$  defined in (3.11). The Born-Oppenheimer approximation (7.1) yields  $\langle\psi, \psi\rangle = 1 + \mathcal{O}(M^{-1})$  thus it remains to estimate the weight function  $|G|^{-2}$ . This weight function satisfies the Hamilton-Jacobi equation

$$(8.19) \quad H_G(\nabla_X \log |G|^{-2}, X) := \nabla_X\theta(X) \cdot \nabla_X \log |G|^{-2} + \Delta_X\theta(X) = 0.$$

The stability of Hamilton-Jacobi equations can then be applied to (8.19), as in (8.13), using now the hitting set

$$(8.20) \quad I := \{X \in \mathbb{T}^{3N} \mid \log |G(X)|^{-2} = \log |\tilde{G}(X)|^{-2}\}$$

and the assumption of bounded hitting times  $\tau$  in the hitting problem, and we obtain

$$(8.21) \quad \|\log |G|^{-2} - \log |\tilde{G}|^{-2}\|_{L^\infty} \leq C\|H_G - H_{\tilde{G}}\|_{L^\infty} = \mathcal{O}(M^{-1+\delta}).$$

In this sense we will use that an  $\mathcal{O}(M^{-1})$  perturbation of the Hamiltonian yields an error estimate of almost the same order for the difference of the corresponding densities  $\rho - \tilde{\rho}$ .



**Remark 8.1.** This remark derives the stability estimate (8.13). The definitions of the value functions imply

$$\begin{aligned}
& \underbrace{\tilde{\theta}(\tilde{X}_{\tilde{\tau}}) - \int_0^{\tilde{\tau}} \tilde{h}(\tilde{P}_t, \tilde{X}_t) dt}_{\tilde{\theta}(\tilde{X}_0)} - \underbrace{\left( \theta(X_\tau) - \int_0^\tau h(P_t, X_t) dt \right)}_{\theta(X_0)} \\
&= - \int_0^{\tilde{\tau}} \tilde{h}(\tilde{P}_t, \tilde{X}_t) dt + \theta(\tilde{X}_{\tilde{\tau}}) - \underbrace{\theta(X_0)}_{\theta(\tilde{X}_0)} + \tilde{\theta}(\tilde{X}_{\tilde{\tau}}) - \theta(\tilde{X}_{\tilde{\tau}}) \\
(8.22) \quad &= - \int_0^{\tilde{\tau}} \tilde{h}(\tilde{P}_t, \tilde{X}_t) dt + \int_0^{\tilde{\tau}} d\theta(\tilde{X}_t) + \tilde{\theta}(\tilde{X}_{\tilde{\tau}}) - \theta(\tilde{X}_{\tilde{\tau}}) \\
&= \int_0^{\tilde{\tau}} \underbrace{-\tilde{h}(\tilde{P}_t, \tilde{X}_t) + \nabla_X \theta(\tilde{X}_t) \cdot \nabla_P \tilde{H}(\tilde{P}_t, \tilde{X}_t)}_{\leq \tilde{H}(\nabla_X \theta(\tilde{X}_t), \tilde{X}_t)} dt + \tilde{\theta}(\tilde{X}_{\tilde{\tau}}) - \theta(\tilde{X}_{\tilde{\tau}}) \\
&\leq \int_0^{\tilde{\tau}} (\tilde{H} - H) \left( \nabla_X \theta(\tilde{X}_t), \tilde{X}_t \right) dt + \tilde{\theta}(\tilde{X}_{\tilde{\tau}}) - \theta(\tilde{X}_{\tilde{\tau}}),
\end{aligned}$$

where the Pontryagin principle (8.11) yields the inequality and we use the Hamilton-Jacobi equation

$$H(\nabla_X \theta(\tilde{X}_t), \tilde{X}_t) = 0.$$

To establish the lower bound we replace  $\theta$  along with  $\tilde{X}_t$  by  $\tilde{\theta}$  and  $X_t$  and repeat the derivation above.

**8.3. The Born-Oppenheimer approximation.** The purpose of this section is to present a case when the Born-Oppenheimer approximation holds in the sense that  $\|\psi - \Psi_{\text{BO}}\|_{L^2(dx)}$  is small.

We know from Section 8.1.1 that the solution  $\psi_t = \mathcal{S}_{t,0}\psi_0$  is bounded in  $L^2(dx dX)$ , since  $\mathcal{S}$  is unitary in the Lagrangian coordinates. This unitary  $\mathcal{S}$  implies that the integral in the Lagrangian coordinates  $\int_{\mathbb{T}^{3N-1}} \langle \psi_t, \psi_t \rangle dX_0$  is constant in time. We consider the co-dimension one set

$$I_\psi := \{X \in \mathbb{R}^{3N} \mid \langle \psi(X), \psi(X) \rangle = \int_{\mathbb{T}^{3N-1}} \langle \psi(t, X_0), \psi(t, X_0) \rangle dX_0 / \int_{\mathbb{T}^{3N-1}} dX_0\},$$

where the point values of  $\langle \psi(X), \psi(X) \rangle$  coincides with its  $L^2$  average. We choose a time  $t$  such that  $X_t \in I_\psi$  and assume that the time  $\tau^*$  it takes to hit  $I_\psi$  the next time is bounded, i.e.,

$$\tau^* := \inf\{\tau \mid X_t \in I_\psi, \tau > 0 \ \& \ X_{t+\tau} \in I_\psi\} = \mathcal{O}(1).$$

We also assume that all functions of  $X$  are smooth.

**Lemma 8.2.** *Assume that  $i\dot{\psi} = M^{1/2}\tilde{\mathcal{V}}\psi$  holds, then there exists initial data for  $\psi$  such that the  $L^2(dx)$  orthogonal decomposition  $\psi = \bar{\psi}_0 \oplus \psi_0^\perp$ , where  $\bar{\psi}_0 = \alpha\Psi_{\text{BO}}$  for some  $\alpha \in \mathbb{C}$  satisfies*

$$\begin{aligned}
(8.23) \quad & \frac{\|\psi_0^\perp(t)\|_{L^2(dx)}}{\|\bar{\psi}_0(t)\|_{L^2(dx)}} = \mathcal{O}(M^{-1/2}) \\
& |\langle \psi_t, \psi_t \rangle - 1| = \mathcal{O}(M^{-1}) \\
& \|\psi_t - \Psi_{\text{BO}}(X_t)\|_{L^2(dx)} = \mathcal{O}(M^{-1/2})
\end{aligned}$$

uniformly in time  $t$ , provided the spectral gap condition (7.2) holds, the smoothness estimate (8.30) is satisfied and the hitting time  $\tau^*$  is bounded.

*Proof.* We consider the decomposition  $\psi = \bar{\psi}_0 \oplus \psi_0^\perp$ , where  $\bar{\psi}_0(\tau)$  is an eigenfunction of  $\mathcal{V}(X_\tau)$  in  $L^2(dx)$ , satisfying  $\mathcal{V}(X_\tau)\bar{\psi}_0(\tau) = \lambda_0(\tau)\bar{\psi}_0(\tau)$  for the eigenvalue  $\lambda_0(\tau) \in \mathbb{R}$ . This *ansatz* is motivated by the zero residual

$$(8.24) \quad \mathcal{R}\psi := \dot{\psi} + iM^{1/2}\tilde{\mathcal{V}}\psi = 0$$

and the small residual for the eigenfunction

$$\begin{aligned}
& \langle \Pi(\dot{\bar{\psi}}_0), \bar{\psi}_0 \rangle = 0 \\
& M^{1/2}\tilde{\mathcal{V}}\bar{\psi}_0 = \mathcal{O}(M^{-1/2}),
\end{aligned}$$

where

$$(8.25) \quad w(X) = \langle \Psi_{\text{BO}}(X), w(X) \rangle \Psi_{\text{BO}}(X) \oplus \Pi w(X)$$

denotes the orthogonal decomposition in the eigenfunction direction  $\Psi_{\text{BO}}$  and its orthogonal complement in  $L^2(dx)$ . We consider first the linear operator  $\mathcal{R}$  in (8.24) with a given function  $V_0$  and then we use a contraction setting to show that  $V_0 = \langle \psi, \mathcal{V}\psi \rangle / \langle \psi, \psi \rangle$  also works since  $\|\bar{\psi}_0^\perp\|_{L^2(dx)}$  is small. The orthogonal splitting  $\psi = \bar{\psi}_0 \oplus \psi_0^\perp$  and the projection  $\Pi(\cdot)$  in (8.25) imply

$$\begin{aligned} 0 &= \Pi(\mathcal{R}(\bar{\psi}_0 + \psi_0^\perp)) \\ &= \Pi(\mathcal{R}(\bar{\psi}_0)) + \Pi(\mathcal{R}(\psi_0^\perp)) \\ &= \Pi(\mathcal{R}\bar{\psi}_0) + \psi_0^\perp + iM^{1/2}(\mathcal{V} - V_0)\psi_0^\perp + i\Pi\left(\frac{GM^{-1/2}}{2}\Delta_X(G^{-1}\psi_0^\perp)\right), \end{aligned}$$

where the last step follows from the orthogonal splitting

$$\Pi((\mathcal{V} - V_0)\psi_0^\perp) = (\mathcal{V} - V_0)\psi_0^\perp$$

together with the second order change in the subspace projection

$$\psi_0^\perp(\tau + \Delta\tau) = \Pi(\tau + \Delta\tau)(\psi_0^\perp(\tau + \Delta\tau)) = \Pi(\tau)(\psi_0^\perp(\tau + \Delta\tau)) + \mathcal{O}(\Delta\tau^2)$$

which yields  $\Pi(\dot{\psi}_0^\perp) = \dot{\psi}_0^\perp$ ; here  $\Pi(\tau)\cdot$  denotes the projection on the orthogonal complement to the eigenvector  $\bar{\psi}_0(\tau)$ . To explain the second order change start with a function  $v$  satisfying  $\langle v, \Psi_{\text{BO}}(X_\tau) \rangle = 0$  and  $\Psi_{\text{BO}}(X_\sigma) = \Psi_{\text{BO}}(X_\tau) + \mathcal{O}(\Delta\tau)$  for  $\sigma \in [\tau, \tau + \Delta\tau]$  to obtain

$$\begin{aligned} \Pi(\sigma)(\Pi(\tau + \Delta\tau)v - \Pi(\tau)v) &= \Pi(\sigma)\left(\langle v, \Psi_{\text{BO}}(X_\tau) \rangle \Psi_{\text{BO}}(X_\tau) - \langle v, \Psi_{\text{BO}}(X_{\tau+\Delta\tau}) \rangle \Psi_{\text{BO}}(X_{\tau+\Delta\tau})\right) \\ &= \Pi(\sigma)\mathcal{O}(\Delta\tau^2) + \Pi(\sigma)\left(\langle v, \mathcal{O}(\Delta\tau) \rangle \Psi_{\text{BO}}(X_\tau)\right) \\ &= \mathcal{O}(\Delta\tau^2) + \mathcal{O}(\Delta\tau)\left(\Psi_{\text{BO}}(X_\tau) - \langle \Psi_{\text{BO}}(X_\tau), \Psi_{\text{BO}}(X_\sigma) \rangle \Psi_{\text{BO}}(X_\sigma)\right) \\ &= \mathcal{O}(\Delta\tau^2). \end{aligned}$$

Let  $\tilde{\mathcal{S}}_{\tau,\sigma}$  be the solution operator from time  $\sigma$  to  $\tau$  for the generator

$$v \mapsto iM^{1/2}(\mathcal{V} - V_0)v + i\Pi\left(\frac{GM^{-1/2}}{2}\Delta_X(G^{-1}v)\right) =: iM^{1/2}\hat{\mathcal{V}}v.$$

Consequently, the perturbation  $\psi_0^\perp$  can be determined from the projected residual

$$\dot{\psi}_0^\perp = -iM^{1/2}\hat{\mathcal{V}}\psi_0^\perp - \Pi(\mathcal{R}\bar{\psi}_0)$$

and we have the solution representation

$$(8.26) \quad \psi_0^\perp(\tau) = \tilde{\mathcal{S}}_{\tau,0}\psi_0^\perp(0) - \int_0^\tau \tilde{\mathcal{S}}_{\tau,\sigma}\Pi(\mathcal{R}\bar{\psi}_0(\sigma)) d\sigma.$$

Integration by parts introduces the factor  $M^{-1/2}$  we seek

$$\begin{aligned} \int_0^\tau \tilde{\mathcal{S}}_{\tau,\sigma}\Pi\mathcal{R}\bar{\psi}_0(\sigma) d\sigma &= \int_0^\tau iM^{-1/2}\frac{d}{d\sigma}(\tilde{\mathcal{S}}_{\tau,\sigma})\hat{\mathcal{V}}^{-1}\Pi\mathcal{R}\bar{\psi}_0(\sigma) d\sigma \\ &= \int_0^\tau iM^{-1/2}\frac{d}{d\sigma}\left(\tilde{\mathcal{S}}_{\tau,\sigma}\hat{\mathcal{V}}^{-1}\Pi\mathcal{R}\bar{\psi}_0(\sigma)\right) d\sigma \\ &\quad - \int_0^\tau iM^{-1/2}\tilde{\mathcal{S}}_{\tau,\sigma}\frac{d}{d\sigma}\left(\hat{\mathcal{V}}^{-1}(X_\sigma)\Pi\mathcal{R}\bar{\psi}_0(\sigma)\right) d\sigma \\ &= iM^{-1/2}\hat{\mathcal{V}}^{-1}\Pi\mathcal{R}\bar{\psi}_0(\tau) - iM^{-1/2}\tilde{\mathcal{S}}_{\tau,0}\hat{\mathcal{V}}^{-1}\Pi\mathcal{R}\bar{\psi}_0(0) \\ &\quad - \int_0^\tau iM^{-1/2}\tilde{\mathcal{S}}_{\tau,\sigma}\frac{d}{d\sigma}\left(\hat{\mathcal{V}}^{-1}(X_\sigma)\Pi\mathcal{R}\bar{\psi}_0(\sigma)\right) d\sigma. \end{aligned} \tag{8.27}$$

To analyze the integral in the right hand side we will use the fact

$$\hat{\mathcal{V}}^{-1} = \left( I + (\mathcal{V} - V_0)^{-1} \left[ \hat{\mathcal{V}} - (\mathcal{V} - V_0) \right] \right)^{-1} (\mathcal{V} - V_0)^{-1},$$

which can be verified by multiplying both sides from the left by  $I + (\mathcal{V} - V_0)^{-1} \left[ \hat{\mathcal{V}} - (\mathcal{V} - V_0) \right]$ . A spectral decomposition in  $L^2(dx)$ , based on the electron eigenpairs  $\{\lambda_k, \bar{\psi}_k\}_{k=1}^\infty$  and satisfying  $\mathcal{V}\bar{\psi}_k = \lambda_k\bar{\psi}_k$ , then implies

$$\begin{aligned} \hat{\mathcal{V}}^{-1}\Pi(\mathcal{R}\bar{\psi}_0) &= \left( I + (\mathcal{V} - V_0)^{-1} \left[ \hat{\mathcal{V}} - (\mathcal{V} - V_0) \right] \right)^{-1} (\mathcal{V} - V_0)^{-1}\Pi(\mathcal{R}\bar{\psi}_0) \\ (8.28) \quad &= \sum_{k \neq 0} \left( I + (\mathcal{V} - V_0)^{-1} \left[ \hat{\mathcal{V}} - (\mathcal{V} - V_0) \right] \right)^{-1} (\lambda_k - V_0)^{-1} \psi_k \langle \Pi(\mathcal{R}\bar{\psi}_0), \psi_k \rangle \\ &= \sum_{k \neq 0} (\lambda_k - V_0)^{-1} \psi_k \langle \Pi(\mathcal{R}\bar{\psi}_0), \psi_k \rangle + \mathcal{O}(M^{-1}) \end{aligned}$$

which applied to the integral in the right hand side of (8.27) shows that  $\|\bar{\psi}_0^\perp\|_{L^2(dx)} = \mathcal{O}(M^{-1/2})$  on a bounded time interval, when the spectral gap condition holds and  $\psi_k$  are smooth.

The evolution on longer times requires an additional idea: one can integrate by parts recursively in (8.27) to obtain

$$\begin{aligned} \int_0^\tau \tilde{\mathcal{S}}_{\tau,\sigma} \Pi \mathcal{R} \bar{\psi}_0(\sigma) d\sigma &= \left[ \tilde{\mathcal{S}}_{\tau,\sigma} \left( \tilde{\mathcal{B}} \tilde{\mathcal{R}} - \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}}) + \tilde{\mathcal{B}} \frac{d}{d\sigma} \left( \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}}) \right) - \dots \right) \right]_{\sigma=0}^{\sigma=\tau}, \\ \tilde{\mathcal{B}} &:= iM^{-1/2} \hat{\mathcal{V}}^{-1}, \quad \tilde{\mathcal{R}} := \Pi \mathcal{R} \bar{\psi}_0(\sigma), \end{aligned}$$

so that by (8.26) we have

$$\psi_0^\perp(\tau) = \tilde{\mathcal{S}}_{\tau,0} \psi_0^\perp(0) - \left[ \tilde{\mathcal{S}}_{\tau,\sigma} \left( \tilde{\mathcal{B}} \tilde{\mathcal{R}} - \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}}) + \tilde{\mathcal{B}} \frac{d}{d\sigma} \left( \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}}) \right) - \dots \right) \right]_{\sigma=0}^{\sigma=\tau}.$$

By choosing

$$\bar{\psi}_0^\perp(\sigma) \Big|_{\sigma=0} = - \left( \tilde{\mathcal{B}} \tilde{\mathcal{R}}(\sigma) - \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}})(\sigma) + \tilde{\mathcal{B}} \frac{d}{d\sigma} \left( \tilde{\mathcal{B}} \frac{d}{d\sigma} (\tilde{\mathcal{B}} \tilde{\mathcal{R}})(\sigma) \right) - \dots \right) \Big|_{\sigma=0}$$

we get

$$(8.29) \quad \bar{\psi}_0^\perp(\tau) = - \sum_{n=0}^{\infty} \tilde{\mathcal{B}}_0^n \mathcal{R}_0(\tau),$$

where  $\tilde{\mathcal{B}}_0 := -iM^{-1/2} \hat{\mathcal{V}}^{-1} \frac{d}{d\tau}$  and  $\mathcal{R}_0 := iM^{-1/2} \hat{\mathcal{V}}^{-1} \tilde{\mathcal{R}}$ . We assume this expansion (8.29) is convergent in  $L^2(dx)$  for each  $\tau$ , which follows from the smoothness estimate

$$(8.30) \quad \|\tilde{\mathcal{B}}_0^n \mathcal{R}_0(\tau)\|_{L^2(dx)} \rightarrow 0 \text{ as } n \rightarrow \infty$$

and (8.28).

The next step, verifying that also the non linear problem for  $V_0$  works, is based on the contraction obtained from

$$V_0 - \lambda_0 = \frac{\langle \psi, (\mathcal{V} - \lambda_0) \psi \rangle}{\langle \psi, \psi \rangle} = \mathcal{O}(\|\psi_0^\perp\|_{L^2(dx)})$$

and that  $\psi_0^\perp$  depends on  $V_0$  in (8.26), (8.27) and (8.28) with a multiplicative factor  $\mathcal{O}(M^{-1/2})$ .

Finally, to conclude that  $|\langle \psi, \psi \rangle - 1| = \mathcal{O}(M^{-1})$ , we use the evolution equation

$$\frac{d}{dt} \langle \psi, \psi \rangle = M^{-1/2} |G|^2 \text{Im} \langle \Delta \frac{\psi}{G}, \frac{\psi}{G} \rangle = \mathcal{O}(M^{-1})$$

where the last equality uses the obtained bound of  $\psi_0^\perp$  in the first part of (8.23). The assumption of a finite hitting time  $\tau^*$  then implies that  $|\langle \psi, \psi \rangle - 1| = \mathcal{O}(\tau^* M^{-1}) = \mathcal{O}(M^{-1})$ , since we may assume that  $\langle \psi, \psi \rangle = 1$  on  $I_\psi$ .  $\square$

**Remark 8.3** (Error estimates for the densities). We have the densities

$$(8.31) \quad \rho_S = G_S^{-2} \langle \psi, \psi \rangle \quad \text{for the Schrödinger equation,}$$

$$(8.32) \quad \rho_{\text{BO}} = G_{\text{BO}}^{-2} \quad \text{for the Born-Oppenheimer dynamics.}$$

From the stability of the Hamilton-Jacobi equation for  $\log(|G|^{-2})$  and the estimate  $\|\partial_{X^i X^j}(\theta - \tilde{\theta})\|_{L^\infty} = \mathcal{O}(M^{-1+\delta})$  in (8.18) we have

$$G_S^{-2} = G_{\text{BO}}^{-2} + \mathcal{O}(M^{-1+\delta}),$$

and Lemma 8.2 implies

$$(8.33) \quad \langle \psi, \psi \rangle = 1 + \mathcal{O}(M^{-1}),$$

which proves

$$\rho_S = \rho_{\text{BO}} + \mathcal{O}(M^{-1+\delta}).$$

## 9. NUMERICAL EXAMPLES

In order to demonstrate the presented theory we consider two different low dimensional Schrödinger problems. For both of these problems we show that there exists a Schrödinger eigenfunction density which converges weakly to the corresponding molecular dynamics density as  $M \rightarrow \infty$  with a convergence rate within the upper bound predicted in the theoretical part of this paper.

**9.1. Example 1: A single WKB state.** The first problem we consider is the time-independent Schrödinger equation

$$(9.1) \quad \mathcal{H}\Phi := \left( -\frac{1}{2M} \partial_{XX} + \bar{\mathcal{V}} \right) \Phi = E\Phi$$

with heavy coordinate  $X \in (-\pi, \pi]$  and two-state light coordinate  $x \in \{x_-, x_+\}$ . Periodicity is assumed over the heavy coordinate,  $\Phi(X, x) = \Phi(X + 2\pi, x)$ , and the potential operator  $\bar{\mathcal{V}}$  is defined by the matrix

$$(9.2) \quad \bar{\mathcal{V}}(X) = \begin{bmatrix} V(X) & \frac{1}{2}V(X)e(X) + c \\ \frac{1}{2}V(X)e(X) + c & 0 \end{bmatrix},$$

where we have chosen  $V(X) = -2\cos(X) + \cos(4X)$ ,  $e(X) = 1 + X^2$  and  $c$  to be a non-negative constant relating to the size of the spectral gap of  $\bar{\mathcal{V}}$ . The action  $\bar{\mathcal{V}}\Phi$  is thus defined by

$$(\bar{\mathcal{V}}\Phi)(X, \cdot) \equiv \bar{\mathcal{V}}(X) \begin{pmatrix} \Phi(X, x_-) \\ \Phi(X, x_+) \end{pmatrix}.$$

For each  $X$  the potential matrix (9.2) gives rise to the eigenvalue problem

$$\bar{\mathcal{V}}(X)v = \lambda_{\pm}(X)v$$

with the eigenvalues

$$\lambda_{\pm}(X) = \frac{1}{2} \left( V(X) \pm \text{Sgn}(X) \sqrt{V(X)^2 + 4(V(X)e(X)/2 + c)^2} \right),$$

where  $\text{Sgn}(X) = \pm 1$  as defined below. When constructing the molecular dynamics density for this problem

$$\rho_{\text{MD}}(X) = \frac{C}{\sqrt{2(E - \lambda(X))}},$$

one has to determine on which of the two eigenfunctions  $\lambda_{\pm}$  to base this density. When  $c = 0$  the difficulty that the eigenvalue functions  $\lambda_+$  and  $\lambda_-$  can cross is added to the problem. In order to determine the continuation of eigenvalue functions at the crossings we introduce a function  $\text{Sgn}(X)$  which is a sign function with  $\text{Sgn}(-\pi) = 1$  that changes sign at points where

$$V(X)^2 + 4 \left( \frac{1}{2}V(X)e(X) + c \right)^2 = 0.$$

Since this situation can only occur when  $c = 0$ , it is possible to set

$$\text{Sgn}(X) := \text{sgn}(V(-\pi))\text{sgn}(V(X)).$$

See Figure 8 for a typical eigenvalue function crossing, which makes the function  $\lambda_{\pm} : \mathbb{R} \rightarrow \mathbb{R}$  smooth (in contrast to the choice  $\text{Sgn} \equiv 1$ ).

To solve (9.1) numerically, we use the finite difference method to discretise the operator  $\mathcal{H}$  on a grid  $\{X_j\}_{j=1}^N \times \{x_-, x_+\}$  with the step-size  $h = 2\pi/N$  and  $X_j = jh$ . The discrete eigenvalue problem

$$\mathcal{H}^{(h)}\Upsilon_j = E_j\Upsilon_j$$

is solved for the 10 eigenvalues being closest to the fixed energy  $E$  and a molecular dynamics approximation of the eigensolution is constructed by

$$\Phi_{\text{MD}}(X, x) := \sqrt{\rho_{\text{MD}}(X)} e^{iM^{1/2}\Theta(X)} v(X, x),$$

where  $v(X, \cdot)$  is one of the eigenvectors of  $\bar{V}(X)$  and

$$(9.3) \quad \Theta(X) := \int_0^X \sqrt{2(E_1 - \lambda(s))} ds$$

is approximated by a trapezoidal quadrature yielding  $\Theta^{(h)}$ . Thereafter a Schrödinger eigensolution  $\Phi^{(h)}$  which is close to the molecular dynamics eigensolution is obtained by projecting  $\Phi_{\text{MD}}$  onto the subspace spanned by  $\{\Upsilon\}_{j=1}^J$  as described in Algorithm 2. By denoting  $\rho_{\Phi^{(h)}}(X) = \langle \Phi^{(h)}, \Phi^{(h)} \rangle$  and  $\rho_{\text{MD}}(X) = \langle \Phi_{\text{MD}}, \Phi_{\text{MD}} \rangle$ , the observables  $g_1(X) = X^2$  and  $g_2(X) = V(X)$  are used to compute the convergence rate of

$$(9.4) \quad \left| \frac{\int_{-\pi}^{\pi} g_i(X) \rho_{\text{MD}}(X) dX - \int_{-\pi}^{\pi} g_i(X) \rho_{\Phi^{(h)}}(X) dX}{\int_{-\pi}^{\pi} g_i(X) \rho_{\text{MD}}(X) dX} \right|,$$

as  $M$  increases. Further details of the numerical solution idea are described in Algorithm 1.

Plots of the results for the test case with the spectral gap  $c = 5$  and  $E = 0$ , and for the test case with crossing eigenvalue functions when  $c = 0$  and  $E = 1.2$  are given below. Most noteworthy is Figure 11, which demonstrates that the obtained convergence rate for (9.4) is  $\mathcal{O}(M^{-1})$  for both scenarios.

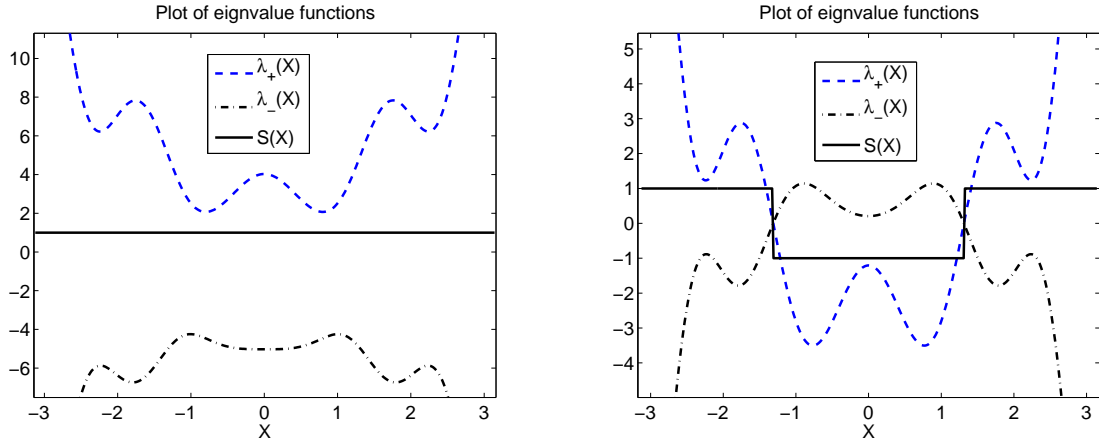


FIGURE 8. Left plot: Eigenvalue functions when  $c = 5$ . There is a spectral gap which makes the sign function  $S = 1$ . Right plot: Eigenvalue functions when  $c = 0$ . The eigenvalue functions exhibit crossing, consequently the function  $S$  changes its sign from  $\pm 1$  to  $\mp 1$  at the crossing points.

**9.2. Example 2: A caustic state.** Next, we consider the one dimensional, time independent, periodic Schrödinger equation

$$(9.7) \quad \left( -\frac{1}{2M} \partial_{XX} + V \right) \Phi = E\Phi, \quad X \in (-2\sqrt{E}, 2\sqrt{E})$$

---

**Algorithm 1** Algorithm for problems in Example 1
 

---

**Input:** Energy  $E$ ; potential functions  $V$ ,  $e$  and  $c$ ; mass  $M$ ; number of grid points  $N$  and grid  $\{X_i\}_{i=1}^N$ .

**Output:** Schrödinger projection density  $\rho_{\Phi^{(h)}}$ .

1. Construct the discrete operator  $\mathcal{H}^{(h)}$  from (9.1) using finite differences and solve the eigenvalue problem

$$\mathcal{H}^{(h)}\Upsilon_i = E_i\Upsilon_i$$

for the 10 eigenvalues being closest to  $E$  by using MATLAB `eigs(H,10,E)`.

2. Sort the eigenvalues and eigenvectors by distance from  $E$  and keep only the  $E_i$ s which are less than  $M^{-1/2}$  away from  $E$ . Let  $\bar{J}$  be the number of kept eigenvalues and  $E_0$  the eigenvalue closest to  $E$ .

3.

**for**  $i = 1$  to  $N$  **do**

Solve the eigenvalue problem

$$\bar{V}(X_i, \cdot)v_{\pm}(X_i, \cdot) = \lambda_{\pm}(X_i)v_{\pm}(X_i, \cdot),$$

where  $\bar{V}$  is the matrix defined in (9.2).

**end for**

4. Construct the molecular dynamics density according to the formula

$$\rho_{\text{MD}}(X) = \frac{(E_0 - \lambda(X))^{-1/2}}{\int_{[0, 2\pi]} (E_0 - \lambda(X))^{-1/2} dX},$$

where we choose  $\lambda(X)$  above from the two eigenvalues  $\lambda_{\pm}(X)$  by the criterion that the eigenvalue chosen must fulfil  $\|\lambda\|_{\infty} < E_0$ .

5. Construct a discrete molecular dynamics approximation to the eigenfunction

$$(9.5) \quad \Phi_{\text{MD}}(X, x) = \sqrt{\rho_{\text{MD}}(X)} e^{iM^{1/2}\Theta(X)} v(X, x),$$

where  $v(X, x)$  is one of the eigenvectors  $v_{\pm}$ ,

$$(9.6) \quad \Theta(X) := \int_0^X \sqrt{2(E_1 - \lambda(s))} ds,$$

and we approximate  $\Theta$  by a trapezoidal quadrature  $\Theta^{(h)}$ .

6. Project the molecular dynamics solution  $\Phi_{\text{MD}}$  onto the eigenspace  $\{\Upsilon_i\}_{i=1}^{\bar{J}}$ ,  $\bar{J} \leq 10$  by Algorithm 2 to obtain a projection solution  $\Phi^{(h)}$ .

7. Derive the Schrödinger projection density by

**for**  $i = 1$  to  $N$  **do**

$$\rho_{\Phi^{(h)}}(X_i) = |\Phi^{(h)}(X_i, x_-)|^2 + |\Phi^{(h)}(X_i, x_+)|^2,$$

**end for**

and scaling  $\rho_{\Phi^{(h)}} = \rho_{\Phi^{(h)}} / \|\rho_{\Phi^{(h)}}\|$ .

---

with  $V(X) = X^2$  and  $E = 1$ . The eikonal equation corresponding to (9.7) is

$$(9.8) \quad \frac{1}{2}P^2 + V(X) = E.$$

As in Example 1, we would like to use the eikonal equation to construct a numerical approximate solution of (9.7) whose density converges weakly as  $M \rightarrow \infty$  to the density generated from a solution of (9.7). The molecular dynamics density corresponding to this eikonal equation becomes by (3.19)  $\rho_{\text{BO}} = C(E - V(X))^{-1/2}$ . The density  $\rho_{\text{BO}}$  goes to infinity at the caustics  $X = V^{-1}(E) = \pm\sqrt{E}$  and the approach in Example 1 does not work directly. We will instead construct the numerical approximate solution using the stationary phase method as outlined below based on the WKB Fourier integral ansatz.

---

**Algorithm 2** Projection algorithm
 

---

**Input:** Mass  $M$ ; wave solution  $\Phi$ ; eigenvalues  $\{E_i\}_{i=1}^{\bar{J}}$  and corresponding eigenvectors  $\{\Upsilon_i\}_{i=1}^{\bar{J}}$ .

**Output:** Schrödinger projection wave solution  $\Phi^{(h)}$ .

1. Organize eigenvalues by multiplicity by a numerical approximation. Construct a  $\bar{J} \times \bar{J}$ , zero matrix  $A$  which keeps track of multiplicity relations as follows:

**for**  $i = 1$  to  $\bar{J}$  **do**

**for**  $j = i$  to  $\bar{J}$  **do**

**if**  $|E_i - E_j| < M^{-3/4}$  **then**

      Consider eigenvalues equal since the expected spectral gap is  $\mathcal{O}(M^{-1/2})$ , and store this relation by

**if**  $A_{kj} = 0$  for all  $k < i$  **then**

        Set  $A_{ij} = 1$ .

**end if**

**end if**

**end for**

**end for**

2. For vectors  $b \in \{0, 1\}^{\bar{J}}$ , define the projection

$$\Phi^{(h,b)} := \sum_{j,k=1}^{\bar{J}} b_k A_{k,j} \langle \Phi, \Upsilon_j \rangle \Upsilon_j$$

and, letting  $\rho$  and  $\rho_{\Phi^{(h,b)}}$  denote the densities generated by  $\Phi$  and  $\Phi^{(h,b)}$  respectively, set

$$b^* = \arg \min_{b \in \{0,1\}^{\bar{J}}} \|\rho - \rho_{\Phi^{(h,b)}}\|.$$

3. Return the projection  $\Phi^{(h)} := \Phi^{(h,b^*)}$ .

---

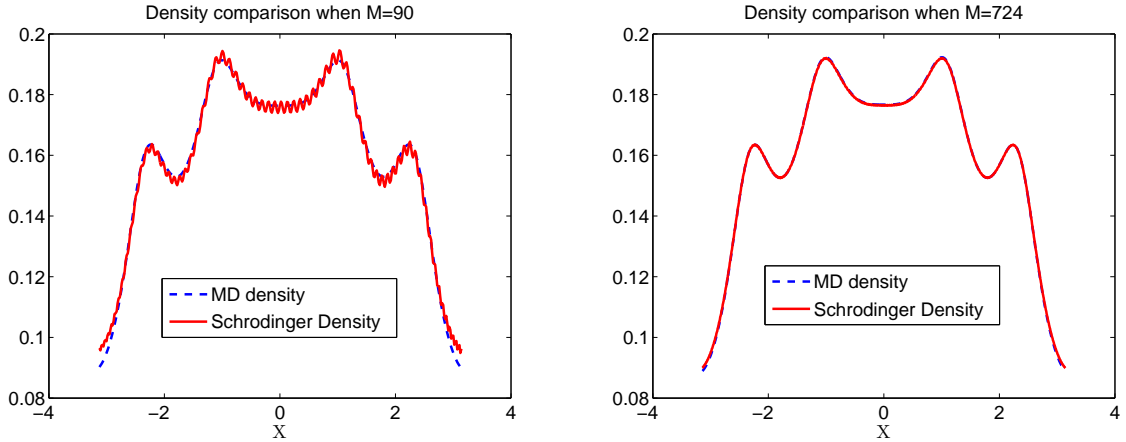


FIGURE 9. Plot of the MD density  $\rho_{\text{MD}}$  and the Schrödinger projection density  $\rho_{\Phi^{(h)}}$  in the case  $c = 5$  and  $E = 0$  for the two different masses  $M = 90$  (left plot) and  $M = 724$  (right plot) illustrating the convergence of the densities.

By the Legendre transform

$$\theta^*(P) = \min_X (XP - \theta(X))$$

an invertible mapping between the momentum and position coordinates fulfilling  $X = \nabla_P \theta^*(P)$  is constructed. Using equation (9.8), one sees that  $\nabla_P \theta^*(P) = V^{-1}(E - P^2/2)$ . Since  $\theta^*(0) = 0$ , one can derive

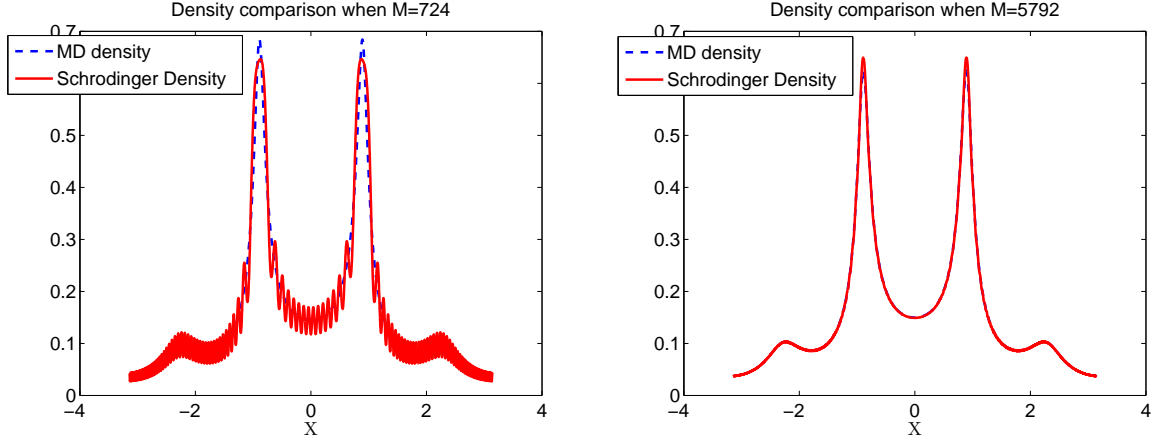


FIGURE 10. Plot of the MD density  $\rho_{\text{MD}}$  and Schrödinger projection density  $\rho_{\Phi^{(h)}}$  in the case  $c = 0$  and  $E = 1.2$  for the two different masses  $M = 724$  (left plot) and  $M = 5792$  (right plot).

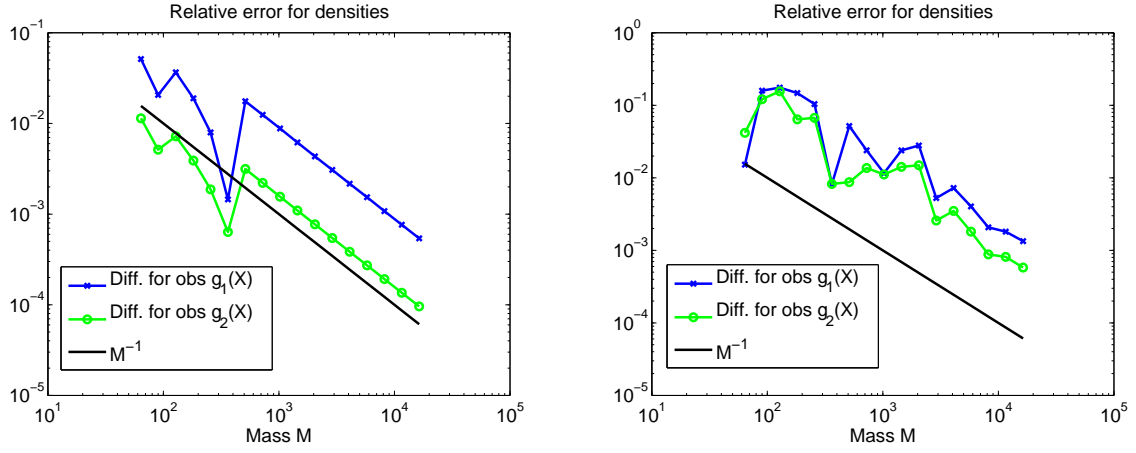


FIGURE 11. Left plot: Plot of the observable density errors given in (9.4) with an eigenvalue gap, when  $c = 5$  and  $E = 0$ . Right plot: Plot of the observable density errors given in (9.4) with an eigenvalue crossing, when  $c = 0$  and  $E = 1.2$ .

that for this particular choice of  $V$

$$\theta^*(P) = \int_0^P \sqrt{E - s^2/2} ds = \frac{E}{\sqrt{2}} \left[ \sin^{-1} \left( \frac{P}{\sqrt{2E}} \right) + \frac{P}{\sqrt{2E}} \sqrt{1 - \frac{P^2}{2E}} \right].$$

In neighbourhoods of the caustics  $[-2E^{1/2}, -X_0]$  and  $(X_0, 2E^{1/2}]$ , we construct the approximate solution by

$$\Phi(X) = \frac{u(X)}{\sqrt{|\nabla_X V(X)|}}$$

where  $u$  is the inverse Fourier transform

$$u(X) := \int_{-2\sqrt{E}}^{2\sqrt{E}} e^{iM^{1/2}(-XP + \theta^*(P))} dP$$



and  $X_0 \in (-V^{-1}(E), V^{-1}(E))$  is a value yet to be chosen. In the region  $(-X_0, X_0)$  the approximate solution is constructed by

$$(9.9) \quad \Phi(X) = C \frac{\bar{u}(X)}{(E - V(X))^{1/4}}.$$

Here

$$(9.10) \quad \bar{u}(X) := e^{-iM^{1/2}\theta(X)}\psi_+ + e^{iM^{1/2}\theta(X)}\psi_-,$$

with, according to the Legendre transform,  $\theta(X) := X\sqrt{2(E - V(X))} - \theta^*(\sqrt{2(E - V(X))})$  and  $\psi_{\pm}$  determined by the stationary phase method:

1. Set  $P(p) = P_0 + p$  with  $P_0 = \sqrt{2(E - V(X_0))}$  and let

$$Y(p) := \text{sgn}(p) \sqrt{2 \frac{-X(P_0 + p) + \theta^*(P_0 + p) + \theta(X_0)}{\partial_{PP}\theta^*(P_0)}},$$

using

$$(9.11) \quad \theta(X) := X\sqrt{2(E_0 - V(X))} - \theta^*(\sqrt{2(E_0 - V(X))}),$$

and determine its inverse  $p(Y)$  in a neighbourhood of  $Y = 0$  by computing  $(p_i, Y(p_i))$  on a grid around  $p = 0$  and, for  $k \geq 3$ , fit a  $3k + 1$ th degree polynomial to the values  $(Y(p_i), p_i)$  using the method of least squares.

2. Evaluate the stationary phase expansion

$$(9.12) \quad u(X_0) = \sum_{p_0 = \pm\sqrt{2(E - V(X_0))}} e^{i\pi \text{sgn}(\partial_{PP}\theta^*(P_0))/4} \left[ \left| \frac{1}{2} \partial_{PP}\theta^*(P_0) \right|^{-1/2} e^{-iM^{1/2}\theta(X_0)} \right. \\ \left. \times \sum_{j=0}^k \frac{M^{-j/2}}{j!} \left( i \left( \frac{1}{2} \partial_{PP}\theta^*(P_0) \right)^{-1} \partial_{YY} \right)^j \left| \partial_Y p \right| \Big|_{Y=0} + \mathcal{O}(M^{-j/2}) \right]$$

to obtain

$$u(X_0^-) = e^{iM^{1/2}\theta(X_0)}(\psi_+ + \mathcal{O}(M^{-k/2})) + e^{-iM^{1/2}\theta(X_0)}(\psi_- + \mathcal{O}(M^{-k/2})),$$

where

$$\psi_{\pm} := e^{i\pi \text{sgn}(\partial_{PP}\theta^*(\pm P_0))/4} \left| \frac{1}{2} \partial_{PP}\theta^*(\pm P_0) \right|^{-1/2} \sum_{k=0}^3 \frac{M^{-k/2}}{k!} \left( i \left( \frac{\partial_{PP}\theta^*(\pm P_0)}{2} \right)^{-1} \partial_{YY} \right)^k \left| \partial_Y p \right| \Big|_{Y=0}.$$

The constant  $C$  in (9.9) is chosen so that the wave solution parts are continuous at the gluing point,  $\Phi(\pm X_0^-) = \Phi(\pm X_0^+)$ . It is most easy to determine  $C$  when  $X_0$  is chosen so that  $|u(X_0)|$  is at a local maximum; see Figure 12 for an illustration of the gluing procedure.

At the end a Schrödinger eigenfunction solution  $\Phi^{(h)}$  is obtained by projecting  $\Phi$  onto the space spanned by a set of eigensolutions to the discretized version of the Schrödinger problem,  $\{\Upsilon_j\}_{j=1}^J$ , as is described in Algorithm 2.

Two convergence results are needed to make the method work. First, the density generated from the stationary phase based on the approximate solution  $\rho(X) := |\Phi|^2(X)/\|\Phi\|_2^2$  must converge weakly to the Schrödinger projection based density  $\rho_{\Phi^{(h)}}(X) := |\Phi^{(h)}|^2(X)/\|\Phi^{(h)}\|_2^2$  as  $M \rightarrow \infty$ ; see Figure 13 for an illustration of how these functions converge. Second,  $\rho_{\Phi^{(h)}}$  must converge to the molecular dynamics density  $\rho_{\text{MD}}(X) := C(E - V(X))^{-1/2}$  as  $M$  increases; see Figure 14.

A numerical test of the convergence rate of

$$(9.13) \quad \left| \frac{\int_{-2\sqrt{E_0}}^{2\sqrt{E_0}} g_1(X) \rho_{\text{MD}}(X) dX}{\int_{-2\sqrt{E_0}}^{2\sqrt{E_0}} g_2(X) \rho_{\text{MD}}(X) dX} - \frac{\int_{-2\sqrt{E_0}}^{2\sqrt{E_0}} g_1(X) \rho_{\Phi^{(h)}}(X) dX}{\int_{-2\sqrt{E_0}}^{2\sqrt{E_0}} g_2(X) \rho_{\Phi^{(h)}}(X) dX} \right|$$

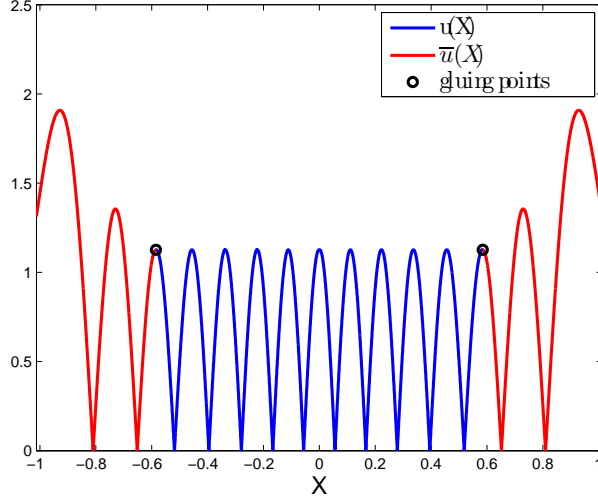


FIGURE 12. Plot illustrating the gluing procedure of the functions  $u(X)$  and  $\bar{u}(X)$  at the points  $\pm X_0$ .

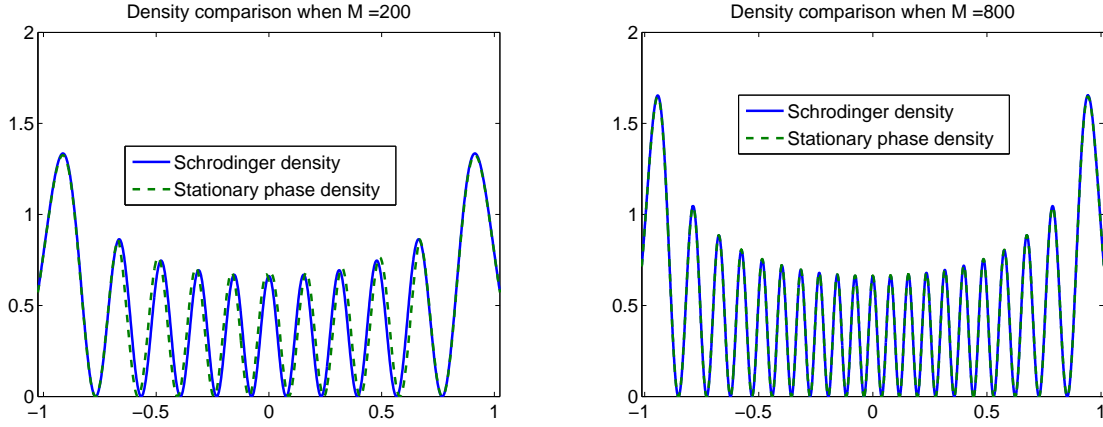


FIGURE 13. Comparison of the approximate solution based density  $\rho$  and the Schrödinger projection based solution  $\rho_{\Phi(h)}$  for  $M = 200$  (left plot) and  $M = 800$  (right plot).

as  $M$  increases is illustrated in Figure 15 for the observables

$$(9.14) \quad g_1(X) = \frac{(1.5 - X)^6(1.5 + X)^6(1 + e^{-X^2})}{1.5^{12}} \quad \text{and} \quad g_2(X) = \frac{(1.5 - X)^6(1.5 + X)^6(1 - X^2 + X^4)}{1.5^{12}}.$$

Further details of the solution procedure in Exampe 2 are given in Algorithm 3.

## 10. THE STATIONARY PHASE EXPANSION

Consider the phase function  $\tilde{X} \cdot \tilde{P} - \theta^*(\tilde{X}, \tilde{P})$  and let  $\tilde{P}_0(\tilde{X})$  be any solution to the stationary phase equation  $\tilde{X} = \nabla_{\tilde{P}} \theta^*(\tilde{X}, \tilde{P}_0)$ . We rewrite the phase function

$$\tilde{X} \cdot \tilde{P} - \theta^*(\tilde{X}, \tilde{P}) = \underbrace{\tilde{X} \cdot \tilde{P}_0 - \theta^*(\tilde{X}, \tilde{P}_0)}_{=\theta(\tilde{X}, \tilde{X})} + (\tilde{P} - \tilde{P}_0) \cdot \int_0^1 (1-t) \partial_{PP} \theta^*(\tilde{P}_0 + t[\tilde{P} - \tilde{P}_0]) dt [\tilde{P} - \tilde{P}_0].$$

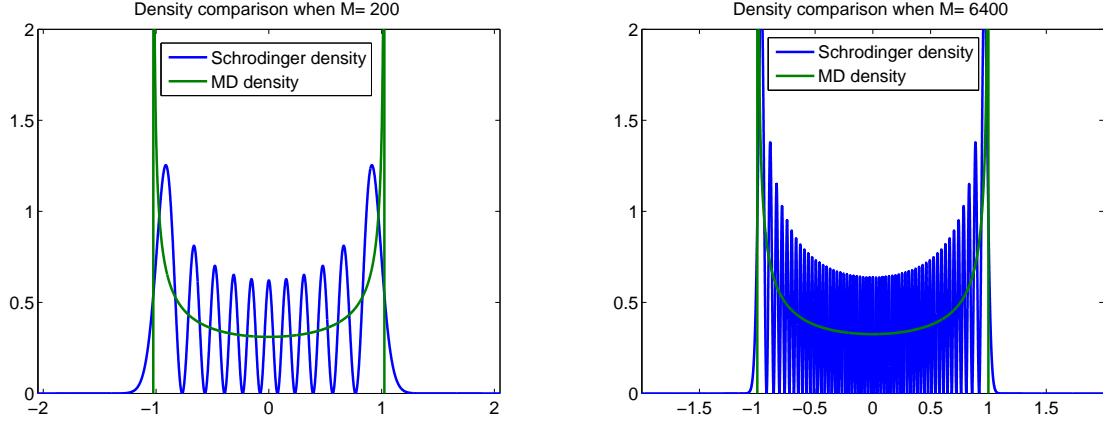


FIGURE 14. Comparison of the Schrödinger projection density  $\rho_{\Phi(\hbar)}$  and the molecular dynamics density  $\rho_{\text{MD}}$  for  $M = 200$  (left plot) and  $M = 6400$  (right plot).

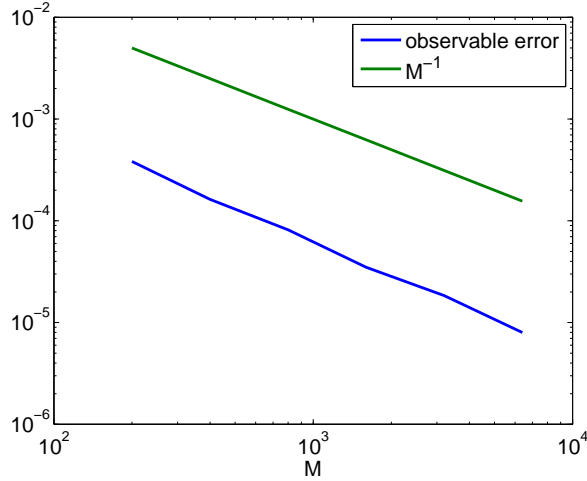


FIGURE 15. Convergence rate of (9.13) for the observables  $g_1$  and  $g_2$  as defined in (9.14).

The relation

$$\frac{1}{2}Y \cdot \partial_{PP} \bar{\theta}(\check{P}_0) Y = (\check{P} - \check{P}_0) \cdot \int_0^1 (1-t) \partial_{PP} \bar{\theta}(\check{P}_0 + t[\check{P} - \check{P}_0]) dt [\check{P} - \check{P}_0]$$

defines the function  $Y(\check{P})$ , and its inverse  $\check{P}(Y)$ , so that the phase is a quadratic function in  $Y$ . The stationary phase expansion of an integral takes the form, see [10],

$$(10.1) \quad \int_{\mathbb{R}^d} w(\check{P}) e^{-iM^{1/2}(\check{X} \cdot \check{P} - \theta^*(\check{X}, \check{P}))} d\check{P} \\ \simeq \sum_{\nabla_P \theta^*(\check{P}_0) = \check{X}} (2\pi M^{-1/2})^{d/2} \left| \det \frac{\partial(\check{P})}{\partial(\check{X})} \right|^{1/2} e^{i\frac{\pi}{4} \text{sgn}(\partial_{PP} \theta^*(\check{P}_0))} e^{-iM^{1/2} \theta(\check{X}, \check{X})} \\ \times \sum_{k=0}^{\infty} \frac{M^{-k/2}}{k!} \left( \sum_{l,j} i(\partial_{P^l P^j} \theta^*)^{-1}(\check{P}_0) \partial_{Y^l Y^j} \right)^k \left( w(\check{P}(Y)) \left| \det \frac{\partial(\check{P})}{\partial(Y)} \right| \right) \Big|_{Y=0}.$$

---

**Algorithm 3** Algorithm for Example 2

---

**Input:** An energy  $E$ , an one-dimensional potential function  $V$ , mass  $M$ , Schrödinger equation (9.7).

**Output:** The Schrödinger projection density  $\rho_{\Phi^{(h)}}$ .

1. Identify the right caustic point  $X_+ > 0$  satisfying  $X_+ = V^{-1}(E)$ . For a fixed  $E \in \mathbb{R}$ , consider the periodic eigenvalue problem. Solve (9.7) numerically by constructing the discretised operator form of  $-(2M)^{-1}\partial_{XX} + V$  using finite differences and denoted  $\mathcal{H}^{(h)}$ , and solve the eigenvalue problem

$$(9.15) \quad \mathcal{H}^{(h)}P_i = E_iP_i$$

for the 10 eigenvalues closest to  $E$  using the Matlab eigenvalue solver **eigs(H,10,E)**. Let  $E_0$  denote the eigenvalue closest to  $E$  and consider from now on solving (9.7) for the energy  $E_0$  and its corresponding eikonal equation  $\frac{1}{2}P^2 + V(X) = E_0$ .

2. Determine  $\theta^*(P)$  by

$$\theta^*(P) = \int_0^P \nabla_P \theta^*(p) dp$$

3. Evaluate the Fourier integral

$$(9.16) \quad u(X) := \int_{-2\sqrt{E}}^{2\sqrt{E}} e^{iM^{1/2}(-XP + \theta^*(P))} dP, \quad |X| > X_0,$$

where  $X_0$  is chosen as the smallest value  $X > X_+/2$  such that  $|u(X)|$  is at a local maximum, and for  $|X| \leq X_0$  compute  $\bar{u}$  by (9.10) using the stationary phase method.

4. Construct the approximate solution

$$\Phi(X) := \begin{cases} C\bar{u}(X)(E_0 - V(X))^{-1/4} & |X| \leq X_0, \\ u(X)/\sqrt{|\nabla_X V(X)|} & |X| \geq X_0, \end{cases}$$

with

$$C = \frac{u(X_0)(E_0 - V(X_0))^{1/4}}{\sqrt{|\nabla_X V(X_0)|\bar{u}(X_0)}}.$$

5.

Project  $\Phi$  onto the eigenspace  $\{\Upsilon_i\}_{i=1}^{\bar{J}}$ ,  $\bar{J} \leq 10$  by Algorithm 2 to obtain a projection solution  $\Phi^{(h)}$  and compute its corresponding approximate density

$$\rho_{\Phi^{(h)}} = \frac{|\Phi^{(h)}|^2(X)}{\|\Phi^{(h)}\|_2^2}.$$

---

#### ACKNOWLEDGMENT

The research of P.P. and A.S. was partially supported by the National Science Foundation under the grant NSF-DMS-0813893 and Swedish Research Council grant 621-2010-5647, respectively. P.P. also thanks KTH and Nordita for their hospitality during his visit when the presented research was initiated.

#### REFERENCES

- [1] F.A. Berezin and M.A. Shubin, *The Schrödinger equation*, Kluwer Academic Publishers, 1991.
- [2] M. Born and R. Oppenheimer, *Zur quantentheorie der molekeln*, Ann. Physik (1927), no. 84, 4571–484.
- [3] F.A. Bornemann, P. Nettesheim, and C. Schütte, *Quantum-classical molecular dynamics as an approximation to full quantum dynamics*, J. Chem. Phys. **105** (1996), 1074–1083.
- [4] A. Bouzounia and D. Robert, *Uniform semiclassical estimates for the propagation of quantum observables*, Duke Math. J. **111** (2002), 223–252.
- [5] J. Briggs, S. Boonchui, and S. Khemmani, *The derivation of the time-dependent Schrödinger equation*, J. Phys. A: Math. Theor. **40** (2007), 1289–1302.
- [6] J. Briggs and J.M. Rost, *On the derivation of the time-dependent equation of Schrödinger*, Foundations of Physics **31** (2001), 693–712.

- [7] E. Cancès, M. Defranceschi, W. Kutzelnigg, C. LeBris, and Y. Maday, *Computational chemistry: a primer*, Handbook of Numerical Analysis, vol. X, North-Holland, 2007.
- [8] E. Cancès, F. Legoll, and G. Stolz, *Theoretical and numerical comparison of some sampling methods for molecular dynamics*, Math. Model. Num. Anal. **41** (2007), 351–389.
- [9] J. Carlsson, M. Sandberg, and A. Szepessy, *Symplectic Pontryagin approximations for optimal design*, Math. Model. Num. Anal. **43** (2009), 3–32.
- [10] J.J. Duistermaat, *Fourier integral operators*, Courant Institute, 1973.
- [11] J.-P. Eckmann and R. S en eor, *The Maslov-WKB Method for the (an-)harmonic oscillator*, Arch. Rat. Mech. Anal. **61** (1976), 153–173.
- [12] L.C. Evans, *Partial differential equation*, American Mathematical Society, Providence, RI, 1998.
- [13] C. Fefferman and L. Seco, *Eigenvalues and eigenfunctions of ordinary differential operators*, Adv. Math. **95** (1992), 145–305.
- [14] D. Frenkel and B. Smith, *Understanding molecular simulation*, Academic Press, 2002.
- [15] G.A. Hagedorn, *High order corrections to the time-independent Born-Oppenheimer approximation II: diatomic Coulomb systems*, Comm. Math. Phys. **116** (1988), 23–44.
- [16] B. Helffer, *Semi-classical analysis for the Schr odinger operator and applications*, Lecture Notes in Mathematics, vol. 1336, Springer Verlag, 1988.
- [17] H. Jeffreys, *On certain approximate solutions of linear differential equations of the second order*, Proc. London Math. Soc. **23** (1924), 428–436.
- [18] J. B. Keller, *Corrected Bohr-Sommerfeld quantum conditions for nonseparable systems*, Ann. Phys. **4** (1958), 180–188.
- [19] M. Klein, A. Martinez, R. Seiler, and X. P. Wang, *On the Born-Oppenheimer expansion for polyatomic molecules*, Comm. Math. Phys. **143** (1992), 607–639.
- [20] C. Lasser and S. R oblitz, *Computing expectations values for molecular quantum dynamics*, SIAM J. Sci. Comput. **32** (2010), 1465–1483.
- [21] C. LeBris, *Computational chemistry from the perspective of numerical analysis*, Acta Numerica, vol. 14, pp. 363–444, CUP, 2005.
- [22] E. Lieb and R. Seiringer, *The stability of matter in quantum mechanics*, CUP, 2010.
- [23] D. Marx and J. Hutter, *Ab initio molecular dynamics: Theory and implementation, modern methods and algorithms of quantum chemistry*, Tech. report, John von Neumann Institute for Computing, J ulich, 2001.
- [24] A. Martinez and V. Sordani, *Twisted pseudodifferential calculus and application to the quantum evolution of molecules*, Memoirs Am. Math. Soc., **200** (2009), n. 936.
- [25] V. P. Maslov and M. V. Fedoriuk, *Semi-classical approximation in quantum mechanics*, D. Reidel Publishing Company, 1981; based on: V. P. Maslov, *Theory of perturbations and asymptotic methods*, Moskov. Gos. Univ.. Moscow 1965 (Russian).
- [26] M. Dimassi and J. Sj ostr and, *Spectral asymptotics in the semiclassical limit*, LMS Lecture Note Series, vol. 268, CUP, 1999.
- [27] N. F. Mott, *On the theory of excitation by collision with heavy particles*, Proc. Camb. Phil. Soc. **27** (1931), 553–560.
- [28] G. Panati, H. Spohn, and S. Teufel, *Space-adiabatic perturbation theory*, Adv. Theor. Math. Phys. **7** (2003), 145–204.
- [29] Rayleigh, *On the propagation of waves through a stratified medium, with special reference to the question of reflection*, Proc. Roy. Soc. (London) Series A **86** (1912), 207–226.
- [30] M. Reed and B. Simon, *Methods of Modern Mathematical Physics, I: Functional Analysis*, Academic Press INC., 1980.
- [31] M. Sandberg and A. Szepessy, *Convergence rates of symplectic Pontryagin approximations in optimal control theory*, Math. Model. Num. Anal. **40** (2006), 149–173.
- [32] L. Schiff, *Quantum mechanics*, McGraw-Hill, 1968.
- [33] E. Schr odinger, *Collected papers on wave mechanics*, Blackie and Son, London, 1928.
- [34] A. Szepessy, *Langevin molecular dynamics derived from Ehrenfest dynamics*, Math. Mod. Meth. Appl. S. **21** (2011), 2289–2334.
- [35] D. J. Tanner, *Introduction to quantum mechanics: A time-dependent perspective*, University Science Books, 2006.
- [36] J. C. Tully, *Mixed quantum-classical dynamics*, Faraday Discuss. **110** (1998), 407–419.
- [37] E. von Schwerin and A. Szepessy, *A stochastic phase-field model determined from molecular dynamics*, Math. Model. Num. Anal. **44** (2010), 627–646.
- [38] S. Zelditch, *Quantum ergodicity and mixing of eigenfunctions*, in Franoise, Jean-Pierre; Naber, Gregory L.; Tsun, Tsou Sheung, Encyclopedia of mathematical physics. Vol. 1, 2, 3, 4, 5, Academic Press/Elsevier Science, Oxford, 2006.

DEPARTMENT OF MATHEMATICS UNIVERSITY OF VIENNA NORDBERGSTRÆ 15 1090 WIEN, AUSTRIA  
*E-mail address:* [christian.bayer@univie.ac.at](mailto:christian.bayer@univie.ac.at)

DEPARTMENT OF NUMERICAL ANALYSIS, KUNGL. TEKNISKA HÖGSKOLAN, 100 44 STOCKHOLM, SWEDEN  
*E-mail address:* [hhoel@kth.se](mailto:hhoel@kth.se)

DEPARTMENT OF NUMERICAL ANALYSIS, KUNGL. TEKNISKA HÖGSKOLAN, 100 44 STOCKHOLM, SWEDEN  
*E-mail address:* [smakadir@csc.kth.se](mailto:smakadir@csc.kth.se)

DEPARTMENT OF MATHEMATICAL SCIENCES, UNIVERSITY OF DELAWARE, NEWARK, DE 19716, USA  
*E-mail address:* [plechac@math.udel.edu](mailto:plechac@math.udel.edu)

DEPARTMENT OF NUMERICAL ANALYSIS, KUNGL. TEKNISKA HÖGSKOLAN, 100 44 STOCKHOLM, SWEDEN  
*E-mail address:* [msandb@kth.se](mailto:msandb@kth.se)

DEPARTMENT OF MATHEMATICS, KUNGL. TEKNISKA HÖGSKOLAN, 100 44 STOCKHOLM, SWEDEN  
*E-mail address:* [szepessy@kth.se](mailto:szepessy@kth.se)

DIVISION OF MATHEMATICS, KING ABDULLAH UNIVERSITY OF SCIENCE AND TECHNOLOGY, THUWAL 23955-6900, KINGDOM  
OF SAUDI ARABIA  
*E-mail address:* [raul.tempone@kaust.edu.sa](mailto:raul.tempone@kaust.edu.sa)