
CHAPTER 14

Wigner Function–Based Device Modeling

Hans Kosina, Mihail Nedjalkov

Institute for Microelectronics, TU Vienna, Vienna, Austria

CONTENTS

1. Introduction	731
1.1. History and State of the Art Review	732
2. The Wigner Function Formalism	733
2.1. The Wigner Function	734
2.2. Marginal Distributions	735
2.3. The Wigner Equation	737
3. Electron–Phonon Interaction	739
3.1. The System Hamiltonian	739
3.2. A Hierarchy of Transport Equations	740
3.3. Integral Form of the Wigner Equation	744
4. The Monte Carlo Method	747
4.1. The General Scheme	747
4.2. Particle Models	749
4.3. The Negative Sign Problem	753
4.4. Particle Annihilation	754
5. Simulation Results	755
5.1. Comparison with Other Numerical Methods	755
5.2. The Effect of Scattering	755
5.3. Inclusion of Extended Contact Regions	759
6. Conclusion	761
References	762

1. INTRODUCTION

Modeling of electronic transport in mesoscopic systems requires a theory that describes open, quantum-statistical systems driven far from thermodynamic equilibrium. Several formulations of quantum transport have been employed practically, such as those based on the density matrix, nonequilibrium Green's functions, and the Wigner function.

ISBN: 1-58883-052-7

Copyright © 2006 by American Scientific Publishers
All rights of reproduction in any form reserved.

Handbook of Theoretical and Computational Nanotechnology
Edited by Michael Rieth and Wolfram Schommers
Volume 10: Pages (731–763)

A quantum-mechanical phase-space distribution was introduced by Eugene Wigner in 1932 [1]. The purpose was the formulation of a quantum correction for the thermodynamic equilibrium of a many-body system by means of a quasiprobability function. In more recent times, the definition of the Wigner function has been generalized as a Fourier transform of a many-body Green's function [2].

The Wigner function is a real-valued but not necessarily positive definite quasidistribution and represents a quantum generalization of Boltzmann's N-particle distribution. The Wigner function formalism is attractive as it allows the expression of quantum dynamics in a phase-space formulation, directly comparable with the classical analogue. A phase-space approach may appear more intuitive compared with the more abstract density matrix and Green's function approaches. The method of quasidistributions has proved especially useful in providing reductions to classical physics and kinetic regimes under suitable conditions.

To discuss the physical interpretation of a quasidistribution, let us consider the simple case of a one-particle distribution. Starting with the classical case, the distribution $f_{cl}(\mathbf{p}, \mathbf{r}, t)$ is proportional to the probability density of finding a particle of momentum \mathbf{p} and position \mathbf{r} in the phase-space volume $d^3p d^3r$. This is a purely classical interpretation, directly conflicting with the uncertainty principle. The quantum mechanical quasidistribution $f_w(\mathbf{p}, \mathbf{r}, t)$, however, is not positive definite and has to be interpreted as a joint density of \mathbf{p} and \mathbf{r} [3]. Only the marginal distributions are positive definite, that is, integrating $f_w(\mathbf{p}, \mathbf{r}, t)$ over momentum space gives the probability density in \mathbf{r} -space, and vice versa.

An excellent review of quantum-mechanical phase-space distributions in scattering theory has been given by Carruthers and Zachariason [4]. This work deals with potential scattering, the two-body problem, and the N-body problem. A coupled hierarchy for reduced distribution functions and its truncation to the Boltzmann-Vlasov equation is presented. Tatarskii [3] concentrates on quantum-mechanical systems in a pure state and investigates the representation of quantum mechanics by phase-space distributions. He points out that not every function that solves the Wigner equation describes a pure state. Therefore, initial conditions for the Wigner equation have to be subjected to a supplementary restriction. Today, phase-space quantization is considered to be a third autonomous and logically complete formulation of quantum mechanics beyond the conventional ones based on operators in Hilbert space or path integrals [5, 6]. This formulation is free of operators and wave functions. Observables and matrix elements are computed through phase-space integrals of c-number functions weighted by a Wigner function.

Important quantum mechanical properties of electronic transport in semiconductor structures are often those associated not with the degeneracy of the Fermi system but rather with quantum interference effects [7]. A wide variety of electronic quantum transport problems of interest are essentially one-particle in nature. In such cases, a full many-body description of the problem is not necessary, and a description of electronic transport that makes use of the one-particle approximation can be used from the very outset. However, even when the electron–electron interaction effects are of interest, certain approximations do exist, allowing their description on a one-particle level [7]. Therefore, we shall consider in the following only electronic systems with one-particle degrees of freedom.

1.1. History and State of the Art Review

Reports on finite-difference solutions of the one-particle Wigner equation for device applications are due to Ravaioli [8], Klusdahl [9], and coworkers, and date back to the mid 1980s. Frensley [10–12] was the first who introduced boundary conditions on the Wigner function to model open quantum systems. Later, self-consistency was added to the Wigner equation solvers [13, 14]. Main and Haddad included a reduced Boltzmann scattering operator in transient Wigner function–based simulations [15]. Research on finite-difference solution methods for the Wigner equation culminated in 1990 when the review articles of Frensley [16] and Buot and Jensen [17] appeared.

The 1990s have seen further extensions and applications of the finite-difference Wigner function method. High-frequency operation of resonant tunneling diodes has been studied by Jensen and Buot [18, 19], and the transient response by Gullapalli [20] and Biegel [21], and later by [22]. A new finite-difference discretization scheme has been proposed in [23].

In 2002, implementations of Monte Carlo methods for solving the Wigner device equation were reported [24, 25]. Although with the finite-difference method, scattering was restricted to the relaxation time approximation and the momentum space to one dimension, the Monte Carlo method allows scattering processes to be included on a more detailed level, assuming a three-dimensional momentum-space [26, 27]. Issues such as choosing proper up-winding schemes, restrictions on matrix size and momentum space resolution are largely relaxed or do not exist when using the Monte Carlo method. Construction of new Monte Carlo algorithms is complicated by the fact that the kernel of the integral equation to solve is not positive semidefinite. As a consequence, the commonly applied Markov chain Monte Carlo method shows a variance exponentially increasing with time, prohibiting its application to realistic structures or larger evolution times [25, 28, 29]. Because of this so-called negative sign problem, the concept of Wigner paths alone [30, 31] is not sufficient to construct a stable Monte Carlo algorithm. Instead, additional measures have to be introduced that prevent a runaway of the particle weights and hence of the variance [26, 32]. Note that in [26], the statistical weights are termed affinities.

Large basic research efforts on the Monte Carlo modeling of electron–phonon interaction based on the Wigner function formalism have been reported in [28, 31, 33–35].

The effect of a spatially varying effective mass in Wigner device simulations has been demonstrated in [36] and [37]. A nonparabolic version of the Wigner equation has been derived by Bufler [38]. Multiband models have been reported in [39–41].

A Wigner equation including a magnetic field has been solved in [42]. The gauge-invariant formulation of the Wigner equation has been given by Levinson [43], and a discussion can be found in various works [4, 44–47]. Two-time and frequency-dependent Wigner functions are considered in [2, 47–49].

Finally, we note that the Wigner function formalism is often used to derive reduced transport models, such as the quantum hydrodynamic model [50, 51–53], or to find quantum corrections to classical models, such as the ensemble Monte Carlo method [54] or the spherical harmonics expansion method [55, 56].

2. THE WIGNER FUNCTION FORMALISM

In the Schrödinger picture, a physical system is quantum-mechanically described by a state vector $|\Psi(t)\rangle$ as function of time t . Often, the precise quantum-mechanical state of a system is not known, but rather some statistical information about the probabilities for the system being in one of a set of states. Suppose that there is a set of orthonormal states $\{|\Psi_1\rangle, |\Psi_2\rangle, \dots\}$, and that the probabilities that the system is in one of these states are $\{p_1, p_2, \dots\}$. Then, the expectation value of operator \hat{A} associated with the observable A is given by

$$\langle A \rangle = \sum_i p_i \langle \Psi_i | \hat{A} | \Psi_i \rangle \quad (1)$$

which is a quantum and statistical average. Introducing the density operator $\hat{\rho}$ as

$$\hat{\rho} = \sum_i p_i |\Psi_i\rangle \langle \Psi_i| \quad (2)$$

the expectation value becomes

$$\langle A \rangle = \text{Tr}(\hat{\rho} \hat{A}) = \text{Tr}(\hat{A} \hat{\rho}) \quad (3)$$

Formulations (1) and (3) require the operator \hat{A} to be self-adjoint. Equation (3) can be easily verified by expressing the trace of some operator \hat{X} in the basis $\{|\Psi_i\rangle\}$.

$$\text{Tr} \langle \hat{X} \rangle = \sum_i \langle \Psi_i | \hat{X} | \Psi_i \rangle \quad (4)$$

The fact that the probabilities sum up to unity, $\sum_i p_i = 1$, is expressed by the fact that the trace of the density operator is also unity, $\text{Tr}(\hat{\rho}) = 1$. If the system is in a pure state $|\Psi_i\rangle$ it

holds $p_i = 1$ and $p_j = 0 \forall j \neq i$, and the density operator is idem-potent, $\hat{\rho}^2 = \hat{\rho}$. Otherwise, the system is in a mixed state, and $\hat{\rho}$ does not obey the idem-potency condition. From the Schrodinger equation for the state vector and the definition of $\hat{\rho}$, we immediately obtain the Liouville-von Neumann equation for the evolution of the density operator.

$$i\hbar \frac{\partial \hat{\rho}}{\partial t} = [\hat{H}, \hat{\rho}] \quad (5)$$

Introducing the one-particle approximation [7] implies that the electron system is modeled as consisting of many, noninteracting electrons. In the next step, one chooses the coordinate representation, where the set of basis vectors is given by the electron position eigenstates $|\mathbf{r}\rangle$. The eigenstates of the system are then represented by the wavefunctions $\Psi_i(\mathbf{r}, t) = \langle \mathbf{r} | \Psi_i(t) \rangle$, and the density operator by the density matrix $\rho(\mathbf{r}_1, \mathbf{r}_2, t)$.

$$\rho(\mathbf{r}_1, \mathbf{r}_2, t) = \langle \mathbf{r}_1 | \hat{\rho}(t) | \mathbf{r}_2 \rangle = \sum_i p_i \Psi_i(\mathbf{r}_1, t) \Psi_i^*(\mathbf{r}_2, t) \quad (6)$$

The Liouville-von Neumann equation in coordinate representation is found as

$$\frac{\partial \rho(\mathbf{r}_1, \mathbf{r}_2, t)}{\partial t} = (H_{\mathbf{r}_1} - H_{\mathbf{r}_2}) \rho(\mathbf{r}_1, \mathbf{r}_2, t) \quad (7)$$

2.1. The Wigner Function

The Wigner function is obtained from the density matrix by means of the Wigner-Weyl transformation. This transformation consists of a change of independent coordinates to diagonal and cross-diagonal coordinates

$$\mathbf{r} = \frac{1}{2}(\mathbf{r}_1 + \mathbf{r}_2), \quad \mathbf{s} = \mathbf{r}_1 - \mathbf{r}_2 \quad (8)$$

followed by a Fourier transformation with respect to \mathbf{s} [16]. The variables \mathbf{r}_1 and \mathbf{r}_2 may be expressed in terms of the new ones.

$$\mathbf{r}_1 = \mathbf{r} + \frac{\mathbf{s}}{2}, \quad \mathbf{r}_2 = \mathbf{r} - \frac{\mathbf{s}}{2} \quad (9)$$

Then, the elementary definition of the Wigner distribution is given by the following transformation of the density matrix.

$$f_w(\mathbf{k}, \mathbf{r}, t) = \int \rho\left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t\right) e^{-i\mathbf{k}\cdot\mathbf{s}} d\mathbf{s} \quad (10)$$

The Wigner function (10) is real-valued, but not positive semidefinite. In terms of the wave functions, the definition (10) becomes

$$f_w(\mathbf{k}, \mathbf{r}, t) = \sum_i p_i \int \Psi_i\left(\mathbf{r} + \frac{\mathbf{s}}{2}, t\right) \Psi_i^*\left(\mathbf{r} - \frac{\mathbf{s}}{2}, t\right) e^{-i\mathbf{k}\cdot\mathbf{s}} d\mathbf{s} \quad (11)$$

The normalization of the Wigner function results from the normalization of the wave functions.

$$\frac{1}{(2\pi)^3} \int d\mathbf{r} \int d\mathbf{k} f_w(\mathbf{k}, \mathbf{r}, t) = 1 \quad (12)$$

Here, the \mathbf{k} -integration can be performed first, giving $\int e^{-i\mathbf{k}\cdot\mathbf{s}} d\mathbf{k} = (2\pi)^3 \delta(\mathbf{s})$. The normalization (12) ensures that the quantity Nf_w , where N is the number of electrons in the system, will approach the classical distribution function f_{cl} in the classical limit [35].

Sometimes it is convenient to use the inverse Fourier transform of (10).

$$\rho\left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t\right) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) e^{i\mathbf{k}\cdot\mathbf{s}} d\mathbf{k} \quad (13)$$

Changing variables gives a transformation that inverts the Wigner–Weyl transformation.

$$\rho(\mathbf{r}_1, \mathbf{r}_2, t) = \frac{1}{(2\pi)^3} \int f_w\left(\mathbf{k}, \frac{\mathbf{r}_1 + \mathbf{r}_2}{2}, t\right) e^{i\mathbf{k} \cdot (\mathbf{r}_1 - \mathbf{r}_2)} d\mathbf{k} \quad (14)$$

An important feature of the phase-space approach is the possibility of expressing quantum-mechanical expectation values in the same way as it is done in classical statistical mechanics, employing integration over the phase-space. The expectation values of operators of the form $A(\hat{\mathbf{r}})$ and $B(\hat{\mathbf{k}})$, where $\hat{\mathbf{k}} = \hat{\mathbf{p}}/\hbar$, are given as follows.

$$\langle A(\hat{\mathbf{r}}) \rangle = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) A(\mathbf{r}) d\mathbf{k} d\mathbf{r} = \sum_i p_i \int A(\mathbf{r}) |\Psi_i(\mathbf{r}, t)|^2 d\mathbf{r} \quad (15)$$

$$\langle B(\hat{\mathbf{k}}) \rangle = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) B(\mathbf{k}) d\mathbf{k} d\mathbf{r} = \sum_i p_i \int B(\mathbf{k}) |\Phi_i(\mathbf{k}, t)|^2 d\mathbf{k} \quad (16)$$

If the classical observable $C(\mathbf{k}, \mathbf{r})$ is a function of both momentum and position, the definition of a corresponding Hermitian operator \hat{C} is not unique. In this case, the Weyl quantization can be applied. Thereby, the function C is expressed through its Fourier transform c .

$$C(\mathbf{k}, \mathbf{r}) = \int c(\mathbf{a}, \mathbf{b}) e^{i(\mathbf{k} \cdot \mathbf{a} + \mathbf{r} \cdot \mathbf{b})} d\mathbf{a} d\mathbf{b} \quad (17)$$

The operator \hat{C} is defined by the following rule of correspondence.

$$\hat{C} = \int c(\mathbf{a}, \mathbf{b}) e^{i(\hat{\mathbf{k}} \cdot \mathbf{a} + \hat{\mathbf{r}} \cdot \mathbf{b})} d\mathbf{a} d\mathbf{b} \quad (18)$$

Then, the expectation value of \hat{C} is given by the phase-space integral.

$$\text{Tr}(\hat{C}\hat{\rho}) = \int C(\mathbf{k}, \mathbf{r}) f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} d\mathbf{r} \quad (19)$$

To proceed with (18), one may employ the Baker–Campbell–Hausdorff formula,

$$e^{\hat{A} + \hat{B}} = e^{\hat{A}} e^{\hat{B}} e^{-\frac{[\hat{A}, \hat{B}]}{2}} \quad (20)$$

which is generally valid when $[\hat{A}, [\hat{A}, \hat{B}]] = [\hat{B}, [\hat{A}, \hat{B}]] = 0$, or in particular when $[\hat{A}, \hat{B}]$ is a c-number.

2.2. Marginal Distributions

The Wigner function (10) can assume negative values. Only the marginal distributions of $f_w(\mathbf{k}, \mathbf{r}, t)$ are positive semidefinite and have the meaning of probability distributions in real space and momentum space, respectively.

$$n(\mathbf{r}) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} = \sum_i p_i |\Psi_i(\mathbf{r}, t)|^2 \quad (21)$$

$$p(\mathbf{k}) = \frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{r} = \sum_i p_i |\Phi_i(\mathbf{k}, t)|^2 \quad (22)$$

Here, $\Phi_i(\mathbf{k}, t)$ denotes the momentum representation of the state vector $|\Psi_i\rangle$. The integration in (22) can easily be carried out after changing variables, using (8).

$$\begin{aligned} & \int d\mathbf{r} \int d\mathbf{s} \Psi_i\left(\mathbf{r} + \frac{\mathbf{s}}{2}, t\right) \Psi_i^*\left(\mathbf{r} - \frac{\mathbf{s}}{2}, t\right) e^{-i\mathbf{k} \cdot \mathbf{s}} \\ &= \int d\mathbf{r}_1 \int d\mathbf{r}_2 \Psi_i(\mathbf{r}_1, t) \Psi_i^*(\mathbf{r}_2, t) e^{-i\mathbf{k} \cdot (\mathbf{r}_1 - \mathbf{r}_2)} = (2\pi)^3 |\Phi_i(\mathbf{k}, t)|^2 \end{aligned} \quad (23)$$

The marginal distributions (21) and (22) can also be expressed as the diagonal elements of the density matrix.

$$\frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} = \langle \mathbf{r} | \hat{\rho} | \mathbf{r} \rangle = \rho(\mathbf{r}, \mathbf{r}) \quad (24)$$

$$\frac{1}{(2\pi)^3} \int f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{r} = \langle \mathbf{k} | \hat{\rho} | \mathbf{k} \rangle = \sigma(\mathbf{k}, \mathbf{k}) \quad (25)$$

Here, $|\mathbf{k}\rangle$ denotes the electron momentum eigenstate with eigenvalue $\hbar\mathbf{k}$ and σ the density matrix in momentum representation. Note that the latter can be used for a dual definition of the Wigner function [28, 57].

$$f_w(\mathbf{k}, \mathbf{r}, t) = \int \sigma\left(\mathbf{k} + \frac{\mathbf{l}}{2}, \mathbf{k} - \frac{\mathbf{l}}{2}, t\right) e^{i\mathbf{r}\cdot\mathbf{l}} d\mathbf{l} \quad (26)$$

This definition follows, for example, from (11), when the Ψ_i are replaced by

$$\Psi_i(\mathbf{r}, t) = (2\pi)^{-3/2} \int \Phi_i(\mathbf{k}', t) e^{i\mathbf{k}'\cdot\mathbf{r}} d\mathbf{k}' \quad (27)$$

Other marginal distributions than the elementary ones, (21) and (22), have to be constructed with care. Only Hermitian operators give real marginal distributions. For the current density, this operator would be $(\hat{\mathbf{k}}\hat{\rho} + \hat{\rho}\hat{\mathbf{k}})/2$. Expressing $\hat{\rho}$ in terms of the wave functions, we get the elementary current definition from wave mechanics.

$$\begin{aligned} \mathbf{j}(\mathbf{r}) &= \frac{\hbar}{2m^*} \langle \mathbf{r} | \hat{\mathbf{k}}\hat{\rho} + \hat{\rho}\hat{\mathbf{k}} | \mathbf{r} \rangle \\ &= \frac{\hbar}{2m^*} \sum_i p_i (\langle \mathbf{r} | \hat{\mathbf{k}} | \Psi_i \rangle \langle \Psi_i | \mathbf{r} \rangle + \langle \mathbf{r} | \Psi_i \rangle \langle \Psi_i | \hat{\mathbf{k}} | \mathbf{r} \rangle) \\ &= \frac{\hbar}{2im^*} \sum_i p_i [\Psi_i^*(\mathbf{r}) \nabla \Psi_i(\mathbf{r}) - \Psi_i(\mathbf{r}) \nabla \Psi_i^*(\mathbf{r})] \end{aligned} \quad (28)$$

Choosing the momentum representation of $\hat{\rho}$, we get the current density expressed in terms of the Wigner function.

$$\begin{aligned} \mathbf{j}(\mathbf{r}) &= \frac{\hbar}{2m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 (\langle \mathbf{r} | \hat{\mathbf{k}} | \mathbf{k}_1 \rangle \langle \mathbf{k}_1 | \hat{\rho} | \mathbf{k}_2 \rangle \langle \mathbf{k}_2 | \mathbf{r} \rangle + \langle \mathbf{r} | \mathbf{k}_1 \rangle \langle \mathbf{k}_1 | \hat{\rho} | \mathbf{k}_2 \rangle \langle \mathbf{k}_2 | \hat{\mathbf{k}} | \mathbf{r} \rangle) \\ &= \frac{\hbar}{2m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 \sigma(\mathbf{k}_1, \mathbf{k}_2) (\mathbf{k}_1 + \mathbf{k}_2) e^{i(\mathbf{k}_1 - \mathbf{k}_2)\cdot\mathbf{r}} = \frac{1}{(2\pi)^3} \int \frac{\hbar}{m^*} \mathbf{k} f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} \end{aligned} \quad (29)$$

Here, the Wigner function has been introduced using (26). The current density is given by the first-order moment of the Wigner function, in full analogy with the classical phase space definition.

For the definition of the energy density we discuss several options. Starting from the trace operation for the statistical average, one would consider the symmetrized operator $(\hat{\mathbf{k}}^2\hat{\rho} + \hat{\rho}\hat{\mathbf{k}}^2)/2$ and derive the marginal distribution.

$$\begin{aligned} w_1(\mathbf{r}) &= \frac{\hbar^2}{4m^*} (\langle \mathbf{r} | \hat{\mathbf{k}}^2 \hat{\rho} | \mathbf{r} \rangle + \langle \mathbf{r} | \hat{\rho} \hat{\mathbf{k}}^2 | \mathbf{r} \rangle) \\ &= -\frac{\hbar^2}{4m^*} \sum_i p_i [\nabla^2 \Psi_i(\mathbf{r}) + \Psi_i(\mathbf{r}) \nabla^2 \Psi_i^*(\mathbf{r})] \\ &= \sum_i p_i [E_i - V(\mathbf{r})] |\Psi_i(\mathbf{r})|^2 \end{aligned} \quad (30)$$

The last expression in (30) is obtained with the help of the stationary Schrodinger equation. Apparently, w_1 describes the kinetic energy density, as the potential energy term $V(\mathbf{r})n(\mathbf{r})$ is subtracted from the total energy term. This energy density can become negative in tunneling

regions, where for one or more states $E_i < V(\mathbf{r})$ holds. In a derivation similar to (29), one finds the Wigner representation of w_1 .

$$\begin{aligned} w_1(\mathbf{r}) &= \frac{\hbar^2}{4m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 \sigma(\mathbf{k}_1, \mathbf{k}_2) (\mathbf{k}_1^2 + \mathbf{k}_2^2) e^{i(\mathbf{k}_1 - \mathbf{k}_2) \cdot \mathbf{r}} \\ &= \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} \left(|\mathbf{k}|^2 - \frac{1}{4} \nabla_{\mathbf{r}}^2 \right) f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} \end{aligned} \quad (31)$$

To ensure positiveness of the energy density, in [58] the Hermitian operator $\hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}}$ is considered. Its marginal distribution can be shown to be positive semidefinite.

$$w_2(\mathbf{r}) = \frac{\hbar^2}{2m^*} \langle \mathbf{r} | \hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}} | \mathbf{r} \rangle = \frac{\hbar^2}{2m^*} \sum_i p_i |\langle \mathbf{r} | \hat{\mathbf{k}} | \Psi_i \rangle|^2 \quad (32)$$

$$= \frac{\hbar^2}{2m^*} \sum_i p_i |\nabla \Psi_i(\mathbf{r})|^2 \geq 0 \quad (33)$$

The Wigner representation of w_2 is obtained as

$$\begin{aligned} w_2(\mathbf{r}) &= \frac{\hbar^2}{4m^*} \int d\mathbf{k}_1 \int d\mathbf{k}_2 \sigma(\mathbf{k}_1, \mathbf{k}_2) (\mathbf{k}_1^2 - \mathbf{k}_2^2) e^{i(\mathbf{k}_1 - \mathbf{k}_2) \cdot \mathbf{r}} \\ &= \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} \left(|\mathbf{k}|^2 + \frac{1}{4} \nabla_{\mathbf{r}}^2 \right) f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} \end{aligned} \quad (34)$$

Conditions for obtaining non-negative marginal distributions are theoretically discussed in [59]. The Weyl correspondence (18) gives the definition of the energy density as the second-order moment of the Wigner function.

$$w_3(\mathbf{r}) = \frac{1}{(2\pi)^3} \int \frac{\hbar^2}{2m^*} |\mathbf{k}|^2 f_w(\mathbf{k}, \mathbf{r}, t) d\mathbf{k} \quad (35)$$

It can be seen that (35) is just the arithmetic mean of (31) and (34), $w_3 = (w_1 + w_2)/2$. Therefore, (35) represents the marginal distribution of the symmetrized operator $(\hat{\mathbf{k}}^2\hat{\rho} + 2\hat{\mathbf{k}}\hat{\rho}\hat{\mathbf{k}} + \hat{\rho}\hat{\mathbf{k}}^2)/4$.

All three definitions of the energy density give the same statistical average $\langle \hat{\epsilon} \rangle = \text{Tr}[\epsilon(\hat{\mathbf{k}})\hat{\rho}]$. The differences among the definitions are in the ∇^2 term, which vanishes after the \mathbf{r} -integration. However, only the density w_1 seems to have a clear physical interpretation as the kinetic energy density.

2.3. The Wigner Equation

In this section, we consider a system consisting of one electron interacting with a potential distribution $V_{\text{tot}}(\mathbf{r})$. This potential is assumed to be a superposition of some potential $V(\mathbf{r})$ and a uniform electric field: $V_{\text{tot}}(\mathbf{r}) = V(\mathbf{r}) - \hbar\mathbf{F} \cdot \mathbf{r}$, with $\hbar\mathbf{F} = -e\mathbf{E}$. Although the existence of a field term is not physically motivated at this point, it is introduced here to demonstrate its treatment in the Wigner function formalism. The potential $V(\mathbf{r})$ comprises the electrostatic potential and the band-edge profile of the semiconductor. A uniform effective mass m^* is assumed. In the usual coordinate representation, the Hamiltonian of the system is then given by

$$H = H_0 + V(\mathbf{r}) - \hbar\mathbf{F} \cdot \mathbf{r} \quad (36)$$

with

$$H_0 = -\frac{\hbar^2}{2m^*} \nabla_{\mathbf{r}}^2 \quad (37)$$

The electron phonon interaction neglected here will be discussed in detail in Section 3. The evolution equation for the Wigner function is found by taking the time derivative of the defining Eq. (10) and substituting the Liouville-von Neumann Eq. (7) on the right-hand side.

$$\frac{\partial}{\partial t} f_w(\mathbf{k}, \mathbf{r}, t) = \frac{1}{i\hbar} \int (H_{r_1} - H_{r_2}) \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (38)$$

In the following, the three parts of the Hamiltonian (36) will be separately transformed. Unlike in Section 2.2, where calculations were done in momentum representation, we choose below the configuration representation to carry out the transformations [33].

The free-electron Hamiltonian is given by H_0 . To calculate the Wigner transform of H_0 , we have to transform the gradients first. Differentiating the density matrix with respect to the new variables r and s

$$\nabla_r \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) = \nabla_{r_1} \rho + \nabla_{r_2} \rho, \quad \nabla_s \rho \left(\mathbf{r} + \frac{\mathbf{s}}{2}, \mathbf{r} - \frac{\mathbf{s}}{2}, t \right) = \frac{1}{2} \nabla_{r_1} \rho - \frac{1}{2} \nabla_{r_2} \rho \quad (39)$$

gives the relations

$$\nabla_{r_1} + \nabla_{r_2} = \nabla_r, \quad \nabla_{r_1} - \nabla_{r_2} = 2\nabla_s, \quad \nabla_{r_1}^2 - \nabla_{r_2}^2 = 2\nabla_r \cdot \nabla_s \quad (40)$$

Now the free-electron term transforms to a diffusion term. For the sake of brevity, we write $\rho_{r,s} = \rho(\mathbf{r} + \mathbf{s}/2, \mathbf{r} - \mathbf{s}/2, t)$ in the following.

$$\frac{1}{i\hbar} \int -\frac{\hbar^2}{2m^*} (\nabla_{r_1}^2 - \nabla_{r_2}^2) \rho_{r,s} e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} = -\frac{\hbar}{im^*} \nabla_r \cdot \int (\nabla_s \rho_{r,s}) e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (41)$$

$$= -\frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r \int \rho_{r,s} e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (42)$$

$$= -\frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r f_w(\mathbf{k}, \mathbf{r}, t) \quad (43)$$

Next, we transform the potential term $V(\mathbf{r})$.

$$\frac{1}{i\hbar} \int \left[V \left(\mathbf{r} + \frac{\mathbf{s}}{2} \right) - V \left(\mathbf{r} - \frac{\mathbf{s}}{2} \right) \right] \rho_{r,s} e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' \quad (44)$$

This transformation is readily found by replacing $p_{r,s}$ on the left-hand side by the inverse Fourier transformation (13). The remaining integral over s is denoted by V_w and referred to as the Wigner potential.

$$V_w(\mathbf{q}, \mathbf{r}) = \frac{1}{(2\pi)^3 i\hbar} \int \left[V \left(\mathbf{r} + \frac{\mathbf{s}}{2} \right) - V \left(\mathbf{r} - \frac{\mathbf{s}}{2} \right) \right] e^{-i\mathbf{q} \cdot \mathbf{s}} d\mathbf{s} \quad (45)$$

Using the simple relation $-(\mathbf{F} \cdot \mathbf{r}_1 - \mathbf{F} \cdot \mathbf{r}_2) = -\mathbf{F} \cdot \mathbf{s}$, the constant-field term transforms as

$$\frac{1}{i\hbar} \int (-\hbar \mathbf{F} \cdot \mathbf{s}) \rho_{r,s} e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} = -\frac{1}{\hbar} \mathbf{F} \cdot \nabla_k f_w(\mathbf{k}, \mathbf{r}, t) \quad (46)$$

Collecting the above results gives the Wigner equation for the system Hamiltonian (36).

$$\left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r + \mathbf{F} \cdot \nabla_k \right) f_w(\mathbf{k}, \mathbf{r}, t) = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' \quad (47)$$

The terms are arranged so to form the classical Liouville operator on the left-hand side. The interaction of the electron with the potential distribution $V(\mathbf{r})$ is described by the potential operator on the right-hand side. As can be seen, the Wigner function in k and r depends in a nonlocal manner on the Wigner function in all other momentum points k' and through V_w also on the potential at all other locations $r \pm s/2$.

3. ELECTRON-PHONON INTERACTION

The Wigner equation has frequently been solved using the finite-difference method [16, 60], assuming the phenomenological relaxation time approximation for dissipative transport. Recently developed Monte Carlo methods allowed phonon scattering to be included semi-classically in quantum device simulations [24, 27]. Use of a Boltzmann scattering operator acting on the Wigner distribution was originally suggested by Frensley [16]. In this section, the Wigner equation with a Boltzmann scattering operator is rigorously derived, using a many-phonon single-electron Wigner function formalism as the starting point.

3.1. The System Hamiltonian

The Hamiltonian (36) is now extended to describe a system consisting of one electron interacting with a many-phonon system and a given potential distribution.

$$H = H_0 + V(\mathbf{r}) - \mathbf{F} \cdot \mathbf{r} + H_p + H_{ep} \quad (48)$$

The additional components of this Hamiltonian are given by [34]

$$H_p = \sum_{\mathbf{q}} \hbar \omega_{\mathbf{q}} b_{\mathbf{q}}^{\dagger} b_{\mathbf{q}} \quad (49)$$

$$H_{ep} = i\hbar \sum_{\mathbf{q}} \mathcal{F}(\mathbf{q}) (b_{\mathbf{q}} e^{i\mathbf{q} \cdot \mathbf{r}} - b_{\mathbf{q}}^{\dagger} e^{-i\mathbf{q} \cdot \mathbf{r}}) \quad (50)$$

Here, H_p is the Hamiltonian of the free phonon-system, H_{ep} the electron-phonon interaction Hamiltonian, $b_{\mathbf{q}}$ and $b_{\mathbf{q}}^{\dagger}$ denote the annihilation and creation operators for a phonon with momentum $\hbar \mathbf{q}$ and energy $\hbar \omega_{\mathbf{q}}$, and $\hbar \mathcal{F}(\mathbf{q})$ is the interaction matrix element.

We introduce a set of basis vectors $|\mathbf{r}, \{n\}\rangle$ in the occupation number representation. A set of occupation numbers is defined as $\{n\} = n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}_q}, \dots$, where $n_{\mathbf{q}}$ is the number of phonons with momentum \mathbf{q} . The Wigner-Weyl transformation of the density matrix $\rho(\mathbf{r}_1, \{n\}, \mathbf{r}_2, \{m\})$ gives the generalized Wigner function $f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$ [28, 33].

$$f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) = \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \right| \hat{\rho}(t) \left| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (51)$$

Note that only the electron coordinates are transformed, such that f_g is a Wigner function on the electron phase-space, but still is the density matrix for the phonon system.

The evolution of the generalized Wigner function is found by taking the time derivative of (51) and using the Liouville-von Neumann equation for the evolution of the density matrix.

$$\frac{\partial}{\partial t} f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) = \frac{1}{i\hbar} \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \right| [\hat{H}, \hat{\rho}(t)] \left| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (52)$$

To continue, one may express the density matrix in the state vectors of the system.

$$\rho(\mathbf{r}_1, \{n\}, \mathbf{r}_2, \{m\}, t) = \sum_i p_i \Psi_i(\mathbf{r}_1, \{n\}, t) \Psi_i^*(\mathbf{r}_2, \{m\}, t) \quad (53)$$

The creation and annihilation operators, and the occupation number operator $b_{\mathbf{q}}^{\dagger} b_{\mathbf{q}}$ satisfy the following well-known eigenvalue equations.

$$\begin{aligned} b_{\mathbf{q}}^{\dagger} \Psi(\mathbf{r}, \{n\}, t) &= \sqrt{n_{\mathbf{q}} + 1} \Psi(\mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}} + 1, \dots\}, t) \\ b_{\mathbf{q}} \Psi(\mathbf{r}, \{n\}, t) &= \sqrt{n_{\mathbf{q}}} \Psi(\mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}} - 1, \dots\}, t) \\ b_{\mathbf{q}}^{\dagger} b_{\mathbf{q}} \Psi(\mathbf{r}, \{n\}, t) &= n_{\mathbf{q}} \Psi(\mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}}, \dots\}, t) \end{aligned} \quad (54)$$

With the help of these equations and the representation (53), the transformation of the free-phonon Hamiltonian is readily found.

$$\begin{aligned} & \frac{1}{i\hbar} \int \left\langle \mathbf{r} + \frac{\mathbf{s}}{2}, \{n\} \left| \left[\widehat{H}_p, \hat{\rho}(t) \right] \right| \mathbf{r} - \frac{\mathbf{s}}{2}, \{m\} \right\rangle e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \\ &= \frac{1}{i\hbar} (\epsilon(\{n\}) - \epsilon(\{m\})) f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) \end{aligned}$$

The energy of the phonon state $|\{n\}\rangle$ is denoted by $\epsilon(\{n\})$.

$$\epsilon(\{n\}) = \sum_{\mathbf{q}} n_{\mathbf{q}} \hbar \omega_{\mathbf{q}} \quad (55)$$

The electron-phonon interaction Hamiltonian is transformed following the same lines [33]. Combining the two terms of the Hamiltonian (50) and the two terms of the commutator in (52) results in four terms related to the electron-phonon interaction. In the equation for the generalized Wigner function shown below, these four terms appear under the sum.

$$\begin{aligned} & \left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_{\mathbf{r}} + \frac{1}{\hbar} \mathbf{F} \cdot \nabla_{\mathbf{k}} \right) f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) \\ &= \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_g(\mathbf{k}', \mathbf{r}, \{n\}, \{m\}, t) d\mathbf{k}' + \frac{1}{i\hbar} (\epsilon(\{n\}) - \epsilon(\{m\})) f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t) \\ &+ \sum_{\mathbf{q}} \mathcal{F}(\mathbf{q}) e^{i\mathbf{q} \cdot \mathbf{r}} \sqrt{n_{\mathbf{q}} + 1} f_g\left(\mathbf{k} - \frac{\mathbf{q}}{2}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}} + 1, \dots\}, \{m\}, t\right) \\ &- e^{-i\mathbf{q} \cdot \mathbf{r}} \sqrt{n_{\mathbf{q}}} f_g\left(\mathbf{k} + \frac{\mathbf{q}}{2}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}} - 1, \dots\}, \{m\}, t\right) \\ &- e^{i\mathbf{q} \cdot \mathbf{r}} \sqrt{m_{\mathbf{q}}} f_g\left(\mathbf{k} + \frac{\mathbf{q}}{2}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \dots, m_{\mathbf{q}} - 1, \dots\}, t\right) \\ &+ e^{-i\mathbf{q} \cdot \mathbf{r}} \sqrt{m_{\mathbf{q}} + 1} f_g\left(\mathbf{k} - \frac{\mathbf{q}}{2}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \dots, m_{\mathbf{q}} + 1, \dots\}, t\right) \end{aligned} \quad (56)$$

Each term under the sum represents a phonon interaction event that changes only one set of phonon variables, increasing or decreasing the occupation number of the single-phonon state $|\mathbf{q}\rangle$ by one and changing the electron momentum by $\pm \mathbf{q}/2$.

3.2. A Hierarchy of Transport Equations

The equation for the generalized Wigner function (56) is too complex for the purpose of mesoscopic device simulation. Several approximations need to be introduced in order to arrive at a more feasible quantum transport equation. In the following, these approximations are discussed.

3.2.1. Weak Scattering Limit

The generalized Wigner equation couples one element of the phonon density matrix, $f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m\}, t)$, with four neighboring elements,

$$f_g(\mathbf{k}, \mathbf{r}, \{n_{\mathbf{q}_1}, n_{\mathbf{q}_2}, \dots, n_{\mathbf{q}} \pm 1, \dots\}, \{m\}, t) \quad (57)$$

$$f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{m_{\mathbf{q}_1}, m_{\mathbf{q}_2}, \dots, m_{\mathbf{q}} \pm 1, \dots\}, t) \quad (58)$$

The equations for the four nearest neighbor elements couple to second nearest neighbors of the element $\{n\}, \{m\}$, and so forth. In the weak scattering limit, all couplings between elements of the first and the second off-diagonals are neglected. Only the main diagonal terms and the first off-diagonal terms remain, as shown in Fig. 1. Higher order electron-phonon interactions are neglected in this way.

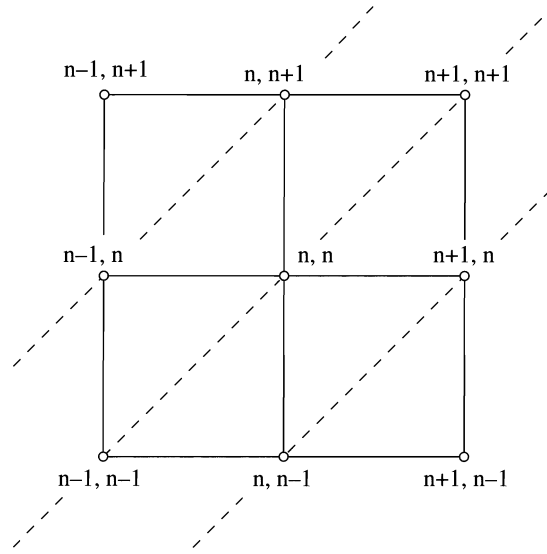


Figure 1. Terms of the phonon density matrix retained in the weak scattering limit.

3.2.2. The Reduced Wigner Function

The reduced Wigner function, $f_w(\mathbf{k}, \mathbf{r}, t)$, is defined as the trace of the generalized Wigner function over all phonon states [28, 61].

$$f_w(\mathbf{k}, \mathbf{r}, t) = \sum_{\{n\}} f_g(\mathbf{k}, \mathbf{r}, \{n\}, \{n\}, t) \quad (59)$$

Further approximations are needed to evaluate this trace and hence to derive a closed equation for the reduced Wigner function [62]. One approximation is to replace any occupation number $n_{\mathbf{q}}$ involved in a transition by the equilibrium phonon number, $N_{\mathbf{q}}$, and to assume that the phonon system stays in equilibrium during the evolution of the electron state. With these assumptions, the trace operation can be performed, and a closed equation set for the reduced Wigner function can be obtained. The set consists of an equation for the reduced Wigner function coupled to two auxiliary equations.

$$\begin{aligned} & \left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_{\mathbf{r}} - \Theta_w \right) f_w(\mathbf{k}, \mathbf{r}, t) \\ &= 2 \text{Re} \sum_{\mathbf{q}} \mathcal{F}^2(\mathbf{q}) \left\{ e^{i\mathbf{q} \cdot \mathbf{r}} f_1 \left(\mathbf{k} - \frac{\mathbf{q}}{2}, \mathbf{r}, t \right) - e^{-i\mathbf{q} \cdot \mathbf{r}} f_2 \left(\mathbf{k} + \frac{\mathbf{q}}{2}, \mathbf{r}, t \right) \right\} \end{aligned} \quad (60)$$

In this equation, we denote the Wigner potential operator by Θ_w and set the classical force to $\mathbf{F} = 0$.

$$\Theta_w[f_w](\mathbf{k}, \mathbf{r}, t) = \int V_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' \quad (61)$$

The auxiliary equations arise from the first off-diagonal terms of the equation for the generalized Wigner function. In the following equation, the lower sign gives f_1 and the upper f_2 :

$$\begin{aligned} & \left(\frac{\partial}{\partial t} + \frac{\hbar}{m^*} \left(\mathbf{k} \pm \frac{\mathbf{q}}{2} \right) \cdot \nabla_{\mathbf{r}} \mp i\omega_{\mathbf{q}} - \Theta_w \right) f_{1,2} \left(\mathbf{k} \pm \frac{\mathbf{q}}{2}, \mathbf{r}, t \right) \\ &= \pm e^{\pm i\mathbf{q} \cdot \mathbf{r}} \left\{ \left(N_{\mathbf{q}} + \frac{1}{2} \mp \frac{1}{2} \right) f_w(\mathbf{k}, \mathbf{r}, t) - \left(N_{\mathbf{q}} + \frac{1}{2} \pm \frac{1}{2} \right) f_w(\mathbf{k} \pm \mathbf{q}, \mathbf{r}, t) \right\} \end{aligned} \quad (62)$$

Although the equation for the reduced Wigner function is real-valued, the two auxiliary equations are complex-valued. Note that f_w depends either on some initial momentum \mathbf{k} or

the momentum after a completed electron-phonon interaction, $k \pm q$. On the other hand, f_1 and f_2 depend on intermediate states $k \pm \mathbf{q}/2$, where only half of the phonon momentum has been transferred.

3.2.3. Mean Field Approximation

To simplify the equation system, one may assume a mean field over the length scale of an electron-phonon interaction. This mean field can be set to the local force field $\hbar \mathbf{F}(\mathbf{r}) = -\nabla V(\mathbf{r})$. Note that this field is kept constant during an electron-phonon interaction event, even though the electron moves on an \mathbf{r} -space trajectory. For a uniform electric-field, the potential operator becomes local, $\Theta_w[f_w] = -\mathbf{F} \cdot \nabla_{\mathbf{k}} f_w$, and the two auxiliary equations (62) can explicitly be solved. The solutions $f_{1,2}$ are expressed as path integrals over the reduced Wigner function. In this way, a single equation for the reduced Wigner function is derived from (60).

$$\left(\frac{d}{dt} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_{\mathbf{r}} - \Theta_w \right) f_w(\mathbf{k}, \mathbf{r}, t) = \int_0^t d\tau \int d\mathbf{k}' [S(\mathbf{k}, \mathbf{k}', \tau) f_w(\mathbf{k}' - \mathbf{F}\tau, \mathbf{R}(\mathbf{k}, \mathbf{k}', \tau), t - \tau) - S(\mathbf{k}', \mathbf{k}, \tau) f_w(\mathbf{k} - \mathbf{F}\tau, \mathbf{R}(\mathbf{k}, \mathbf{k}', \tau), t - \tau)] \quad (63)$$

The scattering kernel is of the form

$$S(\mathbf{k}', \mathbf{k}, \tau) = \frac{2V}{(2\pi)^3} \mathcal{F}^2(\mathbf{q}) \sum_{\nu=\pm 1} \left(N_{\mathbf{q}} + \frac{1}{2} - \frac{\nu}{2} \right) \cos \int_0^\tau \frac{1}{\hbar} \left(\epsilon_{(\mathbf{k}-\mathbf{F}\tau')} - \epsilon_{(\mathbf{k}'-\mathbf{F}\tau')} + \nu \hbar \omega_{\mathbf{q}} \right) dt' \quad (64)$$

and the \mathbf{r} -space trajectory defined as

$$\mathbf{R}(\mathbf{k}, \mathbf{k}', \tau) = \mathbf{r} - \frac{\hbar(\mathbf{k} + \mathbf{k}')}{2m^*} \tau + \frac{\hbar \mathbf{F}}{2m^*} \tau^2 \quad (65)$$

To interpret the above equations, we assume some phase space point \mathbf{k}, \mathbf{r} and some time t to be given. A transition from \mathbf{k} to \mathbf{k}' as described by (64) starts in the past, at time $t - \tau$, where the retarded momentum $\mathbf{k} - \mathbf{F}\tau$ has to be considered [see (63)]. At the beginning of the electron-phonon interaction, half of the phonon momentum is transferred, which determines the initial momentum $\mathbf{k} - \mathbf{F}\tau \pm \mathbf{q}/2$ of a phase space trajectory. With $\mathbf{k}' = \mathbf{k} \pm \mathbf{q}$, the initial momentum becomes

$$\mathbf{k} - \mathbf{F}\tau \pm \frac{\mathbf{q}}{2} = \frac{\mathbf{k} + \mathbf{k}'}{2} - \mathbf{F}\tau \quad (66)$$

During the interaction duration τ , the particle drifts over a phase space trajectory and arrives at \mathbf{r} and $\mathbf{k} \pm \mathbf{q}/2$ at time t . At this time, the electron-phonon interaction is completed by the transfer of another $\pm \mathbf{q}/2$, which produces the final momentum $\mathbf{k} \pm \mathbf{q}$. Also included are virtual phonon emission and absorption processes, where the initial momentum transfer $\pm \mathbf{q}/2$ at $t - \tau$ is compensated by $\mp \mathbf{q}/2$ at t . This model thus includes effects due to a finite collision duration, such as collisional broadening and the intra-collisional field effect. A discussion of the integral form of (63) can be found in [63].

3.2.4. Levinson Equation

For a uniform electric field and an initial condition independent of \mathbf{r} , (63) simplifies to the Levinson equation [43].

$$\left(\frac{\partial}{\partial t} + \mathbf{F} \cdot \nabla_{\mathbf{k}} \right) f_w(\mathbf{k}, t) = \int_0^t d\tau \int d\mathbf{k}' [S(\mathbf{k}, \mathbf{k}', \tau) f_w(\mathbf{k}' - \mathbf{F}\tau, t - \tau) - S(\mathbf{k}', \mathbf{k}, \tau) f_w(\mathbf{k} - \mathbf{F}\tau, t - \tau)] \quad (67)$$

S is given by (64). This equation is equivalent to the Barker-Ferry equation [64] with an infinite electron lifetime. Recently, Monte Carlo methods for the solution of the Levinson equation have been developed, which allow the numerical study of collisional broadening, retardation effects, and the intracollisional field effect [65, 66].

3.2.5. Classical Limit

The classical limit of the scattering operator in (63) is obtained by an asymptotic analysis. For this purpose, the equation is written in a dimensionless form. The primary scaling factors are k_0 for the wave-vector \mathbf{k} and t_0 for the time t . Additional scaling factors to be introduced are s_0 for the scattering rate S , ϵ_0 for the energy ϵ , F_0 for the force \mathbf{F} , r_0 for the real-space vector \mathbf{r} , and $\omega_{\mathcal{F}}$ for the interaction matrix element \mathcal{F} .

The key issue is now to choose an appropriate scale k_0 . Scaling the phonon energy to unity gives $\hbar k_0^2 = m^* \omega_q$. The kinetic equation is now considered on a timescale that is much larger than the timescale of the lattice vibrations. Therefore, one sets $t_0 = (\epsilon \omega_q)^{-1}$, where $\epsilon \ll 1$ denotes a dimensionless parameter. The remaining scaling factors are found as

$$s_0 = \frac{1}{t_0^2 k_0^3}, \quad \epsilon_0 = \frac{\hbar^2 k_0^2}{m^*} \quad (68)$$

$$F_0 = \frac{k_0}{t_0}, \quad r_0 = \frac{\hbar k_0 t_0}{m^*} \quad (69)$$

The frequency scale of the electron–phonon interaction can be chosen as $\omega_{\mathcal{F}} = \mathcal{F}(q_{\text{th}})$, where q_{th} is the wave number of a thermal electron. The scaled Levinson equation has the same form as the unscaled equation (67). The scaled scattering rate varies on a time scale of order ϵ^{-1} . To keep the time integral of order $O(1)$, the amplitude of the scattering rate should be of order ϵ^{-1} as well, which is obtained by setting [67]

$$\epsilon = \frac{2V\omega_{\mathcal{F}}^2}{(2\pi)^3} \sqrt{\frac{m_*^3}{\hbar^3 \omega_q}} \quad (70)$$

This gives a scaled scattering rate of the form

$$S(\mathbf{k}', \mathbf{k}, \tau) = \frac{\mathcal{F}^2(\mathbf{q})}{\epsilon} \sum_{\nu=\pm 1} \left(N_q + \frac{1}{2} - \frac{\nu}{2} \right) \cos \int_0^\tau \frac{1}{\epsilon} \left(\frac{(\mathbf{k} - \mathbf{F}t')^2}{2} - \frac{(\mathbf{k}' - \mathbf{F}t')^2}{2} + \nu \right) dt' \quad (71)$$

The classical limit is valid in the regime where the quantity defined by (70) is small, and thus for timescales $t_0 = (\epsilon \omega_q)^{-1}$ much larger than the inverse phonon frequency. The scattering operator in (67) converges for $\epsilon \rightarrow 0$ to the Fermi golden rule operator in the weak sense. From the asymptotic analysis also a first-order correction to the Fermi golden rule is found [67]. Using parameters for GaAs at room temperature, one computes $\epsilon = 0.011$, which suggests that assuming the asymptotic regime is appropriate.

A heuristic argument for the convergence to the golden rule is as follows. Changing variables in the scattering operator in (67) gives

$$Q[f_w](\mathbf{k}, t) = \int_0^{t/\epsilon} d\tau \int d\mathbf{k}' [\epsilon S(\mathbf{k}, \mathbf{k}', \epsilon\tau) f_w(\mathbf{k}' - \epsilon\mathbf{F}\tau, t - \epsilon\tau) - \epsilon S(\mathbf{k}', \mathbf{k}, \epsilon\tau) f_w(\mathbf{k} - \epsilon\mathbf{F}\tau, t - \epsilon\tau)]$$

and

$$\epsilon S(\mathbf{k}', \mathbf{k}, \epsilon\tau) = \mathcal{F}^2(\mathbf{q}) \sum_{\nu=\pm 1} \left(N_q + \frac{1}{2} - \frac{\nu}{2} \right) \cos \left[(\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu) \tau - \epsilon\mathbf{F} \cdot (\mathbf{k} - \mathbf{k}') \frac{\tau^2}{2} \right]$$

Expanding f_w and ϵS into a Taylor series in ϵ and keeping only terms of zeroth order leads to the integral,

$$\int_0^\infty \cos[(\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu) \tau] d\tau = \pi \delta[\epsilon(\mathbf{k}) - \epsilon(\mathbf{k}') + \nu] \quad (72)$$

which evaluates to the energy-conserving S-function of the golden rule. Undoing the scaling gives the well-known form of the scattering rate.

$$S(\mathbf{k}', \mathbf{k}) = \frac{V}{(2\pi)^3} \sum_{\nu=\pm 1} \frac{2\pi}{\hbar} M^2(\mathbf{q}) \left(N_{\mathbf{q}} + \frac{1}{2} - \frac{\nu}{2} \right) \delta[\epsilon(\mathbf{k}') - \epsilon(\mathbf{k}) + \nu \hbar \omega_{\mathbf{q}}] \quad (73)$$

The interaction matrix element is denoted here by $\mathbf{M} = \hbar \mathcal{F}$. Introducing the total scattering rate $\lambda(\mathbf{k}) = \int S(\mathbf{k}', \mathbf{k}) d\mathbf{k}'$, the Boltzmann scattering operator takes on the following form.

$$\mathcal{Q}[f_w](\mathbf{k}, \mathbf{r}, t) = \int S(\mathbf{k}, \mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' - \lambda(\mathbf{k}, \mathbf{r}) f(\mathbf{k}, \mathbf{r}, t) \quad (74)$$

Finally, we consider the classical limit of the potential operator. Scaling the r-dependent equation (63) gives the scaled form

$$\Theta_w[f](\mathbf{k}, \mathbf{r}, t) = \frac{1}{i(2\pi)^3 \varepsilon} \int d\mathbf{k}' \int d\mathbf{s} f(\mathbf{k}', \mathbf{r}, t) \left[V\left(\mathbf{r} + \frac{\varepsilon \mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\varepsilon \mathbf{s}}{2}\right) \right] e^{-i(\mathbf{k}-\mathbf{k}') \cdot \mathbf{s}}$$

This expression converges for $\varepsilon \rightarrow 0$ to the classical drift term of the Boltzmann equation.

$$\Theta_{cl}[f](\mathbf{k}, \mathbf{r}, t) = \nabla_r V(\mathbf{r}) \cdot \nabla_k f(\mathbf{k}, \mathbf{r}, t) \quad (75)$$

3.2.6. Wigner Equation with Boltzmann Scattering Operator

To obtain a model more suitable for device simulation, the nonlocal potential operator is maintained, whereas in the scattering operator the classical limit is introduced. The result is a Wigner equation with a Boltzmann scattering operator. It is convenient to introduce formally a classical force field $\mathbf{F}(\mathbf{r})$ in this equation to make the form of the Liouville operator equal to that of the Boltzmann equation. This is accomplished by redefining the potential operator.

$$\tilde{V}_w(\mathbf{k}, \mathbf{r}) = \frac{1}{(2\pi)^3 i \hbar} \int \left(V\left(\mathbf{r} + \frac{\mathbf{s}}{2}\right) - V\left(\mathbf{r} - \frac{\mathbf{s}}{2}\right) + \hbar \mathbf{F} \cdot \mathbf{s} \right) e^{-i\mathbf{k} \cdot \mathbf{s}} d\mathbf{s} \quad (76)$$

Substituting $\Theta_w[f] = \tilde{\Theta}_w[f] - \mathbf{F} \cdot \nabla_k f$ into (63) gives the following equation,

$$\left(\frac{\partial}{\partial t} + \frac{\hbar \mathbf{k}}{m^*} \cdot \nabla_r + \mathbf{F}(\mathbf{r}) \cdot \nabla_k \right) f_w = \mathcal{Q}[f_w] + \tilde{\Theta}_w[f_w] \quad (77)$$

From a formal point of view, the classical force field \mathbf{F} can be chosen arbitrarily, as the corresponding terms in (77) cancel each other. Typical choices are the mean electric field in a device region, the local electric field, or, of course, $\mathbf{F} = 0$. Alternatively, an equation of the form (77) can also be obtained by using an approximation. The potential is decomposed as $V = V_{cl} + V_{qm}$, where V_{cl} is a smooth potential such as the electrostatic potential, that can be treated in the classical limit (75), and V_{qm} represents a rapidly varying component that has to be treated quantum mechanically.

3.3. Integral Form of the Wigner Equation

From the integro-differential form of the Wigner equation, a path-integral formulation can be derived. The equation to be transformed reads

$$\begin{aligned} & \left(\frac{\partial}{\partial t} + \mathbf{v}(\mathbf{k}) \cdot \nabla_r + \mathbf{F}(\mathbf{r}) \cdot \nabla_k \right) f_w(\mathbf{k}, \mathbf{r}, t) \\ &= \int [S(\mathbf{k}, \mathbf{k}') + \tilde{V}_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r}) S(\mathbf{k} - \mathbf{k}')] f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' \\ & \quad - [\lambda(\mathbf{k}, \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r})] f_w(\mathbf{k}, \mathbf{r}, t) \end{aligned} \quad (78)$$

At this point, we introduced a fictitious scattering mechanism $\alpha \delta(\mathbf{k} - \mathbf{k}')$, referred to as self-scattering [68]. Because of the δ -function, this mechanism does not change the state of the

electron and hence does not affect the solution of the equation. For the sake of brevity, we define an integral kernel Γ and the symbols μ and U .

$$\mu(\mathbf{k}, \mathbf{r}) = \lambda(\mathbf{k}, \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r}) \quad (79)$$

$$\Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r}) = \frac{S(\mathbf{k}, \mathbf{k}') + \tilde{V}_w(\mathbf{k} - \mathbf{k}', \mathbf{r}) + \alpha(\mathbf{k}, \mathbf{r}) \delta(\mathbf{k} - \mathbf{k}')}{\mu(\mathbf{k}', \mathbf{r})} \quad (80)$$

$$U(\mathbf{k}, \mathbf{r}, t) = \int \Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r}) \mu(\mathbf{k}', \mathbf{r}) f_w(\mathbf{k}', \mathbf{r}, t) d\mathbf{k}' \quad (81)$$

The Liouville operator in (78) is treated by the method of characteristics. One introduces path variables $\mathbf{K}(t)$ and $\mathbf{R}(t)$ and takes the total time derivative of f_w .

$$\frac{d}{dt} f_w(\mathbf{K}(t), \mathbf{R}(t), t) = \left(\frac{\partial}{\partial t} + \frac{d\mathbf{K}(t)}{dt} \cdot \nabla_{\mathbf{k}} + \frac{d\mathbf{R}(t)}{dt} \cdot \nabla_{\mathbf{r}} \right) f_w \quad (82)$$

The right-hand side equals the Liouville operator if the path variables satisfy the following equations of motion.

$$\frac{d}{dt} \mathbf{K}(t) = \mathbf{F}(\mathbf{R}(t)) \quad \frac{d}{dt} \mathbf{R}(t) = \mathbf{v}(\mathbf{K}(t)) \quad (83)$$

Now we assume some phase-space point \mathbf{k}, \mathbf{r} and some time t to be given. A phase-space trajectory with the initial condition $\mathbf{K}(t' = t) = \mathbf{k}$ and $\mathbf{R}(t' = t) = \mathbf{r}$ is obtained by formal integration.

$$\mathbf{K}(t') = \mathbf{k} + \int_t^{t'} \mathbf{F}(\mathbf{R}(y)) dy \quad \mathbf{R}(t') = \mathbf{r} + \int_t^{t'} \mathbf{v}(\mathbf{K}(y)) dy \quad (84)$$

Note that $\mathbf{k}, \mathbf{r}, t$ are treated as constants in the following derivation, only t' is a variable. Introducing the functions

$$\tilde{f}_w(t') = f_w(\mathbf{K}(t'), \mathbf{R}(t'), t'), \quad \tilde{\mu}(t') = \mu(\mathbf{K}(t'), \mathbf{R}(t')), \quad \tilde{U}(t') = U(\mathbf{K}(t'), \mathbf{R}(t'), t') \quad (85)$$

allows (78) to be rewritten as an ordinary differential equation of first order.

$$\frac{d}{dt'} \tilde{f}_w(t') + \tilde{\mu}(t') \tilde{f}_w(t') = \tilde{U}(t') \quad (86)$$

If multiplied by an integrating factor $\exp[\int_0^{t'} \tilde{\mu}(y) dy]$, the equation takes on a form that can be easily integrated in time.

$$\frac{d}{dt'} \exp\left[\int_0^{t'} \tilde{\mu}(y) dy\right] \tilde{f}_w(t') = \exp\left[\int_0^{t'} \tilde{\mu}(y) dy\right] \tilde{U}(t') \quad (87)$$

The choice of the upper and lower bounds of time integration depends on whether the problem under consideration is time-dependent or stationary.

The ordinary differential equation (87), which is the result of treating the Liouville operator by the method of characteristics, has the same structure as the corresponding differential equation for the Boltzmann equation. Therefore, we can refer to the work on the Boltzmann equation regarding the details of the time integration of (87) [69, 70].

3.3.1. The Time-Dependent Equation

The upper bound of the time integration should be $t' = t$ to obtain $\tilde{f}_w(t) = f_w(\mathbf{k}, \mathbf{r}, t)$, the value of the unknown at the given phase space point. At $t' = 0$, an initial distribution $f_i(\mathbf{k}, \mathbf{r})$ is assumed to be given. In analogy with the Boltzmann equation [70], the integral form of the Wigner equation is obtained.

$$\begin{aligned} f_w(\mathbf{k}, \mathbf{r}, t) = & \int_0^t dt' \int d\mathbf{k}' \exp \left\{ - \int_{t'}^t \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \\ & \times \Gamma[\mathbf{K}(t'), \mathbf{k}', \mathbf{R}(t')] \mu[\mathbf{k}', \mathbf{R}(t')] f_w[\mathbf{k}', \mathbf{R}(t'), t'] \\ & + \exp \left\{ - \int_0^t \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} f_i(\mathbf{K}(0), \mathbf{R}(0)) \end{aligned} \quad (88)$$

This equation states that the Wigner function at time t depends on the Wigner function at some previous time t' . Using (88) in an iterative procedure, with each iteration the time variable would move to smaller values. Therefore, another equation is desirable that describes the evolution of the system in forward time direction. Such an equation is given by the adjoint equation of (88).

$$\begin{aligned} g_w(\mathbf{k}', \mathbf{r}', t') = & \int_{t'}^\infty d\tau \int d\mathbf{k} g_w[\mathbf{K}(\tau), \mathbf{R}(\tau), \tau] \exp \left\{ - \int_{t'}^\tau \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \\ & \times \Gamma(\mathbf{k}, \mathbf{k}', \mathbf{r}') \mu(\mathbf{k}', \mathbf{r}') + g_0(\mathbf{k}', \mathbf{r}', t') \end{aligned} \quad (89)$$

The derivation of the adjoint equation (89) is discussed in detail in [69, 70].

3.3.2. The Stationary Equation

In a stationary system, the potential and all material parameters are independent of time. A phase-space trajectory is invariant under time translations. This property can be conveniently used to adjust the time reference of each trajectory [71, 72]. In the stationary case, we assume the phase-space point \mathbf{k}, \mathbf{r} to be given at $t' = 0$. So the initial condition for the phase-space trajectory is $\mathbf{K}(0) = \mathbf{k}$ and $\mathbf{R}(0) = \mathbf{r}$. For the upper bound of time integration of (87), we choose now $t' = 0$ to obtain $\tilde{f}_w(0) = f_w(\mathbf{k}, \mathbf{r})$. The lower time bound has to be chosen such that the functions $\mathbf{K}(t)$ and $\mathbf{R}(t)$ take on values at which the Wigner function is known. In the steady-state, this function is known only at the domain boundary. An appropriate lower time bound is therefore the time when the trajectory enters the simulation domain. This time is denoted by t_b^- and depends on the point \mathbf{k}, \mathbf{r} under consideration. The case that the real space trajectory $\mathbf{R}(t)$ never intersects the domain boundary can occur for a classically bound state. Then the trajectory forms a closed loop and the appropriate choice is $t_b^- = -\infty$. Integration of (87) in the time bounds discussed above results in the integral form of the stationary Wigner equation (cf. [71]).

$$\begin{aligned} f(\mathbf{k}, \mathbf{r}) = & f_0(\mathbf{k}, \mathbf{r}) + \int_{t_b^-(\mathbf{k}, \mathbf{r})}^0 dt' \int d\mathbf{k}' \exp \left\{ - \int_{t'}^0 \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \\ & \times \Gamma[\mathbf{K}(t'), \mathbf{k}', \mathbf{R}(t')] \mu(\mathbf{k}', \mathbf{r}') f_w[\mathbf{k}', \mathbf{R}(t')] \end{aligned} \quad (90)$$

$$f_0(\mathbf{k}, \mathbf{r}) = f_b \{ \mathbf{K}[t_b^-(\mathbf{k}, \mathbf{r})], \mathbf{R}[t_b^-(\mathbf{k}, \mathbf{r})] \} \exp \left\{ - \int_{t_b^-(\mathbf{k}, \mathbf{r})}^0 \lambda[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \quad (91)$$

Here, f_b denotes the boundary distribution. The integral form (90) represents a backward equation. The corresponding forward equation is given by the adjoint equation.

$$\begin{aligned} g_w(\mathbf{k}, \mathbf{r}) = & g_0(\mathbf{k}, \mathbf{r}) + \int_0^{t_b^+(\mathbf{k}', \mathbf{r})} dt \int d\mathbf{k}' g_w[\mathbf{K}(\tau), \mathbf{R}(\tau)] \\ & \times \exp \left\{ - \int_t^\tau \mu[\mathbf{K}(y), \mathbf{R}(y)] dy \right\} \Gamma(\mathbf{k}', \mathbf{k}, \mathbf{r}) \mu(\mathbf{k}, \mathbf{r}) \Theta_D(\mathbf{r}) \end{aligned} \quad (92)$$

Θ_D denotes the indicator function of the simulation domain D . The initial conditions for the phase space trajectory are $\mathbf{K}(t) = \mathbf{k}'$ and $\mathbf{R}(t) = \mathbf{r}$.

4. THE MONTE CARLO METHOD

Monte Carlo is a numerical method that can be applied to solve integral equations. Applying this method to the various integral formulations of the Wigner equation gives rise to a variety of Monte Carlo algorithms, as discussed in the following.

4.1. The General Scheme

This section introduces the general scheme of the Monte Carlo method and outlines its application to the solution of integrals and integral equations. To calculate some unknown value m by the Monte Carlo method, one has to find a random variable ξ whose expectation value equals $E\{\xi\} = m$. The variance of ξ is designated σ^2 , with σ being the standard deviation.

Now consider N independent random variables $\xi_1, \xi_2, \dots, \xi_N$ with distributions identical to that of ξ . Consequently, their expectation values and their variance are equal.

$$E\{\xi_i\} = m, \quad \text{Var}\{\xi_i\} = \sigma^2, \quad i = 1, 2, \dots, N \quad (93)$$

Expectation value and variance of the sum of all these random variables are given by

$$E\{\xi_1 + \xi_2 + \dots + \xi_N\} = E\{\xi_1\} + E\{\xi_2\} + \dots + E\{\xi_N\} = Nm \quad (94)$$

$$\text{Var}\{\xi_1 + \xi_2 + \dots + \xi_N\} = \text{Var}\{\xi_1\} + \text{Var}\{\xi_2\} + \dots + \text{Var}\{\xi_N\} = N\sigma^2 \quad (95)$$

Using the properties $E\{c\xi\} = cE\{\xi\}$ and $\text{Var}\{c\xi\} = c^2\text{Var}\{\xi\}$, one obtains from (94) and (95)

$$E\left\{\frac{1}{N}(\xi_1 + \xi_2 + \dots + \xi_N)\right\} = m \quad (96)$$

$$\text{Var}\left\{\frac{1}{N}(\xi_1 + \xi_2 + \dots + \xi_N)\right\} = \frac{\sigma^2}{N} \quad (97)$$

Therefore, the random variable

$$\bar{\xi} = \frac{1}{N} \sum_{i=1}^N \xi_i \quad (98)$$

has the same expectation value as ξ and an N times reduced variance. A Monte Carlo simulation of the unknown m consists of drawing one random number $\bar{\xi}$. Indeed, this is equivalent to drawing N values of the random variable ξ , and evaluating the sample mean (98).

The Monte Carlo method gives an estimate of both the result and the error. According to the central limit theorem, the sum $\rho_N = \xi_1 + \xi_2 + \dots + \xi_N$ of a large number of identical random variables is approximately normal. For this reason, the following three-sigma rule holds only approximately

$$P\{|\rho_N - Nm| < 3\sqrt{N\sigma^2}\} \approx 0.997 \quad (99)$$

In this equation, the expectation value and the variance of ρ_N are given by (94) and (95), respectively. Dividing the inequality by N and using $\bar{\xi} = \rho_N/N$ we arrive at an equivalent inequality and the probability will not change:

$$P\left\{|\bar{\xi} - m| < 3\frac{\sigma}{\sqrt{N}}\right\} \approx 0.997 \quad (100)$$

This formula indicates that the sample mean $\bar{\xi}$ will be approximately equal to m . The error of this approximation will most probably not exceed the value $3\sigma/\sqrt{N}$. This error evidently approaches zero as N increases [73].

4.1.1. Monte Carlo Integration

We apply the Monte Carlo method to the evaluation of an integral.

$$m = \int_a^b \phi(x) dx \quad (101)$$

For this purpose, the integrand has to be decomposed into a product $\phi = p\psi$, where p is a density function, which means that p is non-negative and satisfies $\int_a^b p(x) dx = 1$. Integral (101) becomes

$$m = \int_a^b p(x)\psi(x) dx \quad (102)$$

and denotes the expectation value $m = E\{\Psi\}$ of some random variable $\Psi = \psi(X)$. Now the general scheme described in the previous section can be applied. First, a sample x_1, \dots, x_N is generated from the density p . Then the sample ψ_1, \dots, ψ_N is obtained by evaluating the function ψ : $\psi_i = \psi(x_i)$. The sample mean

$$m \simeq \bar{\psi} = \frac{1}{N} \sum_{i=1}^N \psi_i \quad (103)$$

approximates the expectation value. To employ the error estimation (100), the variance of Ψ can be approximately evaluated by the sample variance

$$\sigma^2 \simeq \bar{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (\psi_i - \bar{\psi})^2 \quad (104)$$

Because the factorization of the integrand is not unique, different random variables can be introduced depending on the choice of the density p . All of them have the same expectation value but different variance.

4.1.2. Integral Equations

The kinetic equations considered in this work can be formulated as integral equations of the form

$$f(x) = \int K(x, x')f(x') dx' + f_0(x) \quad (105)$$

where the kernel K and the source term f_0 are given functions. Equations of this form are known as Fredholm integral equations of the second kind. In the particular cases of the Boltzmann equation and the Wigner equation, the unknown function f represents the phase-space distribution function. The multidimensional variable x stands for (k, r, t) in the transient case and for (k, r) in the steady state.

Substituting (105) recursively into itself gives the Neumann series, which if convergent, is a formal solution to the integral equation [74].

$$f = f^{(0)} + f^{(1)} + f^{(2)} + \dots \quad (106)$$

The iteration terms are defined recursively beginning with $f^{(0)}(x) = f_0(x)$.

$$f^{(n+1)}(x) = \int K(x, x')f^{(n)}(x') dx', \quad n = 0, 1, 2, \dots \quad (107)$$

The series (106) yields the function value in some given point x . However, in many cases one is interested in mean values of f rather than in a point-wise evaluation. Such a mean value represents a linear functional and can be expressed as an inner product.

$$(f, A) = \int f(x)A(x) dx \quad (108)$$

It is to note that (105) is a backward equation. The corresponding forward equation is given by the adjoint equation,

$$g(x') = \int K^\dagger(x', x)g(x)dx + A(x') \quad (109)$$

where the kernel is defined by $K^\dagger(x', x) = K(x, x')$. Multiplying (105) by $g(x)$ and (109) by $f(x')$, and integrating over x and x' , respectively, results in the equality

$$(f, A) = (g, f_0) \quad (110)$$

By means of (110), one can calculate a statistical mean value not only from f , but also from g , the solution of the adjoint equation. The given function A has to be used as the source term of the adjoint equation. The link with the numerical Monte Carlo method is established by evaluating the terms of the Neumann series by Monte Carlo integration, as pointed out in the previous section.

Note that usage of (110) precludes a point-wise evaluation of the distribution function using a forward algorithm, because $A(x) = \delta(x)$ cannot be treated by the Monte Carlo method. The probability for a continuous random variable x' to assume a given value x is zero. Only the probability of finding x' within a small but finite volume around x is non-zero.

4.2. Particle Models

Each term of the Neumann series of the adjoint equation describes a sequence of alternating free flight and scattering events. A transition consisting of a free flight with initial state \mathbf{k}_i at time t_i and a scattering process to the final state \mathbf{k}_f at time t_f is described by the following expression. For the sake of brevity, the r -dependence of Γ and μ is omitted in the following.

$$P(\mathbf{k}_f, t_f, \mathbf{k}_i, t_i) = \Gamma[\mathbf{k}_f, \mathbf{K}_i(t_f)] \mu[\mathbf{K}_i(t_f)] \exp\left\{-\int_{t_i}^{t_f} \mu(\mathbf{K}_i(\tau)) d\tau\right\} \quad (111)$$

In a Monte Carlo simulation, t_f , the time of the next scattering event, is generated from an exponential distribution, given by the terms $\mu \exp()$ in (111). Then, a transition from the trajectory end point $\mathbf{K}_i(t_f)$ to the final state \mathbf{k}_f is realized using the kernel Γ . In contrast to the classical case, where P would represent a transition probability, such an interpretation is not possible in the case of the Wigner equation because P is not positive semidefinite. The problem originates from the Wigner potential, which assumes positive and negative values. However, because of its antisymmetry with respect to q , the Wigner potential can be reformulated in terms of one positive function V_w^+ [27].

$$V_w^+(\mathbf{q}, \mathbf{r}) = \max(0, V_w(\mathbf{q}, \mathbf{r})) \quad (112)$$

$$V_w(\mathbf{q}, \mathbf{r}) = V_w^+(\mathbf{q}, \mathbf{r}) - V_w^+(-\mathbf{q}, \mathbf{r}) \quad (113)$$

Then, the kernel Γ is rewritten as a sum of the following conditional probability distributions.

$$\Gamma(\mathbf{k}, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}' - \mathbf{k}) + \frac{\gamma}{\mu} [w(\mathbf{k}, \mathbf{k}') - w^*(\mathbf{k}, \mathbf{k}')] \quad (114)$$

$$s(\mathbf{k}, \mathbf{k}') = \frac{S(\mathbf{k}', \mathbf{k})}{\lambda(\mathbf{k}')}, \quad w(\mathbf{k}, \mathbf{k}') = \frac{V_w^+(\mathbf{k} - \mathbf{k}')}{\gamma}, \quad w^*(\mathbf{k}, \mathbf{k}') = w(\mathbf{k}', \mathbf{k}) \quad (115)$$

The normalization factor associated with the Wigner potential is defined as

$$\gamma(\mathbf{r}) = \int V_w^+(\mathbf{q}, \mathbf{r}) d\mathbf{q} \quad (116)$$

In the following, different variants of generating the final state \mathbf{k}_f from the kernel Γ will be discussed.

4.2.1. The Markov Chain Method

In analogy to the simple integral (102), we have now to decompose the kernel P into a transition probability p and the remaining function P/p . More details on the Markov chain method can be found in [75, 76]. With respect to (111), one could use the absolute value of Γ as a transition probability. Practically, it is more convenient to use the absolute values of the components of Γ , giving the following transition probability.

$$p(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\nu} s(\mathbf{k}_f, \mathbf{k}') + \frac{\alpha}{\nu} \delta(\mathbf{k}_f - \mathbf{k}') + \frac{\gamma}{\nu} w(\mathbf{k}_f, \mathbf{k}') + \frac{\gamma}{\nu} w^*(\mathbf{k}_f, \mathbf{k}') \quad (117)$$

The normalization factor is $\nu = \lambda + \alpha + 2\gamma$. In the first method considered here, the free-flight time is generated from the exponential distribution appearing in (111).

$$p_t(t_f, t_i, \mathbf{k}_i) = \mu[\mathbf{K}_i(t_f)] \exp \left\{ - \int_{t_i}^{t_f} \mu[\mathbf{K}_i(\tau)] d\tau \right\} \quad (118)$$

For the sake of brevity, the state at the end of the free flight is labeled $\mathbf{k}' = \mathbf{K}_i(t_f)$ in the following. To generate the final state \mathbf{k}_f , one of the four terms in (117) is selected with the associated probabilities λ/ν , α/ν , γ/ν , and γ/ν , respectively. Apparently, these probabilities sum up to one. If classical scattering is selected, \mathbf{k}_f is generated from s . If self-scattering is selected, the state does not change and $\mathbf{k}_f = \mathbf{k}'$ holds. If the third or fourth term are selected, the particle state is changed by scattering from the Wigner potential and \mathbf{k}_f is selected from w or w^* , respectively. The particle weight has to be multiplied by the ratio

$$\frac{\Gamma}{p} = \pm \left(1 + \frac{2\gamma}{\lambda + \alpha} \right) \quad (119)$$

where the minus sign applies if \mathbf{k}_f has been generated from w^* . For instance, for a quantum mechanical system, where the classical scattering rate λ is less than the Wigner scattering rate γ , the self-scattering rate α can be chosen in such a way that $\lambda + \alpha = \gamma$. Then, the multiplier (119) evaluates to ± 3 . An ensemble of particles would evolve as shown schematically in Fig. 2.

In the second method, we again use the transition rate (117), but now the free flight time is generated with rate ν rather than with μ . In this case, (111) can be rewritten as

$$P(\mathbf{k}_f, t_f, \mathbf{k}_i, t_i) = \frac{\Gamma(\mathbf{k}_f, \mathbf{k}') \mu(\mathbf{k}')}{p(\mathbf{k}_f, \mathbf{k}') \nu(\mathbf{k}')} p(\mathbf{k}_f, \mathbf{k}') \times \nu(\mathbf{k}') \exp \left\{ - \int_{t_i}^{t_f} \nu[\mathbf{K}_i(\tau)] d\tau \right\} \exp \left\{ 2 \int_{t_i}^{t_f} \gamma[\mathbf{R}_i(\tau)] d\tau \right\} \quad (120)$$

The exponential distribution distribution is used to generate t_f and the distribution p to generate \mathbf{k}_f . The remaining terms form the factor by which the particle weight changes

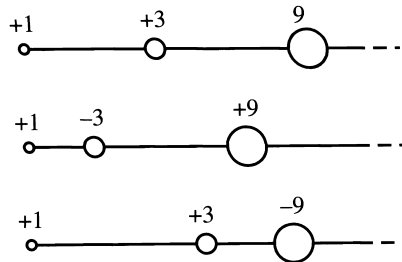


Figure 2. With the Markov chain method, the number of numerical particles is conserved. The magnitude of the particle weight increases with each event, and the sign of the weight changes randomly according to a given probability distribution.

during one free flight. Because of $(\Gamma\mu)/(p\nu) = \pm 1$, the multiplier for the i th free flight evaluates to

$$m_i = \pm \exp \left\{ 2 \int_{t_i}^{t_{i+1}} \gamma[\mathbf{R}_i(\tau)] d\tau \right\} \quad (121)$$

Note that the absolute values of both multipliers, (119) and (121), are always greater than one. With each transition of the Markov chain, the particle weight is multiplied by such factor. Therefore, the absolute value of the particle weight will inevitably grow with the number of transitions on the trajectory. To solve the problem of growing particle weights, one can split particles. In this way, an increase in particle weight is transformed to an increase in particle number.

4.2.2. Pair Generation Methods

The basic idea of splitting is refined so to avoid fractional weights. Different interpretations of the kernel are presented that conserve the magnitude of the particle weight. Choosing the initial weight to be +1, all generated particles will have weight +1 or −1. This is achieved by interpreting the potential operator in (77) as a generation term of positive and negative particles. We consider the kernel (114).

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_f, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}_f - \mathbf{k}') + \frac{\gamma}{\mu} [w(\mathbf{k}_f, \mathbf{k}') - w^*(\mathbf{k}_f, \mathbf{k}')] \quad (122)$$

If the Wigner scattering rate γ is larger than the classical scattering rate λ , the self-scattering rate α has to be chosen large enough to satisfy the inequality $\gamma/\mu \leq 1$. Typical choices are $\mu = \text{Max}(\lambda, \gamma)$ or $\mu = \lambda + \gamma$. These expressions also hold for the less interesting case $\gamma < \lambda$, where quantum interference effects are less important than classical scattering effects.

As in the classical Monte Carlo method, the distribution of the free-flight duration is given by the exponential distribution (118). At the end of a free flight, the complementary probabilities $p_s = \lambda/\mu$ and $1 - p_s = \alpha/\mu$ are considered. With probability p_s , classical scattering is selected. The final state is generated from s . The complementary event is self-scattering. In addition, with probability $p_w = \gamma/\mu$ a pair of particle states is generated from the distributions w and w^* . The multiplier of the weight is +1 for a state generated from one of first three terms and −1 for a state generated from w^* . Therefore, the magnitude of the initial particle weight is conserved, as shown in Fig. 3.

Method G1: In the following, we discuss the case $\gamma > \lambda$, where quantum effects are dominant. We begin with the smallest possible value for μ : $\mu = \text{Max}(\lambda, \gamma) = \gamma$. Because $p_w = \gamma/\mu = 1$, a particle pair is generated after each free flight as shown in Fig. 4. At the same instances, classical or self-scattering events occur. In Fig. 4 and the following figures, only the trajectory of a sample particle is shown and not the whole cascade of trajectories of the generated particles.

Method G2: Choosing the self-scattering rate to be $\alpha = \gamma$, the kernel can be regrouped as

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_f, \mathbf{k}') + \left(1 - \frac{\lambda}{\mu}\right) [\delta(\mathbf{k}_f - \mathbf{k}') + w(\mathbf{k}_f, \mathbf{k}') - w^*(\mathbf{k}_f, \mathbf{k}')] \quad (123)$$

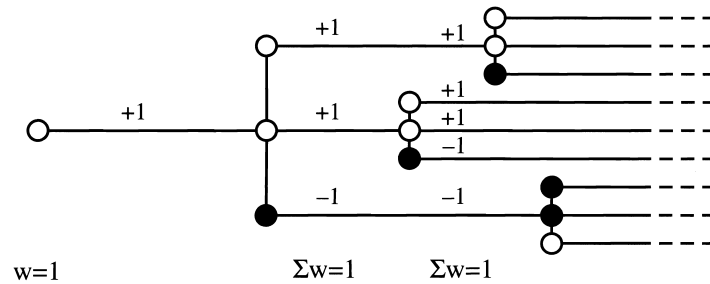


Figure 3. With the pair generation method, the magnitude of the particle weight is conserved, but one initial particle generates a cascade of numerical particles. At all times, mass is exactly conserved.

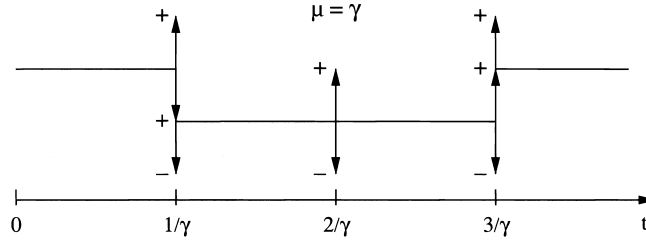


Figure 4. Trajectory of a sample particle resulting from method G1.

With probability $p_s = \lambda/\mu$, classical scattering is selected. Otherwise, a self-scattering event and a pair generation event occur. In this algorithm, classical scattering and pair generation cannot occur at the same time, as shown in Fig. 5. Compared to method G1, the average free flight time is now reduced, because μ has been increased from γ to $\lambda + \gamma$.

4.2.3. Single-Particle Generation Methods

The idea of this method is to further reduce the free-flight time. We rewrite the kernel as

$$\Gamma(\mathbf{k}_f, \mathbf{k}') = \frac{\lambda}{\mu} s(\mathbf{k}_f, \mathbf{k}') + \frac{\alpha}{\mu} \delta(\mathbf{k}_f - \mathbf{k}') + \frac{2\gamma}{\mu} \left[\frac{1}{2} w(\mathbf{k}_f, \mathbf{k}') - \frac{1}{2} w^*(\mathbf{k}_f, \mathbf{k}') \right] \quad (124)$$

In this case, the self-scattering rate α has to be chosen large enough to satisfy the inequality $2\gamma/\mu \leq 1$. Typical choices are $\mu = \text{Max}(\lambda, 2\gamma)$ and $\mu = \lambda + 2\gamma$. As in method G1, classical scattering is selected with probability $p_s = \lambda/\mu$, whereas the complementary event is self-scattering. In addition, with probability $p_w = 2\gamma/\mu$, particle generation is selected. If selected, with equal probability either the distribution w or w^* is chosen to generate the final state \mathbf{k}_f . If w^* has been chosen, the weight is multiplied by -1 .

Method G3: Assuming $\gamma > \lambda/2$ and $\mu = 2\gamma$ gives $p_w = 1$. Therefore, after each free flight, either a positive or negative particle is generated, as depicted in Fig. 6. At the same instances, classical or self-scattering events occur.

Note that in method G3 ($\mu = 2\gamma$), the free-flight time is reduced by a factor of two compared to method G1 ($\mu = \gamma$), which means that now the kernel is applied twice as frequently. In method G3, single particles are generated at a rate of 2γ , whereas in method G1 particle pairs are generated at half of this rate.

Method G4: In this method, we set $\alpha = 2\gamma$ and obtain $\mu = \lambda + 2\gamma$. In analogy with method G2, classical scattering and particle generation are now complementary events. Figure 7 indicates that these two types of events occur at different times. From all methods discussed above, this method uses the shortest free flight time.

From a numerical point of view, method G1 and method G2 have the advantage that they exactly conserve charge as they generate particles pairwise with opposite sign. Method G3 and method G4 generate only one particle each time. Because the sign of the weight is selected randomly, charge is conserved only on average. Simulation experiments, however, have shown that the quality of the pseudo random number generator is good enough to generate almost equally many positive and negative particles even during long simulation

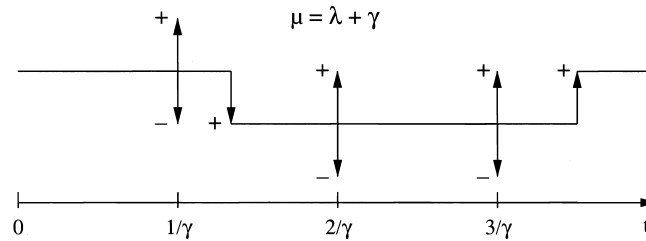


Figure 5. Trajectory of a sample particle resulting from method G2.

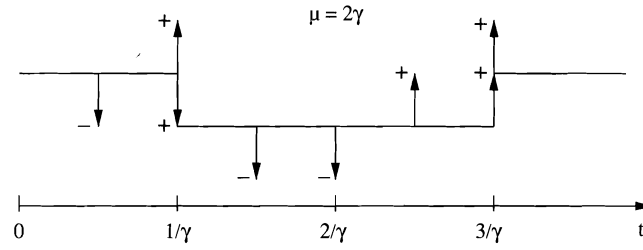


Figure 6. Trajectory of a sample particle resulting from method G3.

times, such that the small difference of net generated particles has no visible effect on the solution.

4.2.4. Other Methods

In method G1 to method G4, the weight of the generated particles is ± 1 , because the generation rate used equals 2γ . If a generation rate larger than 2γ or a fixed time-step less than $1/2\gamma$ were used, the magnitude of the generated weight would be less than one. This approach has been followed in [24], where the resulting fractional weights are termed affinities. On the other hand, a generation rate less than 2γ would result in an under-sampling of the physical process. Then, the magnitude of the generated weights would be generally greater than one.

4.3. The Negative Sign Problem

In the following, we analyze the growth rates of particle weights and particle numbers associated with the different Monte Carlo algorithms. In the first Markov chain method discussed in Section 4.2.1, the weight increases at each scattering event by the multiplier (119). The growth rate of the weight can be estimated for the case of constant coefficients γ and μ . Because free-flight times are generated with rate μ , the mean free-flight time will be $1/\mu$. During a given time interval t , on-average $n = \mu t$ scattering events will occur. The total weight is then estimated asymptotically for $t \gg 1/\mu$.

$$|W(t)| = \left(1 + \frac{2\gamma}{\mu}\right)^n = \left(1 + \frac{2\gamma t}{n}\right)^n \simeq \exp(2\gamma t) \quad (125)$$

This expression shows that the growth rate is determined by the Wigner scattering rate γ independently of the classical and the self-scattering rates.

With the second Markov chain method, one readily obtains that the total weight after n free flights grows as a function of the path integral over $\gamma[\mathbf{R}(\tau)]$.

$$|W(t_n)| = \prod_{i=0}^{n-1} |m_i| = \exp\left\{2 \int_0^{t_n} \gamma(\mathbf{R}(\tau)) d\tau\right\} \quad (126)$$

In this equation, the m_i are given by (121). This result generalizes (125) for a position-dependent γ . The growth rate 2γ is equal to the L_1 norm of the Wigner potential.

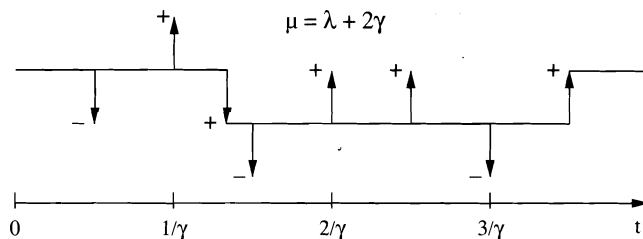


Figure 7. Trajectory of a sample particle resulting from method G4.

In the pair generation methods, the potential operator

$$\Theta_w[f_w](\mathbf{k}) = \int V^+(\mathbf{q})[f_w(\mathbf{k} - \mathbf{q}) - f_w(\mathbf{k} + \mathbf{q})] d\mathbf{q} \quad (127)$$

has been interpreted as a generation term. It describes the creation of two new states, $\mathbf{k} - \mathbf{q}$ and $\mathbf{k} + \mathbf{q}$. The generation rate is equal to γ . When generating the second state, the sign of the statistical weight is changed. It should be noted that the Wigner equation strictly conserves mass, as can be seen by taking the zero-order moment of (77).

$$\frac{\partial n}{\partial t} + \text{div } \mathbf{J} = 0 \quad (128)$$

Looking at the number of particles regardless of their statistical weights, that is, counting each particle as positive, would correspond to using the following potential operator.

$$\Theta_w^*[f_w](\mathbf{k}) = \int V_w^+(\mathbf{q})[f_w(\mathbf{k} - \mathbf{q}) + f_w(\mathbf{k} + \mathbf{q})] d\mathbf{q} \quad (129)$$

Using (129), a continuity equation for numerical particles is obtained.

$$\frac{\partial n^*}{\partial t} + \text{div } \mathbf{J}^* = 2\gamma(\mathbf{r})n^* \quad (130)$$

Assuming a constant γ , the generation rate in this equation will give rise to an exponential increase in the number of numerical particles N^* .

$$N^*(t) = N^*(0) \exp(2\gamma t) \quad (131)$$

This discussion shows that the appearance of an exponential growth rate is independent of the details of the particular Monte Carlo algorithm, and must be considered to be a fundamental consequence of the non-positive kernel.

4.4. Particle Annihilation

The discussed particle models are unstable, because either the particle weight or the particle number grows exponentially in time. Using the Markov chain method, it has been demonstrated that tunneling can be treated numerically by means of a particle model [25]. However, because of the exponentially increasing particle weight at the very short timescale $(2\gamma)^{-1}$, application of this algorithm turned out to be restricted to single-barrier tunneling and small barrier heights only. This method can be useful for devices where quantum effects are weak, and the potential operator is a small correction to the otherwise classical transport equation.

A stable Monte Carlo algorithm can be obtained by combining one of the particle generation methods with a method to control the particle number. One can assume that two particles of opposite weight and a sufficiently small distance in phase space annihilate each other. The reason is that the motions of both particles are governed by the same equation. Therefore, when they come close to each other at some time instant, the two particles have approximately the same initial condition and thus a common probabilistic future. In an ensemble Monte Carlo method, a particle removal step should be performed at given time steps. During the time step, the ensemble is allowed to grow to a certain limit, then particles are removed and the initial size of the ensemble is restored. In this work, the problem has been solved for the stationary transport problem. In the algorithm, the trajectory of only one sample particle is followed, whereas other numerical particles are temporarily stored on a phase space grid. Due to the opposite sign, particle weights annihilate to a large extent in the cells of the grid. The total residual weight in each cell has to be minimized, as it represents a measure for the numerical error of the method [32].

5. SIMULATION RESULTS

Virtually all published results of Wigner function–based device modeling focus on resonant tunneling diodes [77, 78]. In this section, three different devices are discussed. Their parameter values are collected in Table 1, where RTD1 [36] and RTD2 [24] are devices from literature. The semiclassical scattering model includes polar optical, acoustic deformation potential, and ionized-impurity scattering. Parameter values for GaAs have been assumed.

5.1. Comparison with Other Numerical Methods

RTD1 has been used as a benchmark device to compare different numerical approaches to quantum transport. In this device, the potential is assumed to vary linearly only in the double-barrier region and to be constant in the two contact regions. Results of the Monte Carlo method outlined in this work have been compared to nonequilibrium Green's function–based results [79]. The latter have been obtained by NEMO-1D, a one-dimensional nanoelectronic modeling tool [80]. NEMO-1D has served as a quantitatively predictive design and analysis tool for resonant tunneling diodes [81–83].

RTD1 shows a rather large coherent off-resonant valley current. Therefore, phonon scattering has only little effect on the current–voltage characteristics of this device. Both simulators predict only a slight increase in valley current due to inelastic scattering (Fig. 8). The resonance voltages predicted by the two solvers agree very well.

A comparison between finite-difference results and Monte Carlo results is shown in Fig. 9. An important parameter is the cutoff length L_c used in the numerical Wigner transformation. Assuming only one spatial coordinate, L_c is introduced as follows.

$$V_w(k_x, x) = \frac{1}{2\pi i\hbar} \int_{-\infty}^{\infty} \left[V\left(x + \frac{s}{2}\right) - V\left(x - \frac{s}{2}\right) \right] e^{-ik_x s} ds \quad (132)$$

$$\approx -\frac{1}{\pi i\hbar} \int_0^{L_c/2} \left[V\left(x + \frac{s}{2}\right) - V\left(x - \frac{s}{2}\right) \right] \sin(k_x s) ds \quad (133)$$

The cutoff length has to be selected carefully when solving the Wigner equation numerically. The comparison of current–voltage characteristics shown in Fig. 9 demonstrates that only a sufficiently large value for L_c gives a realistic result. A too small value results in an overestimation of the valley current.

5.2. The Effect of Scattering

In RTD2, the potential changes linearly in a region of 40 nm length, starting 10 nm before the emitter barrier and extending 19 nm after the collector barrier, as shown in Fig. 10. The Wigner potential is discretized using $N_k = 640$ equidistant k_x points and $\Delta x = 0.5$ nm spacing in x -direction. Assuming a cutoff length of $L_c = 80$ nm, one would require at least $N_k = L_c/\Delta x = 160$. This minimum value is often used in finite-difference simulations for the Wigner equation, but in the Monte Carlo simulation we use the considerably larger value stated above in order to get a better resolution of the energy domain. The annihilation mesh is three-dimensional. In x -direction, the grid covers the region where the Wigner potential is nonzero. Because of the cylindrical symmetry of the Wigner function, only two momentum coordinates have to be considered. The mesh extends to an energy of 6 eV in both axial and radial k -direction.

Table 1. Parameter values of the simulated resonant tunneling diodes.^a

Device name	Barrier height (eV)	Barrier width (nm)	Well width (nm)	Device length (nm)	Contact doping (cm ⁻³)
RTD1	0.27	2.83 ($5a_0$)	4.52 ($8a_0$)	100.6 ($178a_0$)	2×10^{18}
RTD2	0.3	3.0	5.0	200.0	10^{16}
RTD3	0.47	3.0	4.0	270.0	10^{18}

^aThe lattice constant of GaAs is $a_0 = 0.565$ nm.

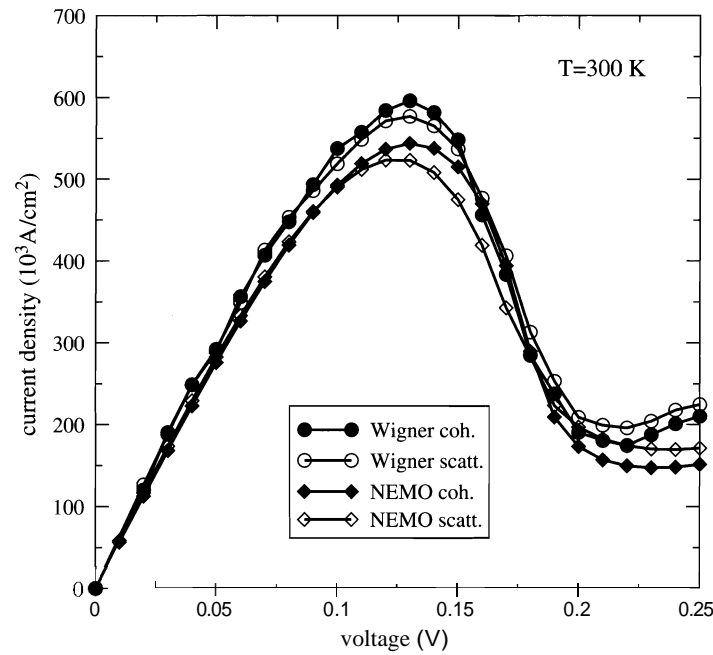


Figure 8. Current-voltage characteristics of RTD1 at 300 K obtained from Wigner Monte Carlo and NEMO-1D. Transport is coherent (coh.) or dissipative (scatt.).

The Wigner generation rate (127) is of the order 10^{15}s^{-1} for RTD2 (Fig. 11). The relation of this rate to the typically much smaller semiclassical scattering rate is a quantitative measure of the fact that quantum interference effects are dominant. The zero-field contact regions have been chosen sufficiently large, such that the Wigner potential drops to zero within these regions.

Figure 12 shows the electron concentration in RTD2 at voltages below the resonance voltage. Classical behavior is observed before and after the double barrier, whereas in the

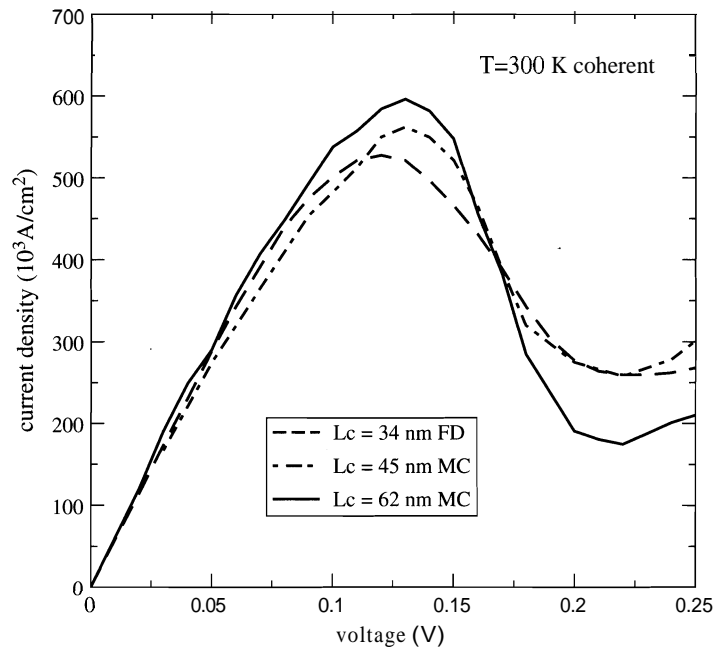


Figure 9. Effect of the cutoff length on the current-voltage characteristics in Wigner simulations. The finite-difference (FD) result is taken from [36].

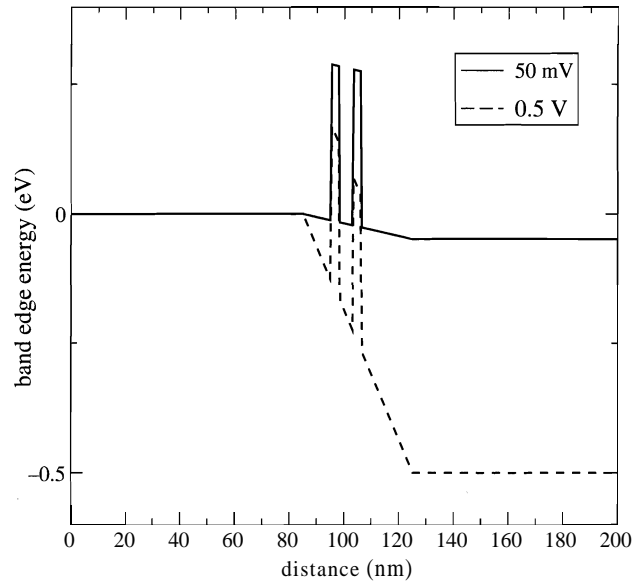


Figure 10. Conduction band edge of RTD2 for different voltages. A linear voltage drop is assumed

quantum well the behavior of the solution is nonclassical. In front of the barrier an accumulation layer forms, with its maximum concentration increasing with the band bending. In the quantum well, the concentration increases as the resonance is approached. After the barrier a depletion layer forms, which grows with applied voltage. In this region, the concentration at 0.15 V varies exponentially in response to the linear potential (see Fig. 10), which is again a classical property.

For voltages above the resonance voltage, the concentration in the well drops, whereas the depletion layer continues to grow (Fig. 13). The mean kinetic energy of the electrons is depicted in Fig. 14. The energy density has been calculated from the second-order moment of the Wigner function (35) and divided by the electron density to get the mean energy per electron. In the zero-field regions, an energy close to the equilibrium energy is obtained, which demonstrates that the energy conservation property of the Wigner potential operator is also satisfied by the numerical Monte Carlo procedure. One has to keep in mind that the Wigner potential can produce a rather large momentum transfer. For the chosen value

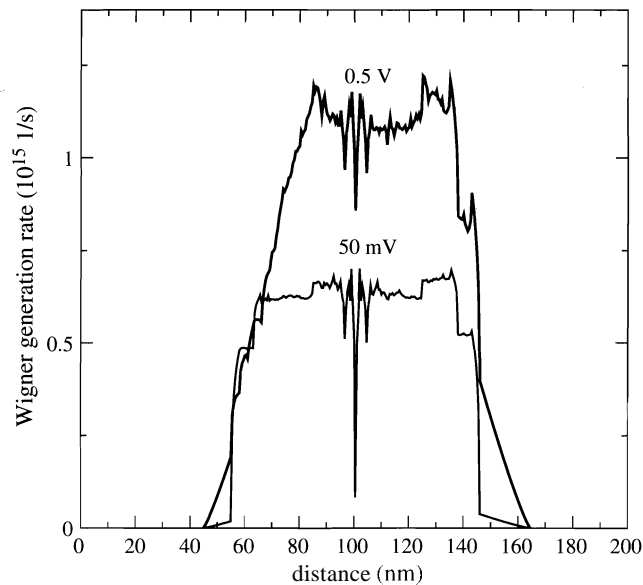


Figure 11. Pair generation rate $\gamma(x)$ in RTD2 caused by the Wigner potential for two different voltages

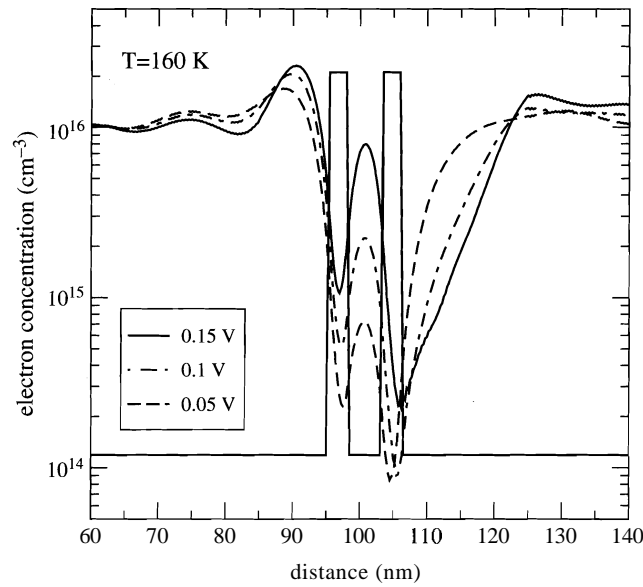


Figure 12. Electron concentration in RTD2 for voltages less than the resonance voltage.

for Ax , the related energy transfer can reach values as large as 5 eV, which shows that a large degree of cancellation occurs in the estimator for the mean energy. Electrons injected from the second barrier into the collector space charge region show initially a high kinetic energy.

Phonon scattering strongly affects the current–voltage characteristic of RTD2 (Fig. 15). As compared to the coherent case, phonon scattering leads to an increase in the valley current and a resonance voltage shift. The large difference in the valley current can be explained by the electron concentration in off-resonance condition (Fig. 16). With phonon scattering included, a significantly higher concentration forms in the emitter notch, and injection in the double barrier is increased. This indicates that a quasi bound state forms in the emitter notch. The population of this state increases when scattering is switched on. On the other hand, in resonance condition where the applied voltage is lower, such a bound state does not form and very similar electron concentrations are observed for the coherent and noncoherent case (Fig. 17).

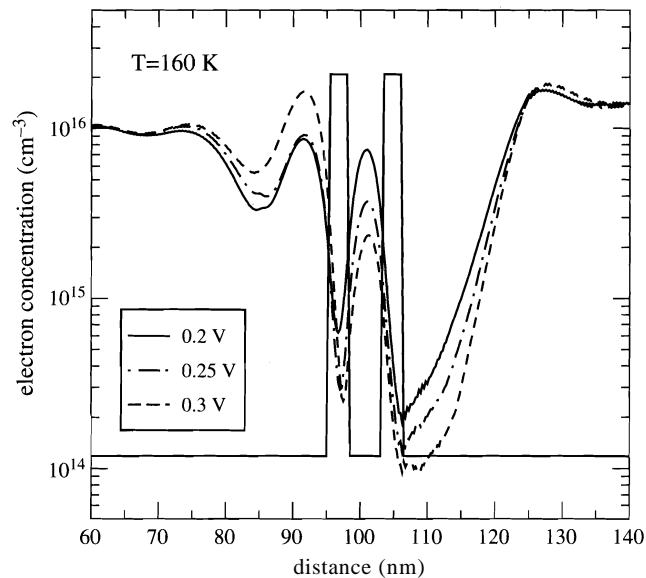


Figure 13. Electron concentration in RTD2 for voltages greater than the resonance voltage.

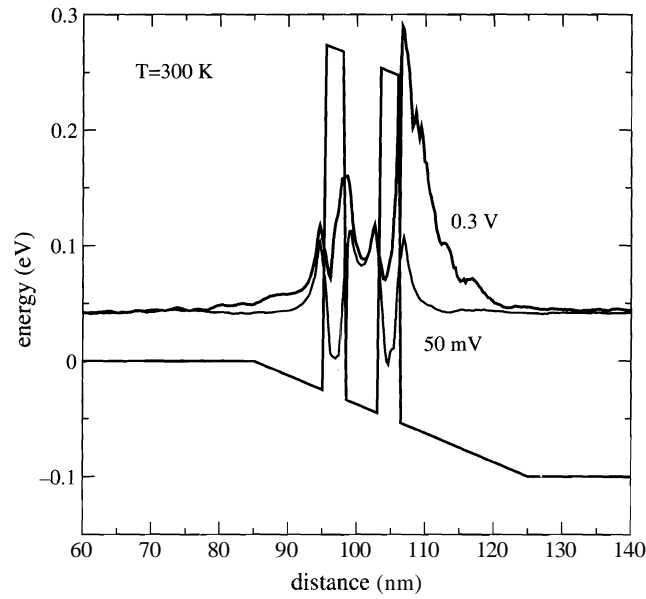


Figure 14. Mean kinetic energy in RTD2 for two different voltages

5.3. Inclusion of Extended Contact Regions

As discussed in Section 3.2.5, the Wigner equation simplifies to the Boltzmann equation when the potential variation is sufficiently smooth. The proposed quantum Monte Carlo method turns into the semiclassical Monte Carlo method for vanishing Wigner potential. Therefore, one can simulate a quantum region embedded in an extended classical region with the interface between the regions correctly treated in an implicit way. By means of the Wigner generation rate γ , the simulation domain can be decomposed into quantum regions ($\gamma > 0$) and classical regions ($\gamma \simeq 0$). In Fig. 18, these regions within RTD2 are marked. The electron concentration and the mean energy are smooth in the extended contact regions and not affected by the strong onset of the Wigner generation rate, as shown in Fig. 11.

In the simulation of RTD3, the Wigner potential $V_w^+(k_x, x)$ is discretized using $N_k = 1200$ equidistant k_x points and $\Delta x = 0.5$ nm spacing in the x -direction. A cutoff length of

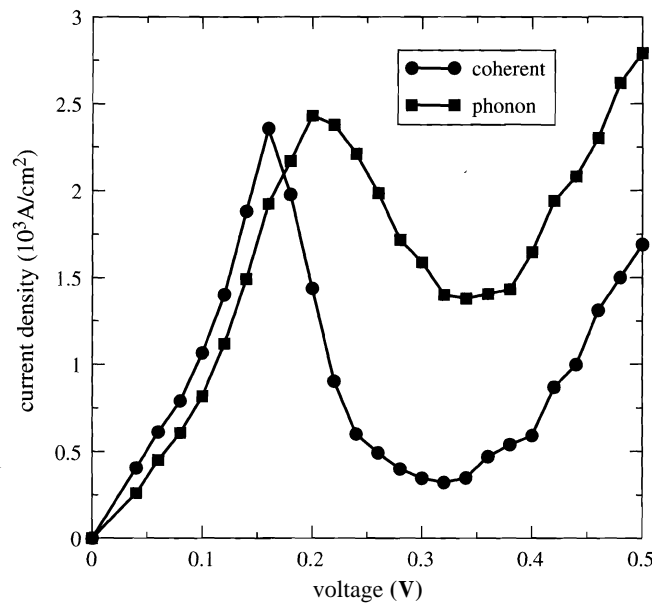


Figure 15. Influence of phonon scattering on the current–voltage characteristics of the RTD2

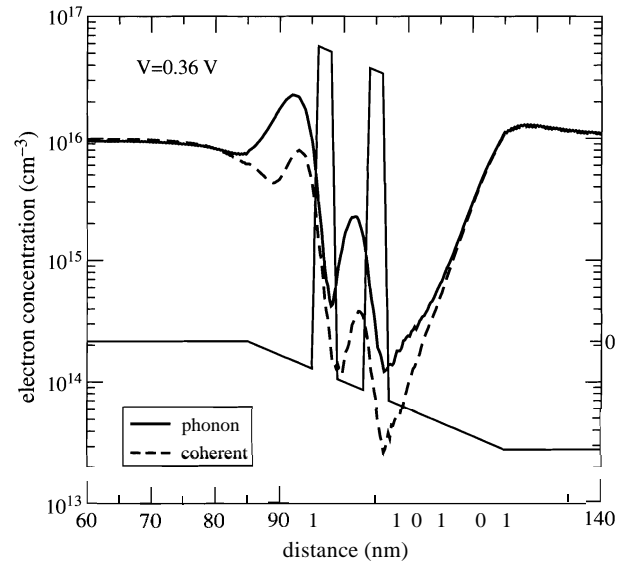


Figure 16. Electron concentration in RTD2 in off-resonance condition.

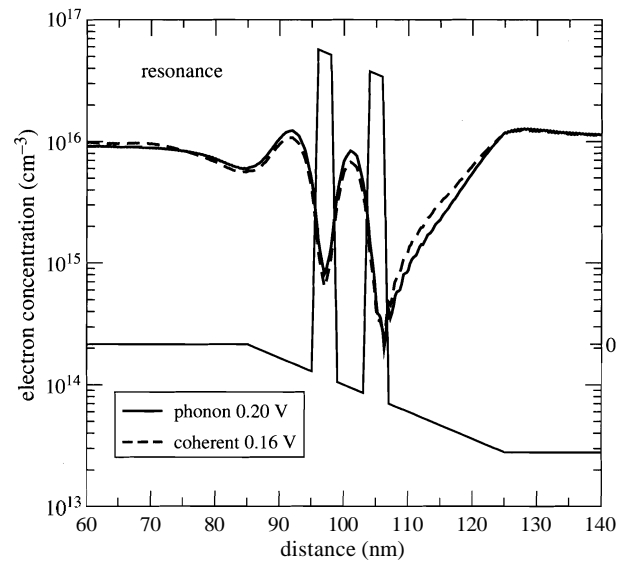


Figure 17. Electron concentration in RTD2 in resonance condition.

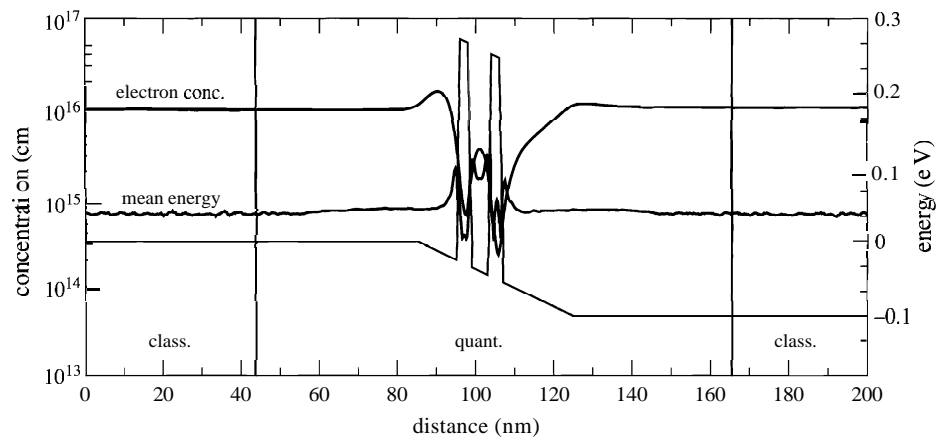


Figure 18. Electron concentration and mean electron energy in RTD2 at $T = 300$ K and 0.1 V applied voltage.

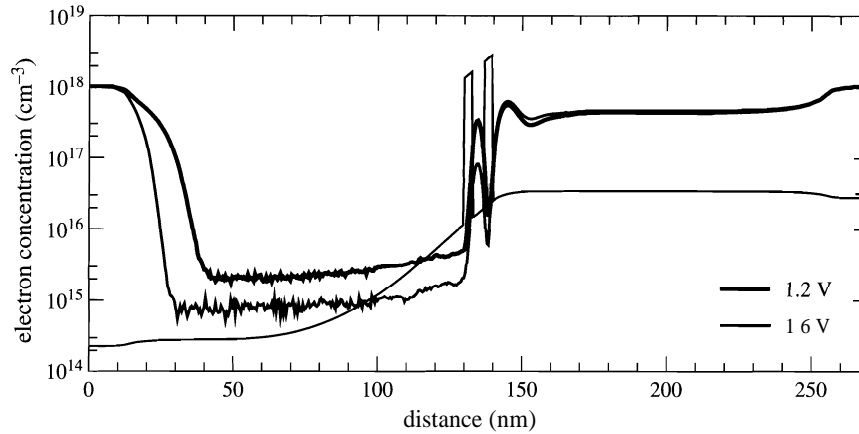


Figure 19. Electron concentration profiles in RTD3.

$L_c = 60$ nm is assumed. The annihilation mesh consists of 480 points in the longitudinal and 120 points in the perpendicular momentum direction, and the real space coordinate is discretized using $\Delta x = 0.5$ nm. The electrostatic potential has been computed using the self-consistent Schrodinger-Poisson solver NANOTCAD-1D [84]. Figure 19 shows the electron concentration profile in the device. At the resonance voltage of 1.2 V, the concentration in the quantum well is considerably higher than in the off-resonance condition at 1.6 V. The concentration in the depletion region left of the barrier depends on the injected current and is thus correlated with the concentration in the well.

6. CONCLUSION

The examples presented in Section 5 demonstrate that a numerical solver for the Wigner equation can provide quantitatively correct results. One requirement is that the cutoff length is chosen sufficiently large. The completeness relation of the discrete Fourier transform reflecting Heisenberg's uncertainty principle, $\Delta k = \pi/L_c$, shows that a small L_c will result in a coarse grid in momentum space, and resonance peaks in the transmission coefficient might not be resolved properly. In the past, the Wigner equation has been solved most frequently by finite-difference methods. Due to the nonlocality of the potential operator, all points in momentum space are coupled, resulting in a very poor sparsity pattern of the matrix. Therefore, increasing the number of grid points in k-space, related to the cutoff length by $N_k = L_c/\Delta x$, is limited by prohibitive memory and computation time requirements. This might be one reason why quantitatively correct solutions were difficult to obtain in the past. We believe that the frequently reported accuracy problems with finite-difference Wigner function–based device simulations result from a too coarse k-space discretization. As this problem occurs already for one-dimensional geometries, higher dimensional simulations using the finite-difference method are probably out of reach. It is interesting to note that Frensky, who pioneered the finite-difference method for the Wigner equation [16], later abandoned this method and developed the quantum-transmitting boundary method to describe coherent transport in open systems [85].

The Monte Carlo method allows the number of k-points to be increased. In this work, the Wigner potential has been discretized using N_k of the order 10^3 . However, high-performance resonant tunneling diodes with very high peak-to-valley current ratio pose still a problem for the Monte Carlo method. In such a device, the density can vary over several orders of magnitude, which often cannot be resolved by the Monte Carlo method. This problem is also well-known from the classical Monte Carlo method. As a solution, one could apply statistical enhancement techniques in such cases. At present, an equidistant k-grid is used for the discretization of the Wigner potential. Because the transmission coefficient of double-barrier structures may show very narrow resonance peaks, using an equidistant k-grid may not be the optimal choice. However, because of the discrete Fourier transform of the potential

involved in the computation of the Wigner potential, the use of a nonequidistant k -grid appears to be problematic.

In a Wigner function–based simulation of one-dimensional heterostructures, fundamental simulation parameters such as the cutoff length are closely linked to physical device parameters such as the spacing from the contacts. This property stems from the choice of plane-wave basis sets in a quantum mechanical regime of broken translational symmetry. Although analytically appealing, this basis set can cause numerical difficulties. Other approaches such as the nonequilibrium Green's function formalism may have the advantage that other basis sets can be used more straightforwardly.

These considerations indicate that from a numerical point of view, the Wigner function formalism might not be the optimal choice for resonant tunneling simulation. However, because the formalism describes quantum effects and scattering effects with equal accuracy, it appears well suited especially when a quasi-ballistic transport condition without energetically sharp resonances is present. One strength of the Wigner function approach is the treatment of contact regions. Nonequilibrium transport can be simulated in the whole device formed by a central quantum region embedded in extended classical regions. The presented Wigner Monte Carlo method can bridge the gap between classical device simulation and pure quantum ballistic simulations.

Development of Monte Carlo methods for the solution of the Wigner equation is still in the beginning. Research efforts are needed especially with respect to the negative sign problem. The particle generation–annihilation algorithm developed by the authors is just one solution to that problem. Improved variants of this algorithm or even new solution strategies are yet to be devised. Extension of the Monte Carlo methods to higher dimensional device geometries is straightforward.

REFERENCES

1. E. Wigner, *Phys. Rev.* 40, 749 (1932).
2. G. Mahan, "Many-Particle Physics." Plenum Press, New York, 1983.
3. V. Tatarskii, *Soviet Physics Uspekhi* 26, 311 (1983).
4. P. Carruthers and F. Zachariasen, *Rev. Mod. Phys.* 55, 245 (1983).
5. T. Cutright and C. Zachos, *Mod. Phys. Lett. A* 16, 2381 (2001).
6. C. Zachos, *Int. J. Mod. Phys. A* 17, 297 (2002).
7. J. Rammer, *Rev. Mod. Phys.* 63, 781 (1991).
8. U. Ravaioli, M. Osman, W. Pötz, N. Kluksdahl, and D. Ferry, *Physica B* 134, 36 (1985).
9. N. Kluksdahl, W. Pötz, U. Ravaioli, and D. Ferry, *Superlattices Microstruct.* 3, 41 (1987).
10. W. Frensley, *Phys. Rev. Lett.* 57, 2853 (1986).
11. W. Frensley, in "International Electron Devices Meeting." p. 571. Institute of Electrical and Electronics Engineers, Los Angeles, CA, 1986.
12. W. Frensley, *Phys. Rev. B* 36, 1570 (1987).
13. N. Kluksdahl, A. Krivan, D. Ferry, and C. Ringhofer, *Phys. Rev. B* 39, 7720 (1989).
14. W. Frensley, *Solid-State Electronics* 32, 1235 (1989).
15. R. Mains and G. Haddad, *J. Appl. Phys.* 64, 5041 (1988).
16. W. Frensley, *Rev. Mod. Phys.* 62, 745 (1990).
17. F. Buot and K. Jensen, *Phys. Rev. B* 42, 9429 (1990).
18. K. Jensen and F. Buot, *Phys. Rev. Lett.* 66, 1078 (1991).
19. F. Buot and K. Jensen, *COMPEL* 10, 241 (1991).
20. K. Gullapalli, D. Miller, and D. Neikirk, *Phys. Rev. B* 49, 2622 (1994).
21. B. Biegel and J. Plummer, *IEEE Trans. Electron Devices* 44, 733 (1997).
22. D. Woolard, P. Zhao, and H. Cui, *Physica B* 314, 108 (2002).
23. R. Mains and G. Haddad, *J. Comput. Phys.* 112, 149 (1994).
24. L. Shifren and D. Ferry, *Physica B* 314, 72 (2002).
25. M. Nedjalkov, R. Kosik, H. Kosina, and S. Selberherr, in "Simulation of Semiconductor Processes and Devices," p. 187. Business Center for Academic Societies Japan, Kobe, Japan, 2002.
26. L. Shifren, C. Ringhofer, and D. Ferry, *Phys. Lett. A* 306, 332 (2003).
27. H. Kosina, M. Nedjalkov, and S. Selberherr, in "Nanotech" (M. Laudon and B. Romanowicz, Eds.), p. 190. Computational Publications, San Francisco, CA, 2003.
28. F. Rossi, C. Jacoboni, and M. Nedjalkov, *Semicond. Sci. Technol.* 9, 934 (1994).
29. M. Nedjalkov, I. Dimov, F. Rossi, and C. Jacoboni, *J. Math. Comp. Modelling* 23, 159 (1996).
30. M. Pascoli, P. Bordone, R. Brunetti, and C. Jacoboni, *Phys. Rev. B* 58, 3503 (1998).
31. P. Bordone, A. Bertoni, R. Brunetti, and C. Jacoboni, *Math. Comp. Simulation* 62, 307 (2003).
32. H. Kosina, M. Nedjalkov, and S. Selberherr, *J. Comput. Electronics* 2, 147 (2003).

33. A. Bertoni, P. Bordone, R. Brunetti, and C. Jacoboni, *J. Phys. Condens. Matter* 11, 5999 (1999).
34. P. Bordone, M. Pascoli, R. Brunetti, A. Bertoni, and C. Jacoboni, *Phys. Rev. B* 59, 3060 (1999).
35. C. Jacoboni, R. Brunetti, P. Bordone, and A. Bertoni, *Int. J. High Speed Electronics Systems* 11, 387 (2001).
36. H. Tsuchiya, M. Ogawa, and T. Miyoshi, *IEEE Trans Electron Devices* 38, 1246 (1991).
37. J.-J. Shih, H.-C. Huang, and G. Wu, *Phys. Rev. B* 50, 2399 (1994).
38. F. Bufler and J. Schlosser, *J. Phys. Condens. Matter* 6, 7445 (1994).
39. D. Miller and D. Neikirk, *Appl. Phys. Lett.* 58, 2803 (1991).
40. L. Demaio, L. Barletti, A. Bertoni, P. Bordone, and C. Jacoboni, *Physica B* 314, 104 (2002).
41. M. B. Unlu, B. Rosen, H.-L. Cui, and P. Zhao, *Phys. Lett. A* 327, 230 (2004).
42. G. Wu and K.-P. Wu, *J. Appl. Phys.* 71, 1259 (1992).
43. I. Levinson, *Soviet Phys. JETP* 30, 362 (1970).
44. P. Holland and K. Kypriandis, *Phys. Rev. A* 33, 4380 (1986).
45. S. Sonogo, *Phys. Rev. A* 44, 5369 (1991).
46. D. Ferry and S. Goodnick, "Transport in Nanostructures." Cambridge University Press, Cambridge, UK, 2001.
47. M. Levanda and V. Fleurov, *Ann. Phys.* 292, 199 (2001).
48. W. Hansch, "The Drift-Diffusion Equation and Its Applications in MOSFET Modeling, Computational Microelectronics." Springer Verlag, Vienna, 1991.
49. R. Brunetti, A. Bertoni, P. Bordone, and C. Jacoboni, in "International Workshop on Computational Electronics" (J. R. Barker and J. R. Watling, Eds.), p. 131. University of Glasgow, Glasgow, Scotland, 2000.
50. J. Zhou and D. Ferry, *IEEE Trans. Electron Devices* 39, 473 (1992).
51. C. Gardner, *SIAM J. Appl. Math.* 54, 409 (1994).
52. C. L. Gardner and C. Ringhofer, *Phys. Rev. E* 53, 157 (1996).
53. P. Degond and C. Ringhofer, *J. Statistical Phys.* 112, 587 (2003).
54. H. Tsuchiya and U. Ravaioli, *J. Appl. Phys.* 89, 4023 (2001).
55. Z. Han, N. Goldsman, and C. Lin, in "International Electron Devices Meeting," p. 62. Institute of Electrical and Electronics Engineers, San Francisco, CA, 2000.
56. N. Goldsman, C.-K. Lin, Z. Han, and C.-K. Huang, *Superlattices Microstruct.* 27, 159 (2000).
57. A. Fannjiang, S. Jin, and G. Papanicolaou, *SUM J. Appl. Math.* 63, 1328 (2002).
58. R. Kosik, Ph.D. Thesis, Vienna University of Technology, 2004.
59. J. J. Wlodarz, *Phys. Lett. A* 264, 18 (1999).
60. B. Biegel and J. Plummer, *Phys. Rev. B* 54, 8070 (1996).
61. R. Brunetti, C. Jacoboni, and F. Rossi, *Phys. Rev. B* 39, 10781 (1989).
62. M. Nedjalkov, R. Kosik, H. Kosina, and S. Selberherr, *Microelectron. Eng.* 63, 199 (2002).
63. M. Nedjalkov, H. Kosina, R. Kosik, and S. Selberherr, *J. Computat. Electronics* 1, 27 (2002).
64. J. Barker and D. Ferry, *Phys. Rev. Lett.* 42, 1779 (1979).
65. M. Nedjalkov, H. Kosina, S. Selberherr, and I. Dimov, *VLSI Design* 13, 405 (2001).
66. T. Gurov, M. Nedjalkov, P. Whitlock, H. Kosina, and S. Selberherr, *Physica B* 314, 301 (2002).
67. C. Ringhofer, M. Nedjalkov, H. Kosina, and S. Selberherr, *SUM J. Appl. Math.* (2004) (in press).
68. C. Jacoboni and L. Reggiani, *Rev. Mod. Phys.* 55, 645 (1983).
69. H. Kosina, M. Nedjalkov, and S. Selberherr, *IEEE Trans. Electron Devices* 47, 1898 (2000).
70. H. Kosina and M. Nedjalkov, *Int. J. High Speed Electronics Systems* 13, 727 (2003).
71. H. Kosina, M. Nedjalkov, and S. Selberherr, *J. Appl. Phys.* 93, 3553 (2003).
72. M. Nedjalkov, H. Kosina, and S. Selberherr, *J. Appl. Phys.* 93, 3564 (2003).
73. I. Sobol, "The Monte Carlo Method." Mir Publishers, Moscow, 1984.
74. F. Byron and R. Fuller, "Mathematics of Classical and Quantum Physics." Dover, New York, 1992.
75. S. Ermakow, "Die Monte-Carlo-Methode und verwandte Fragen." R. Oldenburg Verlag, Munich, 1975.
76. J. Hammersley and D. Handscomb, "Monte Carlo Methods." John Wiley, New York, 1964.
77. J. Sun, G. Haddad, P. Mazumder, and J. Schulman, *Proc. IEEE* 86, 641 (1998).
78. H. Mizuta and T. Tanoue, "The Physics and Applications of Resonant Tunneling Diodes." Cambridge University Press, Cambridge, UK, 1995.
79. H. Kosina, G. Klimeck, M. Nedjalkov, and S. Selberherr, in "Simulation of Semiconductor Processes and Devices," p. 171. Institute of Electrical and Electronics Engineers, Boston, MA, 2003.
80. R. Lake, G. Klimeck, R. Bowen, and D. Jovanovic, *J. Appl. Phys.* 81, 7845 (1997).
81. R. Bowen, G. Klimeck, R. Lake, W. Frensley, and T. Moise, *J. Appl. Phys.* 81, 3207 (1997).
82. G. Klimeck, R. Lake, D. Blanks, C. Fernando, R. C. Bowen, T. Moise, and Y. Kao, *Phys. Status Solidi (b)* 204, 408 (1997).
83. G. Klimeck, R. Lake, and D. Blanks, *Phys. Rev. B* 58, 7279 (1998).
84. G. Iannaccone, M. Macucci, P. Coli, G. Curatola, G. Fiori, M. Gattobigio, and M. Pala, in "IEEE Conference on Nanotechnology," p. 117. IEEE, Maui, Hawaii, 2001.
85. W. Frensley, *Superlattices Microstruct.* 11, 347 (1992).