# The Mathematical Study and Approximation of Semiconductor Models

Joseph W. Jerome*

July, 1992

### Abstract

Emphasis in this paper is upon the development of the mathematical properties of the drift-diffusion semiconductor device model, especially steady-state properties, where the system is incorporated into a fixed point mapping framework. An essential feature of the paper is the description of the Krasnosel'skii calculus, the appropriate extension to nonlinear equations and systems of the famous inf-sup saddle point approximation theory for linear equations. The introductory section describes certain scaling effects and the association with numerical methods constructed via exponential fitting. The fixed point and numerical fixed maps are introduced in the next section, while the finite element approximation theory is described in the third section. The fourth section deals with the evolution system and modular pre-algorithms defined by Newton's method, combined with a fully implicit time semidiscretization. More general moment models are introduced in the fifth section, and some historical comments appear in the final section. Specific algorithms and implementation issues are not discussed. The principal section headings are as follows.
1. The Drift-Diffusion Model, Scaling, and Exponential Fitting.
2. Steady-State: The Fixed Point and Numerical Fixed Point Maps.
3. A Nonlinear Finite Element Convergence Theory.
4. The Evolution System and Newton's Method.
5. More General Moment Models: A Review.
6. Epilogue: Historical Perspective.

---

*Department of Mathematics, Northwestern University, Evanston, Il 60208

# 1 The Drift-Diffusion Model, Scaling, and Exponential Fitting

## 1.1 Overview

One may view classical or semiclassical semiconductor modeling as a pyramid structure, in which the Boltzmann Transport Equation (BTE) forms the summit and the drift-diffusion system the base; intermediate are systems derived from taking moments of the Boltzmann equation. The hydrodynamic model and its momentum relaxation time limit, called the energy model, are examples. These will be discussed more fully in the final technical section of the paper. The hydrodynamic model has the drawback of requiring adequate closure assumptions and accurate representations for the average collision mechanisms. In addition, the model has hyperbolic as well as parabolic components for the moment equations, in the transient case, and thus issues of shock capturing mechanisms arise. On the other hand, the drift-diffusion model has only parabolic components, exclusive of the potential equation, which must be adjoined to any model in the hierarchy.

A derivation of the drift-diffusion model from the Boltzmann equation, making use of the diffusion approximation, is carried out in [15]. This model is by far the best understood of the above models, and it is the primary topic of this paper. It was introduced by Van Roosbroeck in 1950 [50] as a conservation system for electron and hole carriers, with the ambient electric field determined from the potential equation. This model is closely related to the older ionic transport model of Nernst and Planck, which is described in detail in [51]. We shall not have occasion to draw upon the latter model in this work, however.

When appropriate scalings are selected for the quantities appearing, it is seen that the conservation part of the system is highly convective. There is a scaling, viz. , that used to bring the order of the electron and hole concentrations to unity, in which the phenomenon of convection domination is augmented by a singularly perturbed potential equation. The combination of these two facts has made this system quite challenging to solve numerically. Major breakthroughs in the solution of the model occurred during the 1960s, particularly in the one-dimensional steady state case. A major factor was the research of H. K. Gummel, who, in two significant papers [26], [53], the second with D. L. Scharfetter on the transient model, introduced two decisive ideas of considerable impact:

1. System Decoupling.

2. Exponential Fitting for the Continuity Equations.

The system decoupling map is a compact continuous map whose fixed points define solutions of the drift-diffusion model. The study of this map is an important goal of this paper. The Scharfetter-Gummel discretization has been employed widely in the numerical simulations of this model. Two-dimensional versions depend upon area element methods and subsequent reductions to one-dimensional discretizations along boundaries. We describe the one-dimensional version below in terms of generalized spline functions. In this case, the method proceeds by approximating the flux (current) by a piecewise constant function, where the breaks correspond to the nodal points. The concentration values at the nodal points are determined from a well defined linear system, the definition of which depends upon the numerical method selected. The solution of the linear system is equivalent to determining the constants in an exponential

fitting method, arising from the exact solution of the locally constant flux equations. The representation is in terms of $L-$splines.

The first mathematical treatment of the drift-diffusion model is due to Mock [47], who studied the steady state case. Mock made essential use of the decoupling, as originally introduced by Gummel, to define a solution fixed point map. During the 1970s and early 1980s, Newton outer iteration procedures began to complement successive approximation based upon iteration of the Gummel map. Major computational breakthroughs occurred when effective inner iterations were designed in tandem with the outer iterations. Noteworthy was the paper [8], in which damped Newton outer iteration, a form of continuation, was combined with Gauss-Seidel and other forms of inner iteration, thus replacing sparse direct linear algebra methods with iterative methods. This paper also presents a two-dimensional Scharfetter-Gummel discretization. The difficult implementation problem in three dimensions appears to be the geometric partitioning problem, carried out so that the linear systems for the potential equation are of M-matrix type.

Two books have appeared in the mathematical literature on the drift-diffusion model: [48] and [45]. The former also discusses the transient problem. Another will appear shortly: [39]. Three recent conference proceedings related to device and/or circuit modeling are currently available: [7], [6], and [29].

## 1.2   Drift-Diffusion System and Scalings

We shall introduce the drift-diffusion system in this section, as well as a simplified version, the latter for the purpose of explaining, through the choice of well known scalings, the highly convective character of the continuity equations, and the singular perturbation aspect of the potential equation. This includes the determination of the junction width, separating oppositely doped regions, in an accompanying example as proportional to the square root of the singular perturbation parameter.

The drift-diffusion model may be obtained by taking zeroth order moments of the BTE and adjoining the Poisson equation. Thus, one obtains the system for $N$ carriers with electric field $E$, concentration $n_i$, recombination $R_i$, current density $J_i$, (signed) charge $e_i$, $i = 1, \cdots, N$:

$$\frac{e_i \partial n_i}{\partial t} + \nabla \cdot J_i = -e_i R_i, \tag{1.1}$$

$$E = -\nabla \phi, \tag{1.2}$$

$$\nabla \cdot (\epsilon \nabla \phi) = -\sum e_i n_i - k_1. \tag{1.3}$$

Here the dielectric constant is denoted by $\epsilon$, $k_1$ is the net impurity concentration, and $\phi$ is the electrostatic potential. There still remains the issue of determining the constitutive current relations. Classical drift-diffusion theory gives, for $N = 2$, $n_1 = n$, and $n_2 = p$,

$$J_n = -e\mu_n n \nabla \phi + eD_n \nabla n, \tag{1.4}$$

$$J_p = -e\mu_p p \nabla \phi - eD_p \nabla p. \tag{1.5}$$

The electronic charge modulus $e$ is positive here, and $n$ and $p$ denote the electron and hole densities, respectively. The use of the Einstein relations linking the mobilities, $\mu_n, \mu_p$, and the diffusion coefficients, $D_n, D_p$, is common. These relations are specified by

$$D_n = (kT/e)\mu_n, \tag{1.6}$$

$$D_p = (kT/e)\mu_p, \tag{1.7}$$

where $k$ is Boltzmann's constant and $T$ is the temperature. The mobilities are field dependent in realistic modeling situations (cf. (4.4) and (4.5)). It is also possible to derive the constitutive relations (1.4), (1.5), from the first order moment relations developed in §5, under the assumption that the momentum relaxation times tend to zero. The details are given in [52]. In fact, the constitutive relations include a heat flux term as well, which is suppressed at constant temperature. If it is not suppressed, one has what is called the energy model. The introduction of exponential relations for $n$ and $p$ is also common, and will appear in the subsequent sections.

The simplifications employed in the remainder of this subsection are as follows:

1. Carrier recombination is suppressed.

2. The dielectric is assumed constant over the device.

3. The mobility and diffusion coefficients are assumed constant.

We are now ready to introduce a scaling for the system. We shall refer to this as the unit scaling. A convenient parameter is given by the expression $U_T \equiv kT/e$, called the thermal voltage; its value at $T = 300K$ is:
$$U_T = .0259V.$$

We now present the unit scaling.

- Potentials are scaled by $U_T$.

- All concentrations are scaled by $S = \|k_1\|_{max}$.

- Length is scaled by a dimension parameter $l$.

The new system can be written, in the case of constant mobilities and under the assumption that the Einstein relations are valid:

$$-\lambda^2 \nabla^2 u + n - p = k_1, \qquad (1.8)$$
$$-\nabla \cdot (\nabla n - n\nabla u) = 0, \qquad (1.9)$$
$$-\nabla \cdot (\nabla p + p\nabla u) = 0. \qquad (1.10)$$

Here,
$$\lambda^2 = \epsilon U_T/(l^2 eS),$$

where $S$ is the appropriate concentration scale. In the case of the unit scaling, so-called because the doping has maximal value 1 in these units, the potential equation (1.8) can be singularly perturbed, and all concentrations are expected to be of order comparable to unity. Moreover, $\lambda^2 \approx 10^{-1} - 10^{-7}$ in typical situations.

One final remark about scalings is in order. It is not necessary to use any scaling at all, even in micron and submicron devices. What is essential is the proper choice of units in this case. One may view this procedure as self-induced scale, via unit selection (see [20]).

We shall present a very simple example to indicate the boundary layer possible near a p/n junction. This has the simultaneous effect of demonstrating that the continuity equations can be highly convective. The calculation here is an adaptation of that in [58, pp.142-147]. Consider then a one-dimensional p/n metallurgical junction at $x = 0$. The junction is situated so that

the p-region occupies a bounded part of $x < 0$, and the n-region occupies a bounded part of $x > 0$. The junction is assumed to be a so-called step junction:

$$k_1 = \begin{cases} -N_a, & x < 0, \\ N_d, & x > 0. \end{cases}$$

The calculations are based upon the following assumptions:

1. The device is in equilibrium, i. e. , the electron and hole current densities are zero.

2. In particular, the boundary conditions are determined as follows. There are no external voltage biases, and the electron and hole concentrations are specified in the usual way by thermal equilibrium and charge balance, respectively:

$$np = (n_i/S)^2 \quad := \quad c^2, \tag{1.11}$$
$$n - p - k_1 \quad = \quad 0. \tag{1.12}$$

3. There is a carrier depletion region surrounding the junction, written as $[-x_{p0}, x_{n0}]$, in which moment conditions expressing the "smallness" of $n - p$ hold:

$$\int_{-x_{p0}}^{x_{n0}} (n(x) - p(x)) \, dx \approx 0, \tag{1.13}$$

$$\int_{-x_{p0}}^{x_{n0}} x(n(x) - p(x)) \, dx \approx 0. \tag{1.14}$$

4. Adjacent to the depletion region, there is a neutral region, characterized by a zero electric field.

It is immediate from (1.13) that the charge on the p side of the depletion region balances that on the n side. This observation is used in (1.20) below. The boundary conditions, (1.11), (1.12), give the familiar relations for $n$ and $p$:

$$n \quad = \quad (1/2)(k_1 + \sqrt{k_1^2 + 4c^2}), \tag{1.15}$$
$$p \quad = \quad (1/2)(-k_1 + \sqrt{k_1^2 + 4c^2}). \tag{1.16}$$

These relations are valid at the endpoints of the device, but also, because of the neutral region and equilibrium assumptions, in the entirety of the neutral region. This makes use of the solution formulas, which also determine the boundary values of $u$,

$$n = c \exp(u), \quad p = c \exp(-u). \tag{1.17}$$

We define the contact potential by $u_n - u_p := u_0$. Here, $u_n$ and $u_p$ are the values of the potential at $x_{n0}$ and $-x_{p0}$, respectively. Use of (1.15), (1.16), and (1.17), together with a simple approximation based on the size of $N_a$ and $N_d$, yields the expression

$$u_0 = \log N_a N_d + 2 \log(S/n_i). \tag{1.18}$$

Integration of the equation, $E = -\nabla u$, across the depletion region gives the relation, for $W = x_{p0} + x_{n0}$,

$$u_0 = (1/(2\lambda^2)) N_d x_{n0} W. \tag{1.19}$$

This formula, obtained by integration by parts applied to the integrated version of (1.8), makes use of (1.14). The unknown quantities are $x_{n_0}$ and $W$. The two equations specifying these quantities are (1.19) and the charge balance relation:

$$N_a(W - x_{n_0}) = N_d x_{n_0}. \tag{1.20}$$

Note that $u_0$ is specified by (1.18). Solution of these relations gives, finally, the formula for the width of the depletion region:

$$W = \lambda\sqrt{2u_0(N_a + N_d)/(N_a N_d)}. \tag{1.21}$$

## 1.3   The Scharfetter-Gummel Discretization: Exponential Fitting

In this section, we illustrate the class of discretizations of convection-diffusion equations which is known by the name of Scharfetter-Gummel discretizations. As shown below, the class is described in one dimension by the property:

- The associated flux of the numerical approximation is piecewise constant.

Furthermore, the procedures for selecting the nodal values of the approximation differentiate methods within the class.

We shall restrict attention to single equations of the form

$$Ln = -\nabla\cdot J \;=\; F \text{ in } G, \tag{1.22}$$
$$n \;=\; \bar{n} \text{ on } \partial G, \tag{1.23}$$

where $J$ is specified by

$$J = \nabla n - n\nabla u. \tag{1.24}$$

Here, $u$, $F$, and $\bar{n}$ are given functions.

We consider the one-dimensional version of (1.22), (1.23), and (1.24), and require that the scheme be exact for a piecewise constant flux $J$ whose discontinuities coincide with selected grid points; note that our interpretation of solution will be that of a weak solution throughout. For simplicity of exposition only we shall assume that $\bar{n} \equiv 0$. If the exactness requirement is interpreted in terms of approximation theory, we are seeking an approximation of the form

$$n_h = \sum_i \alpha_i M_i, \tag{1.25}$$

where $\{\alpha_i\}$ is a set of nodal values determined by a specified numerical method, and where $\{M_i\}$ is a nodal basis of local support functions defined by $L$. More precisely, given a grid of the interval $G = [0,1]$, of the form $x_i = ih$, $i = 0, \cdots N$, for $Nh = 1$, $M_i$ is associated with the $ith$ grid point, $i \neq 0$ and $i \neq N$, and is specified by the following requirements:

1. $M_i(x_i) = 1$,  support $M_i = [x_{i-1}, x_{i+1}]$.

2. $M_i$ is continuous.

3. On each subinterval determined by the grid, $M_i$ is in the null space of $L$.

We shall refer to discretizations of the form (1.25) as being of the class of Scharfetter-Gummel type. We shall comment later on the critical issue of nodal value determination. The functions $M_i$ are generalizations of the chapeau functions, and are examples of the generalized B-splines, introduced in [30] and studied extensively in [41]. The approximation given by (1.25) is an example of a generalized spline, introduced in [25]. It is quite easy to give explicit formulas for the B-spline functions, $M_i$. For example, when $i = 1$,

$$M_1(x) = \begin{cases} \exp[u(x) - u(x_1)] \int_{x_0}^{x} e^{-u(s)} \, ds / \int_{x_0}^{x_1} e^{-u(s)} \, ds, & x_0 \le x \le x_1, \\ \exp[u(x) - u(x_1)] \int_{x}^{x_2} e^{-u(s)} \, ds / \int_{x_1}^{x_2} e^{-u(s)} \, ds, & x_1 \le x \le x_2. \end{cases}$$

The general formula is obtained via the identifications $0 \to i - 1$, $1 \to i$, and $2 \to i + 1$. In all cases, the definition is completed by the support requirement in the second part of item one above.

It is also quite straightforward to compute the piecewise constant flux $J$. Note that, on the intervals $(x_0, x_1)$ and $(x_{N-1}, x_N)$, the flux is determined by $M_1$ and $M_{N-1}$, respectively, while on all other subintervals, $(x_i, x_{i+1})$, the flux has components from both $M_i$ and $M_{i+1}$. The total flux can then be assembled from the following result. Denoting by $J_{M_i}^-$ the flux due to $M_i$ on $(x_{i-1}, x_i)$, and by $J_{M_i}^+$ the flux due to $M_i$ on $(x_i, x_{i+1})$, we have

$$J_{Mi}^- = e^{-u(x_i)} / \int_{x_{i-1}}^{x_i} e^{-u(s)} \, ds,$$

$$J_{M_i}^+ = -e^{-u(x_i)} / \int_{x_i}^{x_{i+1}} e^{-u(s)} \, ds.$$

In order to evaluate the integrals appearing in the flux representations, it has been common to employ the piecewise linear interpolant of $u$. When this is done, direct flux evaluation gives the following representation for the assembled flux on $(x_i, x_{i+1})$, with $\alpha_i = n_i$:

$$J = (1/h)[B(\triangle u)n_{i+1} - B(-\triangle u)n_i], \tag{1.26}$$

where we have adopted the conventions $\triangle u = u(x_{i+1}) - u(x_i)$ and $B(z) = z/[\exp(z) - 1]$. The latter function, known as the Bernoulli function, must be computed for small values of $\mid z \mid$ by series methods involving the Bernoulli numbers. As noted in [45], an exponential fitting method of this type resolves the currents adequately, even if the mesh allows for substantial variation in the function $u$.

Although the form given in (1.26) is well-known to workers in computational electronics, it is not widely understood that the form holds, irrespective of the numerical method used to characterize the nodal values. The original method of Scharfetter and Gummel was to define these values by the box method, a method tailored for divergence equations. In this case, the matrix entries are given by

$$d_i = B(\triangle u_i) + B(-\triangle u_{i+1}),$$
$$c_i = -B(\triangle u_{i+1}),$$
$$a_i = -B(-\triangle u_{i+1}),$$

where the diagonal subscripts for $d_i$ range from 1 to $N-1$ and the superdiagonal and subdiagonal subscripts for $c_i$ and $a_i$, respectively, range from 1 to $N - 2$.

# 2 Steady-State: The Fixed Point and Numerical Fixed Point Maps

We shall begin with a brief description of the quasi-Fermi level representation of the drift-diffusion system, and use this to introduce the fixed point map in a special case. Subsequent subsections will develop this as well as the discrete analogues.

## 2.1 Alternative System Representation

As defined earlier, the drift-diffusion model is presented in terms of the conduction electron density $n$, the hole density $p$, and the electrostatic potential $u$ as:

$$-\nabla\cdot[\epsilon\nabla u] + n - p - k_1 = 0, \tag{2.27}$$
$$\nabla\cdot[\mu_n n\nabla u - D_n\nabla n] = -R, \tag{2.28}$$
$$-\nabla\cdot[\mu_p p\nabla u + D_p\nabla p] = -R. \tag{2.29}$$

Here, we have employed a common recombination term, $R$, and units in which $kT = 1$ and $e = 1$. As before, the dielectric constant satisfies $\epsilon(x) \geq \epsilon_0 > 0$, and is a function of the material and therefore of the position $x$ only. The system is augmented by boundary conditions of mixed type, including Dirichlet conditions on the contact portions of the device boundary, and homogeneous Neumann boundary conditions on the complement. In terms of the quasi-Fermi levels $v$ and $w$, by which we express $n$ and $p$ via the scaling for which $n = \exp(u - v)$ and $p = \exp(w - u)$, and under the assumption of Einstein's relations (cf. (1.6) and (1.7)), the system (2.27), (2.28), and (2.29) is rewritten as

$$-\nabla\cdot[\epsilon\nabla u] + e^{u-v} - e^{w-u} - k_1 = 0, \tag{2.30}$$
$$-\nabla\cdot[\mu_n e^{u-v}\nabla v] = Q(u,v,w)[e^{w-v} - 1], \tag{2.31}$$
$$-\nabla\cdot[\mu_p e^{w-u}\nabla w] = -Q(u,v,w)[e^{w-v} - 1]. \tag{2.32}$$

Also included are the boundary conditions, specified by functions $\bar{u}$, $\bar{v}$, and $\bar{w}$. Here we make use, for the first time, of the expression

$$R = Q(u,v,w)[e^{w-v} - 1], \tag{2.33}$$

for the generation-recombination term. $Q$ is an appropriate positive rational function. In this case we define the mapping $\mathbf{U}_{qf} : (v,w) \mapsto u$ by solution of (2.30) for $u$ if $v$ and $w$ are given. For zero recombination, the mapping $\mathbf{V}_{qf} : u \mapsto v$ is defined through solution of (2.31) for given electrostatic potential $u$. The mapping $\mathbf{W}_{qf} : u \mapsto w$ is defined similarly. Finally, for (2.30), (2.31), (2.32), and zero recombination, we define $\mathbf{T}_{qf}$ through $\mathbf{T}_{qf} = (\mathbf{V}_{qf}, \mathbf{W}_{qf}) \circ \mathbf{U}_{qf}$.

The ideas and techniques of the following subsections are described in [33] for the map $\mathbf{T}_{qf}$ and in [43] for the map $\mathbf{T}_h$. Some technical difficulties are associated with the transition points on the boundary of $G$, where a transition point is defined to be common to the (relative) closures of the sets $\Sigma_D$ and $\Sigma_N$, corresponding to Dirichlet and Neumann boundary conditions. According to regularity theory, this is a point where the electric field may become unbounded. A second feature is the possible development of a family of fixed point maps, and the study of their properties. Such a study was carried out in [39]. Existence is a by-product of this type of study, which attempts to validate various decoupling approaches to the system. This has led

to allowing coupled current continuity subsystems in the definition of the fixed point map, and associating with the current continuity subsystem a well-defined map $\mathbf{VW}_{qf}$. We shall not take this general approach, but rather shall define a single decoupling map, still denoted $\mathbf{VW}_{qf}$.

## 2.2   The Map $\mathbf{U}_{qf}$

We begin by discussing the definition of the mapping $\mathbf{U}_{qf}$. Specifically, we restrict the domain to be the closed, convex set $K$ in $L_2(G) \times L_2(G)$, where

$$K = \{[v, w] : \alpha \leq v \leq \beta, \alpha \leq w \leq \beta\}, \tag{2.34}$$

with inequalities taken almost everywhere, and where

$$\alpha = \min(\inf_{\Sigma_D} \bar{v}, \ \inf_{\Sigma_D} \bar{w}), \tag{2.35}$$

$$\beta = \max(\sup_{\Sigma_D} \bar{v}, \ \sup_{\Sigma_D} \bar{w}). \tag{2.36}$$

The significance of this set is its relation to maximum principles satisfied by the quasi-Fermi levels $v$ and $w$. We are now prepared to define $\mathbf{U}_{qf}$ in the context of the following lemma.

**Lemma 2.1** *Given a pair $[v, w]$ in $K$, the image under the map $\mathbf{U}_{qf}$ is the unique element $u$ satisfying the weak relation*

$$\langle \epsilon \nabla u, \nabla \phi \rangle \ + \ \langle \exp(u - v) - \exp(w - u) - k_1, \phi \rangle, \tag{2.37}$$

*subject to the boundary condition,*
$$\Gamma(u - \bar{u}) \mid_{\Sigma_D} = 0. \tag{2.38}$$
*The test function space in this relation, denoted $Y_0 := H^1_{0,\Sigma_D}$, consists of $H^1$ functions $\phi$ with zero trace on $\Sigma_D$. Moreover, $u$ satisfies the maximum principle,*

$$u \geq \gamma \ = \ \min(\inf_{\Sigma_D} \bar{u}, \ \gamma'), \tag{2.39}$$

$$u \leq \delta \ = \ \max(\sup_{\Sigma_D} \bar{u}, \ \delta'), \tag{2.40}$$

*where $\gamma'$ and $\delta'$ are uniquely defined by*

$$2\sinh(\gamma' - \alpha) - \inf_G k_1 \ = \ 0, \tag{2.41}$$

$$2\sinh(\delta' - \beta) - \sup_G k_1 \ = \ 0. \tag{2.42}$$

*Proof* We first demonstrate the validity of the bounds (2.39), (2.40). Thus, if $u$ is a solution in $H^1$, with $\exp u$ in $L_2$, we select for the admissible test function $\phi$,

$$\phi = (u - \delta)^+.$$

The restriction of the trace of $\phi$ to $\Sigma_D$ is zero, hence $\phi$ is admissible. Now the substitution of $\phi$ into (2.37) reduces integrations to the set $\{u > \delta\}$. For the term involving the gradients, this

9

uses the chain rule for the composition of the positive part and an $H^1$ function ([24, Lemma 7.6]). Once this reduction is achieved, one uses the nonnegativity of

$$\exp(u - v) - \exp(w - u) - k_1,$$

on the set $\{u > \delta\}$, to conclude $\nabla(u - \delta)^+ = 0$, and hence $(u - \delta)^+ = 0$. The nonegativity above follows from the inequalities

$$\exp(u - v) - \exp(w - u) - k_1 \geq 2\sinh(u - \beta) - \sup_G k_1$$

$$\geq 2\sinh(\delta - \beta) - \sup_G k_1 \geq 2\sinh(\delta' - \beta) - \sup_G k_1 = 0.$$

Note that we have used the definitions of $\beta$, $\delta$, and $\delta'$ above. In a similar way, one shows that $(u - \gamma)^- = \phi$ is the zero function, where $t^- = -(-t)^+$. The use of $\phi$ as a test function, coupled with the inequality,

$$\exp(u - v) - \exp(w - u) - k_1 \leq 0, \quad \text{on } \{u < \gamma\},$$

leads to this result. If we combine these two observations, we obtain the inequalities (2.39), (2.40). We now proceed to the equation itself. It is convenient to consider the problem of convex minimization, later shown to be equivalent,

$$\Phi(u) = \min_{L_2(G)} \Phi(f). \tag{2.43}$$

Here $\Phi$ is the convex functional defined by,

$$\Phi(f) =$$

$$\frac{1}{2} \int_G \epsilon \mid \nabla f \mid^2 \, dx + \int_G F(f(x), x) \, dx - \int_G k_1(x) f(x) \, dx, \tag{2.44}$$

if $f \in \{f : F(f, \cdot) \in L_2(G)\} \cap \{f : (f - \bar{u}) \in Y_0\}$, and is defined to have the value $\infty$, for $f$ not in this convex set. Here, $F$ is an appropriate primitive:

$$F(s, x) =$$

$$\exp(-v(x)) \int_0^s \exp(\sigma) \, d\sigma + \exp(w(x)) \int_0^s - \exp(-\tau) \, d\tau. \tag{2.45}$$

Note that $F$ is convex in $s$ for each fixed $x$. In the terminology of [19], $\Phi$ is a proper convex functional on $L_2(G)$, hence has a finite minimum $u$ under the following two sufficient conditions:

- $\Phi$ is coercive, i. e. ,
$$\Phi(f)/\|f\|_{L_2} \to \infty \text{ as } \|f\|_{L_2} \to \infty;$$

- $\Phi$ is lower semicontinuous.

The verification of these two conditions is carried out in [33] and in [39].

It remains to show that the minimization principle (2.43) implies (2.37). The technique is standard in the sense that one selects $f$ of the form

$$f = u + \eta\phi, \tag{2.46}$$

where $\eta$ is an arbitrary real number, and $\phi$ is constrained to be in $Y_0 \cap L_\infty$. Usual techniques established in the calculus of variations yield (2.37) for $\phi$ so constrained. A density argument employing smooth functions gives the equation for all $\phi$ in $Y_0$. Uniqueness makes use of monotonicity of F in its first argument and standard properties of norms. *Box*

10

## 2.3 The Current Continuity Subsystem: The Map $\mathbf{VW}_{qf}$

It would appear natural to consider the subsystem (2.47) and (2.48) to follow, and attempt to define a joint system map from $u$ to an image $[v, w]$, in terms of a weak solution of the system, and thus determine a map $\mathbf{VW}_{qf}$. It does not appear, however, that such a correspondence need be unique in general, much less continuous. If we wished to pursue the construction of a family of fixed point maps, associated with the solution of the semiconductor system (2.30), (2.31), (2.32), it would be necessary to introduce the concept of admissible lagging of terms in the current continuity subsystem, primarily in the recombination expression, as has been carried out in [39]. In any approach, in the case of nonzero recombination, it is essential to consider dependence upon $[v, w] \in K$, as well as upon $u$. The approach we employ here is the simplest possible lagging, involving system decoupling as described in the next paragraph.

In this section, we assume that $\tilde{v}$, $\tilde{w}$, and $u = \mathbf{U}_{qf}[\tilde{v}, \tilde{w}]$ are prescribed, as well as lagging in the recombination term of each equation of the current continuity subsystem, so that new functions, $R_v$ and $R_w$, respectively, are obtained. This occurs through the substitution of $\tilde{w}$ for $w$ in (2.47) and $\tilde{v}$ for $v$ in (2.48), respectively. Here, we assume, for the form of $R$, that described in (2.33), where $Q \geq 0$ is assumed $C^\infty$ on the closed set $-\infty < u < \infty$, $\alpha \leq v \leq \beta$, $\alpha \leq w \leq \beta$. $R$ is also decreasing in $v$ and increasing in $w$.

Thus, for purposes of eventually defining a map $\mathbf{VW}_{qf} = [\mathbf{V}_{qf}, \mathbf{W}_{qf}]$, we are studying the decoupled subsystem,

$$-\nabla \cdot J_n - R_v \;=\; 0, \tag{2.47}$$
$$-\nabla \cdot J_p + R_w \;=\; 0. \tag{2.48}$$

In order to establish the existence of weak solutions of this subsystem, where boundary conditions are imposed as previously, we can use use arguments based upon the minimization of convex functionals, as in the previous subsection. A critical component of the proof is the derivation of maximum principles: the pair of quasi-Fermi levels solving the subsystem (2.47) and (2.48) is in $K$, so that the inequalities contained in (2.34) are satisfied. Moreover, this occurs precisely because of the exponential factor in the recombination term, (2.33).

We are now prepared to state the result which forms the basis for the definition of the map $\mathbf{VW}_{qf}$. Consider the weak form of (2.47) and (2.48):

$$\langle \mu_n \exp(u - V)\nabla V, \nabla \psi \rangle - \langle R_v(u, V, W), \psi \rangle = 0, \tag{2.49}$$

$$\langle \mu_p \exp(W - u)\nabla W, \nabla \omega \rangle + \langle R_w(u, V, W), \omega \rangle = 0, \tag{2.50}$$

for all $[\psi, \omega]$ in $Y_0$, where the latter test function pair in $\prod_1^2 H^1$ has zero trace on $\Sigma_D$. We have the following result, stated without proof.

**Lemma 2.2** *The map $\mathbf{VW}_{qf}$ is well defined. In particular, for $[\tilde{v}, \tilde{w}] \in K$ and $u = \mathbf{U}_{qf}[\tilde{v}, \tilde{w}]$, there is a unique pair $[V, W] = \mathbf{VW}_{qf}[u, \tilde{v}, \tilde{w}]$ satisfying the system (2.49), (2.50), subject to the stated boundary conditions. Moreover, $[V, W] \in K$.*

## 2.4 Compactness and Continuity of $\mathbf{VW}_{qf}$: Fixed Points of $\mathbf{T}_{qf}$

Compactness estimates are cited in Lemma 2.3 to follow. The proof is omitted for brevity.

**Lemma 2.3** *Given* $[\tilde{v}, \tilde{w}]$ *in* $K$ *and its image* $u$ *under* $\mathbf{U}_{qf}$, *there is an 'a priori' bound* $\lambda$ *on the* $\prod_1^2 H^1(G)$ *norm of the image point of the map* $\mathbf{VW}_{qf}$, *which is independent of* $[\tilde{v}, \tilde{w}]$. *In particular, the range of the map* $\mathbf{VW}_{qf}$ *is relatively compact in* $\prod_1^2 L_2(G)$.

The continuity required of $\mathbf{VW}_{qf}$, for the existence of fixed points of the Gummel map as presented at the end of this subsection, is product $L_2$ continuity in the dependence upon $\tilde{v}$ and $\tilde{w}$. In fact, a significantly stronger continuity result is valid when the range is normed by the energy space norm. Since the weaker result is straightforward, we shall present it and its proof here.

**Lemma 2.4** *The mapping* $\mathbf{VW}_{qf}$ *is continuous, in its* $L_2(G)$ *dependence upon* $\tilde{v}$ *and* $\tilde{w}$, *and in its* $H^1(G)$ *dependence upon* $u$, *as a mapping into* $\prod_1^2 L_2(G)$.

*Proof* Because of the compact embedding described previously, strong sequential continuity of $\mathbf{VW}_{qf}$, as a mapping into $\prod_1^2 L_2$, is implied by weak sequential convergence in $\prod_1^2 H^1$. This is equivalent to showing that the corresponding weak solution systems, as formulated in (2.49) and (2.50), converge, as systems, to the particular system taken at the function pair at which continuity is verified. Moreover, uniqueness allows this weak convergence to be demonstrated subsequentially. The equivalence of weak convergence with the system convergence, as just stated, is a simple consequence of the following fact for the case $q = q' = 2$.

- For $q'$ satisfying $\frac{1}{q} + \frac{1}{q'} = 1$, we have convergence in $L_{q'}$ given by

$$\lim_{k \to \infty} R_v(u, v_k, w_k) = R_v(u, v, w), \tag{2.51}$$

  with a similar relation for $R_w$.

The relation (2.51) is used a second time in demonstrating the convergence (or subsequential convergence) of the systems themselves. This leaves only the convergence of the leading (gradient) parts of (2.49) and (2.50) to be demonstrated. Without loss of generality, we may assume smooth test functions. Given a sequence $\{[\tilde{v}_j, \tilde{w}_j]\} \subset K$, convergent in the product $L_2$ metric, let $u$ denote the image of its limit under $\mathbf{U}_{qf}$. By the usual criterion, we need only show that every subsequence of the image sequence $\{[V_j, W_j]\}$, under the map $\mathbf{VW}_{qf}$, has a further subsequence convergent to an invariant limit. Given such a subsequence, select a further subsequence, conveniently denoted $\{[V_{j_k}, W_{j_k}]\}$, which is weakly convergent in $\prod_1^2 H^1$. Using the well known fact that the corresponding weak convergence is not altered if the integrands are multiplied by a boundedly convergent sequence, and taking a further subsequence if necessary, we conclude that the integral terms,

$$\langle \mu_n(\nabla u) \exp(u_{j_k} - V_{j_k}) \nabla V_{j_k}, \nabla \psi \rangle$$

and

$$\langle \mu_p(\nabla u) \exp(W_{j_k} - u_{j_k}) \nabla W_{j_k}, \nabla \omega \rangle,$$

are convergent to the expected limit. A routine application of the triangle inequality, making use of the term just described, as well as a term estimated by the bounded pointwise convergence of $\mu_n(\nabla u_{j_k})$ to $\mu_n(\nabla u)$ (and similarly for $\mu_p$), yields the convergence of the leading parts as required. In this connection, the local convergence pointwise of $\nabla u_{j_k}$ is converted to global

pointwise convergence by a diagonalization procedure. This concludes the proof of continuity. *Box*

In §2.2, we defined the map $\mathbf{U}_{qf}$, and in §2.3 we defined the map $\mathbf{VW}_{qf}$. By the Gummel map, $\mathbf{T}_{qf}$, we mean the composition, $\mathbf{VW}_{qf} \circ \mathbf{U}_{qf}$, of the maps just cited. According to the theory of this section, $\mathbf{T}_{qf}$ is a mapping of the closed, convex set $K$ into itself. Moreover, $\mathbf{T}_{qf}$ is continuous as a mapping of $\prod_1^2 L_2(G)$ into itself, and has relatively compact range. It follows that $\mathbf{T}_{qf}$ has at least one fixed point in $K$ as a consequence of the Schauder fixed point theorem. Note that we have extended here the usual meaning of composition. We have the following result.

**Theorem 2.1** *There exists a weak solution of the drift-diffusion system. Such a solution may be identified with a fixed point of the map $\mathbf{T}_{qf}$ defined above.*

## 2.5  The Discretized Model and Maximum Principles

In this section we assume that $G$ is a polyhedral domain and introduce the piecewise linear finite element method for the system. We also describe the associated maximum principles. The finite element equations for the potential equation are given by

$$\langle \epsilon(x)\nabla U_h, \nabla \phi_i \rangle + \langle e^{U_h - v} - e^{w - U_h}, \phi_i \rangle - \langle k_1, \phi_i \rangle = 0 \quad \text{for} \ \ i = 1, \cdots, M, \qquad (2.52)$$

where $U_h$ is a finite element function and the $\phi_i$ are appropriate test functions comprising a nodal basis of the piecewise linear finite element subspace $S_h$. We select the piecewise linear interpolant $\bar{u}_I$ of $\bar{u}$ so that $U_h \in \bar{u}_I + S_h$, where the members of $S_h$ vanish on the Dirichlet boundary $\Sigma_D$ of the polyhedral domain $G$. The functions of $S_h$ are continuous and are linear in each simplex, $S$. As usual, $h = \max_S\{\text{diam } S\}$. The fact that $U_h$ satisfies the bounds (2.39) and (2.40) is verified in the paper [43]. Certain mesh restrictions are required, since the proof requires that the matrices corresponding to the Laplacean and the current continuity equations be M-matrices. In order to state these discrete maximum principles in a format applicable both to $\mathbf{U}_h$, and also $\mathbf{V}_h$ and $\mathbf{W}_h$ to be defined subsequently, we consider solutions of the gradient equation, where $a$ is strictly positive:

$$-\nabla \cdot [a(x)\nabla u(x)] + f(x, u(x)) = g(x). \qquad (2.53)$$

Here, $a, g \in L_\infty$, and $f$ is increasing and locally Lipschitz in $u$ for each $x \in G$, with $f^{-1}(x, \cdot)$ the corresponding inverse. On each element $S$ we have the following definition.

DEFINITION. Let $S$ be an $N$-dimensional simplicial finite element such that

- $V$ is the volume;

- $\vec{v}_i$ is a vertex;

- $e_{ij}$ is the edge connecting vertices $\vec{v}_i$ and $\vec{v}_j$;

- $F_k$ is the face opposite the vertex $k$, with measure $|F_k|$;

- $h_i$ is the normal distance of $v_i$ to $F_i$;

- $\gamma_{ij}$ is the angle between the inward normal vectors to the faces $F_i$ and $F_j$;

13

- $\phi_l$ is the piecewise linear nodal basis function which is 1 at vertex $\vec{v}_l$;

- $$\alpha_{ij} \equiv \int_S a(x) \nabla \phi_i \cdot \nabla \phi_j dx$$

  is the $ij$th entry of the *element* stiffness matrix;

- $\langle a(x) \rangle \equiv \int_S a(x) dx / V$, the average of $a(x)$ over the element $S$;

- $a_{ij}$ is the $ij$th element of the assembled stiffness matrix.

*Remark 1.* It was shown in [43] that

$$\alpha_{ij} \equiv \int_S a(x) \nabla \phi_i \cdot \nabla \phi_j dx = \langle a(x) \rangle \cos(\gamma_{ij}) \frac{1}{h_i h_j} V,$$

or

$$\alpha_{ij} = \langle a(x) \rangle \cos(\gamma_{ij}) \frac{|F_i||F_j|}{N^2 V}.$$

In [43] it was also shown that $L_\infty$ stability of $\mathbf{U}_h$, in terms of satisfying the maximum principle, is a consequence of the following assumption.

*Assumption.*

- In $N$ dimensions, where $N \geq 2$, we require that for every edge $jk$ the off-diagonal element $a_{jk}$ in the stiffness matrix satisfies

$$a_{jk} = \sum_{S \text{ adjacent } jk} \langle a(x) \rangle_S \cos(\gamma_{jk}^{(S)}) \frac{V^{(S)}}{h_j h_k} \leq -\frac{\rho}{h_{\max}^2} \sum_{S \text{ adjacent } jk} V^{(S)},$$

  with $\rho > 0$. In two dimensions, the well-known requirement that for every edge $jk$ in the triangulation we have

$$\tfrac{1}{2}[\langle a(x) \rangle_{T_1} \cot(\omega_1) + \langle a(x) \rangle_{T_2} \cot(\omega_2)] \geq \rho > 0,$$

  where the $T_i$ are the two triangles adjacent to edge $jk$, and the $\omega_i$ are the two angles opposite to the edge $jk$, is a slightly more restrictive version of this condition. In higher dimensions, we can impose the sufficient condition that the angle between the vectors normal to any two faces of the same polyhedron in the mesh has to be bounded uniformly from above by $\pi/2 - \eta$.

- For all $c, d \in \mathbf{R}, c < d : |f(x,u) - f(x,v)|/|u - v| \leq D(d,c)$ if $c \leq u, v \leq d$, where $D(\cdot, \cdot)$ is a Lipschitz constant which is a monotonically increasing function of $d$ and a monotonically decreasing function of $c$.

- The numbers $h_i$ satisfy $h_i \geq h_0 h$, where $h_0$ does not depend on $h$.

14

The finite element maps may be defined by the following system.

$$\langle \mu_n \exp(U_h - V_h)\nabla V_h, \nabla \phi_i \rangle - \langle R_v(U_h, V_h, \tilde{w}), \phi_i \rangle = 0, \text{for } i = 1, \cdots, M, \tag{2.54}$$

$$\langle \mu_p \exp(W_h - U_h)\nabla W_h, \nabla \phi_i \rangle + \langle R_w(U_h, \tilde{v}, W_h), \phi_i \rangle = 0, \text{for } i = 1, \cdots, M. \tag{2.55}$$

*Remark 2.* In the application to this paper, the role of $a$ is played for the Poisson equation by $\epsilon(x)$, and for the current continuity equations by the mobility-exponential products, after interpretation of the nonlinear approximations as linear approximations. The first part of the assumption above is satisfied uniformly in all cases if $\langle a \rangle$ is replaced by the lower bounds expressed in terms of the maximum principles. The function $f(x, u)$ in the respective applications is $e^{u-\tilde{v}(x)} - e^{\tilde{w}(x)-u}$, for Poisson's equation, and, for the current continuity equations, is expressed via $-R_v$ and $R_w$, respectively.

The following lemma summarizes the key stability properties and is stated without proof.

**Lemma 2.5** *The range of the mapping $\mathbf{VW}_h = [\mathbf{V}_h, \mathbf{W}_h]$ is now defined, via the auxiliary mapping $\mathbf{U}_h$, as the image pair $[V_h, W_h]$, satisfying the system (2.54) and (2.55), given the domain point $[\tilde{v}, \tilde{w}]$ in $K$. Moreover, $[V_h, W_h] \in K$ under the above assumption. Finally, by the map, $\mathbf{T}_h$, we mean the composition, $\mathbf{VW}_h \circ \mathbf{U}_h$, and $K$ is invariant under this map.*

# 3    A Nonlinear Finite Element Convergence Theory

A very powerful extension of the approximation theory for positive definite self-adjoint operators, in the context of Galerkin approximation, was introduced by Babuška and Aziz in [1]. Continuous bilinear forms, satisfying an inf-sup condition and some auxiliary conditions, were identified as properly defining an invertible operator framework allowing for analysis and approximation theory. Without exploiting it directly, these authors created an underlying fixed point map, via the continuous inverse map. For the reader's benefit, we summarize here the essential features. We shall not present the most general formulation, for simplicity. In [1], it is desired to solve the operator equation, $\mathbf{L}u = f$, approximately. One assumes in this theory the solvability of the direct and adjoint problems. If $B$ denotes the bilinear form of the weak formulation on a Hilbert space $E$, assume:

- continuity:

$$| B(v, w) | \leq C_1 \|v\| \, \|w\|. \tag{3.1}$$

- sup condition:

$$\text{For } w \neq 0, \ \sup_v | B(v, w) | > 0. \tag{3.2}$$

- inf-sup condition:

$$\inf_{\|v\|=1} \sup_{\|w\|\leq 1} | B(v, w) | \geq C_2 > 0. \tag{3.3}$$

Assume also sup and inf-sup conditions on an approximation space, $E_n$:

$$\text{For } \psi \neq 0, \ \sup_\phi | B(\phi, \psi) | > 0, \tag{3.4}$$

$$\inf_{\|\phi\|=1} \sup_{\|\psi\|\leq 1} \mid B(\phi, \psi) \mid \geq c_2 > 0. \tag{3.5}$$

One concludes the finite element approximation, $u_h$, is defined and is within

$$\delta_*\{1 + (C_1/c_2)\} \tag{3.6}$$

of $u$, where

$$\delta_* := \|u - u_*\|, \tag{3.7}$$

and $u_*$ is arbitrary in $E_n$. We note briefly the role of the hypotheses.

1. (3.1) $\Rightarrow$ $\mathbf{L}$ may be identified with a continuous linear map $\mathbf{R}$ on $E$.

2. (3.3) $\Rightarrow$ $\mathbf{R}^{-1}$ exists on a closed domain of $E$.

3. (3.2) $\Rightarrow$ Domain and range of $\mathbf{R}$ are all of $E$.

There is a fixed point formulation if 1 is not in the spectrum:

$$\mathbf{T}u = u, \quad \mathbf{T}v := (\mathbf{I} - \mathbf{R})^{-1}(v - \mathbf{J}f), \quad \mathbf{R} = \mathbf{JL}. \tag{3.8}$$

Here, $\mathbf{J}$ denotes the Riesz map. $\mathbf{T}$ is affine. Its derivative is defined for all $v$ by

$$\mathbf{T}'(v) \equiv (\mathbf{I} - \mathbf{JL})^{-1}. \tag{3.9}$$

Note that 1 is an eigenvalue of $\mathbf{T}'$ if and only if 0 is an eigenvalue of $\mathbf{L}$. The latter is excluded, so the exclusion of 1 as a spectral value is a solvability hypothesis when $\mathbf{T}$ and $\mathbf{T}'$ are compact.

In the discussion to follow, the reader may visualize the solvability conditions as generalized by the uniform invertibility condition on the derivative of the fixed point map at the fixed point. In the present setting, this invertibility is simply an eigenvalue condition via the compactness properties. We now introduce this theory.

Given a fixed point $x_0$ of a smooth mapping $\mathbf{T}$, a numerical approximation map $\mathbf{T}_n$, and a linear projection map $\mathbf{P}_n$, a theory is constructed to estimate $\|x_n - \mathbf{P}_n x_0\|$, where $\mathbf{T}_n x_n = x_n$. In fact, the authors of [44] characterize the map $\mathbf{P}_n\mathbf{T}$ as the "Galerkin" approximate map, and $\mathbf{T}_n$ as a "perturbed Galerkin" map. Since $x_0$ is a fixed point of $\mathbf{T}$, the estimates represent the dispersion between these two methods. The mapping $\mathbf{P}_n\mathbf{T}$, while not actually implemented numerically, has a convergence rate which is readily estimated. Now, the manner in which the 'a priori' estimates are derived is to deduce a zero of the map $\mathbf{I} - \mathbf{T}_n$, in a ball centered at $\mathbf{P}_n x_0$, by constructing an equivalent contraction map: The methodology involves derivative inversion and a mean value calculus. The result is stated as Theorem 3.1 below. A similar approach is employed for the 'a posteriori' estimates. The relevant result is Theorem 3.2. In our application of this theory to the semiconductor model, we shall work with energy norms, augmented by $L_\infty$ norms when appropriate.

## 3.1 The Abstract Calculus

Let $E$ be a Banach space and $\mathbf{T}$ a mapping from an open set $\Omega$ into $E$. We assume the existence of a fixed point $x_0$ for $\mathbf{T}$:

$$\mathbf{T}x_0 = x_0. \tag{3.10}$$

If $\{E_n\}$ denotes a sequence of subspaces of $E$ of dimension $r(n) \geq n$, suppose that $\mathbf{T}_n : \Omega_n \mapsto E_n$, $\Omega_n := \Omega \cap E_n$, has a fixed point:

$$\mathbf{T}_n x_n = x_n. \tag{3.11}$$

Finally, let $\{\mathbf{P}_n\}$ be a family of linear projections onto $E_n$. We shall describe here the framework of the calculus developed by Krasnosel'skii et al. [44] for the convergence of the solutions of discretizations of fixed point equations (3.11) to the solutions of the original fixed point equation (3.10). First, we demonstrate that for sufficiently small meshwidth $h$, a solution $x_n$ to the discretized problem (3.11) exists close to all solutions $x_0$ to the original problem (3.10). Second, we show that for sufficiently small meshwidth $h$, a solution $x_0$ to the original problem (3.10) exists close to all solutions $x_n$ to the discretized problem (3.11).

Our first convergence result follows through Theorem 19.1 in [44], as quoted below.

**Theorem 3.1** *Let the operators $\mathbf{T}$ and $\mathbf{P}_n\mathbf{T}$ be Fréchet-differentiable in $\Omega$, and $\mathbf{T}_n$ Fréchet-differentiable in $\Omega_n$. Assume that (3.10) has a solution $x_0 \in \Omega$ and the linear operator $\mathbf{I} - \mathbf{T}'(x_0)$ is continuously invertible in $E$. Let*

$$\|\mathbf{P}_n(x_0) - x_0\| \to 0,$$

$$\|\mathbf{P}_n\mathbf{T}\mathbf{P}_n x_0 - \mathbf{T}x_0\| \to 0, \quad \|\mathbf{P}_n\mathbf{T}'(\mathbf{P}_n x_0) - \mathbf{T}'(x_0)\| \to 0,$$

$$\|[\mathbf{T}_n - \mathbf{P}_n\mathbf{T}]\mathbf{P}_n x_0\| \to 0, \quad \|[\mathbf{T}_n' - (\mathbf{P}_n\mathbf{T})'](\mathbf{P}_n x_0)\| \to 0,$$

*as $n \to \infty$. Finally, assume that for any $\epsilon > 0$ there exist $n_\epsilon$ and $\delta_\epsilon > 0$ such that*

$$\|\mathbf{T}_n'(x) - \mathbf{T}_n'(\mathbf{P}_n x_0)\| \leq \epsilon \quad for \quad (n \geq n_\epsilon; \ \|x - \mathbf{P}_n x_0\| \leq \delta_\epsilon, \ x \in \Omega_n). \tag{3.12}$$

*Then there exist $n_0$ and $\delta_0 > 0$ such that, when $n \geq n_0$, equation (3.11) has a unique solution $x_n$ in the ball $\|x - x_0\| \leq \delta_0$. Moreover,*

$$\|x_n - x_0\| \leq \|[\mathbf{I} - \mathbf{P}_n]x_0\| + \|x_n - \mathbf{P}_n x_0\| \to 0 \quad as \quad n \to \infty, \tag{3.13}$$

*and $\|x_n - \mathbf{P}_n x_0\|$ satisfies the following two-sided estimate ($c_1, c_2 > 0$):*

$$c_1\|\mathbf{P}_n\mathbf{T}x_0 - \mathbf{T}_n\mathbf{P}_n x_0\| \leq \|x_n - \mathbf{P}_n x_0\| \leq c_2\|\mathbf{P}_n\mathbf{T}x_0 - \mathbf{T}_n\mathbf{P}_n x_0\|. \tag{3.14}$$

*Remark 3.* The result continues to remain valid if (3.12) holds, under the stipulation that the metric distance $\rho(x, P_n x_0)$ is measured by a norm $\|\cdot\|_*$, stronger than $\|\cdot\|$. However, in this case, invariance under the regularization maps, employed in [44] to obtain the fixed points, is required to hold with respect to the stronger norm. Notice that in this theorem the actual rate of convergence depends only on the terms in the two sided estimate (3.14). The additional convergence assumptions need not hold with this same rate.

Next, we assert that for sufficiently fine meshwidth $h$, a solution $x_0$ to (3.10) exists close to the solution $x_n$ to the discretized problem (3.11). From [44], we cite Theorem 19.2 as stated in the following theorem.

**Theorem 3.2** *Let the operators* $\mathbf{T}, \mathbf{P}_n\mathbf{T}$, *and* $\mathbf{T}_n$ *be Fréchet differentiable in some neighborhood of the point* $\tilde{x}_n \in \Omega_n$, *and* $\mathbf{I} - \mathbf{T}'_n(\tilde{x}_n)$ *continuously invertible in* $E_n$,

$$\|[\mathbf{I} - \mathbf{T}'_n(\tilde{x}_n)]^{-1}\| = \kappa_n.$$

*Let*

$$\gamma_n \equiv (1 + \kappa_n\|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|)\|[\mathbf{T}' - \mathbf{P}_n\mathbf{T}'](\tilde{x}_n)\| + \kappa_n\|[\mathbf{T}'_n - \mathbf{P}_n\mathbf{T}'](\tilde{x}_n)\| < 1,$$

*and for some* $\delta_n$ *and* $q_n$ $(\delta_n > 0; 0 < q_n < 1)$,

$$\sup_{\|x - \tilde{x}_n\| \leq \delta_n} \|\mathbf{T}'(x) - \mathbf{T}'(\tilde{x}_n)\| \leq \frac{q_n}{\kappa'_n}, \tag{3.15}$$

$$\|\tilde{x}_n - \mathbf{T}\tilde{x}_n\| \leq \frac{\delta_n(1 - q_n)}{\kappa'_n},$$

*where*

$$\kappa'_n = \frac{1 + \kappa_n\|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|}{1 - \gamma_n}.$$

*Then* (3.10) *has a unique solution* $x_0$ *in the ball* $\|x - \tilde{x}_n\| \leq \delta_n$, *and we have the error estimate*

$$\frac{\alpha_n}{1 + q_n} \leq \|\tilde{x}_n - x_0\| \leq \frac{\alpha_n}{1 - q_n}, \tag{3.16}$$

*where*

$$\alpha_n \equiv \|[\mathbf{I} - \mathbf{T}'(\tilde{x}_n)]^{-1}(\tilde{x}_n - \mathbf{T}\tilde{x}_n)\| \leq \kappa'_n\|\tilde{x}_n - \mathbf{T}\tilde{x}_n\|.$$

*Remark 4.* As in the preceding theorem, this result continues to be valid if (3.15) holds under the stipulation that $\rho(x, \tilde{x}_n)$ is measured by a stronger norm, $\|\cdot\|_*$. The regularization maps must act invariantly with respect to the stronger norm. Again, the actual rate of convergence depends only on the terms in the two sided estimate (3.16), while the additional convergence assumptions need not hold with this same rate. Also, $\tilde{x}_n = x_n$ is not required.

## 3.2   The Application to the Semiconductor Problem: Setting

Throughout the application subsections, we assume that the recombination term $R$ is zero, that the mobility and diffusion coefficients are constant, and that the dielectric constant $\epsilon$ is normalized to unity. The corresponding reduced form of the system will be evident as we proceed. We set $E = \prod_1^2 H^1(G)$ and $E_n = $ linear span $\{\bar{v}_I, S_h\}\otimes$ linear span $\{\bar{w}_I, S_h\}$, with $\mathbf{P}_n$ the orthogonal projection onto $E_n$. We also assume that $G$ is a polyhedral domain.

The map $\mathbf{T}$, required to apply the operator calculus of [44], must be defined on an open set in function space. In this context the suitable space is $\prod_1^2 H^1(G)$. However, a number of the results require 'a priori' bounds on the extrema of the functions $u, v$, and $w$, similar to the maximum principles derived in §2. Thus, $\mathbf{T}$ will be an extension map of $\mathbf{T}_{qf}$.

Because the set $K$ is not open, we modify the definition of $\mathbf{T}$ such that the assumption that the preimage $[v, w]$ lies in $K$ can be removed. To achieve this, we compose a $\mathbf{T}$-like map with a truncation operator $\mathbf{Tr}$, which leaves $[v, w]$ unaffected within $K$ (where the solution lies), but which restricts the range to a set $K_1$ (cf. (3.17)) which is only slightly larger than $K$. By carefully selecting $K_1$ we achieve that the intermediate function $u$ in the definition of $\mathbf{T}$ satisfies

'a priori' $L_\infty$ bounds which are only slightly wider than those for $u$ in §2. However, $u$ satisfies these slightly wider bounds as the range of a map defined for all $[v, w]$ *in an appropriate open subset of* $\prod_1^2 H^1$, and not just on the set $K$. We introduce $\mathbf{h}_i \in C_0^\infty(\mathbf{R})$, $i = 1, 2$, such that *support* $\mathbf{h}_i = [\alpha_i, \beta_i]$, and

$$
\begin{aligned}
h_1(t) &= t, \qquad \inf_{\Sigma_D} \bar{v} \leq t \leq \sup_{\Sigma_D} \bar{v}, \\
h_2(t) &= t, \qquad \inf_{\Sigma_D} \bar{w} \leq t \leq \sup_{\Sigma_D} \bar{w}.
\end{aligned}
$$

Below we will define an open ball $\Omega$, centered at zero in $\prod_1^2 H^1$, on which

$$
\mathbf{Tr}[v, w] := [\mathbf{h}_1(v), \mathbf{h}_2(w)], \qquad [v, w] \in \Omega.
$$

Note that the range of $\mathbf{Tr}$ is contained in $K_1 \subset \prod_1^2 L_\infty$, where

$$
K_1 = \{[v, w] \in \prod_1^2 L_\infty : \alpha_1 \leq v \leq \beta_1, \alpha_2 \leq w \leq \beta_2\}. \tag{3.17}
$$

$\mathbf{Tr}$ is Lipschitz continuously differentiable when restricted to $\Omega \cap \prod_1^2 L_\infty$, but is simply Fréchet differentiable on $\Omega$. This fact is what necessitates the use of a more delicate norm in the development and application of the theory of the preceding subsection, as suggested in the remarks following the theorems. The reader is referred to [39] for complete details. Here, we simply outline the general features of the theory. We consider the extension maps $\mathbf{U}$ of $\mathbf{U}_{qf}$, $\mathbf{V}$ of $\mathbf{V}_{qf}$, and $\mathbf{W}$ of $\mathbf{W}_{qf}$ defined as previously, with elements in the domain of $\mathbf{U}$ now taken from $K_1 \supset K$. In terms of these quantities, $\mathbf{T}$ may be defined by

$$
\mathbf{T} = [\mathbf{V} \circ \mathbf{U} \circ \mathbf{Tr}, \mathbf{W} \circ \mathbf{U} \circ \mathbf{Tr}]. \tag{3.18}
$$

The domain $\Omega$ of the map $\mathbf{T}$ is defined in tandem with the composition maps defining $\mathbf{T}$ in such a way as to ensure that $\Omega$ is invariant under $\mathbf{T}$ and contains a fixed point. $\mathbf{T}_n$ may be defined analogously; for consistency with the previous subsection we may wish to consider $\mathbf{T}_n$ as restricted to $\Omega_n$, but this is unimportant. Note that $\mathbf{T}_n$ may be viewed as an extension of the mapping $\mathbf{T}_h$, introduced in an earlier section. Arguments similar to those of the previous section yield fixed points of $\mathbf{T}_n$, via an application of Brouwer's fixed point theorem applied to $\bar{\Omega} \cap E_n$.

An important approximation property of $\mathbf{P}_n$ on the union of the convex hull, $co\, R_\mathbf{T}$, of the range of $\mathbf{T}$, with $\prod_1^2 H^{1+\theta}(G) \cap H_{0,\Sigma_D}^1(G)$, is

$$
\|\mathbf{P}_n \tau - \tau\|_{\prod H^1} \leq ch^\theta, \qquad \|\tau\|_{\prod H^{1+\theta}} \leq 1. \tag{3.19}
$$

This is a consequence of standard approximation theory. Here, $1 < \theta \leq 1$ depends upon Euclidean dimension $N$, and reflects the transition point boundary condition singularities. Note that, in (3.19), $\tau$ is a member of the set, $co\, R_\mathbf{T} \cup \prod_1^2 (H^{1+\theta} \cap H_{0,\Sigma_D}^1)$. Similar approximation results hold for the approximation of $\mathbf{T}$ by $\mathbf{T}_n$.

## 3.3 Differentiability Properties

We state these results for completeness, without proof. We begin with $\mathbf{U}$.

**Lemma 3.1** *Let* $\mathbf{U} : (v, w) \mapsto u$ *be the mapping defined implicitly through the solution of the boundary value problem,*

$$\langle \nabla u, \nabla \phi \rangle + \langle e^{u-v} - e^{w-u} - k_1, \phi \rangle = 0, \tag{3.20}$$

*where* $\phi \in H^1_{0,\Sigma_D}$, *subject to suitable mixed boundary conditions in* $N$ *dimensions. Then the derivative* $D_{(v,w)}\mathbf{U}(v, w) : (\sigma, \tau) \mapsto \mu$ *is defined through the solution of the boundary value problem,*

$$\langle \nabla \mu, \nabla \phi \rangle + \langle e^{u-v}[\mu - \sigma] + e^{w-u}[\mu - \tau], \phi \rangle = 0, \tag{3.21}$$

*where* $\mu|_{\Sigma_D} \equiv 0$, $\phi|_{\Sigma_D} \equiv 0$, *and for all* $[v, w]$ *is a uniformly bounded linear mapping from* $\prod_1^2 L_2$ *to* $H^1_{0,\Sigma_D}$. *The mapping is also uniformly bounded from* $\prod_1^2 H^1$ *to* $L_\infty$ *if* $N \leq 3$. *Moreover, the mapping* $(v, w) \mapsto D_{(v,w)}\mathbf{U}(v, w)$ *is Lipschitz continuous from* $\prod_1^2 H^1$ *to the mappings from* $\prod_1^2 L_2$ *to* $H^1_{0,\Sigma_D}$ *if* $N \leq 4$.

We continue with the mappings $\mathbf{V}$ and $\mathbf{W}$.

*Remark 5.* The derivative $D\mathbf{V}(u) : \mu \mapsto \sigma$, is defined through solution of the boundary value problem,

$$\langle e^{u-v}[(\mu - \sigma)\nabla v + \nabla \sigma], \nabla \phi \rangle = 0, \tag{3.22}$$

where $\mu \in H^1$, and $\sigma \in H^1_{0,\Sigma_D}$. Here, $\phi$ is a test function in $H^1_{0,\Sigma_D}$. For smooth $v$, the standard existence theory (cf. [24, Chap. 8]) yields a solution $\sigma$ for $N \leq 3$. The proof of the lemma to follow, as given in [39], in which $H^1$ bounds for $\sigma$ are determined in terms of $\mu$, allows limits of $v$ and $\sigma$, and hence solutions for $v \in H^1 \cap L_\infty$. Finally, the Moser iteration theory (cf. [24]) gives 'a priori' $L_\infty$ bounds on $\sigma$ in terms of $H^1 \cap L_\infty$ bounds on $\mu$.

**Lemma 3.2** *Let* $u \in H^1 \cap L_\infty$ *be given and, for* $N \leq 3$, *let* $v$ *be the solution to the weak formulation of the mixed boundary value problem,*

$$\nabla \cdot (e^{u-v}\nabla v) = 0,$$

*on* $G$. *Then the derivative* $D\mathbf{V}$ *of the mapping* $\mathbf{V}$ *from* $u$ *to* $v$, *defined through this equation, is uniformly bounded from* $H^1_{0,\Sigma_D}$ *to itself. The derivative* $D\mathbf{V}$ *is a locally Lipschitz continuous mapping from* $H^1$ *to the mappings from* $H^1 \cap L_\infty$ *to* $H^1_{0,\Sigma_D}$, *for Euclidean dimension* $N \leq 3$. *A similar statement holds for* $\mathbf{W}$.

We consider now $\mathbf{U}_h$ and $\mathbf{V}_h$; for simplicity, the same symbols are used for the extensions.

**Lemma 3.3** *The derivative* $D_{(v,w)}\mathbf{U}_h(v, w) : (\sigma, \tau) \mapsto \mu_h$ *is defined through the solution of the projection relation,*

$$\langle \nabla \mu_h, \nabla \phi \rangle + \langle e^{U_h-v}[\mu_h - \sigma] + e^{w-U_h}[\mu_h - \tau], \phi \rangle = 0, \tag{3.23}$$

*where* $\mu_h$ *and* $\phi$ *are in* $S_h$. *The derivative* $D\mathbf{V}_h(u) : \mu \mapsto \sigma_h$ *may be defined by*

$$\langle e^{u-v_h}[(\mu - \sigma_h)\nabla v_h + \nabla \sigma_h], \nabla \phi \rangle = 0, \tag{3.24}$$

*where* $\sigma_h, \phi \in S_h$.

We close this subsection by mentioning the approximation properties of the derivative maps.

**Lemma 3.4** *The solutions $D_{(v,w)}\mathbf{U}(v,w)(\sigma,\tau) := \mu$ of (3.21) and $D_{(v,w)}\mathbf{U}_h(v,w)(\sigma,\tau) := \mu_h$ of (3.23) satisfy an estimate of the form*

$$\|\mu - \mu_h\| \le Ch^\theta \|[\sigma,\tau]\|, \tag{3.25}$$

*where $C$ does not depend upon $h, v,$ or $w$. The norm used here is the $H^1$ norm or product $H^1$ norm. Similar statements apply to the finite element approximations defined by the derivative maps associated with the quasi-Fermi levels.*

## 3.4 Verification of the 'A Priori' Estimates

The following lemma affords an analysis of the hypotheses of Theorem 3.1 for the semiconductor application. The properties concerning $\mathbf{T}'$ and $\mathbf{T}'_n$ are, for the most part, consequences of earlier discussed facts concerning the composition mappings. It is assumed throughout that the Euclidean dimension $N$ satisfies $N \le 3$. Unsubscripted norms denote energy norms.

**Lemma 3.5** *For the mapping $\mathbf{T}$ and the piecewise linear finite element projection $\mathbf{P}_n$,*

$$\|\mathbf{P}_n\mathbf{T}\mathbf{P}_n x_0 - \mathbf{T}x_0\| \rightarrow 0, \tag{3.26}$$

$$\|\mathbf{P}_n\mathbf{T}'(\mathbf{P}_n x_0) - \mathbf{T}'(x_0)\| \rightarrow 0, \quad \text{if } \|x_0 - \mathbf{P}_n x_0\|_* \rightarrow 0, \tag{3.27}$$

$$\|[\mathbf{T}_n - \mathbf{P}_n\mathbf{T}]\mathbf{P}_n x_0\| \rightarrow 0, \tag{3.28}$$

$$\|[\mathbf{T}'_n - (\mathbf{P}_n\mathbf{T})'](\mathbf{P}_n x_0)\| \rightarrow 0, \tag{3.29}$$

*while for any $\epsilon > 0$ there exist $n_\epsilon$ and $\delta_\epsilon > 0$ such that*

$$\|\mathbf{T}'_n(x) - \mathbf{T}'_n(\mathbf{P}_n x_0)\| \le \epsilon \quad for \quad (n \ge n_\epsilon; \|x - \mathbf{P}_n x_0\|_* \le \delta_\epsilon, x \in \Omega_n), \tag{3.30}$$

*where $\|\cdot\|_*$ is the maximum of $\|\cdot\|$ and the product $L_\infty$ norm.*

m Proof We shall verify only (3.27) and (3.30); the reader is referred to [38] or [39] for greater detail. For the derivative mapping $\mathbf{T}'(x) : f \mapsto g$ we have

$$
\begin{aligned}
\|[\mathbf{P}_n\mathbf{T}'(\mathbf{P}_n x_0) - \mathbf{T}'(x_0)]f\| &\le \|[\mathbf{P}_n - \mathbf{I}]\mathbf{T}'(\mathbf{P}_n x_0)f\| + \|[\mathbf{T}'(\mathbf{P}_n x_0) - \mathbf{T}'(x_0)]f\| \\
&\le ch^\theta \|\mathbf{T}'(\mathbf{P}_n x_0)f\|_{H^{1+\theta}} + L_{\mathbf{T}'}\|[\mathbf{P}_n - \mathbf{I}]x_0\|_* \|f\|.
\end{aligned}
$$

Here we have used (3.19) and the Lipschitz continuity and uniform boundedness in $x$ of the derivative mapping $\mathbf{T}'(x)$. Relation (3.27) follows from the $H^1 \cap L_\infty$ convergence of $\mathbf{P}_n$ to $\mathbf{I}$ on $co\ R_{\mathbf{T}}$ under proper regularity hypotheses. Inequality (3.30) is a restatement of the uniform continuity of $\mathbf{T}'_n$ in $n$.   Box

## 3.5 Verification of the 'A Posteriori' Estimates

The following lemma addresses the hypotheses of Theorem 3.2 for the semiconductor application. It is assumed throughout that the Euclidean dimension $N$ satisfies $N \le 3$. Unsubscripted norms denote energy norms.

**Lemma 3.6** *Let the operators* $\mathbf{T}, \mathbf{P}_n\mathbf{T},$ *and* $\mathbf{T}_n$ *be as defined before. Let* $\mathbf{I} - \mathbf{T}'_n(\tilde{x}_n)$ *be continuously invertible in* $E_n$ *at the approximate solution* $\tilde{x}_n$ *to (3.11), and let*

$$\|[\mathbf{I} - \mathbf{T}'_n(\tilde{x}_n)]^{-1}\| = \kappa_n \le \kappa. \tag{3.31}$$

*Then for a sufficiently small meshwidth* $h$,

$$\gamma_n \equiv (1 + \kappa_n \|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|)\|[\mathbf{T}' - \mathbf{P}_n\mathbf{T}'](\tilde{x}_n)\| + \kappa_n \|[\mathbf{T}_n - \mathbf{P}_n\mathbf{T}]'(\tilde{x}_n)\| < 1. \tag{3.32}$$

*If* $\|\tilde{x}_n - x_n\| \le Ch^\theta$, *where* $C$ *does not depend on* $h$, *then there exist* $\delta_n$ *and* $q_n$ *($\delta_n > 0$; $0 < q_n < 1$) such that*

$$\sup_{\|x - \tilde{x}_n\|_* \le \delta_n} \|\mathbf{T}'(x) - \mathbf{T}'(\tilde{x}_n)\| \le \frac{q_n}{\kappa'_n}, \tag{3.33}$$

$$\|\tilde{x}_n - \mathbf{T}\tilde{x}_n\| \le \frac{\delta_n(1 - q_n)}{\kappa'_n} = ch^\theta, \tag{3.34}$$

*where*

$$\kappa'_n = \frac{1 + \kappa_n \|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|}{1 - \gamma_n}.$$

m Proof The bound on $\gamma_n$ as stated in (3.32) is proven through

$$
\begin{aligned}
\gamma_n &\equiv (1 + \kappa_n \|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|)\|[\mathbf{T}' - \mathbf{P}_n\mathbf{T}'](\tilde{x}_n)\| + \kappa_n \|[\mathbf{T}_n - \mathbf{P}_n\mathbf{T}]'(\tilde{x}_n)\| \\
&\le (1 + \kappa_n \|\mathbf{P}_n\mathbf{T}'(\tilde{x}_n)\|)Ch^\theta \|\mathbf{T}'(\tilde{x}_n)\|_{H^1, H^{1+\theta}} + \kappa_n \|[\mathbf{T}_n - \mathbf{P}_n\mathbf{T}]'(\tilde{x}_n)\| \\
&\le Ch^\theta.
\end{aligned}
$$

Here $C$ is a generic constant.

The existence of a finite $\delta_n$ and $q_n$ ($\delta_n > 0$; $0 < q_n < 1$), such that (3.33) and (3.34) hold, follows because

$$\sup_{\|x - \tilde{x}_n\|_* \le \delta_n} \|\mathbf{T}'(x) - \mathbf{T}'(\tilde{x}_n)\| \le L_{\mathbf{T}'}\|x - \tilde{x}_n\|_* \le L_{\mathbf{T}'}\delta_n,$$

so that we can choose $q_n = L_{\mathbf{T}'}\delta_n\kappa'_n$. On the other hand, $x_n$ is a fixed point of $\mathbf{T}_n$ and therefore

$$
\begin{aligned}
\|\tilde{x}_n - \mathbf{T}\tilde{x}_n\| &\le \|\tilde{x}_n - x_n\| + \|\mathbf{T}_n x_n - \mathbf{T}x_n\| + \|\mathbf{T}x_n - \mathbf{T}\tilde{x}_n\| \\
&\le (1 + L_{\mathbf{T}})\|\tilde{x}_n - x_n\| + \|\mathbf{T}_n - \mathbf{T}\|\|x_n\| \\
&\le Ch^\theta.
\end{aligned}
$$

This implies that (3.34) holds provided that

$$\delta_n = Ch^\theta \frac{\kappa'_n}{(1 - q_n)}.$$

The proof is concluded if we can show that the choices

$$q_n = L_{\mathbf{T}'}\kappa'_n\delta_n, \qquad \delta_n = Ch^\theta\kappa'_n/(1 - q_n)$$

are compatible with the requirement $q_n < 1$. Thus, set $\alpha = L_{\mathbf{T}'}\kappa'_n$, $\beta = Ch^\theta\kappa'_n$. A simple argument shows, provided $\beta$ is sufficiently small so that $\alpha\beta \le \frac{1}{4}$, we may choose $\delta_n \le 1/(2\alpha)$ with $q_n \le \frac{1}{2}$. Note that, the bounds for $\kappa'_n$ and $\gamma_n$ show that $\kappa'_n$ remains bounded as $h$ decreases so that the bound $\alpha\beta \le \frac{1}{4}$ does not depend on $h$. *Box*

## 3.6   Summary of Results

The following corollary expresses a summary of the major results in conjunction with the employed hypotheses. The norms denote energy norms.

**Corollary 3.1** *Under appropriate regularity hypotheses, let $x_0$ be a fixed point of $\mathbf{T}$ and suppose that $\mathbf{T}'(x_0)$ does not possess 1 as an eigenvalue. Suppose also that the $L_\infty$ estimates required for the regularization mapping invariance of Theorem 3.1 hold. Then there exist an index $n_0$ and a neighborhood of $x_0$ containing fixed points $x_n$ of $\mathbf{T}_n, n \geq n_0$, satisfying*

$$\|x_0 - x_n\| \leq Ch^\theta \quad \forall n. \tag{3.35}$$

*Here $C$ is a constant independent of $n$ and $h$. Conversely, suppose $\{\tilde{x}_n\}$ is a sequence of approximate fixed points of $\mathbf{T}_n$ satisfying $\|\tilde{x}_n - x_n\| \leq ch^\theta$ and (3.31). Then, under the regularization invariance of Theorem 3.2 and the regularity hypothesis, there exists a fixed point $x_0$ of $\mathbf{T}$ such that*

$$\|x_0 - \tilde{x}_n\| \leq Ch^\theta \quad \forall n. \tag{3.36}$$

*Here $C$ is a constant independent of $n$ and $h$. In all cases $N \leq 3$.*

m Proof The hypotheses of the corollary absorb those of Theorems 3.1 and 3.2, due to the compactness of $\mathbf{T}'(x_0)$. Estimates (3.35) and (3.36) follow from the conjunction of (3.13) with (3.14), and from (3.16), respectively. The former two inequalities must be appropriately combined with the triangle inequality,

$$\|\mathbf{P}_n\mathbf{T}x_0 - \mathbf{T}_n\mathbf{P}_nx_0\| \leq \|(\mathbf{P}_n - \mathbf{I})\mathbf{T}x_0\| + \|\mathbf{T}x_0 - \mathbf{T}\mathbf{P}_nx_0\| + \|(\mathbf{T} - \mathbf{T}_n)\mathbf{P}_nx_0\|,$$

which yields order $h^\theta$ convergence. For the 'a posteriori' result, inequality (3.16) must be supplemented by

$$\frac{\alpha_n}{1 - q_n} \leq (\frac{\kappa'_n}{1 - q_n})\|\tilde{x}_n - \mathbf{T}\tilde{x}_n\| \leq \delta_n.$$

The choice, $\delta_n = Ch^\theta \kappa'_n/(1 - q_n)$ and the boundedness of $\kappa'_n$ and $1/(1 - q_n)$ were discussed earlier.     *Box Remark 6.* Estimates of order $h^\theta$ follow from Corollary 3.1 for $\|v - v_h\|$ and $\|w - w_h\|$, if $[v, w]$ and $[v_h, w_h]$ are fixed points of $\mathbf{T}$ and $\mathbf{T}_n$, respectively. Estimates for $\|U - U_h\|$ follow immediately from

$$\|\mathbf{U}(v, w) - U_h\| \leq \|\mathbf{U}(v, w) - \mathbf{U}(v_h, w_h)\| + \|\mathbf{U}(v_h, w_h) - U_h\|.$$

*Remark 7.* The verification of the $L_\infty$ estimates, assumed in Corollary 3.1, can be achieved by elliptic regularity theory in the case of $\mathbf{T}$. For $\mathbf{T}_n$, it is achieved by a careful study of pointwise convergence properties of finite element approximations. Details are given in [39].

## 4   The Evolution System and Newton's Method

A general outer iteration, based upon linearization, is introduced at discrete time steps for the one-dimensional semiconductor device model. The iteration depends upon solving the semidiscrete device equations approximately, specifically, in such a way that the residual is of order $\Delta t$ in an appropriate norm. This maintains the order of the backward Euler method. A

monitoring of the constants, including time-step requirements for solvability of the semidiscrete systems, as well as smoothness and bounded invertibility for the maps defining the Newton approximations, is possible. An invariant-region principle provides an important theoretical basis. Full details, including proofs of the theorems presented in this section, are given in [14].

The devices studied here are one-dimensional, and the dependent variables selected are the electrostatic potential and the carrier concentrations. The model is described more fully in §4.1. The mobility coefficients are selected to have a form similar to that used in computations for silicon devices so as to ensure the physically essential property of saturation; the Einstein relations are assumed only so that physically realistic representation of diffusion is possible. For the dependent variables selected here, the Einstein relations provide no equation simplification. The model assumes zero recombination. As described here, the model is a special case of the spatially multidimensional model considered in [35]; in particular, we deduce that the initial/boundary-value problem possesses a unique solution, globally in time. Uniqueness need not persist for the semidiscrete solutions without additional time-step restrictions, however, so that the linearized time-stepping must have an inherent tracking mechanism; this is provided by continuation.

In this paper, we do not analyze second order time-stepping methods. We note, however, that one such method, possessing the property of $L$-stability, has been introduced in [3].

## 4.1  Transient Semiconductor Equations and Semidiscretization

Our interest here is in a transient initial/boundary-value problem, in one spatial dimension, that models the behavior of a simple semiconductor device; hence, the appropriate space-time domain is $\bar{G} \times [0, T_0]$ where $G = (a, b) \subset \mathbf{R}$ and $T_0$ is the final time of interest. The transient semiconductor equations, which hold for $(x, t) \in G \times (0, T_0]$, were given in §1.2. We repeat them here in the form desired. For simplicity, recombination is excluded.

$$-\nabla \cdot (\epsilon \nabla u) + e(n - p - k_1) = 0, \tag{4.1}$$

$$e\frac{\partial n}{\partial t} - \nabla \cdot J_n = 0, \tag{4.2}$$

$$e\frac{\partial p}{\partial t} + \nabla \cdot J_p = 0. \tag{4.3}$$

Note that (4.1) is the stationary Maxwell equation governing the electrostatic potential $u$, while (4.2) and (4.3) are typical continuity equations governing $n$ and $p$, respectively. In addition, we assume that appropriate initial data $n_0(x), p_0(x)$ and boundary data $\bar{u}(x, t), \bar{n}(x), \bar{p}(x)$ are given.

We write the current densities in the traditional drift-diffusion form introduced in (1.4) and (1.5). One simple mobility model that includes velocity saturation has the form

$$\mu(\nabla u) = \mu_0 \left[ 1 + \left( \frac{\mu_0 |\nabla u|}{v_{sat}} \right)^{\gamma} \right]^{-1/\gamma},$$

where $\mu_0$ and $v_{sat}$ are the low-field mobility value and temperature-dependent saturation velocity, respectively (cf. Caughey and Thomas [11] and Thornber [60]). We will assume $\gamma \equiv 2$ for both carriers and write the mobility functions as

$$\mu_n = \mu_{0n} \left[ 1 + \left( \frac{\mu_{0n} |\nabla u|}{v_{sn}} \right)^2 \right]^{-1/2}, \tag{4.4}$$

24

$$\mu_p = \mu_{0p} \left[ 1 + \left( \frac{\mu_{0p}|\nabla u|}{v_{sp}} \right)^2 \right]^{-1/2}. \tag{4.5}$$

Of course, there are other possibilities for the mobilities, which are less smooth or are more complicated functions of the dependent variables. It is important to note, however, that the mobility models are largely phenomenological expressions, which attempt to incorporate a number of experimentally observed phenomena.

We include the Einstein relations as before,

$$D_n = \frac{kT}{e} \mu_n, \quad D_p = \frac{kT}{e} \mu_p. \tag{4.6}$$

To simplify matters, we select the natural units so that

$$\frac{kT}{e} \equiv 1; \tag{4.7}$$

in fact, we scale all of the equations (cf. [17, 16]), except in our treatment of $\epsilon$ (see also the book by Markowich [45]).

We continue by discussing a solvability result for the semidiscrete version of the problem (4.1)–(4.7), which is based upon an invariant-region principle. The discussion here is related to earlier work described in [4] and [33]. A convergence result, where the residual is controlled in $L_2$, is presented in § 4.2.

We shall find it convenient to expand the mobility terms in (4.2) and (4.3) by the product rule, substitute the second derivative terms via the potential equation (4.1), and make use of the Einstein relations (4.6) and the scaling (4.7). Thus, throughout much of this section, we shall write the semidiscrete system via a fully implicit time discretization as

$$-\epsilon \nabla^2 u_k + n_k - p_k = k_1, \tag{4.8}$$

$$\frac{n_k - n_{k-1}}{\Delta t_k} - \nabla \cdot (\mu_n \nabla n_k) + \mu_n \nabla u_k \nabla n_k + U_{n,k} = 0, \tag{4.9}$$

$$\frac{p_k - p_{k-1}}{\Delta t_k} - \nabla \cdot (\mu_p \nabla p_k) - \mu_p \nabla u_k \nabla p_k + U_{p,k} = 0, \tag{4.10}$$

where

$$U_{n,k} = \epsilon^{-1}(\mu_n' \nabla u_k + \mu_n) n_k (n_k - p_k - k_1), \tag{4.11}$$

$$U_{p,k} = -\epsilon^{-1}(\mu_p' \nabla u_k + \mu_p) p_k (n_k - p_k - k_1), \tag{4.12}$$

and $\mu_\star$ and $\mu_\star'$ are evaluated at $\nabla u_k$.

The invariant-region principle referred to above represents a slight weakening of the usual such principle, since it includes the spatially dependent doping. Specifically, we define a number

$$\lambda_0 = \max(\|\bar{n}\|_{L_\infty}, \|\bar{p}\|_{L_\infty}, \|n_0\|_{L_\infty}, \|p_0\|_{L_\infty}), \tag{4.13}$$

and *functions* $n_k^{max}$ and $p_k^{max}$, via the relations,

$$n_0^{max} = \lambda_0 + \text{Step}(k_1^+), \tag{4.14}$$

$$p_0^{max} = \lambda_0 + \text{Step}(k_1^-), \tag{4.15}$$

$$n_k^{max} - n_{k-1}^{max} = p_k^{max} - p_{k-1}^{max} = 0 \text{ for } 1 \le k \le L, \tag{4.16}$$

25

where

$$\text{Step}(k_1^+) = \begin{cases} \sup k_1 & \text{if } k_1 > 0, \\ 0 & \text{otherwise.} \end{cases} \tag{4.17}$$

It follows that, if

$$n^{max} = \lambda_0 + \sup k_1^+, \ p^{max} = \lambda_0 + \sup k_1^-, \tag{4.18}$$

then $n^{max} \geq n_k^{max}$ and $p^{max} \geq p_k^{max}$. Here, $\bar{n}$ and $\bar{p}$ are linear extensions of the boundary data, and $k_1 = k_1^+ - k_1^-$. Note the different usage of the negative part in this definition, with respect to that in §2.2. We are now prepared for the first result. The invariant region principle is stronger than that described in [14], due to the absence of recombination here.

**Theorem 4.1** *Under hypotheses on the time step (cf. (4.23), and (4.24) stated at the conclusion of the theorem), there is a solution of the Dirichlet boundary-value problem (4.8)–(4.10) with boundary values*

$$u_k(x_e) = \bar{u}(x_e, t_k), \ n_k(x_e) = \bar{n}(x_e), \ p_k(x_e) = \bar{p}(x_e) \ for \ x_e = a, b. \tag{4.19}$$

*The solution triple satisfies an invariant-region principle for the carrier concentrations and a generalized maximum principle for the potential:*

$$0 \leq n_k \leq n_k^{max} \leq n^{max}, \tag{4.20}$$

$$0 \leq p_k \leq p_k^{max} \leq p^{max}, \tag{4.21}$$

$$|u_k| \leq \|\bar{u}(\cdot, t_k)\|_{L_\infty} + \epsilon^{-1}(e^{b-a} - 1)(n^{max} + p^{max}) \stackrel{\text{def}}{=} u_k^{max}. \tag{4.22}$$

*The time step restrictions are given as follows; the second guarantees that the system fixed point map is Lipschitz continuous, while the first is related to the invariant region:*

$$8\epsilon^{-1}\Delta t_k \max(\mu_{0n}, \mu_{0p})[\max(n^{max}, p^{max}) + \|k_1\|_{L_\infty}] \leq 1, \tag{4.23}$$

$$\Delta t_k < \frac{\min(\inf \mu_n, \inf \mu_p)}{\max(v_{sn}^2, v_{sp}^2)}. \tag{4.24}$$

## 4.2 Approximate Solvability and $L_2$ Residual Control

A typical computing procedure involves solving the system (4.8)–(4.10) only approximately. In the next major result, we present a criterion for approximate solvability in terms of the residuals of the individual systems at successive time steps. The computed triple, $[u_k, n_k, p_k]$, need not satisfy the invariant-region property as presented in Theorem 4.1. Since pointwise bounds are essential to the theory, however, they are built into the hypothesis structure directly at the outset.

**Theorem 4.2** *Suppose that the triple $[u_k, n_k, p_k]$ approximately solves (4.8)–(4.10), i.e., the boundary conditions (4.19) are satisfied exactly and*

$$-\epsilon\nabla^2 u_k + n_k - p_k - k_1 = r_1, \tag{4.25}$$

$$\frac{n_k - n_{k-1}}{\Delta t_k} - \nabla \cdot (\mu_n \nabla n_k) + \mu_n \nabla u_k \nabla n_k + U_{n,k} = r_2, \tag{4.26}$$

$$\frac{p_k - p_{k-1}}{\Delta t_k} - \nabla \cdot (\mu_p \nabla p_k) - \mu_p \nabla u_k \nabla p_k + U_{p,k} = r_3, \tag{4.27}$$

where the residuals $r_1$, $r_2$, and $r_3$ satisfy

$$\|r_1\|_{L_2} \le C_u \Delta t_k, \ \ \|r_2\|_{L_2} \le C_n \Delta t_k, \ \ \|r_3\|_{L_2} \le C_p \Delta t_k, \tag{4.28}$$

for $k = 1, \ldots, L$. Define the errors,

$$v_k \stackrel{\text{def}}{=} u(\cdot, t_k) - u_k, \tag{4.29}$$

$$e_k \stackrel{\text{def}}{=} n(\cdot, t_k) - n_k, \tag{4.30}$$

$$q_k \stackrel{\text{def}}{=} p(\cdot, t_k) - p_k. \tag{4.31}$$

Suppose that pointwise bounds,

$$\|u_k\|_{L_\infty} \le c_u, \ \ \|n_k\|_{L_\infty} \le c_n, \ \ \|p_k\|_{L_\infty} \le c_p, \tag{4.32}$$

exist for the computed semidiscrete approximations. Define

$$\Delta t \stackrel{\text{def}}{=} \max_k \Delta t_k. \tag{4.33}$$

If $\Delta t$ is sufficiently small, then these errors converge with optimal order in $L_2$:

$$\sup_{1 \le M \le L} (\|e_M\|_{L_2}^2 + \|q_M\|_{L_2}^2) \tag{4.34}$$

$$+ \min(\inf \mu_n, \inf \mu_p) \sum_{k=1}^{L} (\|\nabla e_k\|_{L_2}^2 + \|\nabla q_k\|_{L_2}^2) \Delta t_k \le C_1 (\Delta t)^2,$$

$$\sup_{1 \le M \le L} \|\nabla v_M\|_{L_2}^2 \le C_2 (\Delta t)^2, \tag{4.35}$$

where $C_1$ and $C_2$ are certain positive constants which can be estimated explicitly. The result remains valid if $n^2$ is replaced by $(n_+)^2$ and $p^2$ by $(p_+)^2$ in $U_{n,k}$ and $U_{p,k}$, respectively.

A natural approximation, giving rise to approximate solvability in the language of the previous theorem, is defined by Newton's method. In the following subsections, we shall explicitly define this method for the semidiscrete system (4.8)-(4.10). Moreover, the classical properties shall be reviewed for effectiveness. For the moment, we indicate here the result of that discussion in relation to Theorem 4.2. The following result specifies the size limitation of the initial residual to achieve the desired $O(\Delta t)$ accuracy.

**Theorem 4.3** *Suppose that at a given time step, $t = t_k$, $m$ iterations of Newton's method are implemented to compute the residual triple, denoted $r^m$, for the system (4.8)-(4.10). If, as motivated by the previous theorem, we desire*

$$\|r^m\|_{L_2} \le c \Delta t_k, \tag{4.36}$$

*for specified $c$, then the initial residual must satisfy*

$$c \Delta t_k \ge (2 M_1^2 M_3)^m \|r^0\|_{L_2}^{2^m}. \tag{4.37}$$

*Here, $2 M_3$ is a Lipschitz constant of the derivative of the map defining the differential system, and $M_1$ is a local bound for the norms of approximate inverses of those derivative maps.*

## 4.3   Approximate Newton Methods for the Semidiscrete System

We shall describe a rather general approximate Newton method, in which the approximate linear inversion is based upon an arbitrary and unspecified inner iteration, arranged in such a way that the algorithm is quadratically convergent. Though the description is general, we shall aim for the development of an algorithm with a small number of outer iterations. The reader will recall from § 4.2 that it is the residual of the nonlinear map, defining the semidiscrete system, which is to be controlled with order $\Delta t_k$ (cf. Theorem 4.2).

The plan of the following subsections is as follows. An operator tableau for the linearized problem is presented in § 4.4, while the properties required of an approximate Newton method are presented in § 4.5. A study of these properties, as they affect the semiconductor model, is carried out in [14]. This reference also discusses appropriate continuation issues.

## 4.4   The Linearized Problem

This is primarily a formal subsection, in which we briefly present the operator for the linearized problem. At the conclusion, we mention considerations of operator domain and range.

Suppose the system (4.8)–(4.10) is denoted by

$$F(u_k, n_k, p_k) = 0, \tag{4.38}$$

and the $m$th linear increment, i.e., the difference between the $m$th and the $(m-1)$th Newton iterates, is denoted by $[\phi_k^m, \psi_k^m, \omega_k^m]$. This increment satisfies

$$\begin{bmatrix} -\epsilon \nabla^2 \star & \star & -\star \\ (F')_{21} & (F')_{22} & (\partial U_p/\partial p)\star \\ (F')_{31} & (\partial U_n/\partial n)\star & (F')_{33} \end{bmatrix} \begin{pmatrix} \phi_k^m \\ \psi_k^m \\ \omega_k^m \end{pmatrix} = - \begin{pmatrix} r_1^{m-1} \\ r_2^{m-1} \\ r_3^{m-1} \end{pmatrix} \tag{4.39}$$

where

$$(F')_{21} \overset{\text{def}}{=} -\nabla\left[\nabla n_k^{m-1}\mu_n'(\nabla u_k^{m-1})\nabla\star\right] + \nabla n_k^{m-1}\nabla u_k^{m-1}\mu_n'(\nabla u_k^{m-1})\nabla\star \tag{4.40}$$

$$+\nabla n_k^{m-1}\mu_n(\nabla u_k^{m-1})\nabla\star + \frac{\partial U_n}{\partial u}\nabla\star$$

$$(F')_{22} \overset{\text{def}}{=} \frac{\star}{\Delta t_k} - \nabla\left[\mu_n(\nabla u_k^{m-1})\nabla\star\right] + \mu_n(\nabla u_k^{m-1})\nabla u_k^{m-1}\nabla\star + \frac{\partial U_n}{\partial n}\star \tag{4.41}$$

$$(F')_{31} \overset{\text{def}}{=} -\nabla\left[\nabla p_k^{m-1}\mu_p'(\nabla u_k^{m-1})\nabla\star\right] - \nabla p_k^{m-1}\nabla u_k^{m-1}\mu_p'(\nabla u_k^{m-1})\nabla\star \tag{4.42}$$

$$-\nabla p_k^{m-1}\mu_p(\nabla u_k^{m-1})\nabla\star + \frac{\partial U_p}{\partial u}\nabla\star$$

$$(F')_{33} \overset{\text{def}}{=} \frac{\star}{\Delta t_k} - \nabla\left[\mu_p(\nabla u_k^{m-1})\nabla\star\right] - \mu_p(\nabla u_k^{m-1})\nabla u_k^{m-1}\nabla\star + \frac{\partial U_p}{\partial p}\star . \tag{4.43}$$

Moreover, $r_1^{m-1}$, $r_2^{m-1}$, and $r_3^{m-1}$ are the residuals of $F(u_k^{m-1}, n_k^{m-1}, p_k^{m-1})$, given explicitly in (4.8)–(4.10); $[u_k^{m-1}, n_k^{m-1}, p_k^{m-1}]$ is the Newton iterate.

It can be advantageous to consider two distinct application frameworks for the mappings $F$ and $F'$. Note that (4.39)–(4.43) is a delineation of the equation,

$$F'(u_k^{m-1}, n_k^{m-1}, p_k^{m-1}) \begin{pmatrix} \phi_k^m \\ \psi_k^m \\ \omega_k^m \end{pmatrix} = - \begin{pmatrix} r_1^{m-1} \\ r_2^{m-1} \\ r_3^{m-1} \end{pmatrix}. \tag{4.44}$$

(Also note that (4.39) represents a considerable notational economy; a similar tableau was introduced in [8].) We have in mind both an $L_2$ and $H^{-1}$ residual measurement framework because of the variety of applications. Thus, $F$ will be defined from triples in the affine subspace of $H^s(G)$ given by

$$X_s = [\bar{u}(\cdot, t_k), \bar{n}, \bar{p}] + Y_s, \quad s = 1, 2, \tag{4.45}$$

where

$$Y_s = \prod_1^3 H^s(G) \cap H_0^1(G), \quad s = 1, 2, \tag{4.46}$$

to the triple product of copies of $H^{s-2}(G)$. $F'(z)$, for fixed $z$, is defined from $Y_s$ to the same space. Any approximate inverse, denoted $S(z)$ in the next subsection, reverses domain and range. Finally, we shall find it convenient, and, in fact, necessary to restrict $F$ to have local domain of definition.

## 4.5   A Class of Approximate Newton Methods

We close this section with a subsection devoted to general principles. By doing so we clarify some results of the previous subsection. A more complete development may be found in [32]. We begin by recalling the three fundamental properties required of an approximate Newton method, based on a family of linear maps, $S(z^m)$, where $\{z^m\}$ comprises the Newton sequence, defined by

$$z^m - z^{m-1} \stackrel{\text{def}}{=} -S(z^{m-1})F(z^{m-1}), \tag{4.47}$$

and where $z$, $F(z) = 0$, is sought as $z = \lim_{m \to \infty} z^m$. The three properties are:

$$\|S(z^{m-1})F(z^{m-1})\| \le M_1 \|F(z^{m-1})\|, \tag{4.48}$$
$$\|[F'(z^{m-1})S(z^{m-1}) - I]F(z^{m-1})\| \le M_2 \|F(z^{m-1})\|^2, \tag{4.49}$$
$$\|F'(x) - F'(y)\| \le 2M_3 \|x - y\|, \tag{4.50}$$

together with a mechanism insuring that successive iterates lie within the domain of definition of $F$. Here, the constants $M_1$, $M_2$, and $M_3$ are independent of the elements in the domain of definition of the respective maps. This framework is reminiscent of that described by Bank and Rose [5].

In order to estimate $\|F(z^m)\|$, we write

$$F(z^m) = -[F'(z^{m-1})S(z^{m-1}) - I]F(z^{m-1}) + R(z^{m-1}, z^m), \tag{4.51}$$

where

$$R(z^{m-1}, z^m) = F(z^m) - F(z^{m-1}) - F'(z^{m-1})(z^m - z^{m-1}). \tag{4.52}$$

Notice that the first term in (4.51) may be dominated by $M_2 \|F(z^{m-1})\|^2$, via (4.49), and the second term by

$$\|R(z^{m-1}, z^m)\| = \left\| \int_0^1 \left[ F'\left(z^{m-1} + s(z^m - z^{m-1})\right) - F'(z^{m-1}) \right] (z^m - z^{m-1}) \, ds \right\|$$

, on use of a standard Taylor expansion. From these remarks we have the estimate,

$$\|F(z^m)\| \le (M_1^2 M_3 + M_2)\|F(z^{m-1})\|^2. \tag{4.53}$$

We have proved the bulk of the following lemma.

**Lemma 4.1** *Let $F$ be a map defined on a closed ball $B_r$ in an affine subspace $x_0 + X_0$, contained in a Banach space $X$, with Fréchet derivative*

$$F'(z) : X_0 \mapsto W, \quad z \in B_r.$$

*Here $W$ is a Banach space containing the range of $F$. Suppose that $S(z) : W \mapsto X_0$ is defined for each $z \in B_r$, and that $F$, $F'$, and $S$ satisfy (4.48), (4.49), and (4.50). If $z^0$ satisfies*

$$M_1 \|F(z^0)\| \leq (1 - \alpha)r, \tag{4.54}$$

*where $0 \leq \alpha < 1$ is such that $z^0 \in B_{\alpha r}$, then $z^1 \in B_r$ and the residual $F(z^1)$ satisfies*

$$\|F(z^1)\| \leq (M_1^2 M_3 + M_2)\|F(z^0)\|^2. \tag{4.55}$$

*Remark 8.* The hypothesis (4.50) can be weakened from $x, y \in B_r$ to $x, y$ on the line segment between $z^{m-1}$ and $z^m$. Also, the estimate (4.55) is related to (4.37), $m = 1$, if the choice $M_2 = M_3$ is made.

*Remark 9.* This section is quite different in spirit from the previous two sections, since it does not depend upon fixed point formulations. The advantage of such fixed point frameworks is more than cosmetic; it avoids the necessity of smoothing in passing to numerical discretization as a way of defining the approximate Newton method via the inner iterations. For a more complete explanation, the reader is referred to [31] and [37].

# 5 More General Moment Models: A Review

## 5.1 Mass, Momentum and Energy Transport Equations

The equations as presented here are discussed in references [10], [52], and [12]. They are, respectively, derived as zeroth, first, and second order moments of the Boltzmann equation, with the latter written for an electron species moving in an electric field without magnetic effects as,

$$\frac{\partial f}{\partial t} + u \cdot \nabla_x f - \frac{e}{m} E \cdot \nabla_u f = C. \tag{5.56}$$

Here, $f = f(x, u, t)$ is the numerical distribution function of a carrier species, $x$ is the position vector, $u$ is the species' group velocity vector, $E = E(x, t)$ is the electric field, $e$ is the electron charge modulus, $m$ is the effective electron mass, and $C$ is the time rate of change of $f$ due to collisions, typically represented by

$$C = \int S(u', u) f(x, u', t)(1 - f(x, u, t)) \, du' - \int S(u, u') f(x, u, t)(1 - f(x, u', t)) \, du', \tag{5.57}$$

in terms of a scattering kernel $S$. The moment equations are expressed in terms of certain dependent variables, where $n$ is the electron concentration, $v$ is the velocity, $p$ is the momentum density, $P$ is the symmetric pressure tensor, $q$ is the the heat flux, $e_I$ is the internal energy, and $C_n$, $C_p$, and $C_W$ represent moments of $C$, taken with respect to the functions

$$\begin{aligned} g_0(u) &\equiv 1, \\ g_1(u) &= mu, \\ g_2(u) &= \frac{m}{2} |u|^2 . \end{aligned}$$

The moment equations are given by:

$$\frac{\partial n}{\partial t} + \nabla \cdot (nv) = C_n, \qquad (5.58)$$

$$\frac{\partial p}{\partial t} + v(\nabla \cdot p) + (p \cdot \nabla)v = -enE - \nabla \cdot P + C_p, \qquad (5.59)$$

$$\frac{\partial}{\partial t}(\frac{mn}{2} \mid v \mid^2 + mne_I) + \nabla \cdot (v \,[\frac{mn}{2} \mid v \mid^2 + mne_I]) =$$
$$-env \cdot E - \nabla \cdot (vP) - \nabla \cdot q + C_W. \qquad (5.60)$$

The first Maxwell equation for the electrostatic potential must be adjoined. The reader should realize that each species contributes a corresponding moment subsystem, with appropriately signed charge. In the subsystem (5.58), (5.59), and (5.60) above, it is understood that the tensor $P$ acts like a matrix; $\nabla$ operates like a row vector in relation to it and $v$ acts like a row vector. Moreover, the following definitions in terms of the moment quantities hold.

1. The concentration is given by
$$n := \int f \; du.$$

2. The average velocity is given by

$$v := \frac{1}{n} \int uf \; du.$$

3. The momentum is given by
$$p := mnv.$$

4. The random velocity is given by
$$c := u - v.$$

5. The pressure tensor is given by

$$P_{ij} := m \int c_i c_j f \; du.$$

6. The internal energy density is given by

$$e_I := \frac{1}{2n} \int \mid c \mid^2 f \; du.$$

   This function represents energy/unit mass/unit concentration.

7. The heat flux is given by
$$q_i := \frac{m}{2} \int c_i \mid c \mid^2 f \; du.$$

The assumptions on the distribution function $f$ are now stated.

The function $f$ decreases sufficiently rapidly at infinity:

$$\lim_{|u|\to\infty} g_i(u)f(u) = 0, \ i = 0, 1, 2.$$

Finally, we make certain observations about integral evaluations which follow from the definitions and assumptions after integration by parts.

**Integral Identities**

1. $\int c_i f \, du = 0.$

2. $\int \frac{\partial f}{\partial u_i} \, du = 0.$

3. $\int u_j \frac{\partial f}{\partial u_i} \, du = -\delta_{ij} n.$

4. $\int |u|^2 \frac{\partial f}{\partial u_i} \, du = -2nv_i.$

The derivation of (5.58), (5.59), and (5.60) now continues as follows. We multiply the Boltzmann equation (5.56) by $g_0, g_1$, and $g_2$, respectively, and integrate over group velocity space. The mass equation (5.58) is immediate from the definitions of $v$ and $C_n$, and from the second identity above. In order to derive the momentum equation (5.59), we begin with the identity,

$$\int u_i u_j f \, du = nv_i v_j + \int c_i c_j f \, du, \tag{5.61}$$

which makes use of the first identity. If (5.61) and the definitions of $P$ and $C_p$ are applied to the integrated product of (5.56) and $g_1$, one obtains (5.59) after an application of the third identity above. The derivation of (5.60) begins with the identity

$$(1/2) \int u_i |u|^2 f \, du =$$

$$\frac{nv_i}{2} |v|^2 + v_i ne_I + \sum_j v_j P_{ij}/m + (1/2) \int c_i |c|^2 f \, du, \tag{5.62}$$

which makes use of the first identity above. An application of (5.62), the identity (5.61), and the definitions of $q$ and $C_W$, to the integrated product of (5.56) and $g_2$, gives (5.60) after an application of the fourth identity above. This completes the mass/momentum/energy system derivation. In addition to these transport equations, we have Poisson's equation for the electric field, where $k_1 :=$ doping and $\epsilon :=$ dielectric :

$$E = -\nabla\phi, \tag{5.63}$$
$$\nabla\cdot(\epsilon\nabla\phi) = -\sum e_i n_i - k_1. \tag{5.64}$$

Here, we have used the convention that there are different species, each of concentration $n_i$ and (signed) charge $e_i$. The entire system consists of equations (5.58), (5.59), and 5.60), repeated according to species, and (5.63), (5.64).

## 5.2 Moment Closure and Relaxation Relations

The system derived in the preceding subsection has fifteen dependent variables in physical space in the case of one species, determined by $\phi$, $n$, $v$, $P$, $e_I$, and $q$. By moment closure is meant the selection of compatible relations among these variables, so that the number of equations is equal in number to the remaining primitive variables selected. One way of proceeding is to introduce a new tensor variable $T$, the carrier temperature, defined by

$$P_{ij} = nkT_{ij},$$

where $k$ is Boltzmann's constant, and a scalar variable $W$, the total carrier energy. A program of reduction to the basic variables $n$, $v$, $W$, and $\phi$ can be implemented by the following assumptions:

1. The pressure tensor is isotropic, with diagonal entries $P_s$ and off-diagonal entries zero, for a suitable scalar function $P_s$. By previous relations, $P_s$ is related to $e_I$ via $mne_I = \frac{3}{2}P_s$.

2. It follows from the previous assumption that temperature may be represented by a scalar quantity $T$ and that the internal energy is represented in terms of $T$ by

$$me_I = \frac{3}{2}kT.$$

3. The total energy density (per unit concentration) $w$ is given by combining internal energy and parabolic energy bands:

$$w = me_I + \frac{1}{2}m \mid v \mid^2,$$

and the total energy (per unit volume) $W$ is the product, $W = nw$.

4. The heat flux is obtained by a differential expression involving the temperature. A popular such choice is

$$q = -\kappa\nabla T.$$

Here, $\kappa$ is the thermal conductivity governed by the Wiedemann-Franz law (cf. [9]).

In the case of $N$ species, the closure relations determine $(d+2)N + 1$ variables in $d$ spatial dimensions. It is possible to rewrite the system (5.58), (5.59), and (5.60) with the closure assumptions incorporated. We have the following:

$$\frac{\partial n}{\partial t} + \nabla\cdot(nv) = C_n, \tag{5.65}$$

$$\frac{\partial p}{\partial t} + v(\nabla\cdot p) + (p \cdot \nabla)v = -enE - \nabla(nkT) + C_p, \tag{5.66}$$

$$\frac{\partial}{\partial t}W + \nabla\cdot(v\,W) = -env\cdot E - \nabla\cdot(vnkT)$$
$$+\nabla\cdot(\kappa\nabla T) + C_W. \tag{5.67}$$

The final step deals with the replacemant of the collision moments. Motivated by the approach of [49], [2], [52], and [27], we define the recombination rate $R$ and the momentum and energy relaxation times, $\tau_p$ and $\tau_w$, respectively, in terms of averaged collision moments as follows.

1. The particle recombination rate $R$ is given by

$$R := -C_n := -\int C\ du.$$

2. The momentum relaxation time $\tau_p$ is given via

$$\frac{p}{\tau_p} := -\int muC\ du := -C_p.$$

3. The energy relaxation time $\tau_w$ is given via

$$-\frac{W - W_0}{\tau_w} := \frac{m}{2}\int \mid u \mid^2 C\ du := C_W.$$

Here, $W_0$ denotes the rest energy, $\frac{3}{2}kT_0$, where $T_0$ is the lattice temperature.

The forms for the relaxation times used in [2] and retained by subsequent authors are:

$$\begin{aligned}
\tau_p &= c_p/T, \\
\tau_w &= c_w\frac{T}{T + T_0} + \frac{1}{2}\tau_p.
\end{aligned}$$

Here, $c_p$ and $c_w$ are physical constants. A comprehensive discussion of these issues is outside the scope of the current exposition. We note, however, that the representations of the relaxation times given in [27] can be used to obtain mobility derivations for representations such as (4.4) and (4.5). It is of interest that numerical simulations can be based upon the hydrodynamic model; for the $n^+ - n - n^+$ diode, see [23] and [20].

# 6   Epilogue: Historical Perspective

The drift-diffusion model discussed herein is an example of a convection dominated system. These models emerge in a number of very different applications, including oil reservoir simulation and that of biological contamination. Professor Jim Douglas, Jr. of Purdue University and Professor Mary Wheeler of Rice University have contributed substantially to these latter two application areas, respectively. The Douglas-Russell transport diffusion algorithm (cf. [18]) was developed to handle time dependent systems of this type. The rationale is that longer time steps are permitted in following characteristics; accurate information about the flow rate, determined from the potential equation, is calculated via a mixed method of finite elements in this approach. A careful study of the well-posedness of the flow map is given in [36]. Perhaps due to scale differences, this approach has not had the same effectiveness in semiconductor device modeling as in the above mentioned areas. Nonetheless, as a general approach to convective systems, this program appears promising, particularly when combined with sophisticated techniques such as domain decomposition. Professor Richard Ewing of Texas A&M University is a proponent of large scale codes based upon this set of ideas.

The writer first became aware of the semiconductor drift-diffusion model in 1980, during a visit to Yale University, through the courtesy of Martin Schultz. Shortly afterwards, the essential features of the Gummel map were identified, and later related in [33] to earlier work of

Mock, mentioned previously, and Seidman (cf. [54]). A long term commitment to the study of this model was facilitated during a year leave at AT&T Bell Laboratories in 1982-83, hosted by Donald Rose. This special relationship with Bell Laboratories continued until 1987. Both Rose and William Coughran raised during this time challenging issues concerning the successive approximation properties of the Gummel map. This was analyzed globally in [34], via an estimation of the Lipschitz constant; a local study of Gummel's method in one dimension was carried out by Kerkhoven in [42]. The fundamental singularity hypothesis employed in [34] was verified later in [21]. The convergence of Newton's method was analyzed in [14]. It was, in fact, an attempt to provide a rigorous framework for approximate Newton methods based upon numerical approximation, that led the writer to formulate approaches employing the numerical fixed point map. For an amplification of this discussion, and the associated issues of smoothing or regularization, the reader is referred to [37].

In §3, we developed an approximation theory valid for drift-diffusion steady-state systems, for transport in a field induced by an underlying potential. Although an actual implementation would employ some form of upwinding, the error analysis for piecewise linear finite elements is instructive. In fact, it has been noticed that quadrature rules for such elements can lead to appropriate discretization of Scharfetter-Gummel type in two dimensions (cf. [59]). The fact that the Krasnosel'skii calculus is the proper extension of Babuška's inf-sup theory appears not to have been noticed before.

The applied mathematician who first appeared to grasp the significance of the general hydrodynamic model was the late Farouk Odeh, who passed away recently, on May 3, 1992. Curiously, the writer heard Odeh's presentation during a 1986 SIAM meeting in Boston; at the same meeting, Stanley Osher presented an effective class of shock capturing methods, and the writer was motivated to bring these to bear on the ballistic diode, in a collaboration with Osher and Emad Fatemi. Some theoretical results were obtained in [22] and in [23]. It is the writer's belief that moment methods, based upon the Boltzmann equation, define preferred routes to effective simulation of present conventional devices. The reason for this is the nonlocal dependence of critical parameters such as mobility; it appears that energy dependence is essential to capture significant effects such as velocity overshoot. The latter effect is related to device switching times, for example. Two dimensional simulations have recently been reported in [40] for both the hydrodynamic model, and a recent energy transport model developed by electrical engineers at the University of Illinois at Urbana (cf. [13]). The methods used in [40] depend upon essentially nonoscillatory shock capturing techniques introduced in [28] and developed in [56] and [57]; applications to fluid dynamics were studied in [55].

We shall close this survey with a look toward the future. It seems very likely that quantum devices will play a decisive role in the development of microelectronics. For example, resonant tunneling diodes are already under significant study. There are various reasons for this, including the possibility of multiple state devices, as realized in the so-called $I - V$ curves for a particular device. From the viewpoint of modeling, this entails the replacement of the Boltzmann equation by its quantum mechanical equivalent, involving the Wigner function. This formalism is amenable to derived moment systems as well. The reader is invited to consult [46] for further elaboration.

# References

[1] Ivo Babuška and A.K. Aziz. Survey lectures on the mathematical foundations of the finite element method. In A.K. Aziz, editor, *The Mathematical Foundations of the Finite Element Method with Applications to Partial Differential Equations*, pages 5–359. Academic Press, 1972.

[2] G. Baccarani and M.R. Wordeman. An investigation of steady-state velocity overshoot effects in Si and GaAs devices. *Solid State Electr.*, 28:407–416, 1985.

[3] Randolph E. Bank, William M. Coughran, Wolfgang Fichtner, Eric H. Grosse, Donald J. Rose, and R. Kent Smith. Transient simulation of silicon devices and circuits. *IEEE Transactions Electron Devices*, ED-32(10):1992–2007, October 1985.

[4] Randolph E. Bank, Joseph W. Jerome, and Donald J. Rose. Analytical and numerical aspects of semiconductor device modeling. In R. Glowinski and J. Lions, editors, *Proc. Fifth International Conference on Computing Methods in Applied Science and Engineering*. North Holland, Amsterdam, 1982.

[5] Randolph E. Bank and Donald J. Rose. Global approximate Newton methods. *Numerische Mathematik*, 37:279–295, 1981.

[6] R.E. Bank, editor. *Computational Aspects of VLSI Design with an Emphasis on Semiconductor Device Simulation*, volume 25 of *Lectures in Applied Mathematics*. American Mathematical Society, Providence, R.I., 1990.

[7] R.E. Bank, R. Bulirsch, and K. Merten, editors. *Mathematical Modeling and Simulation of Electrical Circuits and Semiconductor Devices*. Birkhäuser Verlag, Basel, 1990.

[8] R.E. Bank, D.J. Rose, and W. Fichtner. Numerical methods for semiconductor device simulation. *IEEE Transactions Electron Devices*, 30:1031–1041, 1983.

[9] F.J. Blatt. *Physics of Electric Conduction in Solids*. McGraw Hill, New York, 1968.

[10] K. Blotekjaer. Transport equations for electrons in two-valley semiconductors. *IEEE Transactions Electron Devices*, 17:38–47, 1970.

[11] D.M. Caughey and R.E. Thomas. Carrier mobilities in silicon empirically related to doping and field. *Proc. IEEE*, 55:2192–2193, 1967.

[12] C. Cercignani. *The Boltzmann Equation and its Application*. Springer -Verlag, New York, 1987.

[13] D. Chen, E. Kan, U. Ravaioli, K. Hess, and R. Dutton. Steady-state macroscopic transport equations and coefficients for submicron device modeling. *IEEE Transactions Electron Devices*, submitted.

[14] W.M. Coughran and J.W. Jerome. Modular algorithms for transient semiconductor device simulation, Part I: Analysis of the outer iteration. In R.E. Bank, editor, *Computational Aspects of VLSI Design with an Emphasis on Semiconductor Device Simulation*, pages 107–149. American Mathematical Society, Lectures in Applied Mathematics 25, 1990.

[15] P. Degond, F. Guyot-Delaurens, F.J. Mustieles, and F. Nier. Semiconductor modelling via the Boltzmann equation. In R.E. Bank, R. Bulirsch, and K. Merten, editors, *Mathematical Modelling and Simulation of Electrical Circuits and Semiconductor Devices*, pages 153–167. Birkhäuser Verlag, 1990.

[16] A. DeMari. An accurate numerical one-dimensional solution of the P-N junction under arbitrary transient conditions. *Solid-State Electron.*, 11:1021–1053, 1968.

[17] A. DeMari. An accurate numerical steady-state one-dimensional solution of the P-N junction. *Solid-State Electron.*, 11:33–58, 1968.

[18] J. Douglas and T.J. Russell. Numerical methods for convection-dominated diffusion problems based on combining the method of characteristics with finite element or finite difference procedures. *SIAM J. Numer. Anal.*, 19:871–885, 1982.

[19] I. Ekeland and R. Temam. *Convex Analysis and Variational Problems.* North-Holland, Amsterdam, New York, 1976.

[20] E. Fatemi, J. Jerome, and S. Osher. Solution of the hydrodynamic device model using high-order nonoscillatory shock capturing algorithms. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, CAD-10:232–244, 1991.

[21] Irene M. Gamba. Asymptotic behavior at the boundary of a semiconductor device in two dimensions. Technical Report 740, Instituto di Analisi Numerica, University of Pavia, 1989.

[22] Irene M. Gamba. Stationary transonic solutions for a one-dimensional hydrodynamic model for semiconductors. *Communications in Partial Differential Equations*, 17:553–577, 1992.

[23] C.L. Gardner, J.W. Jerome, and D.J. Rose. Numerical methods for the hydrodynamic device model: Subsonic flow. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, CAD-8:501–507, 1989.

[24] D. Gilbarg and N. Trudinger. *Elliptic Partial Differential Equations of Second Order.* Springer -Verlag, New York, 1977.

[25] T.N.E. Greville. Interpolation by generalized spline functions. Technical Report 476, Mathematics Research Center, University of Wisconsin, 1964.

[26] H.K. Gummel. A self-consistent iterative scheme for one-dimensional steady state transistor calculations. *IEEE Transactions Electron Devices*, 11:455–465, 1964.

[27] W. Hänsch and M. Miura-Mattausch. The hot-electron problem in small semiconductor devices. *J. Appl. Physics*, 60:650–656, 1986.

[28] A. Harten, B. Engquist, S. Osher and S. Chakravarthy. Uniformly high order accurate essentially non-oscillatory schemes, III. *J. Comp. Phys.*, 71:231–303, 1987.

[29] K. Hess, J.P. Leburton, and U. Ravaioli, editors. *Computational Electronics.* Kluwer Academic Publishers, Boston, 1991.

[30] Joseph W. Jerome. Uniform approximation by certain generalized spline functions. *J. Approximation Theory*, 7:143–154, 1973.

[31] Joseph W. Jerome. An adaptive Newton algorithm based on numerical inversion: Regularization as postconditioner. *Numerische Mathematik*, 47:123–138, 1985.

[32] Joseph W. Jerome. Approximate Newton methods and homotopy for stationary operator equations. *Constructive Approximation*, 1:271–285, 1985.

[33] Joseph W. Jerome. Consistency of semiconductor modelling: An existence/stability analysis for the stationary van Roosbroeck system. *SIAM J. Appl. Math.*, 45(4):565–590, August 1985.

[34] Joseph W. Jerome. The role of semiconductor device diameter and energy-band bending in convergence of Picard iteration for Gummel's map. *IEEE Transactions Electron Devices*, ED-32(10):2045–2051, October 1985.

[35] Joseph W. Jerome. Evolution systems in semiconductor device modeling: A cyclic uncoupled line analysis for the Gummel map. *Mathematical Methods in the Applied Sciences*, 9(4):455–492, 1987.

[36] Joseph W. Jerome. Convection-dominated nonlinear systems: Analysis of the Douglas-Russell transport diffusion algorithm based on approximate characteristics and invariant regions. *SIAM J. Numer. Anal.*, 25:815–836, 1988.

[37] Joseph W. Jerome. Numerical approximation of PDE system fixed-point maps via Newton's method. *J. Comp. Appl. Math.*, 38:211–230, 1991.

[38] Joseph W. Jerome and Thomas Kerkhoven. A finite element approximation theory for the drift-diffusion semiconductor model. *SIAM J. Numer. Anal.*, 28:403–422, 1991.

[39] Joseph W. Jerome and Thomas Kerkhoven. *Steady State Drift Diffusion Semiconductor Models.* SIAM, 1993.

[40] Joseph W. Jerome and Chi-Wang Shu. Energy models for one-carrier transport in semiconductor devices. In *Proceedings, IMA Workshop on Semiconductors*. Springer-Verlag, 1992.

[41] J.W. Jerome and L.L. Schumaker. Local support bases for a class of spline functions. *J. Approximation Theory*, 16:16–27, 1976.

[42] Thomas Kerkhoven. On the effectiveness of Gummel's method. *SIAM J. Sci. & Stat. Comp.*, 9:48–60, January 1988.

[43] Thomas Kerkhoven and Joseph W. Jerome. $L_\infty$ stability of finite element approximations to elliptic gradient equations. *Numerische Mathematik*, 57:561–575, 1990.

[44] M.A. Krasnosel'skii, G.M. Vainikko, P.P. Zabreiko, Ya.B. Rititskii, and V.Ya. Stetsenko. *Approximate Solution of Operator Equations.* Wolters-Noordhoff, Groningen, 1972.

[45] P.A. Markowich. *The Stationary Semiconductor Device Equations*. Springer -Verlag, Vienna and New York, 1986.

[46] P.A. Markowich, C.A. Ringhofer, and C. Schmeiser. *Semiconductor Equations*. Springer - Verlag, Vienna and New York, 1990.

[47] M. Mock. On equations describing steady-state carrier distributions in a semiconductor device. *Comm. Pure Appl. Math.*, 25:781–792, 1972.

[48] M.S. Mock. *Analysis of Mathematical Models of Semiconductor Devices*. Boole Press, Dublin, 1983.

[49] J.P. Nougier, J. Vaissiere, D. Gasquet, J. Zimmermann, and E. Constant. Determination of the transient regime in semiconductor devices using relaxation time approximations. *J. Appl. Phys.*, 52:825–832, 1981.

[50] W. Van Roosbroeck. Theory of flow of electrons and holes in germanium and other semiconductors. *Bell System Tech. J.*, 29:560–607, 1950.

[51] Isaak Rubenstein. *Electro-Diffusion of Ions*. SIAM Studies in Applied Mathematics, 1990.

[52] M. Rudan and F. Odeh. Multi-dimensional discretization scheme for the hydrodynamic model of semiconductor devices. *COMPEL*, 5:149–183, 1986.

[53] D.L. Scharfetter and H.K. Gummel. Large signal analysis of a silicon Read diode oscillator. *IEEE Transactions Electron Devices*, 16:64–77, 1969.

[54] T. Seidman. Steady state solutions of diffusion reaction systems with electrostatic convection. *Nonlinear Anal.*, 4:623–637, 1980.

[55] C.-W. Shu, G. Erlebacher, T. Zang, D. Whitaker, and S. Osher. High-order ENO schemes applied to two- and three-dimensional compressible flow. *Applied Numer. Math.*, 9:45–71, 1992.

[56] C.-W. Shu and S.J. Osher. Efficient implementation of essentially non-oscillatory shock capturing algorithms. *J. Comp. Phys.*, 83:32–78, 1989.

[57] C.-W. Shu and S.J. Osher. Efficient implementation of essentially non-oscillatory shock capturing schemes, II. *J. Comp. Phys.*, 83:32–78, 1989.

[58] B.G Streetman. *Solid State Electronic Devices*. Prentice-Hall, Englewood Cliffs, NJ, 1980.

[59] G.-L. Tan, X.-L. Yuan, Q.-M. Zhang, W.-H. Tu, and A.-J. Shey. Two-dimensional semiconductor device analysis based on new finite-element discretization employing the S-G scheme. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, CAD-8:468–478, 1989.

[60] K.K. Thornber. Relation of drift velocity to low-field mobility and high-field saturation velocity. *J. Appl. Phys.*, 51:2127–2136, 1980.