



Facoltà di Ingegneria Elettronica

Dottorato di Ricerca in Ingegneria delle Telecomunicazioni e
Microelettronica

XX Ciclo del Corso di Dottorato

A Multiscale Simulation Environment for Electronic and Optoelectronic Devices

Matthias Auf der Maur

Tutor

Prof. Aldo Di Carlo

Coordinator

Prof. Nicola Blefari Melazzi

Acknowledgement

This work would not have been possible without the direct or indirect help of many people. First of all, I would like to thank my supervisor, Prof. Aldo Di Carlo, that he accepted me as a PhD student in his group and for his assistance in scientific and bureaucratic issues. Due to his dedication and unremitting work he is heading a group of international reputation, which to be part of is a great satisfaction. I'm very grateful to my colleagues, especially Michael Povolotskyi, with whom I had many valuable discussions about scientific questions and solved many software problems, Fabio Sacconi and Alessandro Pecchia. They accepted me with friendliness and provided a pleasant company during the last three years. Without their scientific and moral help and their hard work, TIBERCAD would not be what it is now. I have to thank also Gabriele Penazzi and Giuseppe Romano who joined the project last year, the latter especially for several interesting discussions and for thoroughly reading my thesis. I thank all the group members, in particular Andrea Reale, Francesca Brunetti, Thomas Brown, Stefano Penna, Riccardo Riccitelli, Eleonora Petrolati and Stefano Bellocchio for a cordial company.

Very special thanks are due to Marianna. If we hadn't met four years ago I would never have ended up in Rome. Without having her nearby I would not have survived the many abstract submission deadlines. I thank also her parents who accepted me with sympathy and made me feel at home although being far away from my family and friends.

Finally I would like to thank my parents to have given me together with my brothers a happy childhood, a decent education and who supported us in all our different decisions.

Rome, 2008

Matthias Auf der Maur

Contents

Acknowledgement	i
Contents	ii
Abstract	iv
Riassunto	v
1 Introduction	1
1.1 Historical review	1
1.2 The need for multiscale simulation	2
1.3 The multiscale simulation approach in TIBERCAD	5
2 Physical Models	8
2.1 Strain and related phenomena	8
2.1.1 Elasticity theory of heterostructures	9
2.1.2 Strain related effects	11
2.2 Semiclassical particle transport	14
2.2.1 The drift-diffusion model	16
2.2.1.1 Electron and hole transport	21
2.2.1.2 Exciton transport	22
2.2.1.3 Coupling of exciton and electron/hole transport	24
2.3 Heat transport	25
2.4 Quantum mechanical models	28
2.4.1 Envelope function approximation	28
2.4.2 Atomistic models	34
2.4.3 Quantum transport	36
3 Numerical Implementation of the Drift-Diffusion Model	44
3.1 The stationary drift-diffusion equations	44
3.2 Scaling and the choice of the dependent variables	45
3.3 The drift-diffusion equations in finite element formulation	48
3.3.1 The finite element method	49
3.3.2 The drift-diffusion equations in weak form	53

3.3.3	Application of FEM to the drift-diffusion equations	60
3.3.4	Conditioning of the linearized system	67
4	The TIBERCAD Software	72
4.1	Introduction	72
4.2	Software structure	72
4.2.1	Mesh handling	74
4.2.2	Model hierarchy	75
4.3	User interface	75
5	Simulation Examples	81
5.1	Piezoresistivity effects of HEMT structures	81
5.2	The influence of gate tunneling in MOSFETs	85
5.3	GaAs-based pin-diodes for polariton LASER	89
5.4	Structures for polariton LASERS and LEDs based on GaN	93
5.5	AlGaAs/GaAs/AlGaAs Quantum well	98
A	Numerical Evaluation of Terminal Currents	100
B	Implementation of Cylindrical Symmetry	103
C	A comparison principle for quasi-linear elliptic equations	105
	Bibliography	109
	List of Figures	118
	List of Tables	120

Abstract

Driven by the need for high integration density of integrated circuits and high performance of single devices – especially with respect to operation frequency – the dimensions of conventional electronic and optoelectronic devices have been subject to downscaling since the early days of semiconductor technology. The enormous progress in device technology allows for the fabrication of nanometer-scale structures which are of great interest especially for optoelectronic applications. Classical and semi-classical approaches to the simulation of such structures partially or completely break down, on the one hand because the typical device dimensions get comparable to the particle mean free path and on the other hand because atomistic details of the device structure begins to play a fundamental role. Moreover, carrier confinement is a typical effect in nanostructured devices leading to new physical phenomena with vast application possibilities. Besides conventional semiconductor technology molecular-based electronics has gained much attention for nanometric electronics. Transport in organic molecules and carbon nanotubes needs a strict quantum mechanical treatment based on methods with atomistic resolution.

The consideration of quantum-mechanical effects in the simulation of nanoscale devices is essential for a reliable description of structural, electronic and optical properties and particle transport. Several approaches for such a detailed description exist and are widely used. However, they are usually computationally very intensive and therefore restricted to rather small system.

Normally, what we called nanoscale device up to here represents only the active part of an electronic or optoelectronic device. It does not include surrounding parts such as contact access regions, substrates or similar. However the overall device behaviour can be influenced in a non-trivial way by these “non-active” device parts. Therefore a reliable, quantitative simulation has to take them into account. The surroundings can usually be described using semi-classical models. This situation can be handled only by a *multiscale* simulation, that is able to couple self-consistently the scale of semi-classical, continuous media approaches with microscale quantum-mechanical simulations.

The goal of the TIBERCAD project is to provide a multiscale simulation environment which meets the requirements for the simulation of emerging and future devices. It is designed to capture all the important aspects of modern devices such as strain, heat transport and electronic transport on different scales.

Riassunto

La richiesta di dispositivi elettronici di elevate prestazione e di circuiti ad alta integrazione ha condotto ad una continua riduzione delle loro dimensioni sin dall'avvento dell'era della tecnologia dei semiconduttori. Gli enormi progressi tecnologici consentono attualmente di produrre strutture nanometriche di grande interesse in particolare per i dispositivi optoelettronici. Approcci classici e semi-classici alla simulazione di tali strutture non sono adeguati per sistemi di questo tipo. Da un lato le dimensioni tipiche dei dispositivi diventano comparabili con il cammino libero medio dei portatori e dall'altro i dettagli della struttura atomica non sono più trascurabili. Inoltre il confinamento dei portatori, che è un tipico effetto nei dispositivi nanometrici, può essere sfruttato in diversi modi. Oltre alle tecnologie convenzionali basate su semiconduttori anche i dispositivi formati da molecole e nanotubi a carbonio suscitano un elevato interesse. La descrizione del trasporto in tali strutture richiede un trattamento quantistico che tenga conto della struttura atomica.

Per ottenere una simulazione affidabile delle proprietà strutturali, elettroniche ed optoelettroniche dei dispositivi nanometrici è essenziale considerare gli effetti quantistici. Diversi approcci sono stati sviluppati a tale scopo, che però richiedono risorse computazionali eccessive per il calcolo numerico e che pertanto possono essere utilizzati solo per sistemi piccoli.

In genere le strutture nanometriche sono la parte attiva di un dispositivo che comprende contatti elettrici, regioni di accesso, substrato e altro. Tuttavia il comportamento di tutto il dispositivo può essere influenzato in modo non banale da queste parti, di cui una simulazione affidabile deve quindi tener conto e che possono essere descritte con modelli semi-classici. L'approccio corretto è quindi una simulazione *multiscala* che sia in grado di accoppiare modelli semi-classici e modelli quantistici o atomistici in modo autoconsistente.

Lo scopo del progetto TIBERCAD è di fornire un ambiente di simulazione multiscala che soddisfi le esigenze di una simulazione di dispositivi elettronici avanzati. TIBERCAD è progettato per includere su scale diverse gli aspetti più importanti riscontrati in dispositivi moderni come la tensione meccanica, il trasporto di particelle e di calore.

Chapter 1

Introduction

1.1 Historical review

In 1960, one year after the first integrated circuits were demonstrated, the first working MOSFET was realized at Bell Labs by Kahng and Atalla [51]. This was a milestone in semiconductor technology as the MOSFET became the main ingredient for computer technology. Due to the following enormous progress in semiconductor technology, integrated circuits — especially microprocessors and memory — underwent a tremendous increase in performance, for single transistor devices could be made smaller, and more transistors could be integrated on one chip [62]. This growth is described by Moore’s law, stated in 1965 [75] and illustrated in Fig. 1.1

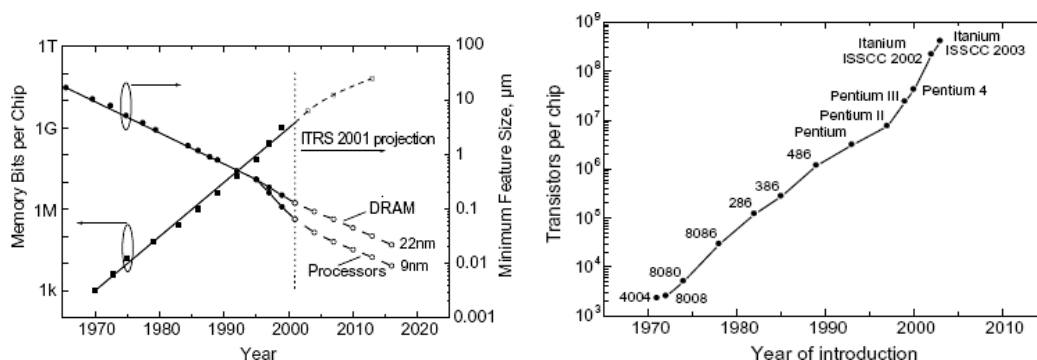


Figure 1.1: Moore’s law: evolution trends for feature size (left) and transistors per chip (right) [62].

The construction of more and more powerful computers provided the opportunity for the implementation of numerical simulation methods of increasing complexity, which in turn aided technological progress.

Modeling of semiconductor devices began in the nineteen-fifties after the formulation of the basic semiconductor equations by Van Roosbroeck [65]. These equations,

which are usually called *drift-diffusion equations* in the engineering and *Van Roosbroeck equations* in the mathematical community, build a system of three coupled partial differential equations: the poisson equation and continuity equations for electrons and holes.

A major breakthrough for the use of the drift-diffusion equations for numerical device simulation is due to Scharfetter and Gummel, which in 1969 proposed a discretisation scheme that overcomes the main numerical problems connected to the stiffness of the semiconductor equations [93]. A rigorous mathematical treatment of the drift-diffusion equations started only in the nineteen-seventies and led to a good understanding of the analytic and numerical properties of the single equations.

In 1966 the Monte Carlo method was applied the first time for the calculation of transport in semiconductors and since then evolved to a powerful and reliable simulation tool [47]. It allows a direct solution of the Boltzmann transport equation (BTE) without the need of a priori assumptions on the particle distribution functions, which in contrary is required when expanding the BTE to get a simpler model such as the drift-diffusion equations. Nevertheless the drift-diffusion model has probably been the most important transport model for the simulation of semiconductor devices. Especially in industrial environments it has been the workhorse for device development and optimisation. This is mainly due to the much lower computational cost of the drift-diffusion model with respect to other approaches.

1.2 The need for multiscale simulation

The downscaling trend of conventional devices as shown in Fig. 1.1 leads to devices with feature sizes in the nanometer scale, for which the assumptions of the models based on a continuous media approximation break down. On the one hand, quantum mechanical behaviour gets important and cannot be neglected anymore, and on the other hand details of the atomistic structure can gain a high impact on device performance. In particular, as devices get smaller bulk properties usually get less important with respect to surface effects, e.g. at the contacts.

There are also emerging new devices where the active parts are based on nanostructures such as nanowires, quantum dots, carbon nanotubes or even molecules. The transport behaviour of such systems cannot be modeled without considering the quantum mechanical properties of the carriers, and in many cases even the use of approximations such as the envelope-function approximation (EFA) get questionable and atomistic approaches have to be used [79, 26].

In a real system the aforementioned structures usually just represent the active part of a device and are therefore not isolated but embedded in a larger environment used e.g. as support or for contacting. As for practical applications the overall device behaviour has to be known, a reliable simulation should consider not only the active part itself, but also its environment.

Quantum mechanical and especially atomistic models require, however, very high computational power and are therefore limited to rather small structures as

molecules, carbon nanotubes or semiconductor structures with symmetries that reduce the dimension of the mathematical description (i.e. which needs a low number of basis functions for the numerical treatment). As such, they can be used to calculate the properties of the small active regions, but they definitely cannot include all the surroundings. The only way to overcome this problem is to perform different simulations on the same device, coupling them together selfconsistently.

Such an approach has been used for a long time in the simulation of MOSFETs and HEMTs to include the effect of quantum confinement of the carriers in the channel (see [104] and references therein). In this case, a stationary 1D Schrödinger-like equation including the Hartree potential is solved on slices perpendicular to the channel and the result is used to describe the carrier distribution along these slices, whereas transport along the channel is described in a classical way.

The approach as described before is well known and widely used, but has a few drawbacks. First, an “adiabatic” behaviour along the channel is assumed in the sense, that the electro-chemical potential is assumed to be slowly varying. Second, transport is calculated classically and therefore does include neither source-drain tunneling nor gate tunneling, which can be both important in short-gate transistors. Third, the approach is mainly useful for structures like MOSFETs or HEMTs or in general when carrier transport in the plane of confinement is of interest. Finally, the approach is usually based on the envelope function approximation which mostly neglects atomic details of the underlying structure. When applied to very small heterostructures this can be inadequate and a truly atomistic description of the active part of the structure is desirable.

The way of selfconsistent coupling between atomistic and continuous models is not evident and is therefore currently an important research topic. Undoubtedly the coupling of nano- and microscale simulations in form of *multiscale* simulations will in future be the preferred approach for the study of electronic and optoelectronic devices. In fact, the multiscale simulation approach is already an established and indispensable tool in fields such as computational materials science (e.g. [41]) or biotechnology and drug design/drug release modeling (e.g. [2]).

Fig. 1.2 shows the different scales that are generally involved in the simulation of electronic and optoelectronic devices. On the highest level, which corresponds to the biggest length scale, a merely architectural point of view is adopted. This is the world of circuit simulation where compact models, simple mathematical models or equivalent circuits built from lumped elements are used for the description of single devices. The model parameters are extracted either from measurements or from lower scale physical device simulations.

The next scale is the scale of microstructures, that is single electronic devices (also MEMS etc. in other engineering fields) which can be characterized based on continuous media approaches, i.e. which do not need a deeper knowledge of atomic details of the structure. This is equivalent to saying that device behaviour is governed by bulk effects rather than surface or interface properties. Examples are hydrodynamic and drift-diffusion models, elasticity theory and fluidodynamics. On this scale quantum mechanical effects do not enter explicitly into the description,

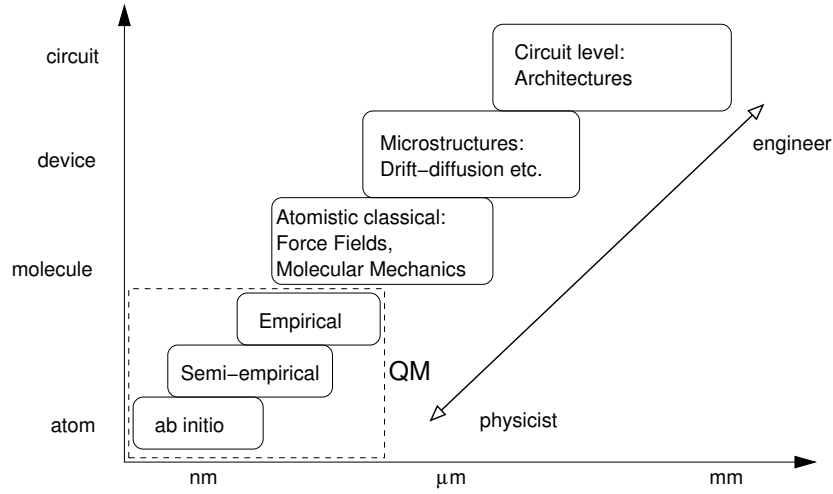


Figure 1.2: The hierarchy of scales important for device simulation.

but they can be present implicitly in model parameters, e.g. in the band parameters of semiconductors.

In the sub-micrometer scale approaches can be used that consider the atomistic structure rather than assuming a continuous media, but still use classical models to describe the properties of the system. An example are valence force methods, where the interatomic forces are parametrized in a classical way [53].

Arriving at the nanometer scale one finally enters the reign of quantum mechanics, where the device behaviour is controlled to a great deal or entirely by quantum effects. Often, approaches such as envelope function or the related effective mass approximations cannot be used or are at least questionable. This is the case especially when dealing with molecules, carbon nanotubes or atomic clusters, but also for nanostructured devices, when the exact atomic structure of e.g. interfaces or surfaces begins to play an important role. The behaviour in the latter case tends to be dominated by surface and interface properties (not least by the interface between contact and active region) rather than by bulk properties. The quantum mechanical approaches can be divided into empirical, semi-empirical and *ab initio* methods, where the computational burden usually grows quickly when moving towards first principles methods.

The device simulators that are mostly used at present emerged during the nineteen-eighties (Silvaco [98]) and nineteen-nineties (DESSIS of ISE-TCAD, acquired by Synopsis and now called Sentaurus Device [101]) and are based on a continuous media simulation approach. Although they incorporate quantum mechanical models in form of quantum-corrections in the semi-classical transport equations and for selfconsistent Schrödinger-Poisson solvers, they do not treat them on an atomistic level, which is clearly a limitation for the simulation of emerging and future nanoelectronic devices.

The TIBERCAD project [85, 7, 9] was launched to provide a simulation environment that adopts the new approach of multiscale simulations to overcome the limitations of the established classic device simulation framework.

It should be noted that a multiscale simulation is usually at the same time also a *multiphysics* simulation. This is due to the evident fact, that the descriptions adopted on different scales are normally connected to different physical or mathematical models. On the contrary, multiphysics simulations are not intrinsically connected to multiscale simulations but nevertheless of crucial importance for a reliable simulation of modern electronic and optoelectronic devices. For this reason TIBERCAD is designed to be a multiphysics *and* multiscale simulation tool, without a priori preferring one aspect over the other, although the multiscale aspect is the more important one from a scientific point of view, especially for future devices.

1.3 The multiscale simulation approach in TIBERCAD

In this section we illustrate the multiscale (and multiphysics) simulation approach of TIBERCAD. The multiscale procedure is visualized schematically in Fig. 1.3. The

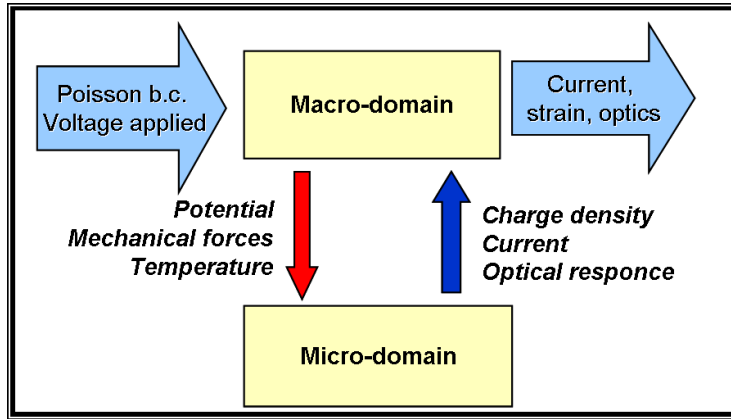


Figure 1.3: Schematic flowchart for a multiscale simulation.

fact of being a multiscale simulation is hidden from the user as far as possible. The user has to define geometry, materials and their physical properties (if not taken from a database) along with appropriate boundary conditions. From this information the simulator constructs the domain to be simulated on different domains and prepares all needed data structures. Normally, the user defined boundary conditions refer to the classical or *macro*-domain, whereas boundary conditions between macro-domain and quantum mechanical or atomistic domains, in the figure denoted by *micro*-domain, are constructed internally.

The solutions of the different macro-domain calculations – assuming the macro-domain problem to be itself a multiphysics problem – represent mean field data such as the (Hartree-) potential, the elastic strain from local deformation or the

local lattice temperature. These data are used in the micro-domain problem for the calculation of the Hamiltonian. The results of the micro-domain calculations in turn represent boundary conditions on the macro-domain (e.g. particle currents) or input data to the corresponding equations such as particle densities or optical responses. The simulations on the two domains have to be solved iteratively to attain a self-consistent global solution of the multiscale/multiphysics problem.

As different models communicate by means of boundary conditions or by point-wise data inside the simulation domain, they have to be able to exchange data. This usually involves an interpolation because the different models are not necessarily solved on the same grid. TIBERCAD implements this aspect in a particular way. The meshes for the different models are all derived from one single *parent mesh* by selective refinement. In this way groups of elements used in different models will have a common parent element, i.e. they belong to the same *element family*. This allows to find data from another model in a given element in an efficient way. Fig. 1.4 illustrates this approach graphically.

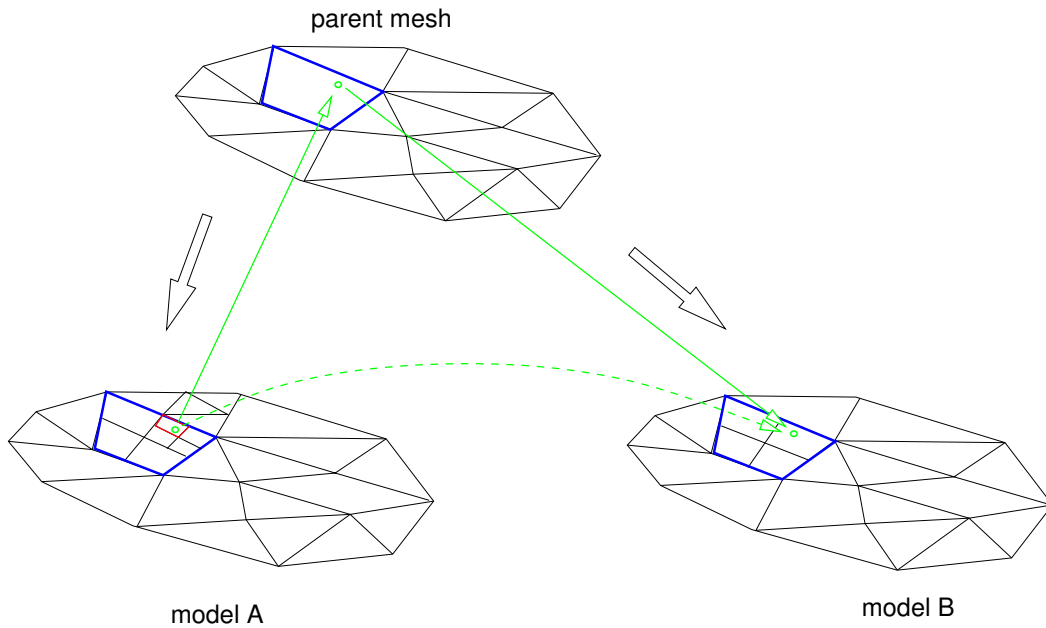


Figure 1.4: Data interpolation between different meshes. Model A can find data of model B in a certain point of an element as it knows that the corresponding element in model B has to be a child of the common parent element.

Whenever possible, the numerical implementations in TIBERCAD are based on the finite element method (FEM). FEM has been widely used especially in structural mechanics, but in the last decades it has been applied to many other engineering fields. Today it is often the method of choice for the numerical treatment of computational problems as it is able to treat complex geometries rather easily and as it

has a sound mathematical basis which allows for well founded approximation and convergence analysis.

Treating all models in the framework of FEM provides a clean and coherent way of handling complex multiphysical problems. As the solutions are expanded in basis functions that are well defined on the whole simulation domain, interpolation of data between mesh points is a well-defined operation. For this reason data exchange between models can be done in a transparent way, without introducing non-quantifiable or obscure sources of errors. Moreover, the possibility of quantifying approximation errors allows for efficient adaptive mesh refinement.

Chapter 2

Physical Models

In this chapter the different physical models implemented in TIBERCAD are described. First, the classical and semi-classical continuum mechanical approaches for the calculation of strain (2.1), particle transport (2.2) and heat transport (2.3), and finally the quantum mechanical models (2.4) are presented. The latter are subdivided into envelope function approximations and atomistic methods. All models that are based on partial differential equations are discretised using the finite element method (FEM) which leads to a consistent description of a device across the different models. The numerical implementation of the drift-diffusion equations derived in sec. 2.2 will be treated in more detail in chapter 3 as this represents the main part of this work.

2.1 Strain and related phenomena

Mechanical strain is present in systems in which external or internal forces lead to a mechanical deformation, i.e. to a change of the interatomic distances of the constituent materials. This is the case in a homogeneous material under an external pressure, but also in a lattice mismatched heterostructure, where the substrate material imposes its lattice constant to the material grown on top of it.

Strain got an important issue in semiconductor devices due to the tremendous progress in semiconductor technology, especially growth technology. This seems paradoxical, and the reason is that only thanks to very high quality growing techniques it is possible to grow strained heterostructures that can be used to build devices with. Consequently one began to take advantage of the properties of strained semiconductor devices [87], and therefore the treatment of strain in the simulation of these systems is of crucial importance.

Strain can be calculated either atomistically or using continuum mechanics. Using an atomistic approach we would have to write down the Hamiltonian H of the system and minimizing the total energy we could find the equilibrium positions \mathbf{R}_k of the ions. Doing this without too much simplifications is computationally very expensive and feasible only for systems with a small number of atoms, e.g. molecules.

Systems with millions of atoms can be treated by using empirical approaches as e.g. *valence force methods* where the inter-ionic forces are parametrised and the other terms in the Hamiltonian are discarded [53]. This is still a computationally heavy task and applicable only to nanostructures.

For the study of “big” micrometer scale structures a continuous media approach has to be used. The implementation in TIBERCAD is based on elasticity theory, assuming cristallographically perfect heterointerfaces (coherent growth) and defect-free structures [84, 83, 82, 60].¹

2.1.1 Elasticity theory of heterostructures

Let us assume a structure that is deformed due to some force such that each point \mathbf{r} moves to a new point \mathbf{r}' . The displacement \mathbf{u} of \mathbf{r} is given by

$$\mathbf{u} = \mathbf{r}' - \mathbf{r}, \quad \text{or} \quad u_i = x_i' - x_i \quad (2.1)$$

and completely characterises the deformation of the system. Generally the deformation is not homogeneous and therefore the displacement is position dependent: $\mathbf{u} = \mathbf{u}(\mathbf{r})$. The distance between two infinitesimally adjacent points in the deformed system reads

$$dl'^2 = (d\mathbf{r} + d\mathbf{u})^2 \quad (2.2)$$

$d\mathbf{u}$ can be written as (in components) $du_i = (\partial u_i / \partial x_k) dx_k$ and with this the above distance gets

$$\begin{aligned} dl'^2 &= dl^2 + 2 \frac{\partial u_i}{\partial x_k} dx_i dx_k + \frac{\partial u_l}{\partial x_k} \frac{\partial u_l}{\partial x_i} dx_k dx_i \\ &= dl^2 + \left(\frac{\partial u_i}{\partial x_k} + \frac{\partial u_k}{\partial x_i} \right) dx_i dx_k + \frac{\partial u_l}{\partial x_k} \frac{\partial u_l}{\partial x_i} dx_i dx_k \end{aligned} \quad (2.3)$$

where dl is the distance before the deformation.² Eq. (2.3) can be rewritten as

$$dl'^2 = dl^2 + 2\varepsilon_{ik} dx_i dx_k \quad (2.4)$$

The symmetric tensor ε_{ik} in the last expression is called *strain tensor*. For small deformations where we can neglect terms quadratic in $\partial u_i / \partial x_k$ it is given by

$$\varepsilon_{ik} = \frac{1}{2} \left(\frac{\partial u_i}{\partial x_k} + \frac{\partial u_k}{\partial x_i} \right) \quad (2.5)$$

Its diagonal components describe stretching or shrinking, and the off-diagonal ones are connected to shear deformations. The trace of ε gives the relative volumic change due to the deformation.

¹The software implementation of this model was done by Dr. M. Povolotskyi.

²Here and in the following we use the Einstein summing convention (without differentiating co- and contravariant indices), i.e. we sum over indices appearing twice in the same expression.

In a strained structure internal forces try to restore locally the interatomic distances to the equilibrium values. These forces can be written as divergence of a second-rank tensor σ_{ik} :

$$F_i = \frac{\partial \sigma_{ik}}{\partial x_k} \quad (2.6)$$

With this the i -th component of the force F acting on a small volume V can be written as

$$\int_V F_i dV = \int_V \frac{\partial \sigma_{ik}}{\partial x_k} dV = \int_{\partial V} \sigma_{ik} ds_k \quad (2.7)$$

where s_k are the components of the surface element. $\sigma_{ik} ds_k$ is therefore the force acting on the surface element ds_k . The tensor σ_{ik} is called *stress tensor*. In equilibrium the internal forces in each small volume V have to compensate such that $F_i = 0$, i.e.

$$\frac{\partial \sigma_{ik}}{\partial x_k} = 0 \quad (2.8)$$

The elastic energy of a deformed system is given by [60]

$$E = \frac{1}{2} \int_V \sigma_{ik} \varepsilon_{ik} dV \quad (2.9)$$

The equilibrium state can be found by minimizing E which is equivalent to (2.8).

As we are interested in small deformations we can use Hooke's law which linearly relates strain to stress:

$$\sigma_{ik} = C_{iklm} \varepsilon_{lm} \quad (2.10)$$

C_{iklm} is called *elasticity tensor*. It is symmetric in ik and lm and with respect to the interchange of the two index pairs. Therefore it can have at most 21 different components, called *elasticity moduli*. The number of independent components is given by the symmetry of the crystal. The for semiconductor devices most important crystalline structures have cubic or hexagonal symmetry. In cubic systems there are only three different elasticity moduli, whereas in hexagonal systems there are five of them.

The system of equations to be solved to get the deformation and strain of semiconductor structures finally reads

$$\frac{\partial}{\partial x_k} (C_{iklm} \varepsilon_{lm}) = \frac{1}{2} \frac{\partial}{\partial x_k} \left[C_{iklm} \left(\frac{\partial u_l}{\partial x_m} + \frac{\partial u_m}{\partial x_l} \right) \right] = f_i \quad (2.11)$$

where f_i is an externally applied mechanical force.

For the numerical calculation we proceed in the following way [82]:

1. Define a simulation mesh covering the simulation domain and its coordinate system. If the axes of the simulation system do not coincide with the crystallographic axes of the constituent materials, calculate the corresponding rotation matrices.

2. Define a reference lattice. The lattice of one of the constituent materials, usually the substrate material, is used. We can then define a *lattice-matching strain* ε_{ij}^0 which is caused by the deformation of the unit cells of the unstrained materials needed to match the reference lattice.

In order to define the reference lattice constants, a conventional unit cell of a minimal size has to be chosen with its faces parallel to the heterointerfaces. Then the reference lattice constants are the lengths of translation vectors of this conventional cell. In the case of hexagonal crystals grown in [0001]-direction and with the axes chosen as $x \parallel [10\bar{1}0]$, $y \parallel [\bar{1}2\bar{1}0]$, $z \parallel [0001]$, the non-zero components of ε_{ij}^0 become: $\varepsilon_{xx}^0 = \varepsilon_{yy}^0 = (a_0 - a)/a$, $\varepsilon_{zz}^0 = (c_0 - c)/c$, where a and c are the lattice constants along the $[10\bar{1}0]$ and $[0001]$ directions, respectively. In a zincblende system grown along $[001]$ we get $\varepsilon_{ij}^0 = \delta_{ij}(a_0 - a)/a$.

3. The total strain tensor is now given by

$$\varepsilon_{ij}(\mathbf{r}) = \tilde{\varepsilon}_{ij}(\mathbf{r}) + \varepsilon_{ij}^0(\mathbf{r}) \quad (2.12)$$

where $\tilde{\varepsilon}_{ij}(\mathbf{r})$ is the strain according to eq. (2.5) due to the displacement $\mathbf{u}(\mathbf{r})$ with respect to the reference lattice.

When applying the strain $\varepsilon_{ij}^0(\mathbf{r})$ to the system, all heterointerfaces will be perfectly matched, but the structure will not be in equilibrium. The equilibrium can be found by minimizing the elastic energy, maintaining the interface matching.

Following the procedure as delineated before, we rewrite eq. (2.11) in the following way, and using the symmetry $C_{ijkl} = C_{ijlk}$ we finally get

$$\frac{\partial}{\partial x_i} C_{ijkl}(\mathbf{r}) \frac{\partial u_k(\mathbf{r})}{\partial x_l} = - \frac{\partial}{\partial x_i} C_{ijkl}(\mathbf{r}) \varepsilon_{ij}^0(\mathbf{r}) + f_i \quad (2.13)$$

which has to be solved under appropriate boundary conditions.

In order to get the equilibrium shape of a strained structure, the deformation $\mathbf{u}(\mathbf{r})$ has to be applied to the discretization mesh, and then the strain in the deformed system has to be recalculated, leading to a new deformation. After a few iterations the final shape can be obtained. The algorithm is illustrated in Fig. 2.1.

Fig. 2.2 shows an example of the application of elasticity theory to a freestanding GaN/AlGaN heterostructure as shown in part (a) of the figure. Due to the fact that AlGaN has a smaller lattice constant than GaN, the structure deforms such as to build a sort of bowl as shown in (b).

2.1.2 Strain related effects

Strain has important effects on the electrical and optical behaviour of semiconductor devices.

On the one hand it is clear that strain provokes a change of the band structure due to the deformation of the primitive cells. In particular this results in a splitting

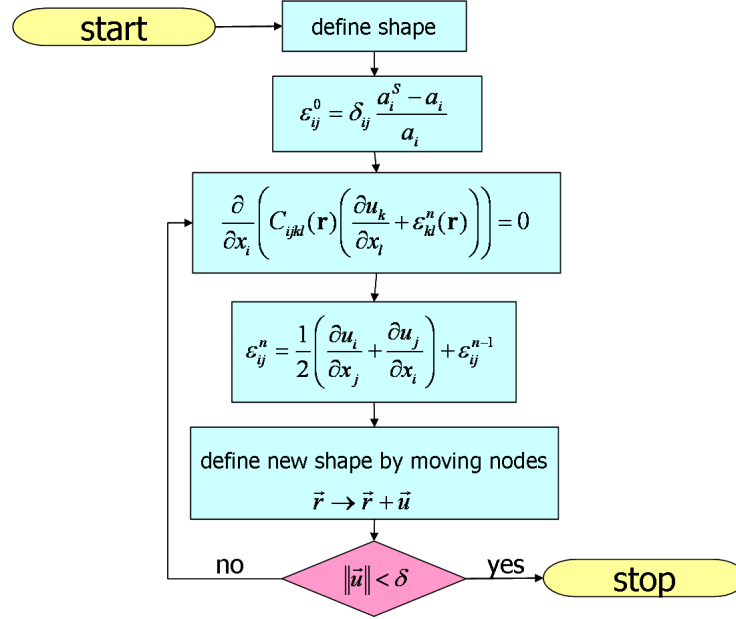


Figure 2.1: Iterative procedure for the calculation of the deformed equilibrium shape.

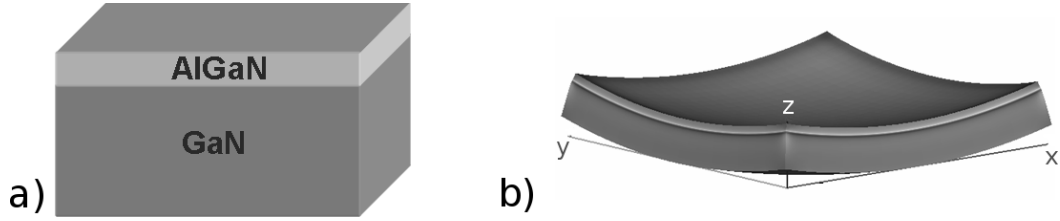


Figure 2.2: Deformed free standing GaN/AlGaN heterostructure. The structure is indicated schematically in (a), (b) shows the convex final shape.

of the degenerate hole bands, in a change of the band extrema and therefore of the band gap and in changes of the quasi-particle masses. These together locally change the effective density of states and result in space dependent band gaps.

On the other hand strain can change the charge distribution inside the conventional cell and lead to a macroscopic electric polarization \mathbf{P} . This is called *piezoelectric effect*.

The band structure effects will be described in section 2.4. Only the piezoelectric effect shall be treated here.

The electric displacement \mathbf{D} in a piezoelectric crystal can be written in general as [59]

$$D_i = D_i^{(0)} + \kappa_{ik} E_k + e_{ikl} \varepsilon_{kl} \quad (2.14)$$

To avoid confusion with the strain tensor we denote the permittivity tensor with

the symbol κ_{ik} . The term $D_i^{(0)}$ represents a spontaneous electric polarization. It is called *pyroelectric polarization*, and we will identify it with the symbol \mathbf{P}^{py} . This effect can only be found in crystals with certain symmetries as the properties of the crystal have to remain unchanged under the symmetry transformation. In particular, as the polarization vector is a polar vector, it changes sign under a parity transformation \mathcal{P} : $\mathcal{P}(\mathbf{P}) = -\mathbf{P}$. Therefore only crystals which break parity symmetry can have a spontaneous electric polarization. This is the case for semiconductors with wurtzite crystal structure, e.g. nitride based materials (GaN, AlN, InN and their alloys). The pyropolarization in this case points along the symmetry axis c as illustrated in Fig. 2.3. The pyroelectric polarization is strongly temperature

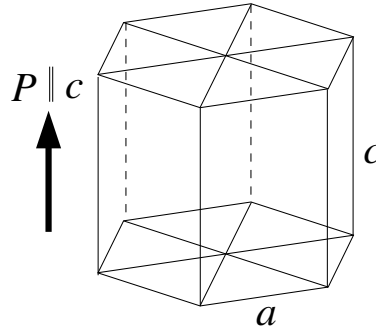


Figure 2.3: Pyropolarization in a hexagonal lattice.

dependent.

The second term on the right hand side of eq. (2.14) is the usual expression for the displacement in function of the electric field, with the difference that the permittivity is written as a tensor as it can be anisotropic.

The third term is called *piezoelectric polarization* \mathbf{P}^{pz} . The third-rank tensor e_{ikl} is the piezoelectric tensor. It is symmetric in the index pair kl and its components depend on the symmetry of the crystal. In a cubic crystal e.g. it has only one independent non-zero component associated with the off-diagonal components of the strain (representing shear deformation), and the piezoelectric polarization can be written as

$$\mathbf{P}^{pz} = 2e_{xyz} \begin{pmatrix} \varepsilon_{yz} \\ \varepsilon_{xz} \\ \varepsilon_{xy} \end{pmatrix} \quad (2.15)$$

In a material with the crystal structure of wurtzite e_{ikl} has the three independent non-zero components, denoted by e_{15} , e_{31} and e_{33} and the piezoelectric polarization reads (using the Voigt-notation for the indices, which contracts the index pair kl into a single index [42])

$$\mathbf{P}^{pz} = \begin{pmatrix} 2e_{15}\varepsilon_{xz} \\ 2e_{15}\varepsilon_{yz} \\ e_{31}\varepsilon_{xx} + e_{31}\varepsilon_{yy} + e_{33}\varepsilon_{zz} \end{pmatrix} \quad (2.16)$$

2.2 Semiclassical particle transport

Semiclassical particle transport is usually based on the Boltzmann transport equation (BTE), formulated first by Boltzmann in 1872:

$$\left[\frac{\partial}{\partial t} + \mathbf{v} \cdot \nabla_{\mathbf{r}} + \frac{1}{m} \mathbf{F}_{eff} \cdot \nabla_{\mathbf{v}} \right] f(\mathbf{r}, \mathbf{v}, t) = \left(\frac{\partial f}{\partial t} \right)_{coll} \quad (2.17)$$

where \mathbf{v} means the particle velocity, $\mathbf{F}_{eff} = m d\mathbf{v}/dt$ is an effective force acting on the particles, m is the effective mass of the particles and the term on the right hand side of the equation is the collision integral which describes the scattering between different states. The effective force can also be written as gradient of an effective potential $\mathbf{F}_{eff} = -\nabla_{\mathbf{r}} U_{eff}$. $f(\mathbf{r}, \mathbf{v}, t)$ represents the particle density in phase space such that a volume element $d^3\mathbf{r} d^3\mathbf{v}$ at (\mathbf{r}, \mathbf{v}) contains $f(\mathbf{r}, \mathbf{v}, t) d^3\mathbf{r} d^3\mathbf{v}$ particles. The forces responsible for particle scattering shall be assumed to be short ranged and the scattering events to take place on a short time scale, in compatibility with the picture of a classical “collision”.

Note that the left hand side of eq. (2.17) is just the expansion of $df(\mathbf{r}, \mathbf{v}, t)/dt$ (without the collision term). This rate of change of the particle density can be interpreted as convection due to the effective force \mathbf{F}_{eff} so that (2.17) can be seen as a balance between the rates of change due to convection and collisions [66]:

$$\left(\frac{\partial f}{\partial t} \right)_{conv} = \left(\frac{\partial f}{\partial t} \right)_{coll} \quad (2.18)$$

The Boltzmann transport equation can be derived either by a phenomenological approach [5, 66], or by a Green’s function description of transport [50].

An expression for the collision integral can be found as follows (assuming particles following the Fermi-Dirac statistics). Let $P(\mathbf{r}, \mathbf{v}' \rightarrow \mathbf{v}, t)$ denote the rate of scattering of a particle at position \mathbf{r} from a velocity \mathbf{v}' to the velocity \mathbf{v} . We assume this rate to be proportional to the occupation of the initial state $f(\mathbf{r}, \mathbf{v}', t)$ and to the number of available states $1 - f(\mathbf{r}, \mathbf{v}, t)$:

$$P(\mathbf{r}, \mathbf{v}' \rightarrow \mathbf{v}, t) = S(\mathbf{r}, \mathbf{v}', \mathbf{v}) f(\mathbf{r}, \mathbf{v}', t) [1 - f(\mathbf{r}, \mathbf{v}, t)] \quad (2.19)$$

where $S(\mathbf{r}, \mathbf{v}', \mathbf{v})$ is the scattering rate from $(\mathbf{x}, \mathbf{v}')$ to (\mathbf{x}, \mathbf{v}) . The total inscattering for a state $(\mathbf{x}, \mathbf{v}, t)$ is then the integral of the above equation over all possible initial states $(\mathbf{x}, \mathbf{v}')$. Analogously a total outscattering rate can be calculated and we finally get for the collision integral

$$\begin{aligned} \left(\frac{\partial f}{\partial t} \right)_{coll} &= \int [P(\mathbf{r}, \mathbf{v}' \rightarrow \mathbf{v}, t) - P(\mathbf{r}, \mathbf{v} \rightarrow \mathbf{v}', t)] d\mathbf{v}' \\ &= \int [S(\mathbf{r}, \mathbf{v}', \mathbf{v}) f'(1 - f) - S(\mathbf{r}, \mathbf{v}, \mathbf{v}') f(1 - f')] d\mathbf{v}' \end{aligned} \quad (2.20)$$

where we abbreviated $f = f(\mathbf{r}, \mathbf{v}, t)$ and $f' = f(\mathbf{r}, \mathbf{v}', t)$. The scattering rates S have to be calculated quantum mechanically for each scattering mechanism.

Several interesting (and measurable) microscopic quantities can be calculated from the moments $\int \mathbf{v}^m f \, d\mathbf{v}$ of the distribution function $f(\mathbf{r}, \mathbf{v}, t)$. The real space carrier density $n(\mathbf{r}, t)$ is given by the zeroth moment

$$n(\mathbf{r}, t) = \int f(\mathbf{r}, \mathbf{v}, t) \, d^3\mathbf{v} \quad (2.21)$$

whereas the particle flux $\mathbf{j}(\mathbf{r}, t)$ is found from the first moment

$$\mathbf{j}(\mathbf{r}, t) = \int \mathbf{v} f(\mathbf{r}, \mathbf{v}, t) \, d^3\mathbf{v} \quad (2.22)$$

The moments of even order m can be considered generalized densities (particles, energy etc.) and the ones of odd order $m + 1$ the corresponding fluxes. Calculating the moments of the Boltzmann equation itself, one can get conservation equations for the corresponding moments of the distribution function. This will be shown in section 2.2.1.

The validity of the Boltzmann equation is connected to the following assumptions:

- All scattering events are assumed to be local and instantaneous.
- Carrier-carrier interaction can be neglected.
- The potential $U_{eff}(\mathbf{r})$ is assumed to vary slowly in \mathbf{r} with respect to the extension of a particle wave packet.

Under certain assumptions the collision integral (2.20) can be simplified. In the *low density approximation* we can use assume that $f \ll 1$ and therefore the collision integral becomes a linear operator. In the *relaxation time approximation* we presume a small driving force $-\nabla U_{eff}$ such that the distribution f can be linearised as $f = f_{eq} + f^{(1)}$, where f_{eq} denotes the equilibrium distribution. Then we can find for the collision integral

$$\left(\frac{\partial f}{\partial t} \right)_{coll} = -\frac{f - f_{eq}}{\tau(\mathbf{r}, \mathbf{v})} \quad (2.23)$$

with the relaxation time given by

$$\tau(\mathbf{r}, \mathbf{v})^{-1} = \int S(\mathbf{r}, \mathbf{v}, \mathbf{v}') \, d\mathbf{v}' \quad (2.24)$$

For further simplification $\tau(\mathbf{r}, \mathbf{v})$ can be substituted by a constant τ .

To study transport of quasi-particles in a semiconductor, e.g. electrons or holes, the quantum effects of the crystal lattice have to be included in the Boltzmann equation. In this case we write the distribution f in terms of the crystal momentum \mathbf{k} , $f = f(\mathbf{r}, \mathbf{k}, t)$, and the Boltzmann equation gets

$$\left[\frac{\partial}{\partial t} + \mathbf{v}(\mathbf{k}) \cdot \nabla_{\mathbf{r}} + \frac{1}{\hbar} \mathbf{F}_{eff} \cdot \nabla_{\mathbf{k}} \right] f(\mathbf{r}, \mathbf{k}, t) = \left(\frac{\partial f}{\partial t} \right)_{coll} \quad (2.25)$$

where $\mathbf{v}(\mathbf{k}) = d\mathbf{r}/dt$ denotes the group velocity of the particles and the change of momentum is given by $\hbar d\mathbf{k}/dt = \mathbf{F}_{eff}$. Using the dispersion relation $E_{eff}(\mathbf{r}, \mathbf{k})$ for the particle under consideration we can write

$$d\mathbf{k}/dt = \mathbf{F}_{eff} = -\frac{1}{\hbar} \nabla_{\mathbf{r}} E_{eff}(\mathbf{r}, \mathbf{k}) \quad (2.26a)$$

$$d\mathbf{r}/dt = \mathbf{v}(\mathbf{k}) = \frac{1}{\hbar} \nabla_{\mathbf{k}} E_{eff}(\mathbf{r}, \mathbf{k}) \quad (2.26b)$$

Some remarks regarding the effective energies U_{eff} (or E_{eff}) used in the above expressions are appropriate at this point. Although not stated explicitly, it is clear that the Boltzmann equation is a single particle equation, and as such it does not inherently take into account particle interactions (apart from short ranged collisions). Let us consider a system of charged particles at low density in an external potential U_{ext} , assuming no collisions. If we would take $U_{eff} = U_{ext}$ in this case, we would completely neglect the Coulomb interaction, which is a long ranged force. In first order the latter can be accounted for by including in $U_{eff}(\mathbf{r})$ the potential generated by all other particles located in positions \mathbf{r}' . Given the real space particle density $n(\mathbf{r})$ from eq. (2.21) this means

$$U_{eff}(\mathbf{r}) = U_{ext}(\mathbf{r}) + \int U_{p-p}(\mathbf{r}, \mathbf{r}') n(\mathbf{r}') d\mathbf{r}' \quad (2.27)$$

where $U_{p-p}(\mathbf{r}, \mathbf{r}')$ is the potential generated at \mathbf{r} by a point source located at \mathbf{r}' . This phenomenological argumentation is compatible with the Hartree approximation in quantum mechanics, and the same result can also be obtained when deriving the Boltzmann equation starting from the Liouville equation and using the Bogoliubov-Born-Green-Kirkwood-Yvon (BBGKY) hierarchy [66]. The important point here to keep in mind is that for charged particles the Boltzmann and Poisson equation intrinsically belong together.

Several approaches exist to solve the integro-differential equation (2.17) (or, equivalently, (2.25)). The most successful direct solution method is based on the Monte Carlo method [47]. Simulations based on this method however are very time consuming. Faster methods can be constructed by formulating kinetic equations for some mean values of the BTE, reducing thereby the number of dependent variables and resulting in a set of pure partial differential equations. The procedure for this approach is drafted in the next section.

2.2.1 The drift-diffusion model

To derive a simpler, computationally less demanding transport model from the BTE described in the last section, usually the *method of moments* [66, 96, 94] is used, although other methods are described in literature [66]. Here we shall follow the derivation exposed in Ref. [94].

The basic idea of the method of moments is to derive kinetic equations for the mean values with respect to the particle velocity of a certain function $\Phi(\mathbf{v})$:

$$\langle \Phi \rangle = \frac{\int \Phi f(\mathbf{r}, \mathbf{v}, t) d\mathbf{v}}{\int f(\mathbf{r}, \mathbf{v}, t) d\mathbf{v}} = \frac{1}{n(\mathbf{r}, t)} \int \Phi f(\mathbf{r}, \mathbf{v}, t) d\mathbf{v} \quad (2.28)$$

where we used eq. (2.21). Obviously $\langle \Phi \rangle$ will be a function of \mathbf{r} and t .

Consider now the partial derivative of the product of $\langle \Phi \rangle$ and $n(\mathbf{r}, t)$ with respect to the time

$$\frac{\partial}{\partial t} (n \langle \Phi \rangle) = \frac{\partial}{\partial t} \int \Phi(\mathbf{v}) f(\mathbf{r}, \mathbf{v}, t) d\mathbf{v} = \int \Phi(\mathbf{v}) \frac{\partial}{\partial t} f(\mathbf{r}, \mathbf{v}, t) d\mathbf{v} \quad (2.29)$$

The partial derivative of the distribution function can be replaced using the Boltzmann equation (2.17), leading to

$$\begin{aligned} \frac{\partial}{\partial t} (n \langle \Phi \rangle) &= \int \Phi \left[-\mathbf{v} \cdot \nabla_{\mathbf{r}} f - \frac{\mathbf{F}}{m} \cdot \nabla_{\mathbf{v}} f + \left(\frac{\partial f}{\partial t} \right)_{coll} \right] d\mathbf{v} \\ &= -\nabla_{\mathbf{r}} \int \Phi \mathbf{v} f d\mathbf{v} - \frac{\mathbf{F}}{m} \int \Phi \nabla_{\mathbf{v}} f d\mathbf{v} + \left(\frac{\partial}{\partial t} \int \Phi f d\mathbf{v} \right)_{coll} \end{aligned} \quad (2.30)$$

To write this last equation in terms of mean values we have to eliminate all derivatives of f inside the integrals by partial integration (using $\lim_{|\mathbf{v}| \rightarrow \infty} f = 0$). The resulting expression reads

$$\frac{\partial}{\partial t} (n \langle \Phi \rangle) + \nabla_{\mathbf{r}} (n \langle \mathbf{v} \Phi \rangle) - n \frac{\mathbf{F}}{m} \langle \nabla_{\mathbf{v}} \Phi \rangle = \left(\frac{\partial}{\partial t} n \langle \Phi \rangle \right)_{coll} \quad (2.31)$$

The last equation reads in symbolical form

$$\frac{\partial}{\partial t} (n \langle \Phi \rangle) + \nabla_{\mathbf{r}} \mathbf{j}_{\langle \Phi \rangle} - n F_{\langle \Phi \rangle} = \left(\frac{\partial}{\partial t} n \langle \Phi \rangle \right)_{coll} \quad (2.32)$$

allowing the interpretation as a conservation law for the *generalized density* $n \langle \Phi \rangle$. $\mathbf{j}_{\langle \Phi \rangle} = n \langle \mathbf{v} \Phi \rangle$ is the associated *generalized flux* and $F_{\langle \Phi \rangle} = \mathbf{F} \langle \nabla_{\mathbf{v}} \Phi \rangle / m$ a *generalized driving force*.

For the special choice $\Phi(\mathbf{v}) = \mathbf{v}^m$ the equation (2.31) corresponds to the moment of order m of the Boltzmann equation, and $n \langle \Phi \rangle$ is the m -th moment of the distribution function f . Note that the moment of order m in this case is a tensor of rank m , which can easily be seen using index notation:

$$\Phi_{i_1 i_2 \dots i_m}^{(m)}(\mathbf{v}) = v_{i_1} v_{i_2} \dots v_{i_m} \quad (2.33)$$

The set of equations for the moments $m = 1 \dots \infty$ build an infinite hierarchy of conservation laws, where the equation for the m -th moment depends on the moment of order $m + 1$. The lowest order moments are the particle density ($m = 0$), the particle flux ($m = 1$), the energy density ($m = 2$) and the energy flux ($m = 3$). In the following the expressions for the first few moments will be explicitly calculated, writing $\Phi^{(0)} = 1$, $\Phi^{(1)} = \mathbf{v}$, $\Phi^{(2)} = \mathbf{v} \otimes \mathbf{v}$ (\otimes being the tensor product).

$$\boxed{m = 0}$$

The zeroth moment is given by $n\langle\Phi^{(0)}\rangle = n\langle\mathbf{v}^0\rangle = n$. The corresponding flux (which turns out to be the first moment) and driving force are $\mathbf{j}_{\langle\Phi^{(0)}\rangle} = n\langle\Phi^{(1)}\rangle = n\langle\mathbf{v}\rangle =: \mathbf{j}_n$ and $F_{\langle\Phi\rangle} = 0$, respectively. Equation (2.32) leads to the well known continuity equation for the particle density n

$$\frac{\partial n}{\partial t} + \nabla_{\mathbf{r}} \mathbf{j}_n = \left(\frac{\partial n}{\partial t} \right)_{coll} \quad (2.34)$$

\mathbf{j}_n has to be calculated from the first moment of the BTE.

$$\boxed{m = 1}$$

The first moment is defined by $n\langle\Phi^{(1)}\rangle = n\langle\mathbf{v}\rangle =: \mathbf{j}_n$. The corresponding flux (connected to the second moment) and the driving force are $\mathbf{j}_{\langle\Phi^{(1)}\rangle} = n\langle\Phi^{(2)}\rangle = n\langle\mathbf{v} \otimes \mathbf{v}\rangle$ and $F_{\langle\Phi\rangle} = \mathbf{F}/m$, respectively. The velocity can be written as the sum of its mean value and the deviation from it, i.e.

$$\mathbf{v} = \langle\mathbf{v}\rangle + \delta\mathbf{v} \quad (2.35)$$

with $\langle\delta\mathbf{v}\rangle = 0$. $\langle\mathbf{v}\rangle$ can be interpreted as mean drift velocity and $\delta\mathbf{v}$ as its statistical (thermal) fluctuation. Substituting we get

$$\begin{aligned} \mathbf{j}_{\langle\Phi^{(1)}\rangle} &= n\langle\mathbf{v} \otimes \mathbf{v}\rangle = n\langle(\langle\mathbf{v}\rangle + \delta\mathbf{v}) \otimes (\langle\mathbf{v}\rangle + \delta\mathbf{v})\rangle \\ &= n\langle\mathbf{v}\rangle \otimes \langle\mathbf{v}\rangle + n\langle\delta\mathbf{v} \otimes \delta\mathbf{v}\rangle \end{aligned} \quad (2.36)$$

The divergence of \mathbf{j} can be written in the following form

$$\nabla_{\mathbf{r}} \mathbf{j}_{\langle\Phi\rangle} = \langle\mathbf{v}\rangle \nabla_{\mathbf{r}} (n\langle\mathbf{v}\rangle) + (n\langle\mathbf{v}\rangle \cdot \nabla_{\mathbf{r}}) \langle\mathbf{v}\rangle + \nabla_{\mathbf{r}} (n\langle\delta\mathbf{v} \otimes \delta\mathbf{v}\rangle) \quad (2.37)$$

To get the last equation we used

$$\begin{aligned} \underbrace{\nabla_{\mathbf{r}} (n\mathbf{v} \otimes \mathbf{v})}_{\partial_j (nv_i v_j)} &= (nv_j \partial_j) v_i + v_i \partial_j (nv_j) \\ &= \underbrace{(n\mathbf{v} \cdot \nabla_{\mathbf{r}}) \mathbf{v} + \mathbf{v} \nabla_{\mathbf{r}} \cdot (n\mathbf{v})}_{\partial_j (nv_i v_j)} \end{aligned} \quad (2.38)$$

With eq. (2.37) the first moment of the boltzmann equation becomes

$$\frac{\partial}{\partial t} \mathbf{j}_n + \langle\mathbf{v}\rangle \nabla_{\mathbf{r}} \cdot \mathbf{j}_n + (\mathbf{j}_n \cdot \nabla_{\mathbf{r}}) \langle\mathbf{v}\rangle + \nabla_{\mathbf{r}} (n\langle\delta\mathbf{v} \otimes \delta\mathbf{v}\rangle) - n \frac{\mathbf{F}}{m} = \left(\frac{\partial \mathbf{j}_n}{\partial t} \right)_{coll} \quad (2.39)$$

and represents the conservation law for the particle flux. The fourth term on the left hand side is abbreviated by defining a mean particle temperature tensor \mathbf{T} from the mean kinetic energy of the fluctuation $\delta\mathbf{v}$

$$\frac{1}{2} k_B \mathbf{T} = \frac{m}{2} \langle\delta\mathbf{v} \otimes \delta\mathbf{v}\rangle \quad (2.40)$$

such that

$$\nabla_{\mathbf{r}}(n\langle\delta\mathbf{v} \otimes \delta\mathbf{v}\rangle) = \frac{1}{m}\nabla_{\mathbf{r}} \cdot (nk_B\mathbf{T}) \quad (2.41)$$

We truncate the hierarchy at this point as we will not introduce the hydrodynamic and energy balance models but rather derive the drift-diffusion model directly. Before doing this, we have to find expressions for the collision terms in eqns. (2.34) and (2.39). In general they do not depend in a simple way on the moments, but making an ansatz for the distribution function and thus fixing its dependence on velocity we could find expressions for them [66]. Here we use a more phenomenological ansatz [94, 20]. For the general collision term we write

$$\left(\frac{\partial}{\partial t}n\langle\Phi\rangle\right)_{coll} = \left(n\frac{\partial\langle\Phi\rangle}{\partial t}\right)_{coll} + \langle\Phi\rangle\left(\frac{\partial n}{\partial t}\right)_{coll} \quad (2.42)$$

The first term represents intra-band processes that do not change the particle density, whereas the second term describes inter-band collisions that change the number of particles in a band. We replace the rate of change of the particle density by the difference of some generation and recombination rates, i.e.

$$\left(\frac{\partial n}{\partial t}\right)_{coll} = G - R \quad (2.43)$$

For the intra-band term we use a relaxation time approximation and write

$$\left(n\frac{\partial\langle\Phi\rangle}{\partial t}\right)_{coll} = -n\frac{\langle\Phi\rangle - \langle\Phi\rangle_{eq}}{\tau_{\langle\Phi\rangle}} \quad (2.44)$$

The relaxation time $\tau_{\langle\Phi\rangle}$ describes how fast a system in non-equilibrium returns to its equilibrium state $\langle\Phi\rangle_{eq}$. Using these expressions we get for the collision terms of the first two moments

$$\left(\frac{\partial}{\partial t}n\langle\Phi^{(0)}\rangle\right)_{coll} = G - R, \quad \text{as } \langle\Phi^{(0)}\rangle = 1 \quad (2.45)$$

and

$$\left(\frac{\partial}{\partial t}n\langle\Phi^{(1)}\rangle\right)_{coll} = -n\frac{\langle\mathbf{v}\rangle}{\tau} + \langle\mathbf{v}\rangle(G - R), \quad \text{as } \langle\mathbf{v}\rangle_{eq} = 0 \quad (2.46)$$

τ in the last equation is called *momentum relaxation time*.

In order to get a closed transport model based on the particle continuity equation an expression for the particle flux has to be found. For this purpose we first expand $\partial\mathbf{j}_n/\partial t$ to get

$$\frac{\partial\mathbf{j}_n}{\partial t} = \frac{\partial}{\partial t}(n\langle\mathbf{v}\rangle) = n\frac{\partial\langle\mathbf{v}\rangle}{\partial t} + \frac{\partial n}{\partial t}\langle\mathbf{v}\rangle \quad (2.47)$$

Then we substitute (2.34) into (2.39) and use (2.45) and (2.46). This leads to

$$n\frac{\partial\langle\mathbf{v}\rangle}{\partial t} + (\mathbf{j}_n \cdot \nabla_{\mathbf{r}})\langle\mathbf{v}\rangle + \frac{1}{m}\nabla_{\mathbf{r}} \cdot (nk_B\mathbf{T}) - n\frac{\mathbf{F}}{m} = -n\frac{\langle\mathbf{v}\rangle}{\tau} \quad (2.48)$$

where τ is the momentum relaxation time introduced in (2.46).

The third term on the left hand side of eq. (2.48) containing the particle temperature tensor \mathbf{T} should be computed from the second moment of the Boltzmann equation. As we truncated after the first moment it has to be approximated in some way. First, we assume the particles to be in local thermal equilibrium with their environment. That is, we do not consider hot carriers. It can then be assumed that there is no correlation between different directions of the thermal motion and that the thermal energy is equally distributed on all directions (theorem of equipartition [57]). Therefore we write

$$T_{ij} = \frac{m}{k_B} \langle \delta v_i \delta v_j \rangle = T_L \delta_{ij} \quad (2.49)$$

where T_L stands for the environmental temperature (the lattice temperature in a crystal) and δ_{ij} is the usual Kronecker delta.

With the above expression for the particle temperature and writing the mean velocity $\langle \mathbf{v} \rangle$ in terms of the particle flux, we get from eq. (2.48)

$$\mathbf{j}_n + n\tau \frac{\partial}{\partial t} \left(\frac{\mathbf{j}_n}{n} \right) + \tau (\mathbf{j}_n \cdot \nabla_{\mathbf{r}}) \left(\frac{\mathbf{j}_n}{n} \right) = -\frac{k_B T_L \tau}{m} \nabla_{\mathbf{r}} n - n \frac{\tau}{m} \nabla_{\mathbf{r}} (k_B T_L) + n\tau \frac{\mathbf{F}}{m} \quad (2.50)$$

To get an explicit and simple formula for the particle flux, we discard the second and third term on the left hand side, assuming them to be small compared to the flux itself. This approximation has the following physical interpretation:

- $\tau \left| \frac{\partial \langle \mathbf{v} \rangle}{\partial t} \right| \ll |\langle \mathbf{v} \rangle|$: this means that the momentum relaxation has to be faster than the variation of $\langle \mathbf{v} \rangle$ in time induced by extrinsic perturbations.
- $\tau |\nabla_{\mathbf{r}} \langle \mathbf{v} \rangle| \ll 1$: this means that the spatial variation of $\langle \mathbf{v} \rangle$ has to be small compared to the scattering rate.

Whereas the first point isn't too restrictive (the momentum relaxation time is typically in the order of 10^{-12} to 10^{-13} s), the second one can be violated in very small devices like short-channel MOSFETs.

The final form of the drift-diffusion model is obtained by defining a mobility μ and a diffusion coefficient D by

$$\mu = \frac{e\tau}{m}, \quad D = k_B T_L \frac{\tau}{m} \quad (2.51)$$

and by writing the driving force \mathbf{F} as gradient of a potential U , $\mathbf{F} = -\nabla_{\mathbf{r}} U$. Note that mobility and diffusion coefficients are connected by the *Einstein relation* $D = (k_B T_L / e) \mu$. The final formulation of the drift-diffusion model thus reads

$$\frac{\partial n}{\partial t} + \nabla_{\mathbf{r}} \mathbf{j}_n = G - R \quad (2.52a)$$

$$\mathbf{j}_n = -D \nabla_{\mathbf{r}} n - \mu \nabla_{\mathbf{r}} (U/e + k_B T_L / e) \quad (2.52b)$$

The first equation is the familiar continuity equation for the particle density. The second equation is called *constitutive equation*. In the case of charged particles the above system has to be completed by the Poisson equation.

2.2.1.1 Electron and hole transport

For the description of electron and hole transport in a semiconductor in the drift-diffusion approximation we start from a set of two semi-classical Boltzmann equations as given in eq. (2.25):

$$\left[\frac{\partial}{\partial t} + \mathbf{v}_n \cdot \nabla_{\mathbf{r}} - \frac{e}{\hbar} \mathbf{E} \cdot \nabla_{\mathbf{k}} \right] f_n(\mathbf{r}, \mathbf{k}, t) = \left(\frac{\partial f_n}{\partial t} \right)_{coll} \quad (2.53a)$$

$$\left[\frac{\partial}{\partial t} + \mathbf{v}_p \cdot \nabla_{\mathbf{r}} + \frac{e}{\hbar} \mathbf{E} \cdot \nabla_{\mathbf{k}} \right] f_p(\mathbf{r}, \mathbf{k}, t) = \left(\frac{\partial f_p}{\partial t} \right)_{coll} \quad (2.53b)$$

Here $\mathbf{E} = -\nabla\varphi$ denotes the electric field. As electrons and holes are charged particles their effect on the electric field has to be included, i.e. the total electric field has to satisfy the Poisson equation

$$-\nabla(\epsilon \nabla \varphi) = -e(n - p + C) \quad (2.54)$$

where C is the total density of the fixed charges in the system, e.g. ionized dopants. The electron and hole densities n and p are given by the zeroth moments $\int f_n d\mathbf{k}$ and $\int f_p d\mathbf{k}$, respectively.

By applying the method of moments to this system as described in the last section, we end up with the following continuity equations for electrons and holes, which together with the Poisson equation (2.54) form the drift-diffusion equations, in the mathematical community referred to as the Van Roosbroeck equations:³

$$\frac{\partial n}{\partial t} + \nabla \cdot \mathbf{j}_n = -R + G \quad (2.55a)$$

$$\frac{\partial p}{\partial t} + \nabla \cdot \mathbf{j}_p = -R + G \quad (2.55b)$$

with the constitutive equations (assuming constant temperature)

$$\mathbf{j}_n = -D_n \nabla n + \mu_n n \nabla \varphi \quad (2.56a)$$

$$\mathbf{j}_p = -D_p \nabla p - \mu_p p \nabla \varphi \quad (2.56b)$$

The electrical current densities are given by $\mathbf{J}_n = -e\mathbf{j}_n$ and $\mathbf{J}_p = e\mathbf{j}_p$.

From statistical physics we know expressions for the non-degenerate particle densities in local equilibrium in terms of electro-chemical potentials ϕ_n and ϕ_p (see e.g. [57, 56])

$$n = N_c \exp\left(\frac{e\varphi - e\phi_n - E_c}{k_B T}\right) \quad (2.57a)$$

$$p = N_v \exp\left(\frac{E_v - e\varphi + e\phi_p}{k_B T}\right) \quad (2.57b)$$

³It is interesting to note, that the continuity equations can easily be derived in an independent way from the first Maxwell equation (see [96])

$$\nabla \times \mathbf{H} = \mathbf{J} + \frac{\partial \mathbf{D}}{\partial t}$$

The effective densities of states N_c and N_v for electrons and holes, respectively, are given by

$$N_c = 2 \left(\frac{k_B T m_n^*}{2\pi \hbar^2} \right)^{3/2} \quad \text{and} \quad N_v = 2 \left(\frac{k_B T m_p^*}{2\pi \hbar^2} \right)^{3/2} \quad (2.58)$$

where $m_{n,p}^*$ are the density of states effective masses of electrons and holes. The expressions for the densities generalize in the degenerate case to [5]

$$n = N_c F_{1/2} \left(\frac{e\varphi - e\phi_n - E_c}{k_B T} \right) \quad (2.59a)$$

$$p = N_v F_{1/2} \left(\frac{E_v - e\varphi + e\phi_p}{k_B T} \right) \quad (2.59b)$$

$F_{1/2}(x)$ is the Fermi integral of order 1/2

$$F_{1/2}(x) = \frac{1}{\pi^{1/2}} \int_0^\infty \frac{y^{1/2}}{1 + e^{y-x}} dy$$

We will always assume the Einstein relations to hold locally, which connect the diffusion coefficients to the carrier mobilities:

$$D_n = \frac{k_B T}{e} \mu_n, \quad D_p = \frac{k_B T}{e} \mu_p \quad (2.60)$$

where T is the local lattice temperature. The factor $U_T = k_B T / e$ is called *thermal voltage*. The relations (2.60) are valid for the non-degenerate case only. In the degenerate case they can be generalized to [5]

$$\mu_n = e D_n \frac{1}{n} \frac{\partial n}{\partial \varphi}, \quad \mu_p = -e D_p \frac{1}{p} \frac{\partial p}{\partial \varphi} \quad (2.61)$$

Using the above expressions for the carrier densities and the Einstein relations the electron and hole flux from eqns. (2.56) can be rewritten as

$$\mathbf{j}_n = \mu_n n \nabla \phi_n, \quad \mathbf{j}_p = -\mu_p p \nabla \phi_p \quad (2.62)$$

In this form the currents can be interpreted as pure diffusion currents with the gradient of the electro-chemical potential as driving force.

2.2.1.2 Exciton transport

Excitons are quasi-particles formed by a bound electron-hole pair. Different types of excitons are distinguished. When the distance between the electron and the hole is big compared to the lattice constant it is called *Wannier-Mott exciton*. When the electron-hole distance is comparable to the lattice constant it is called *Frenkel exciton*. In this latter case the notion of a bound electron-hole pair has to be understood in a rather formal way [61].

A very simple model of a Wannier-Mott exciton can be obtained by adopting the following simplifying assumptions: Let the crystal have cubic symmetry and a direct bandgap, assume two parabolic bands with extrema at $\mathbf{k} = 0$ and let the electron-hole interaction be the Coulomb attraction scaled by the dielectric constant of the semiconductor. The latter assumption obviously is only valid if the exciton radius is large enough. We can then write down a Schrödinger type equation for this two-particle system and obtain for the exciton energy levels i [61]

$$E_x(i, \mathbf{k}) = \frac{\hbar^2 \mathbf{k}^2}{2m_x} + E_g - \frac{m' e^4}{2\epsilon^2 \hbar^2 i^2} \quad (2.63)$$

where $m_x = m_n + m_p$, $E_g = E_c - E_v$, $m' = m_n m_p / (m_n + m_p)$ and ϵ are the exciton mass, the bandgap, the reduced mass of the system and the dielectric constant, respectively. The last term on the right hand side represents the exciton binding energy R . Fig. 2.4 illustrates graphically the above equation.

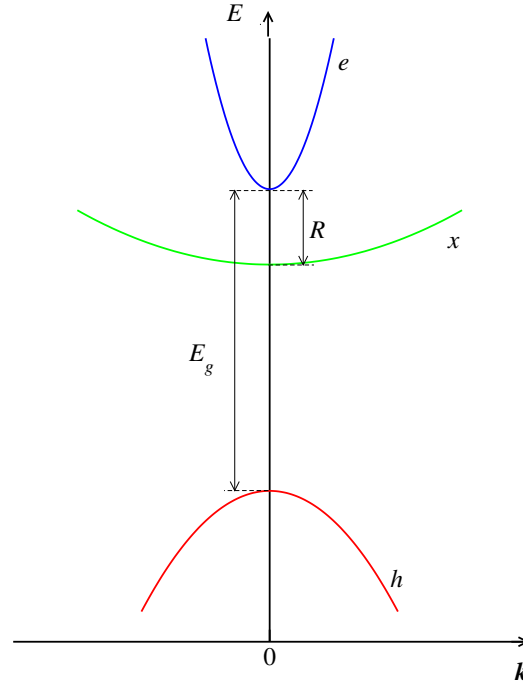


Figure 2.4: Electron, hole and exciton dispersion in a semiconductor.

For the description of the excitons in a drift-diffusion picture we will in the following ignore the electric dipole moment due to the spatial separation of the electron and hole in a Wannier-Mott exciton, despite its interaction with the electric field. As the exciton is composed by two quasi-particles of half-integer spin it is a composite boson and has to be described by Bose-Einstein statistics. But assuming

low density we can use Boltzmann statistics and get for the exciton density

$$x = (2J + 1) \left(\frac{k_B T m_x}{2\pi \hbar^2} \right)^{3/2} e^{(e\phi_x - E_g + R)/k_B T} \quad (2.64)$$

Here, J is the exciton spin and ϕ_x represents the analogue of the electro-chemical potential of holes and electrons. Note that in this case $e\phi_x$ coincides with the chemical potential μ of the exciton gas. The expression in front of the exponential is the effective density of states N_x of the exciton band.

To get the drift-diffusion transport equations of the excitons we can use the results of the previous sections and immediately get

$$\nabla j_x = -\nabla(\mu_x x \nabla \phi_x) = -R_x + G_x \quad (2.65)$$

where we defined an exciton mobility μ_x and exciton recombination/generation rates R_x and G_x .

The following generation/recombination processes are considered in the implementation:

- Dissociation of an exciton into a free electron hole pair, modeled by a relaxation time τ_{diss}
- Nonradiative recombination, modeled by a relaxation time τ_{nr}
- Radiative recombination, modeled by a relaxation time τ_r

With this the total recombination gets $R_x = x(1/\tau_{diss} + 1/\tau_{nr} + 1/\tau_r) = x/\tau_x$, where τ_x is defined as total exciton relaxation time. The main source of exciton generation is the formation of bound pairs from the populations of free electrons and holes. Together with dissociation it couples the exciton system to the electron/hole system. This is described in some more detail in the next section.

2.2.1.3 Coupling of exciton and electron/hole transport

In the drift-diffusion picture adopted in TIBERCAD the coupling between the exciton, electron and hole populations is achieved by means of generation and recombination mechanisms. We assume that electrons and holes form excitons with a rate that is proportional to the electron and hole densities such that the exciton generation rate can be written as

$$G_x = \gamma np \quad (2.66)$$

Furthermore we consider the dissociation of excitons into free electron hole pairs as the reverse process of exciton generation and model it by a dissociation time τ_{diss}

$$R_x = \frac{x}{\tau_{diss}} \quad (2.67)$$

As the generation of excitons “destroys” free electrons and holes, the rate G_x takes the role of a recombination rate in the electron and hole continuity equations. Writing down all three continuity equations we get the system

$$\nabla \mathbf{j}_n = \nabla \mathbf{j}_p = G_{e-h} - R_{e-h} - G_x + R_x \quad (2.68a)$$

$$\nabla \mathbf{j}_x = G_x - R_x \quad (2.68b)$$

where we included a general electron-hole recombination and generation term. Considering thermodynamic equilibrium, i.e. $\mathbf{j}_n = \mathbf{j}_p = \mathbf{j}_x = 0$, we immediately get

$$G_x = R_x \quad \Rightarrow \quad \gamma np = \frac{x}{\tau_{diss}}, \quad (2.69)$$

and using $n_{eq}p_{eq} = n_i^2$ resulting from eq. (2.68a)

$$x_{eq} = \gamma \tau_{diss} n_i^2 \quad (2.70)$$

As thermodynamic equilibrium implies chemical equilibrium we can get an equivalent expression from the law of mass action by viewing the exciton formation and dissociation as a chemical reaction [57]

$$n + p \rightleftharpoons x \quad \Rightarrow \quad \frac{n_{eq}p_{eq}}{x_{eq}} = \chi(T) \quad (2.71)$$

where $\chi(T)$ is some function of temperature. By recalling the expressions for the equilibrium densities given in the last sections it can be found to be (in the non-degenerate case)

$$\chi(T) = \frac{N_c N_v}{N_x} e^{-(R + e\phi_{x,eq})/k_B T} \quad (2.72)$$

This allows to calculate the equilibrium chemical potential of the exciton gas when the effective density of exciton states N_x is known.

Note that eq. (2.70) could have been derived equally well even including the other excitonic recombination processes mentioned in the last section by invoking the principle of detailed balance [19].

For the simulations done with TIBERCAD the values of γ and the exciton life times were estimated from the results of Monte Carlo calculations.

2.3 Heat transport

Due to miniaturization and high density integration of semiconductor devices thermal effects gain crucial importance for their functionality. Self-heating due to power dissipation is in fact considered as one of the major limits for device integration [108]. Therefore an accurate description of heat generation and transport is essential for the simulation of highly miniaturized and integrated devices.

The implementation of macroscopic heat transport in TIBERCAD is based on irreversible thermodynamics [110, 108, 59].⁴ We assume that the system under consideration can be characterized by a set of thermodynamic state variables. In the electron/hole system for example we choose the set [94]

$$(\phi_n, T_n) \quad (\phi_p, T_p) \quad T_L$$

which has the property that at thermodynamic equilibrium $\phi_n = \phi_p = \text{const}$ and $T_n = T_p = T_L = \text{const}$. For simplicity we will neglect hot carrier effects and thus we write $T_n = T_p = T_L$, and we will restrict the considerations to electrons. Our aim is to find expressions that connect the electron and heat fluxes \mathbf{j} and $\mathbf{j}^{(Q)}$ to appropriate source terms (the generalized forces) by means of some kinetic coefficients. By examining the rate of change in entropy we find that the source terms are given by $\nabla\phi/T$ and $\nabla T/T^2$ [59], where $\phi = \phi_n$, and we write therefore (assuming an anisotropic material)

$$j_i = n\mu_{ik}(\partial_k\phi + P\partial_kT) \quad (2.73a)$$

$$j_i^{(Q)} = \beta_{ik}\partial_k\phi - \gamma_{ik}\partial_kT \quad (2.73b)$$

$$j_i^{(Q,L)} = -\kappa_{ik}^L\partial_kT \quad (2.73c)$$

The last equation represents the diffusive heat flux of the lattice which is independent of carrier transport. P is called the thermoelectric power of the electrons. It can be derived starting from (2.52b) by expanding ∇n and using (2.57), (2.58) and the Einstein relation:

$$\begin{aligned} j_i &= -D_{ik}\partial_k n + n\mu_{ik}\partial_k\left(\varphi - \frac{k_B T}{e}\right) \\ &= -D_{ik}\left[\underbrace{\frac{\partial_k N_c}{N_c}}_{\frac{1}{N_c}\frac{\partial N_c}{\partial T}} n + n\partial_k\left(\frac{e\varphi - e\phi - E_c}{k_B T}\right)\right] + n\mu_{ik}\partial_k\left(\varphi - \frac{k_B T}{e}\right) \\ &= -nD_{ik}\left[\frac{3}{2}\frac{1}{T}\partial_k T + \frac{e\partial_k(\varphi - \phi)}{k_B T} - \underbrace{\left(\frac{e\varphi - e\phi - E_c}{k_B T}\right)}_{\ln\left(\frac{n}{N_c}\right)}\frac{\partial_k T}{T}\right] + n\mu_{ik}\partial_k\left(\varphi - \frac{k_B T}{e}\right) \\ &= n\mu_{ik}\left\{\partial_k\phi + \frac{k_B}{e}\left[\ln\left(\frac{n}{N_c}\right) - \frac{5}{2}\right]\partial_k T\right\} \end{aligned} \quad (2.74)$$

Therefore we can identify the electron thermoelectric power in the non-degenerate case as the scalar quantity⁵

$$P = \frac{k_B}{e}\left[\ln\left(\frac{n}{N_c}\right) - \frac{5}{2}\right] \quad (2.75)$$

⁴This model has been implemented in software by Michael Povolotskyi and Giuseppe Romano.

⁵This expression is only correct under the hypothesis that the effective mass is temperature-independent. If the latter is additionally considered anisotropic, P becomes a tensor.

Due to the Onsager reciprocity relation [77] not all coefficients in eqns. (2.73a) and (2.73b) are independent, but $n\mu_{ik}P$ and β_{ik} are connected by the relation [59]

$$T\beta_{ik} = -enPT^2\mu_{ik} \quad (2.76)$$

Substituting $\partial_k\phi = \mu_{ki}^{-1}j_i/n - P\partial_kT$ in (2.73b) we get for the electronic heat flux

$$j_i^Q = -eTPj_i - \underbrace{(\gamma_{ik} - enTP^2\mu_{ik})}_{\kappa_{ik}^n} \partial_kT \quad (2.77)$$

In the above equation we defined the heat conductivity of the electron gas, κ_{ik}^n . The total energy flux can now be written as

$$\begin{aligned} j_i^u &= j_i^{Q,L} + j_i^Q - e\phi j_i \\ &= -\underbrace{(\kappa_{ik}^L + \kappa_{ik}^n)}_{\kappa_{ik}^{tot}} \partial_kT - e(\phi + TP)j_i \end{aligned} \quad (2.78)$$

The last term represents the energy carried by the electron flow. The energy flux satisfies the continuity equation, u being the internal energy density

$$\frac{\partial u}{\partial t} - \partial_i j_i^u = \left(\frac{\partial u}{\partial t} \right)_{rad} \quad (2.79)$$

The above continuity equation can be transformed to [108]

$$c^{tot} \frac{\partial T}{\partial t} + \partial_i \kappa_{ik}^{tot} \partial_k T = H \quad (2.80)$$

where c and κ_{ik} are the total heat capacity and thermal conductivity of the system and H contains all thermal sources and sinks. Here we neglect radiative contributions and we are only interested in the stationary case. Comparing (2.78) and (2.80) we get

$$\begin{aligned} H &= e\partial_i[(\phi + TP)j_i] \\ &= ej_i(\partial_i\phi + P\partial_iT) + eTj_i\partial_iP + e(\phi + PT)\partial_i j_i \\ &= \frac{e}{n}\mu_{ik}^{-1}j_i j_k + eTj_i\partial_iP - eR(\phi + PT) \end{aligned} \quad (2.81)$$

The first term in the last equation represents the well known Joule heat, which in the scalar case simplifies to the more familiar expression $H_{Joule} = e|\mathbf{j}|^2/\mu n$. The second term can be decomposed to obtain the Peltier and Thomson effects. The former is due to spatial variations in the carrier density or changes of the thermoelectric power at material interfaces, the latter is assigned to the change of thermoelectric power due to temperature variations. We finally note that the thermal conductivity of the carriers are usually some orders of magnitude lower with respect to the lattice thermal conductivity and can therefore be neglected [108].

2.4 Quantum mechanical models

The quantum mechanical models implemented in TIBERCAD are based on a single particle Schrödinger-like equation of the form

$$H(\mathbf{r})\psi(\mathbf{r}) = E\psi(\mathbf{r}) \quad (2.82)$$

where H , E and $\psi(\mathbf{r})$ are the Hamiltonian of the system, the eigenenergy and the wavefunction for this eigenenergy, respectively. The Hamiltonian can generally be written in the form

$$H = -\frac{\hbar^2}{2}\nabla_r \left(\frac{1}{m}\nabla_r \right) + V_{eff}(\mathbf{r}) \quad (2.83)$$

where m is a possibly position dependent effective mass. The first term represents the electron kinetic energy. The second term is an effective potential comprising the ionic potential of the atoms forming the structure, the Hartree potential and other contributions.

There are many different solution approaches for the eigenvalue problem (2.82), ranging from ab initio (first principles) to empirical methods and differing mainly in the level of approximation of V_{eff} and the choice of the basis in which to expand the wavefunction $\psi(\mathbf{r})$ [67, 5, 100, 26]. For the study of microstructures most often $\mathbf{k} \cdot \mathbf{p}$ approaches in the framework of the envelope function approximation are employed [26].

The next sections briefly describe the approaches as implemented in TIBERCAD, 2.4.3 being a short introduction to quantum mechanical transport models. They are not intended as exhaustive presentations but rather as an introduction to the basic features.

2.4.1 Envelope function approximation

The basic idea of the *envelope function approximation* (EFA) is to expand the single-particle wavefunction $\psi_i(\mathbf{r})$ of a heterostructure in a basis of periodic functions $U_n(\mathbf{r}) = U_n(\mathbf{r} + \mathbf{R})$, where \mathbf{R} is a lattice vector of the underlying Bravais lattice [64, 24]

$$\psi(\mathbf{r}) = \sum_n F_n(\mathbf{r})U_n(\mathbf{r}) \quad (2.84)$$

Usually, the bulk Bloch functions of one of the constituent materials are chosen as basis [82, 24, 25] such that

$$\psi(\mathbf{r}) = \sum_{n,\mathbf{k}} C_{n,\mathbf{k}} u_{n,\mathbf{k}}(\mathbf{r}) e^{i\mathbf{k}\mathbf{r}} \quad (2.85)$$

Next, only the Bloch function for a certain $\mathbf{k} = \mathbf{k}_0$ are kept in the expansion. This leads to

$$\begin{aligned} \psi(\mathbf{r}) &= \sum_{n,\mathbf{k}} \tilde{C}_{n,\mathbf{k}} u_{n,\mathbf{k}_0}(\mathbf{r}) e^{i\mathbf{k}\mathbf{r}} \\ &= \sum_n F_n(\mathbf{r}) u_{n,\mathbf{k}_0}(\mathbf{r}) \end{aligned} \quad (2.86)$$

where $F_n(\mathbf{r}) = \sum_{\mathbf{k}} \tilde{C}_{n,\mathbf{k}} e^{i\mathbf{k}\mathbf{r}}$ is the *envelope function set*. Usually \mathbf{k}_0 is chosen to be a band extremum.

The expansion (2.86) for the wavefunction can now be inserted into eq. (2.82) to get a Schrödinger-like equation for the envelope functions $F_n(\mathbf{r})$ (omitting the argument in $F_n(\mathbf{r})$ and $U_n(\mathbf{r})$):

$$\begin{aligned} \sum_n \left(-\frac{\hbar^2}{2m} \nabla^2 + V \right) F_n U_n = \\ \sum_n \left[-\frac{\hbar^2}{2m} (\nabla^2 F_n) U_n - \frac{\hbar^2}{m} \nabla F_n \nabla U_n - \frac{\hbar^2}{2m} (\nabla^2 U_n) F_n + V F_n U_n \right] \\ = E \sum_n F_n U_n \quad (2.87) \end{aligned}$$

This equation has to be brought to envelope-function expansion form. This can be achieved by means of the matrix elements of the different operators acting on U_n . In detail, and neglecting any nonlocalities [24]:

$$(a) \quad \nabla U_n = \frac{i}{\hbar} \mathbf{p} U_n = \frac{i}{\hbar} \mathbf{p}_{n'n} U_{n'}, \quad \text{with} \quad \mathbf{p}_{n'n} = \frac{1}{\Omega_c} \int U_{n'}^* \mathbf{p} U_n \, dx$$

$$(b) \quad \nabla^2 U_n = -\frac{2m}{\hbar^2} T U_n = -\frac{2m}{\hbar^2} T_{n'n} U_{n'}, \quad \text{with} \quad T_{n'n} = \frac{1}{\Omega_c} \int U_{n'}^* T U_n \, dx$$

$$(c) \quad V U_n = V_{n'n} U_{n'}, \quad \text{with} \quad V_{n'n} = \frac{1}{\Omega_c} \int U_{n'}^* V U_n \, dx$$

where $\frac{1}{\Omega_c} \int U_{n'}^* U_n \, dx = \delta_{n'n}$.

With the above expressions we recast (2.87) to get (after some changes of indices)

$$\sum_n \left[-\frac{\hbar^2}{2m} \nabla^2 F_n - i \frac{\hbar}{m} \sum_{n'} \mathbf{p}_{nn'} \nabla F_{n'} + \sum_{n'} \underbrace{(T_{nn'} + V_{nn'})}_{H_{nn'}} F_{n'} \right] U_n = E \sum_n F_n U_n \quad (2.88)$$

Equating the coefficients on both sides we finally get the equation for the envelope function

$$-\frac{\hbar^2}{2m} \nabla^2 F_n(\mathbf{r}) - i \frac{\hbar}{m} \sum_{n'} \mathbf{p}_{nn'} \nabla F_{n'}(\mathbf{r}) + \sum_{n'} H_{nn'}(\mathbf{r}) F_{n'}(\mathbf{r}) = E F_n(\mathbf{r}) \quad (2.89)$$

Note that to get the above equation we did not do any approximation other than neglecting the nonlocal part in $H_{nn'}$, which in any case is only important near the heterointerfaces [24].

The solution of (2.89) is still computationally very demanding and further simplification is needed. For this purpose the different bands associated with the different envelope functions F_n are divided into two groups of bands denoted by S and R [82, 24]. The different bands F_s of the S group are the dominant bands whereas the bands F_r contained in the R group are considered as remote bands. The latter are approximated such as to eliminate them from the system of equations to get equations only for the dominant bands. The small F_r are written as as

$$F_r \approx (E - H_{rr})^{-1} \sum_{s'} \left(-i \frac{\hbar}{m} \mathbf{p}_{rs'} \cdot \nabla F_{s'} + H_{rs'} F_{s'} \right) \quad (2.90)$$

For the dominant bands we get, setting $n = s$ in (2.89)

$$\begin{aligned} -\frac{\hbar^2}{2m} \nabla^2 F_s - i \frac{\hbar}{m} \sum_{s'} \mathbf{p}_{ss'} \cdot \nabla F_{s'} + \sum_{s'} H_{ss'} F_{s'} \\ - i \frac{\hbar}{m} \sum_r \mathbf{p}_{sr} \cdot \nabla F_r + \sum_r H_{sr} F_r = E F_s \end{aligned} \quad (2.91)$$

Substituting (2.90) into the above equation we finally get

$$\begin{aligned} -\frac{\hbar^2}{2m} \sum_{s'} \nabla \cdot \left[\gamma_{ss'}^{(r)}(E, \mathbf{r}) \cdot \nabla F_{s'}(\mathbf{r}) \right] + \sum_{s'} \frac{-i\hbar}{m} \mathbf{p}_{ss'} \cdot \nabla F_{s'}(\mathbf{r}) + \\ \sum_{s'} H_{ss'}^{(2)}(E, \mathbf{r}) F_{s'}(\mathbf{r}) + \sum_{s', r} \frac{-i\hbar}{m} \mathbf{p}_{sr} \cdot \nabla \left[(E - H_{rr}(\mathbf{r}))^{-1} H_{rs'} \right] F_{s'}(\mathbf{r}) + \\ \sum_{s', r} \frac{-i\hbar}{m} \frac{\mathbf{p}_{sr} H_{rs'} + H_{sr} \mathbf{p}_{rs'}}{E - H_{rr}(\mathbf{r})} \cdot \nabla F_{s'}(\mathbf{r}) = E F_s(\mathbf{r}) \end{aligned} \quad (2.92)$$

The second rank tensor $\gamma_{ss'}^{(r)}$ is given by

$$\gamma_{ss'}^{(r)}(E, \mathbf{r}) = \mathbf{I} \delta_{ss'} + \frac{2}{m} \sum_r \frac{\mathbf{p}_{sr} \otimes \mathbf{p}_{rs'}}{E - H_{rr}(\mathbf{r})} \quad (2.93)$$

where \mathbf{I} is the second rank unity tensor. The second term in the above expression can be interpreted as a renormalization of the electron mass due to the influence of the remote bands, giving rise to the name *effective mass approximation* often used for this method.

The term $H_{ss'}^{(2)}$ in (2.92) is defined by

$$H_{ss'}^{(2)}(E, \mathbf{r}) = H_{ss'}(\mathbf{r}) + \sum_r \frac{H_{sr}(\mathbf{r}) H_{rs'}(\mathbf{r})}{E - H_{rr}(\mathbf{r})} \quad (2.94)$$

and can be seen as a modified or effective band energy.

Eq. (2.92) still represents a nonlinear eigenvalue problem due to the energy dependent terms. It can be simplified further by discarding small terms. Especially the fourth and fifth term on the left-hand side of (2.92) are non-zero only near the interfaces and scaled by a large energy denominator and can therefore be neglected [24].

In the simplest case, i.e. considering only one dominant band with minimum in Γ , an effective mass approximation for the conduction band can be obtained. It reads, writing $r = c$,

$$-\frac{\hbar^2}{2m} \nabla \left(\frac{1}{m_c(E, \mathbf{r})} \nabla F_c(\mathbf{r}) \right) + H_{cc}^{(2)}(E, \mathbf{r}) F_c(\mathbf{r}) = E F_c(\mathbf{r}) \quad (2.95)$$

where $1/m_c(E, \mathbf{r})$ and $H_{cc}^{(2)}(E, \mathbf{r})$ are given by (2.93) and (2.94), respectively. The last approximation is done by assuming a large band gap such that $E - H_{rr}(\mathbf{r}) \approx E_g(\mathbf{r})$. This removes the energy dependences and leads to the final eigenvalue equation for the conduction band

$$-\frac{\hbar^2}{2m} \nabla \left(\frac{1}{m_c(\mathbf{r})} \nabla F_c(\mathbf{r}) \right) + E_c(\mathbf{r}) F_c(\mathbf{r}) = E F_c(\mathbf{r}) \quad (2.96)$$

m_c and E_c can be replaced by experimental values. We note that the basis functions U_n are the same in the whole structure and therefore the envelope functions are continuous [82, 24].

Equations for the valence bands can be obtained in a similar way, but all the three top valence bands having p -type angular symmetry have to be treated as dominant ones. The latter are usually denoted as $|X\rangle$, $|Y\rangle$ and $|Z\rangle$. The remote bands that have to be included have s -like (Γ_1 , conduction band) and d -like (Γ_{12} , Γ_{15}) symmetry.

We introduce a new operator $\mathbf{k} = -i\nabla$, which basically transforms (2.92) to \mathbf{k} -space, and choose a coordinate system such that $X \parallel [100]$, $y \parallel [010]$ and $z \parallel [001]$. First we write down the effective mass approximation for the valence band states without considering spin-orbit coupling, leading to a three-fold degeneracy in the Γ -point. We formulate the eigenvalue problem in a short form as

$$H_{vv}(\mathbf{r}) \vec{F}_v(\mathbf{r}) = E \vec{F}_v(\mathbf{r}) \quad (2.97)$$

with $\vec{F}_v = (F_X \ F_Y \ F_Z)^T$. A calculation of all the relevant terms of (2.92) considering the aforementioned remote bands leads to the following 3×3 valence band Hamiltonian for a wurtzite crystal [42]

$$H_{vv} = \begin{pmatrix} E_v + \frac{\hbar^2}{2m} k^2 & 0 & 0 \\ 0 & E_v + \frac{\hbar^2}{2m} k^2 & 0 \\ 0 & 0 & E_v + \frac{\hbar^2}{2m} k^2 \end{pmatrix} + \begin{pmatrix} k_x L_1 k_x + k_y M_1 k_y + k_z M_2 k_z & k_x C_1 k_y + k_y M_1 k_x & k_x C_2 k_z + k_z M_2 k_x \\ k_y C_1 k_x + k_x M_1 k_y & k_x M_1 k_x + k_y L_1 k_y + k_z M_2 k_z & k_y C_2 k_z + k_z M_2 k_y \\ k_z C_2 k_x + k_x M_2 k_z & k_z C_2 k_y + k_y M_2 k_z & k_x M_3 k_x + k_y M_3 k_y + k_z L_2 k_z \end{pmatrix} \quad (2.98)$$

The different parameters in the Hamiltonian are given by the matrix elements in (2.92), considering the symmetry properties of the chosen basis. Explicit expressions can be found e.g. in [82] and [42] and references therein.

Formally the same Hamiltonian is found for zinc blende crystals. However, due to the cubic symmetry, the L_i , M_i and C_i in (2.98) are all equal to L , M and C , respectively.

Up to now we did not consider spin at all. Spin-orbit interaction, which is a relativistic effect, is taken into account approximately by extending the Hamiltonian (2.83) in the following way [82]:

$$H = \frac{\mathbf{p}^2}{2m} + V_{eff}(\mathbf{r}) + \frac{\hbar}{4m^2c^2} [\boldsymbol{\sigma} \times \nabla V_{eff}(\mathbf{r})] \cdot \mathbf{p} \quad (2.99)$$

To include spin in the calculation we extend the basis to the six kets $|X \uparrow\rangle$, $|Y \uparrow\rangle$, $|Z \uparrow\rangle$, $|X \downarrow\rangle$, $|Y \downarrow\rangle$ and $|Z \downarrow\rangle$, where the arrows symbolically stand for the direction of the projection of spin onto the z -axis (spin up, spin down). We denote the resulting 6×6 EFA Hamiltonian by $H_{6 \times 6}$. It reads

$$H_{6 \times 6} = \begin{pmatrix} H_{vv} & 0 \\ 0 & H_{vv} \end{pmatrix} + H_{6 \times 6}^{s/o} \quad (2.100)$$

where the spin-orbit coupling Hamiltonian is given by [82]

$$H_{6 \times 6}^{s/o} = \frac{\Delta}{3} \begin{pmatrix} 0 & -i & 0 & 0 & 0 & 1 \\ i & 0 & 0 & 0 & 0 & -i \\ 0 & 0 & 0 & -1 & i & 0 \\ 0 & 0 & -1 & 0 & i & 0 \\ 0 & 0 & -i & -i & 0 & 0 \\ 1 & i & 0 & 0 & 0 & 0 \end{pmatrix} \quad (2.101)$$

The Δ in the above expression is the spin-orbit coupling energy. The basis which diagonalizes the spin-orbit Hamiltonian defines the states known as *light*, *heavy* and *split off* holes.

In some cases, e.g. for low band gap semiconductors, treating the conduction band as a remote band is a too poor approximation. In this situation the top valence bands and the lowest conduction band are considered as dominant ones and a 8×8 EFA can be formulated, using as basis $|S \uparrow\rangle$, $|X \uparrow\rangle$, $|Y \uparrow\rangle$, $|Z \uparrow\rangle$, $|S \downarrow\rangle$, $|X \downarrow\rangle$, $|Y \downarrow\rangle$, $|Z \downarrow\rangle$. The total EFA Hamiltonian is written as

$$H_{8 \times 8} = \begin{pmatrix} H_{4 \times 4} + H_{4 \times 4}^{strain} & 0 \\ 0 & H_{4 \times 4} + H_{4 \times 4}^{strain} \end{pmatrix} + H_{8 \times 8}^{s/o} \quad (2.102)$$

The spin-orbit contribution $H_{8 \times 8}^{s/o}$ can be obtained from (2.101) by adding all-zero rows and columns for the $|S \uparrow\rangle$ and $|S \downarrow\rangle$ states, as s -states do not show spin-orbit interaction. $H_{4 \times 4}$ is given by

$$H_{4 \times 4} = \begin{pmatrix} H_{cc} & H_{cv} \\ H_{cv}^\dagger & H_{vv} \end{pmatrix} \quad (2.103)$$

where H_{cc} is the 1×1 conduction band Hamiltonian and H_{cv} is the 1×3 coupling between the conduction band and the three valence bands. H_{vv} is formally given by (2.98), however the parameters are different as now the conduction band is not anymore considered as remote band. A similar statement holds for H_{cc} , whereas H_{cv} are new terms.

The contribution $H_{4 \times 4}^{strain}$ represents the first order correction due to strain [34]. For a zinc blende crystal it reads

$$H_{4 \times 4}^{strain} = \begin{pmatrix} a_c \text{Tr}(\varepsilon) & 0 & 0 & 0 \\ 0 & l\varepsilon_{xx} + m\varepsilon_{yy} + m\varepsilon_{zz} & n\varepsilon_{xy} & n\varepsilon_{xz} \\ 0 & n\varepsilon_{xy} & m\varepsilon_{xx} + l\varepsilon_{yy} + m\varepsilon_{zz} & n\varepsilon_{yz} \\ 0 & n\varepsilon_{xz} & n\varepsilon_{yz} & m\varepsilon_{xx} + m\varepsilon_{yy} + l\varepsilon_{zz} \end{pmatrix} \quad (2.104)$$

a_c is the absolute deformation potential for the conduction band. l , m and n are given by [42]

$$l = a_v + 2b \quad (2.105a)$$

$$m = a_v - b \quad (2.105b)$$

$$n = \sqrt{3}d \quad (2.105c)$$

where a_v is the absolute valence band deformation potential and b and d are shear deformation potentials. We can note that the relative volumic change of the unit cell, given by the trace of the strain tensor, leads to an absolute shift of the band energies by means of a_c and a_v , whereas b and d lead to an additional individual shift of each valence band. Thus strain can lift the degeneracy of the light and heavy hole bands. We finally note that the \mathbf{k} -vector in $H_{4 \times 4}$ (c.f. (2.98)) has to be substituted by $k_i = (\delta_{ij} - \zeta_{ij})k_j$, where δ_{ij} is the Kronecker delta and $\zeta_{ij} = \partial u_i / \partial x_j$ (c.f. section 2.1) [82, 34].

We will shortly illustrate the application of the EFA as described above to a heterostructure.⁶ Assume an AlGaAs/GaAs/AlGaAs quantum well grown in z -direction ($z \parallel [001]$), thus maintaining translational symmetry in the xy -plane. The parameters in the Hamiltonian, e.g. (2.98), are in this case constant in the xy -plane but vary with z . We separate the \mathbf{k} -space into a parallel and an orthogonal component \mathbf{k}_{\parallel} and \mathbf{k}_{\perp} , where \mathbf{k}_{\parallel} lies in the xy -plane. In z -direction we transform back to real space, using $\mathbf{k}_z = -i\partial_z$. Writing the envelope function as $F_{\mathbf{k}_{\parallel}}(z)$ we then get an eigenvalue problem for each \mathbf{k}_{\parallel} reading

$$H_{\mathbf{k}_{\parallel}}(z, \partial_z)F_{\mathbf{k}_{\parallel}}(z) = E_{\mathbf{k}_{\parallel}}F_{\mathbf{k}_{\parallel}}(z) \quad (2.106)$$

For the actual calculation the above equation has to be discretized on a mesh along the z -axis, and appropriate boundary conditions for $F_{\mathbf{k}_{\parallel}}(z)$ have to be provided.

⁶The single-band, 6×6 and 8×8 EFA has been implemented in TIBERCAD by Michael Polotskiy.

In TIBERCAD the discretization is done using the standard Galerkin finite element method (c.f. section 3.3.1), i.e. we use piecewise linear basis functions $\varphi_i(z)$, expand $F_{\mathbf{k}_{\parallel}}(z) = \sum_i f_i \varphi_i(z)$ and write (2.106) in weak form in such a way as to get a finite-dimensional generalized eigenvalue problem $H_{ij} f_j = E S_{ij} f_j$

$$\underbrace{\int \varphi_i H_{\mathbf{k}_{\parallel}}(z, \partial_z) \varphi_j dz}_{H_{ij}} f_j = E f_j \underbrace{\int \varphi_i \varphi_j dz}_{S_{ij}} \quad (2.107)$$

Obviously the same approach can be used for other structures. In a quantum wire along x e.g. symmetry is broken in y and z direction and \mathbf{k}_{\parallel} will be a 1D space along x .

The result of an EFA calculus can be used to calculate the quantum mechanical particle density and thus the charge density entering into the Poisson equation. As the Hamiltonian depends on the electric potential V , the Schrödinger and Poisson equations have to be solved self-consistently. The electron density can be written based on the conduction band envelope functions and the corresponding energy levels as

$$n_Q(\mathbf{r}) = \sum_s \frac{1}{(2\pi)^d} \int_{BZ_{\parallel}} |F_s(\mathbf{r}, \mathbf{k}_{\parallel})|^2 f \left(\frac{E_s(\mathbf{k}_{\parallel}) + e\bar{\phi}_{n,s}}{k_B T} \right) d\mathbf{k}_{\parallel} \quad (2.108)$$

where $f(x) = 1/(1 + \exp(x))$ denotes the Fermi-function. d denotes the dimension of the parallel Brillouin zone. The envelope functions are normalized by

$$\int_{\Omega} |F_s|^2 dx = 1 \quad (2.109)$$

For the quantized hole states we get a similar expression

$$p_Q(\mathbf{r}) = \sum_s \frac{1}{(2\pi)^d} \int_{BZ_{\parallel}} |F_s(\mathbf{r}, \mathbf{k}_{\parallel})|^2 f \left(-\frac{E_s(\mathbf{k}_{\parallel}) + e\bar{\phi}_{p,s}}{k_B T} \right) d\mathbf{k}_{\parallel} \quad (2.110)$$

The mean electro-chemical potentials $\bar{\phi}_{n/p,s}$ are calculated as mean values $\langle F_s | \phi | F_s \rangle$

$$\bar{\phi}_{n/p,s} = \int_{\Omega} \phi_{n/p}(\mathbf{r}) |F_s(\mathbf{r})|^2 d\mathbf{r} \quad (2.111)$$

The approach to calculate quantum density as described here is in principle only correct in equilibrium. The electrochemical potentials are, however, often approximately constant in regions of quantized states such that the above approach leads to reasonable results.

2.4.2 Atomistic models

In the *tight-binding* scheme (TB) the wavefunction is expanded as a sum of atomic orbitals $|n\alpha\rangle = \psi_{\alpha}(\mathbf{r} - \mathbf{R}_n)$ [26, 79].⁷ The n indexes the atom located at position

⁷The implementation of atomistic models in TIBERCAD is carried out by Alessandro Pecchia and Gabriele Penazzi.

\mathbf{R}_n and the index α specifies the type of the orbital with respect to symmetry and spin quantum numbers. The wave function $|\Psi\rangle$ of the system then reads

$$|\Psi\rangle = \sum_{n,\alpha} C_{n\alpha} |n\alpha\rangle \quad (2.112)$$

This is called *linear combination of atomic orbitals* (LCAO). The use of atomic orbitals reflects (or induces) the assumption that the electronic states in the crystal are not too much disturbed with respect to the states in the isolated atoms, in contrast to free electron models, where the electrons are assumed to be essentially unbound in the crystal and therefore plain waves are used as basis functions. This explains the name of the tight-binding method and lets expect that it produces better results for lower lying (more bound) states, i.e. the valence bands.

Putting the above LCAO expansion into the Schrödinger equation and projecting onto the LCAO basis we get

$$\sum_{n\alpha} [H_{n'\alpha',n\alpha} - ES_{n'\alpha',n\alpha}] C_{n\alpha} = 0 \quad (2.113)$$

The solution of this generalized eigenvalue problem is a set of expansion coefficients $C_{n\alpha}$ and corresponding eigenenergies E . $H_{n'\alpha',n\alpha}$ and $S_{n'\alpha',n\alpha}$ are the Hamiltonian and overlap matrix elements, respectively. They are given by

$$H_{n'\alpha',n\alpha} = \langle n'\alpha' | H | n\alpha \rangle \quad (2.114a)$$

$$S_{n'\alpha',n\alpha} = \langle n'\alpha' | n\alpha \rangle \quad (2.114b)$$

Often orbitals are used that are orthogonalized by a procedure introduced by Löwdin [63]. It has the important property that it preserves the symmetry of the orbitals. However, the orthogonal orbitals have a longer range than the original ones, which can be a disadvantage in practical implementations [79].

Let us assume the effective one-particle Hamiltonian H to be of the form

$$H = \frac{\mathbf{p}^2}{2m} + \sum_n V_n(\mathbf{r} - \mathbf{R}_n) \quad (2.115)$$

where each V_n is the (spherically symmetric) potential produced by the atom at position \mathbf{R}_n . With this, the Hamiltonian matrix elements can be written as

$$H_{n'\alpha',n\alpha} = \int \psi_{\alpha'}^*(\mathbf{r} - \mathbf{R}_{n'}) \left[\frac{\mathbf{p}^2}{2m} + \sum_{n''} V_{n''}(\mathbf{r} - \mathbf{R}_{n''}) \right] \psi_{\alpha}(\mathbf{r} - \mathbf{R}_n) d\mathbf{r} \quad (2.116)$$

Thus we can identify four categories of integrals involving a potential:

- (i) on-site integrals, where the orbitals and the potential are located on the same atomic site, i.e. $n = n' = n''$
- (ii) two-centre integrals, where one orbital and the potential are on the same site, i.e. $n \neq n' = n''$ or $n' \neq n = n''$

- (iii) three-center integrals, where all the orbitals and the potential are located on different sites, i.e $n \neq n' \neq n''$
- (iv) the two orbitals are on the same site, but the potential on a different one, i.e. $n = n' \neq n''$

Usually not all the terms appearing in (2.116) are included in the Hamiltonian. Different implementations of tight-binding can therefore be distinguished on the one hand by the number and type of terms that are included in the matrix elements and on the other hand by the way they are calculated. Often a two-center approximation is used and only nearest or second-nearest neighbours are considered. Another characteristic of different implementations is the size of the basis used in the expansion (2.112).

Very successful are empirical approaches (empirical tight-binding, ETB) where the matrix elements are considered as fitting parameters of experimentally accessible quantities such as band gaps at high symmetry points, effective masses and so on. This approach leads to parameter sets that can very reliably reproduce true band structures.

In ab initio approaches the matrix elements are evaluated from first principles using e.g. density functional theory or pseudo-potential methods.

When the system under consideration exhibits translational symmetry in one or more dimensions, it is possible to decrease the dimension of the eigenvalue problem by writing the orbital basis as a bloch sum, e.g. for a bulk crystal

$$|n\alpha\rangle = \sum_{\mathbf{R}} e^{i\mathbf{k}\cdot(\mathbf{R}+\boldsymbol{\nu}_n)} |\mathbf{R} + \boldsymbol{\nu}_n, \alpha\rangle \quad (2.117)$$

where \mathbf{R} is a lattice vector and $\boldsymbol{\nu}_n$ denotes the position of the n -th atom inside the primitive cell.

As an example we apply the empirical tight-binding scheme for the calculation of the GaAs bulk band structure, choosing a basis consisting of ten orbitals per atom and using the two-center nearest neighbour approximation. This parametrization is called $\text{sp}^3\text{d}^5\text{s}^*$ -parametrization, where s, p and d denote the symmetry of the atomic orbitals. Including spin-orbit interaction we get a 40×40 eigenvalue problem that has to be solved for each \mathbf{k} -point to get the band structure shown in Fig. 2.5. The parametrization of the matrix elements was taken from [48].

2.4.3 Quantum transport

Particle transport can be described quantum mechanically by identifying the particle flux with the probability density flux calculated from the wavefunction $\psi(\mathbf{r})$ of the particle [31, 58]:⁸

$$\mathbf{j} = \frac{1}{2m} (\psi \mathbf{p}^* \psi^* + \psi^* \mathbf{p} \psi) = \frac{i\hbar}{2m} (\psi \nabla \psi^* - \psi^* \nabla \psi) \quad (2.118)$$

⁸Quantum transport in TIBERCAD is worked on by Alessandro Pecchia and Fabio Sacconi.

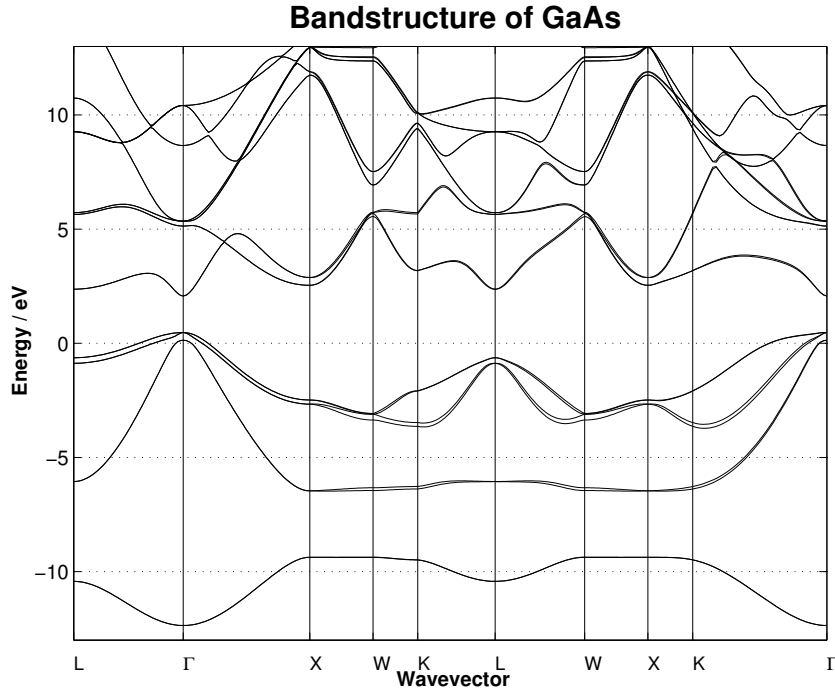


Figure 2.5: Band structure of GaAs calculated with $sp^3d^5s^*$ parametrization given in [48].

When examining the above equation we note that whenever ψ is a solution of the stationary Schrödinger equation for a closed system and therefore is a real function, the flux vanishes. This means that the system has to be treated as an open system by imposing open boundary conditions.

Eq. (2.118) is obtained from the time-dependent Schrödinger equation when deriving the continuity equation for the probability density

$$\frac{\partial \rho}{\partial t} + \nabla \cdot \mathbf{j} = 0 \quad (2.119)$$

The derivation is based on the hermiticity of H , that is on the reality of V . This however means that the number of particles in the device is constant and especially that transport is non-dissipative [92].

Usually eq. (2.118) is not applied in this form to calculate currents. Several different approaches are used, depending on the properties of the system under consideration. In the case of coherent (and non-coherent elastic) transport the Landauer-Büttiker formalism can appropriately describe transport, using the notion of *transmission functions* [31, 18]. It is especially useful for the study of mesoscopic systems at low temperatures. Consider a system as shown in Fig. 2.6. It consists of three terminals (the contacts) connected to a conductor. The contacts i are assumed to be in equilibrium, each being characterized by an electro-chemical potential μ_i and an equilibrium distribution $f_i(E)$. Transmission coefficients are defined that

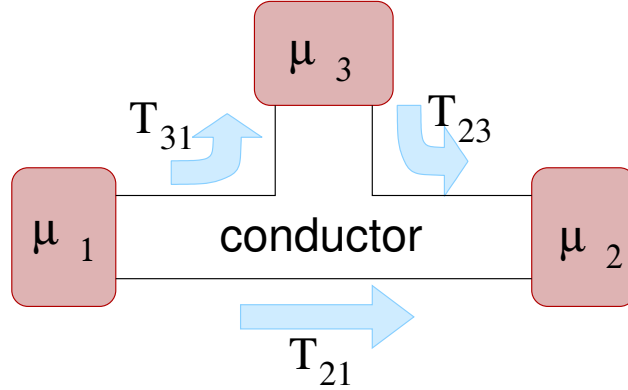


Figure 2.6: A three-terminal device. The three contacts, connected to a conductor, are assumed to be in equilibrium.

describe the total transmission from contact q to contact p as $\bar{T}_{pq}(E)$. The total net current from contact p into the conductor is then given by

$$I_p = \frac{2e}{h} \int \sum_q [\bar{T}_{qp}(E) f_p(E) - \bar{T}_{pq}(E) f_q(E)] dE \quad (2.120)$$

In equilibrium and when there is no inelastic scattering in the conductor, the transmission functions satisfy the relation $\sum_q \bar{T}_{qp}(E) = \sum_q \bar{T}_{pq}(E)$ and the above equation simplifies to

$$I_p = \frac{2e}{h} \int \sum_q \bar{T}_{pq}(E) [f_p(E) - f_q(E)] dE \quad (2.121)$$

If the particle flux is carried by different modes in the different contacts (e.g. given by different Landau-levels in a hall bar), then the total transmission is given by

$$\bar{T}_{pq}(E) = \sum_{m \in p} \sum_{n \in q} T_{mn} \quad (2.122)$$

where T_{mn} is the transmission probability from mode n in contact q to mode m in contact p .

The transmission functions as described before are closely related to the scattering matrix (S -matrix). The latter is defined as the (unitary) matrix that relates incoming wave amplitudes to the outgoing ones, in a similar way as for the S -parameters in microwave engineering. We denote the incoming waves as \mathbf{a} and the outgoing waves as \mathbf{b} such that $\mathbf{b} = \mathbf{S}\mathbf{a}$. The total number of propagating modes is the sum of the modes in each contact (or lead). The connection with the transmission function is then given by the relation $T_{mn} = |s_{mn}|^2$.

The S -matrix can be calculated from the Schrödinger equation, using e.g. the envelope function approximation or a tight-binding description of the conductor as introduced in the last two sections.

Although the approach based on the Landauer-Büttiker formalism in principle cannot treat non-coherent transport, it can be accounted for in a phenomenological way by introducing virtual floating voltage probes that will simulate an incoherent component in the current flow between the ordinary contacts [31].

To go a step beyond, a description based on the density matrix can be adopted, based on the quantum Liouville equation [42]

$$\frac{\partial \rho}{\partial t} = -\frac{i}{\hbar}[H, \rho] + \left(\frac{\partial \rho}{\partial t}\right)_{int} \quad (2.123)$$

where ρ is the *density matrix* of the system under consideration and the second term on the right hand side describes the interaction with the outside world, in analogy to the scattering term in the Boltzmann equation. In a pure state, the density matrix can be defined as $\rho = |\psi\rangle\langle\psi|$. The mean value of any observable A can be expressed using ρ as

$$\langle A \rangle = \text{Tr}(\rho A), \quad (2.124)$$

$\text{Tr}(\bullet)$ being the trace operator. The above expression is independent of the choice of the basis.

The quantum current in terms of the density matrix is given as the trace

$$\mathbf{j}(r) = -i\hbar \left[\frac{\nabla_r - \nabla_{r'}}{2m} \rho(r, r', t) \right]_{r'=r} \quad (2.125)$$

Starting from the density matrix and the quantum Liouville equation it is possible to derive a kinetic transport equation for a “quasi distribution function” $w(\mathbf{r}, \mathbf{v}, t)$, termed Wigner transport equation and Wigner function, respectively [66]. They form a quantum mechanical analogue of the Boltzmann transport equation and can be used for numerical simulation [86].

We now slightly change the picture from Fig. 2.6 to the one given in Fig. 2.7. The

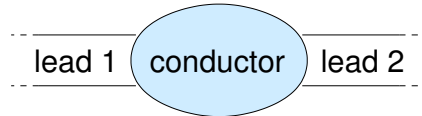


Figure 2.7: A two-lead device. The leads connected to the conductor are semi-infinite and assumed to be homogeneous such that a translational symmetry can be exploited to collapse the dimensionality of the numerical description.

main difference consists in substituting the contact reservoirs with homogeneous semi-infinite leads. To keep the description simple we have drawn only two leads. Assume we can describe the system shown in Fig. 2.7 by three lead and conductor Hamiltonians H_1 , H_2 and H_C and by coupling terms H_{1C} and H_{C2} . Furthermore there shall be no direct coupling between the two leads. We can then write down a

Schrödinger equation in matrix form as

$$\begin{pmatrix} H_1 & H_{1C} & 0 \\ H_{1C}^\dagger & H_C & H_{C2} \\ 0 & H_{C2}^\dagger & H_2 \end{pmatrix} \begin{pmatrix} \Psi_1 \\ \Psi_C \\ \Psi_2 \end{pmatrix} = E \begin{pmatrix} \Psi_1 \\ \Psi_C \\ \Psi_2 \end{pmatrix} \quad (2.126)$$

To simplify the description we have assumed that the wavefunctions assigned to the leads and the conductor are orthogonal (for a detailed discussion see [26, 79]). Note that due to the semi-infinite leads the Hamiltonian in the above eigenvalue problem is in principle always infinite dimensional. However, we still can symbolically transform the equation and obtain a modified eigenproblem which apparently only involves the conductor. As the leads do not interact, we are able to write e.g. for lead 1

$$H_1 \Psi_1 + H_{1C} \Psi_C = E \Psi_1 \quad \Rightarrow \quad \Psi_1 = (E - H_1)^{-1} H_{1C} \Psi_C \quad (2.127)$$

We put this and the equivalent expression for lead 2 into the second row of the system (2.126) to get

$$\left[\underbrace{H_{1C}^\dagger (E - H_1)^{-1} H_{1C}}_{\Sigma_1} + H_C + \underbrace{H_{C2} (E - H_2)^{-1} H_{C2}^\dagger}_{\Sigma_2} \right] \Psi_C = E \Psi_C \quad (2.128)$$

The effect of the semi-infinite leads (or, the fact of being an open system) enters into the eigenvalue problem by means of the operators Σ_i , called *self-energies*, which we will discuss later on. We note that a discretization of (2.128) would lead to a finite-dimensional problem, provided that expressions for the self-energies can be found.

A more direct solution strategy for the eigenvalue problem (2.126) makes use of so called *transfer matrices* [26, 79]. It is most easily formulated within the framework of a nearest-neighbour tight-binding approach. Assume a laterally homogeneous system such that we can divide real and reciprocal space into orthogonal and parallel components \mathbf{R}_\parallel , R_\perp and \mathbf{k}_\parallel , k_\perp , and that we can define a parallel Brillouin zone BZ_\parallel . We then cut the device into parallel layers labeled by an index m along R_\perp such that the interaction between different layers has nearest-neighbour character. We can then define a transfer matrix Γ_m for each layer describing the connection between the layers $m - 1 \leftrightarrow m \leftrightarrow m + 1$:

$$\begin{aligned} \mathbf{H}_{m,m-1} \mathbf{C}_{m-1} + \mathbf{H}_{m,m} \mathbf{C}_m + \mathbf{H}_{m,m+1} \mathbf{C}_{m+1} &= E \mathbf{C}_m \\ \Downarrow \\ \begin{pmatrix} \mathbf{C}_{m+1} \\ \mathbf{C}_m \end{pmatrix} &= \Gamma_m \begin{pmatrix} \mathbf{C}_m \\ \mathbf{C}_{m-1} \end{pmatrix} \end{aligned} \quad (2.129)$$

The transfer matrix of the m -th layer can easily be found and reads

$$\Gamma_m = \begin{pmatrix} \mathbf{H}_{m,m+1}^{-1} (\mathbf{H}_{m,m} - E \mathbf{I}) & \mathbf{H}_{m,m+1}^{-1} \mathbf{H}_{m,m-1} \\ \mathbf{I} & \mathbf{0} \end{pmatrix} \quad (2.130)$$

The \mathbf{C}_m in the above expressions is a vector containing the expansion coefficients. Its dimension N is given by the number of orbitals per parallel unit cell that are used as basis in which the wavefunction is expanded (cf. eq. (2.112)).

The leads are treated in a similar way. In this case the layers are called *principle layers* (PL) or *superlayers* (SLAY). Additionally, we require the expansion coefficients to satisfy a Bloch condition of the form

$$\mathbf{C}_m = e^{ik_\perp d_\perp} \mathbf{C}_{m-1} \quad (2.131)$$

where d_\perp is the width of the PL. This allows the formulation of an eigenvalue problem in the leads of the form

$$\begin{aligned} \mathbf{H}_{m,m-1} \mathbf{C}_m e^{-ik_\perp d_\perp} + \mathbf{H}_{m,m} \mathbf{C}_m + \mathbf{H}_{m,m+1} \mathbf{C}_{m+1} &= E \mathbf{C}_m \\ \mathbf{C}_{m+1} &= e^{-ik_\perp d_\perp} \mathbf{C}_m \\ \Downarrow \\ \begin{pmatrix} \mathbf{H}_{m,m} - E & \mathbf{H}_{m,m+1} \\ \mathbf{I} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{C}_m \\ \mathbf{C}_{m+1} \end{pmatrix} &= e^{-ik_\perp d_\perp} \begin{pmatrix} -\mathbf{H}_{m,m-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{pmatrix} \begin{pmatrix} \mathbf{C}_m \\ \mathbf{C}_{m+1} \end{pmatrix} \end{aligned} \quad (2.132)$$

The $2N$ eigenvalues define N pairs of in general complex k_\perp , corresponding to two sets of complex eigenfunctions, one defining incoming waves and the other outgoing waves. The states with complex k_\perp are called *evanescent states* and describe states that cannot freely propagate.

Let the device extend from layer $m = 1$ to layer $m = N$ such that the first two PL of the right contact correspond to $m = N + 1$, $m = N + 2$ and of the left contact to $m = 0$, $m = -1$. By means of the device transfer matrices we can relate the right contact to the left contact by writing

$$\begin{pmatrix} \mathbf{C}_{N+2} \\ \mathbf{C}_{N+1} \end{pmatrix} = \prod_m \mathbf{\Gamma}_m \begin{pmatrix} \mathbf{C}_0 \\ \mathbf{C}_{-1} \end{pmatrix} \quad (2.133)$$

By applying the boundary conditions obtained from (2.132), the transmission coefficients for incoming waves at the left contact to outgoing waves at the right contact can be found. The sum over all possible states and integration over the parallel Brillouin zone BZ_\parallel and the energy leads then to the total current [26]

$$I = \frac{-e}{(2\pi)^3 \hbar} \int_{BZ_\parallel} d\mathbf{k}_\parallel \int_{-\infty}^{+\infty} dE \sum_{k_{\perp,j}^-, k_{\perp,i}^+} T(E, \mathbf{k}_\parallel, k_{\perp,i}^+ \rightarrow k_{\perp,j}^-) [f_R(E) - f_L(E)] \quad (2.134)$$

where we denoted the incoming and outgoing states on the left and the right lead as $k_{\perp,i}^+$ and $k_{\perp,j}^-$, respectively. Note the formal resemblance with (2.121).

The most powerful description of quantum transport phenomena is given by the Green's function approach [26, 79, 31]. The mathematical concept of the Green's function is the following. Consider a differential operator \mathcal{L} and response function u due to an excitation f such that

$$\mathcal{L}u = f, \quad (2.135)$$

then the Green's function G can be seen as the inverse of \mathcal{L} such that (in operatorial sense)

$$u = Gf = \mathcal{L}^{-1}f \quad (2.136)$$

As $\mathcal{L}G = 1$, G can be obtained as the solution of (2.135) for a point source

$$\mathcal{L}G(x, x') = \delta(x - x') \quad \Rightarrow \quad u(x) = \int G(x, x')f(x') dx' \quad (2.137)$$

Therefore, $G(x, x')$ relates a response in x to an excitation in x' . Applying this concept to the Schrödinger equation by writing $\mathcal{L} = E - H$ and thus $G = (E - H)^{-1}$, we note that the Green's function is not uniquely defined without specifying boundary conditions. These can be incorporated directly into the equation by adding an infinitesimally small positive or negative imaginary part and thus defining the so called *retarded* and *advanced* Green's functions (letting $\eta \rightarrow 0^+$)

$$G^R(x, x') = (E - H + i\eta)^{-1} \quad (2.138a)$$

$$G^A(x, x') = (E - H - i\eta)^{-1} \quad (2.138b)$$

Using these expressions we can rewrite the eigenvalue problem (2.128) for the two-lead device considered above as

$$[H_C + \Sigma_1^R + \Sigma_2^R] \Psi_C = E \Psi_C \quad (2.139)$$

where $\Sigma_n^R = H_{nC}^\dagger g_n^R H_{nC}$, being g_n^R the Green's function of the isolated n -th lead. The latter can usually be calculated easily without the need of actually inverting an infinite-dimensional matrix. We can now define the Green's function of the finite-dimensional conductor including the effects of the contacts by means of the total self-energy Σ^R as

$$G_C^R = (E - H_C - \Sigma^R)^{-1} \quad (2.140)$$

with $\Sigma^R = \sum_n \Sigma_n^R$ for non-interacting leads.

Making use of the relation between Green's functions and the S -matrix (Fisher-Lee relation, see [31]), a compact expression for the transmission function in terms of the Green's function and self-energies can be found. It reads

$$\bar{T}_{pq} = \text{Tr}[\Gamma_p G^R \Gamma_q G^A] \quad (2.141)$$

where we used the identity

$$\Gamma_p = i [\Sigma_p^R - \Sigma_p^A], \quad \Sigma_p^A = (\Sigma_p^R)^\dagger \quad (2.142)$$

Γ_p essentially describes the coupling of lead p with the conductor.

We mention that experimentally accessible quantities can be obtained from the Green's functions, such as the local density of states which is given as [31]

$$\rho(r, E) = -\frac{1}{\pi} \text{Im}[G^R(r, r; E)] \quad (2.143)$$

The self-energy in (2.140) is generally a complex quantity, therefore leading to complex eigenenergies as the Hamiltonian is no longer Hermitian. The physical meaning of this is that the states in the conductor have a finite lifetime, thus a particle in the conductor will eventually escape into one of the leads, leading to a current flow.

Although scattering inside the conductor can be taken into account to some extent in the above approach by means of additional self-energy terms [31], a correct general treatment of transport in interacting systems has to be done in the framework of non-equilibrium Green's functions (NEGF). A description of the involved theory would, however, go beyond the scope of this chapter and the reader is referred to e.g. [31, 79, 50].

Chapter 3

Numerical Implementation of the Drift-Diffusion Model

This chapter discusses the numerical implementation of the drift-diffusion model as derived in 2.2.1 using the finite element method (FEM). Only the stationary case is implemented so far in TIBERCAD, as in many cases the time dependence is of minor interest [65].

3.1 The stationary drift-diffusion equations

The stationary form of the drift-diffusion equations is found by discarding any time dependence in eqns. (2.54) and (2.55):

$$-\nabla (\epsilon \nabla \varphi - \mathbf{P}) = e (p - n + N_d^+ - N_a^-) \quad (3.1a)$$

$$\nabla \mathbf{j}_n = -R \quad (3.1b)$$

$$\nabla \mathbf{j}_p = -R \quad (3.1c)$$

The poisson equation includes spontaneous and strain induced polarization by means of the electric polarization \mathbf{P} as introduced in section 2.1.2. N_d^+ and N_a^- are the ionized donor and acceptor densities, respectively. R denotes the net recombination rate, i.e. the difference between recombination and generation rate.

Together with the constitutive equations (2.56) for the electron and hole flux

$$\mathbf{j}_n = -D_n \nabla n + \mu_n n \nabla \varphi, \quad \mathbf{j}_p = -D_p \nabla p - \mu_p p \nabla \varphi \quad (3.2)$$

the equations (3.1) form a system of three coupled partial differential equations of second order, which have to be solved for some adequate boundary conditions. In eq. (3.2) we assumed the vector potential A to be independent of time such that the electric field can be written as $\mathbf{E} = -\nabla \varphi$.

3.2 Scaling and the choice of the dependent variables

Both for analytical analysis and for numerical implementation it is necessary to scale the equations to dimensionless quantities. Although different approaches to scaling exist, the most appropriate method results from singular perturbation analysis of the semiconductor equations [96, 65]. Let Ω and $C = N_d - N_a$ denote the simulation domain and the net doping density, respectively. Then the basic parameters of the *singular perturbation* or *unit scaling* are the following:

- the characteristic device dimension $x_0 = \text{diam}(\Omega)$
- the thermal voltage $\varphi_0 = U_T = \frac{k_B T}{e}$
- the maximum doping density $C_0 = \sup_{\Omega}(C)$
- the maximum mobility $\mu_0 = \max \left(\sup_{\Omega}(\mu_n), \sup_{\Omega}(\mu_p) \right)$

Table 3.1 lists some physical quantities together with their scaling factors. The scaled quantities, denoted by a tilde over the symbol, are given by the ratio of the unscaled physical quantity and the corresponding scaling factor, e.g. the scaled electric potential reads $\tilde{\varphi} = \varphi/\varphi_0$.

<i>Symbol</i>	<i>Description</i>	<i>Scaling factor</i>
\mathbf{x}	position vector	x_0
φ, ϕ_n, ϕ_p	electric and electro-chemical potentials	φ_0
n, p	electron and hole density	C_0
μ_n, μ_p	electron and hole mobilities	μ_0
N_d, N_a	doping densities	C_0
R	recombination-generation rate	$\frac{\varphi_0 \mu_0 C_0}{x_0^2}$
J_n, J_p, J	current densities	$\frac{e \varphi_0 \mu_0 C_0}{x_0}$
\mathbf{P}	electric polarization	$e x_0 C_0$

Table 3.1: The scaling factors

Applying this scaling scheme to the equations (3.1) and using (3.2) and (2.62) we finally get the scaled system of equations to be solved numerically:

$$-\tilde{\nabla} \left(\lambda^2 \epsilon_r \tilde{\nabla} \tilde{\varphi} - \tilde{\mathbf{P}} \right) = \left(\tilde{p} - \tilde{n} + \tilde{N}_d^+ - \tilde{N}_a^- \right) \quad (3.3a)$$

$$\tilde{\nabla} \left(-\tilde{D}_n \tilde{\nabla} \tilde{n} + \tilde{\mu}_n \tilde{n} \tilde{\nabla} \tilde{\varphi} \right) = \tilde{\nabla} \left(\tilde{\mu}_n \tilde{n} \tilde{\nabla} \tilde{\phi}_n \right) = -\tilde{R} \quad (3.3b)$$

$$\tilde{\nabla} \left(-\tilde{D}_p \tilde{\nabla} \tilde{p} - \tilde{\mu}_p \tilde{p} \tilde{\nabla} \tilde{\varphi} \right) = -\tilde{\nabla} \left(\tilde{\mu}_p \tilde{p} \tilde{\nabla} \tilde{\phi}_p \right) = -\tilde{R} \quad (3.3c)$$

The ϵ_r in the Poisson equation above is the relative dielectric constant. It explicitly appears in the equation because it can be position dependent as is the case for example in heterostructures. We could introduce one more scaling factor $\epsilon_{r,0} = \sup_{\Omega}(\epsilon_r)$, but this was not done in this work. In the following we will consider only the scaled system and denote the scaled quantities with the symbol of the unscaled ones without the tilde.

The parameter λ appearing in the factor before the Laplacian of the electric potential is given by

$$\lambda = \frac{1}{x_0} \sqrt{\frac{\epsilon_0 \varphi_0}{e C_0}} \quad (3.4)$$

This value is tightly connected to the Debye length of a semiconductor which reads

$$\lambda_D = \sqrt{\frac{\epsilon_0 \epsilon_r U_T}{e C}} \quad (3.5)$$

λ acts as a singular perturbation parameter in the Poisson equation (3.3a) [66].

To numerically solve the eqns. (3.3) a set of dependent variables has to be chosen. The following three possible choices of variables for the continuity equations are usually described in literature:

- “natural” variables n, p
- Slotboom variables $v = e^{-\phi_n}, \omega = e^{\phi_p}$
- the electrochemical potentials ϕ_n, ϕ_p

In all cases the electric potential φ is used as third variable.

There is no clear answer to the question, widely discussed in literature (see e.g. [96, 16, 36, 44]), whether one set is better than the others. Often the “natural” variables n and p are used. In this case (3.3a) is linear in φ and the operators in eqns. (3.3b) and (3.3c) are linear in n and p , respectively. However, n and p usually cover ranges of many orders of magnitude. Moreover, due to the terms proportional to the gradient of the electric potential, the continuity equations are of convection-diffusion type and a maximum principle cannot directly be applied to them. This means that special care has to be taken when discretizing the equations by using e.g. some Scharfetter-Gummel type approach or upwinding techniques [93, 70, 76]. Nevertheless, n and p are most often used as variables and there is a lot of literature treating the numerics of the system (see e.g. [74, 65, 66, 54, 36, 16]).

The Slotboom variables are mostly useful for mathematical analysis as they lead to selfadjoint equations. But they are cumbersome to use in the degenerate case. Moreover, their variation is even bigger than that for n and p which makes them difficult to handle numerically.

In this work the electrochemical potentials were chosen as dependent variables for the following reasons:

1. All dependent variables are potentials and of the same order of magnitude
2. The current equations (2.62) can be used also in heterostructures, whereas the eqns. (3.2) would need an explicit modification due to the position dependent effective density of states and band edges. Moreover the generalisation to any kind of quasi-particle seems more intuitive.
3. Whereas the use of n and p implies necessarily some sort of exponential interpolation, we find a linear interpolation of ϕ_n and ϕ_p (as for the electric potential) to be a good approximation. The reasons are as follows.
 - (a) in a homogenous material and far away from interfaces where material properties change, chemical equilibrium will be reached between electrons and holes such that $np = n_i^2$, where n_i is the intrinsic density. The recombination-generation hence vanishes, and the electro-chemical potential follows the electric potential (the chemical potential is constant) which is interpolated linearly.
 - (b) in the case of minority carrier injection e.g. in pn-junctions into a doped material, the electric potential in the doped region (away from the depletion region) is approximately constant and the most important recombination process is due to trap-assisted two-particle transitions (Shockley-Read-Hall recombination) which is given by the minority carrier density divided by the carrier lifetime. An analysis of the (1D-)continuity equation in this case shows, that the minority density varies exponentially in space:

$$\begin{aligned} \frac{d}{dx} \left(D_n \frac{dn}{dx} - \overbrace{\mu_n n \frac{d\phi}{dx}}^{\approx 0} \right) &= \frac{n - n_{eq}}{\tau_n}, \quad n \propto e^{-\phi_n} \\ D_n \frac{d^2 n}{dx^2} - \frac{n - n_{eq}}{\tau_n} &= 0, \quad n(0) = n_0, \quad n(\infty) = n_{eq} \\ \Rightarrow n(x) &= (n_0 - n_{eq})e^{-x/L_n} + n_{eq}, \quad L_n = \sqrt{D_n \tau_n} \end{aligned}$$

That is for $n_0 \gg n_{eq}$ a linear variation of the electrochemical potential is an adequate approximation. Its gradient is controlled by the diffusion length L_n of the carriers.

4. The goal of the TIBERCAD project is to couple atomistic/quantum mechanical to classical drift-diffusion calculations. An elegant way to couple the two

worlds could be the use of the electrochemical potentials as boundary conditions, as the contacts of quantum regions are usually assumed to be reservoirs in equilibrium with a certain Fermi level. To use this type of coupling local equilibrium of the particle populations in the interfacing region of the two models would have to be assumed. The applicability of this assumption has to be studied yet in detail. [78]

5. The particle fluxes are described as a pure drift driven by the gradient of the electro-chemical potentials so that the continuity equations aren't anymore of convection-diffusion type. Due to this, a stable discretization of the nonlinear operators can be found, leading to an M -matrix for the discretized system (see Def. 3.3). The stability of an iterative solution algorithm involving only the system matrix, given that such an algorithm exists and is in principle stable, would then not be affected by the discretization. For example, one may use an explicit time-stepping algorithm to solve the system of equations. As the system matrix is diagonally dominant for any value of the dependent variables, it should be possible to find the minimal time step applying the usual convergence analysis (e.g. [76]).

A drawback of using the electrochemical potentials as variables is the fact that all the nonlinearities become of exponential type and the differential operators in the continuity equations get nonlinear in the electro-chemical potentials. Moreover, the latter operators are expected to be somewhat ill-conditioned due to their dependence on the particle densities.

For the discretization of the equations it will be assumed that the electrochemical potentials are continuous functions in space. This is compatible with the assumption of local equilibrium. However this becomes problematic in the case of heterostructures, where material properties change abruptly. In this case a nearly discontinuous behaviour of the electro-chemical potentials across the hetero-interfaces can be observed (depending on the device structure and operating condition), consequently leading to high gradients which can give rise to numerical complications. Similar results can be obtained also by classical kinetic emission models [80], and generally a discontinuity of the electro-chemical potentials across heterojunctions should be expected in presence of thermionic emission or tunneling [38].

3.3 The drift-diffusion equations in finite element formulation

This section treats the numerical implementation for the solution of the scaled, stationary drift-diffusion equations (3.3) in the framework of the *finite element method (FEM)*. After a short introduction to the FEM, its application to the drift-diffusion model will be described.

3.3.1 The finite element method

Consider a (linear) partial differential equation of second order $\mathcal{L}(u) = \sum_i \partial_i^2 u + b_j \partial_j u + cu + f = 0$ on $\Omega \subset \mathbb{R}^n$ with non-constant coefficients and boundary conditions on $\partial\Omega$ that has a unique solution $u \in U(\Omega)$, where $\partial_i = \partial/\partial x_i$ and $U(\Omega)$ some vector space (we will usually write $U = U(\Omega)$ in the following). Generally we cannot expect to find an analytic, closed form for u . Instead, the equation has to be solved numerically. For this purpose, as the vector space U is infinite-dimensional, the original problem needs to be restated in a finite-dimensional space $U_h = \text{span}\{\psi_i\}$. In many cases U_h is a subspace of U , i.e. $U_h \subset U$. The problem is then reduced to solving a linear system in the coefficients c_i such that $U_h \ni u_h = \sum_i c_i \psi_i$ approximates the true solution $u \in U$. This procedure can be understood as the *discretization* of the problem. In all practical methods this is accompanied with a discretization of the simulation domain Ω , i.e. the definition of a mesh or *triangulation* \mathcal{T}_h . The subscript h in \mathcal{T}_h , U_h and V_h reminds of this fact and stands at the same time for the characteristic mesh spacing of the triangulation. This notation is useful as for convergence analysis families of triangulations $\{\mathcal{T}_{h_k}\}$ with decreasing mesh spacing h_k are considered such that $h_1 > h_2 > \dots > h_k$.

There are essentially three methods that are widely used for the numerical solution of differential equations [76]:

1. The *Finite Difference Method (FDM)* approximates the differential operators on a usually rectangular grid by finite differences.
2. In the *Finite Box* or *Box Integration Method (FBM or BIM)*, the differential equation is integrated over non-overlapping regions $\mathcal{V}_i \subset \Omega$, $\bigcup \mathcal{V}_i = \Omega$ around the mesh points \mathbf{x}_i , leading to a system of equations

$$\int_{\mathcal{V}_i} (\sum_j \partial_j^2 u + b_j \partial_j u + cu + f) dx = 0$$

3. In the *Finite Element Method (FEM)* the differential equation is multiplied by a *test function* $v_h \in V_h$ and then integrated over Ω to get a system

$$\int_{\Omega} (\sum_j \partial_j^2 u + b_j \partial_j u + cu + f) v dx = 0$$

The first two methods can be regarded as special cases of the finite element method. A discretization based on the FDM is usually formulated starting from a classical boundary value problem as shown at the beginning of this section. BIM and FEM are particularly suited for equations in divergence form, and FEM is obtained in a natural way for problems based on a variational formulation.

A short introduction to the FEM shall be given in the following (for a detailed introduction and analysis see e.g. [29, 35, 27, 114], for a more pragmatic introduction [73]). For this purpose we start by stating an abstract linear variational prob-

lem, which we shall assume to have a unique solution:

$$\begin{aligned} &\text{find } u \in U \text{ such that} \\ &a(u, v) = f(v), \quad \forall v \in V \end{aligned} \tag{3.6}$$

The bilinear form $a(u, v) : U \times V \rightarrow \mathbb{R}$ and the linear form $f(v) : V \rightarrow \mathbb{R}$ are assumed to be continuous.

The *discretization* in a mathematical sense of the problem (3.6) consists in approximating it in adequate finite-dimensional spaces U_h and V_h , such that the discretized problem reads

$$\begin{aligned} &\text{find } u_h \in U_h \text{ such that} \\ &a_h(u_h, v_h) = f_h(v_h), \quad \forall v_h \in V_h \end{aligned} \tag{3.7}$$

The finite element method is essentially the construction of the latter spaces, then called *finite element spaces*. Ususally, U_h is called *solution space* and V_h *trial space*. Based on the different choices of U_h and V_h , different methods can be distinguished, given in the following definitions.

Definition 3.1. When $U_h = V_h$, the method is called standard Galerkin method. Otherwise it is called non-standard Galerkin or Petrov-Galerkin method.

Definition 3.2 (Conformity). When the finite-dimensional spaces U_h and V_h are subspaces of the respective infinite-dimensional spaces, i.e. $U_h \subset U$ and $V_h \subset V$, the method is called conformal. Otherwise it is called non-conformal.

Remark 3.1. If the bilinear form $a(u, v)$ is symmetric and positive, i.e. $a(u, v) = a(v, u) > 0$, the mathematical treatment of the problem is somewhat easier. Problem (3.6) is in this case associated to a minimization problem

$$J : v \rightarrow J(v) = \frac{1}{2}a(v, v) - f(v)$$

Problem (3.6) is then termed *variational formulation*. □

Remark 3.2. Strictly speaking, only methods where the discretized bilinear form is equal to the original one, i.e. $a_h(u, v) = a(u, v)$, are considered as conformal methods. Therefore, when $a(u, v)$ has to be approximated e.g. by numerical integration, the resulting method is often regarded as non-conformal. □

Remark 3.3. Consider the case of a conformal standard Galerkin method for a linear problem with symmetric bilinear form. Then a Hilbert space H can be defined such that the bilinear form defines a scalar product on H . Choosing $U_h = V_h \subset H$ it can be seen that $u_h \in V_h$ is the projection of u onto the finite dimensional subspace V_h and the approximation error $u - u_h$ lies in the orthogonal complement of V_h in H . In other words, the discrete solution u_h is the best-approximating element in V_h to u . □

As a practical example for the abstract problem (3.6) we consider the following forms for $a(u, v)$ and $f(v)$

$$a(u, v) = \int_{\Omega} \nabla u \cdot \nabla v \, dx, \quad f(v) = \int_{\Omega} f(x) v \, dx \quad (3.8)$$

such that problem (3.6) becomes

$$\begin{aligned} &\text{find } u \in U \text{ such that} \\ &\int_{\Omega} \nabla u \cdot \nabla v \, dx = \int_{\Omega} f(x) v \, dx, \quad \forall v \in V \end{aligned} \quad (3.9)$$

Without considering the mathematical intricacies, we formally apply Green's formula to the above equation to get

$$-\int_{\Omega} (\Delta u + f) v \, dx + \int_{\partial\Omega} (\nabla u \cdot \boldsymbol{\nu}) v \, dx = 0, \quad \forall v \in V \quad (3.10)$$

where $\boldsymbol{\nu}$ denotes the outward normal on $\partial\Omega$. Thus we are led to the conclusion that by solving the problem (3.9) one is formally solving the associated boundary value problem

$$\begin{cases} -\Delta u = f(x) & \text{on } \Omega \\ \nabla u \cdot \boldsymbol{\nu} = 0 & \text{on } \partial\Omega \end{cases} \quad (3.11)$$

We can also invert the argumentation and start from this boundary value problem [40]. Then, applying Green's formula, we are led to (3.9) which is then called the *weak form* of (3.11). The solution u is accordingly called *weak* or *generalized solution*. The approach as described before is especially suited for problems where the principal part of the differential operator is in divergence form, such as

$$\mathcal{L}u = -\partial_i (a^{ij} \partial_j u + b^i u) + c^i \partial_i u + du \quad (3.12)$$

the weak form of it reading (under rather weak smoothness assumptions for the coefficients)

$$\mathcal{L}(u, v) = \int_{\Omega} [(a^{ij} \partial_j u + b^i u) \partial_i v + (c^i \partial_i u + du) v] \, dx, \quad \forall v \in C_0^1(\Omega) \quad (3.13)$$

where $C_0^1(\Omega)$ is the space of continuously differentiable functions vanishing on $\partial\Omega$.

A boundary value problem in weak form allows for solutions from a broader class of functions than its classical ("strong") counterpart. In particular, the differentiations in (3.12) can be understood in a weak (distributional) sense. The general theory for weakly differentiable functions is naturally formulated in the Sobolev spaces $W^{m,p}$ and $W_0^{m,p}$, and the most useful spaces for the treatment of second-order partial differential equations are the Hilbert spaces $H^1 = W^{1,2}$ and $H_0^1 = W_0^{1,2}$ (cf. [40, 29, 35], [1] for the theory on Sobolev spaces).

We are now able to write down a “recipe” for the construction of a finite element formulation for some (linear) boundary value problem

$$\begin{cases} \mathcal{L}u = 0 & \text{on } \Omega \\ \nabla u \cdot \boldsymbol{\nu} = g(x) & \text{on } \Gamma_N \subset \partial\Omega \\ u = f(x) & \text{on } \Gamma_D \subset \partial\Omega, \Gamma_D \cap \Gamma_N = \emptyset \end{cases} \quad (3.14)$$

1. Based on the form of \mathcal{L} and on the boundary conditions, choose appropriate spaces U and V
2. Choose appropriate finite element spaces $U_h = \text{span}\{\varphi_i\}$ and $V_h = \text{span}\{\psi_i\}$
3. Multiply (3.14) by a trial function $v_h = \psi_i \in V_h$ and integrate over Ω
4. Apply Green’s formula to get rid of second derivatives
5. Expand u_h as $c^j \varphi_j$ and write down the resulting linear algebraic system

Applying this procedure to our model problem (3.11) leads to

$$c^j \int_{\Omega} \nabla \psi_i \cdot \nabla \varphi_j \, dx = \int_{\Omega} f(x) \psi_i \, dx \quad (3.15)$$

which can be written as a linear algebraic equation for the expansion coefficients c^j

$$\begin{aligned} K_{ij} c^j &= f_i \\ \text{with } K_{ij} &= \int_{\Omega} \nabla \psi_i \cdot \nabla \varphi_j \, dx, \quad f_i = \int_{\Omega} f(x) \psi_i \, dx \end{aligned} \quad (3.16)$$

This is formally nothing else than a generalized Fourier expansion of (3.14).

The construction of finite element spaces always includes the following three basic aspects, which distinguish the finite element method from other approaches and which are fundamental for the numerical analysis and implementation [29]:

(FEM 1) A triangulation \mathcal{T}_h is defined on the set $\bar{\Omega}$, i.e. $\bar{\Omega}$ is subdivided into a finite number of subsets $\mathcal{K} \in \mathcal{T}_h$ such that

- each \mathcal{K} is closed and its interior $\mathring{\mathcal{K}}$ non empty
- $\bigcup_{\mathcal{K}} \mathcal{K} = \bar{\Omega}$, i.e. the subdivision \mathcal{K} completely covers the set $\bar{\Omega}$
- for $\mathcal{K}_1 \neq \mathcal{K}_2$, $\mathring{\mathcal{K}}_1 \cap \mathring{\mathcal{K}}_2 = \emptyset$
- the boundary $\partial\mathcal{K}$ of each \mathcal{K} is Lipschitz-continuous (see [29])

The subsets \mathcal{K} are called *finite elements*.

(FEM 2) Let $P_{\mathcal{K}}$ denote the space spanned by the restriction of the basis $v \in V_h$ to \mathcal{K} , i.e. $P_{\mathcal{K}} = \text{span}\{v_h|_{\mathcal{K}}; v_h \in V_h\}$. The spaces $P_{\mathcal{K}}$ shall contain polynomials or functions “close” to polynomials.

(FEM 3) There exists a basis in V_h and U_h formed by functions that can be easily described and have supports as small as possible.

The aspect (FEM 1) is simply the mathematical formulation for the fact, that the simulation domain Ω is discretized in a mesh formed by non-overlapping elements covering the whole domain and having certain regularity properties.

(FEM 2) assures a simple form for the integrals that have to be calculated and is used for convergence analysis.

(FEM 3) finally assures that the finite element basis has near orthogonal properties, i.e. two basis functions have only overlapping support when they are associated to neighbouring nodes. Usually, “neighbouring” in this context means “located on the same element”. This is comparable to the nearest neighbour approximation in localized basis methods such as the tight-binding method (cf. section 2.4.2) and assures that the resulting matrices are sparse.

The main advantages of FEM over other methods are on the one hand it’s ability to treat almost arbitrarily complex geometries due to the first aspect (FEM 1) of finite element space construction (although the same is valid for BIM). On the other hand, it has a very sound mathematical foundation which allows for clear convergence and error analysis. The latter in particular opens the possibility for use of adaptive mesh refinement schemes. In addition, the integrals over the domain Ω can be decomposed into a sum of integrals over the finite elements \mathcal{K}_l

$$K_{ij} = \sum_{\mathcal{K}_l \in \mathcal{T}_h} \int_{\mathcal{K}_l} \nabla \psi_i \cdot \nabla \varphi_j \, dx, \quad f_i = \sum_{\mathcal{K}_l \in \mathcal{T}_h} \int_{\mathcal{K}_l} f(x) \psi_i \, dx \quad (3.17)$$

and the near-orthogonality assures that only a small number of basis functions (usually only the ones that are associated to a node of the element to integrate over) lead to non-vanishing integrals.

The integrations (3.17) are usually not calculated directly in the original coordinates x . Each finite element $\mathcal{K} \in \mathcal{T}_h$ is mapped onto an equivalent *reference element* instead, on which the calculations can be done in an easier way. The linear mapping for an element \mathcal{K} is defined by

$$\mathcal{F}_{\mathcal{K}} : \xi \in \mathbb{R}^n \mapsto x \in \mathcal{K} \quad (3.18)$$

Fig. 3.1 illustrates the triangular case. In many cases the above transformation is affine, e.g. for triangles or rectangles. It then can be written as $x = B_{\mathcal{K}}\xi + b_{\mathcal{K}}$, where $B_{\mathcal{K}}$ and $b_{\mathcal{K}}$ are an invertible $n \times n$ matrix and a vector in \mathbb{R}^n , respectively.

General details about the numerical implementation of FEM will be given in section 3.3.3, if needed for our application. In the next section we will rewrite the drift-diffusion equations in weak form, analyze some basic mathematical properties and describe the solution approach adopted in TIBERCAD.

3.3.2 The drift-diffusion equations in weak form

The aim of this section is to reformulate the system of equations (3.3), restated below for clarity, in weak form to analyze some of its properties, to illustrate it’s

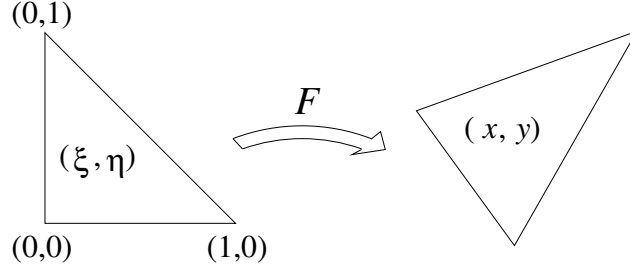


Figure 3.1: The reference element for a triangular finite element.

solution using a Newton-like approach and to prepare the application of the finite element method:

$$-\nabla [\lambda^2 \epsilon_r(\mathbf{r}) \nabla \varphi - \mathbf{P}(\mathbf{r})] = \rho(\mathbf{r}, \varphi, \phi_n, \phi_p) \quad (3.19a)$$

$$-\nabla [\mu_n(\mathbf{r}) n(\mathbf{r}, \varphi, \phi_n) \nabla \phi_n] = R(\mathbf{r}, \varphi, \phi_n, \phi_p) \quad (3.19b)$$

$$-\nabla [\mu_p(\mathbf{r}) p(\mathbf{r}, \varphi, \phi_p) \nabla \phi_p] = -R(\mathbf{r}, \varphi, \phi_n, \phi_p) \quad (3.19c)$$

where $\rho = p - n + N_d^+ - N_a^-$ is the scaled total charge density. In the above equations we evidenced the dependence of the different parameters on the variables to point out the nonlinear character of the system. In the following we will write x instead of \mathbf{r} for the space coordinates.

Remark 3.4. All equations in (3.19) are in divergence form. The Poisson equation (3.19a) is *semilinear* in φ and equations (3.19b, 3.19c) are *quasi-linear* in ϕ_n and ϕ_p , respectively. \square

We now restate the system (3.19) in weak form according to what has been introduced in section 3.3.1 by multiplying each of the equations by a test function $v \in H^1(\Omega)$ and integrating over the domain Ω . To not overload the equations, we will omit from now on the arguments in the different functions. We get, after applying Green's formula

$$\int_{\Omega} (\lambda^2 \epsilon_r \nabla \varphi - \mathbf{P}) \cdot \nabla v \, dx - \int_{\partial\Omega} (\lambda^2 \epsilon_r \nabla \varphi - \mathbf{P}) \cdot \boldsymbol{\nu} v \, dx = \int_{\Omega} \rho v \, dx \quad (3.20a)$$

$$\int_{\Omega} (\mu_n n \nabla \phi_n) \cdot \nabla v \, dx - \int_{\partial\Omega} (\mu_n n \nabla \phi_n) \cdot \boldsymbol{\nu} v \, dx = \int_{\Omega} R v \, dx \quad (3.20b)$$

$$\int_{\Omega} (\mu_p p \nabla \phi_p) \cdot \nabla v \, dx - \int_{\partial\Omega} (\mu_p p \nabla \phi_p) \cdot \boldsymbol{\nu} v \, dx = - \int_{\Omega} R v \, dx \quad (3.20c)$$

where $\boldsymbol{\nu}$ is the outer normal on $\partial\Omega$ as in the last section. The boundary integrals will serve to impose von Neumann or mixed-type boundary conditions. For the following, we assume boundary conditions such that the boundary integrals vanish, corresponding to homogeneous von Neumann boundary conditions. Other boundary conditions will be treated later on. We note, however, that at least one Dirichlet

boundary condition has to be imposed on a part of the boundary $\Gamma_D \subset \partial\Omega$ with strictly positive measure. So we complete the system (3.20) with some Dirichlet boundary conditions

$$\varphi|_{\Gamma_D} = \varphi^D, \quad \phi_n|_{\Gamma'_D} = \phi_n^D, \quad \phi_p|_{\Gamma''_D} = \phi_p^D \quad (3.21)$$

From the system (3.20) we can identify the forms $a(u, v)$ and $f(v)$ introduced in the last section as listed here:

$$\begin{aligned} a^\varphi(\varphi, v) &= \lambda^2 \int_{\Omega} (\epsilon_r \nabla \varphi) \cdot \nabla v \, dx, & f^\varphi(v) &= \int_{\Omega} (\rho v + \mathbf{P} \cdot \nabla v) \, dx \doteq f(v) \\ a^{\phi_n}(\phi_n; \phi_n, v) &= \int_{\Omega} (\mu_n n \nabla \phi_n) \cdot \nabla v \, dx, & f^{\phi_n}(v) &= \int_{\Omega} R v \, dx \doteq g(v) \\ a^{\phi_p}(\phi_p; \phi_p, v) &= \int_{\Omega} (\mu_p p \nabla \phi_p) \cdot \nabla v \, dx, & f^{\phi_p}(v) &= -f^{\phi_n}(v) = -g(v) \end{aligned}$$

It has to be noted that only $a^\varphi(\varphi, v)$ is a symmetric bilinear form, whereas the forms $a^{\phi_n}(\phi_n; \phi_n, v)$ and $a^{\phi_p}(\phi_p; \phi_p, v)$ are non-symmetric and nonlinear, which is evidenced using a special notation. The $f^\bullet(v)$ are obviously all nonlinear in φ , ϕ_n and ϕ_p .

Before considering possible solution methods we shall state some properties of the system (3.20).¹ From (3.19) we can easily see that we are faced with a (quasi-linear) system of equations in divergence form. The single equations can be written in the form

$$Q(u, v) \doteq \int_{\Omega} [\mathbf{A}(x, u, \nabla u) \nabla v - B(x, u) v] \, dx = 0 \quad (3.22)$$

We make the following

Assumption 3.1.

(i) We assume for the moment that the mobilities depend only on position, excluding especially velocity saturation models which would introduce some mathematical problems [65]. Furthermore, μ_n and μ_p shall be essentially bounded away from zero, i.e.

$$0 < \underline{\mu}_{n,p} \leq \mu_{n,p} \leq \bar{\mu}_{n,p} < \infty \quad (3.23)$$

(ii) The function B shall not depend on the gradients of the solutions, assuming $\rho = \rho(x, \varphi, \phi_n, \phi_p)$ and $R = R(x, n, p) = g(x, n, p)(np - n_i^2)$ with $g \geq 0$. Note that the latter excludes generation by impact ionization from our mathematical analysis. ρ and R shall be continuously differentiable with respect to φ, ϕ_n, ϕ_p

(iii) The Dirichlet boundary conditions (3.21) shall be essentially bounded, i.e.

$$(\varphi^D, \phi_n^D, \phi_p^D) \in (L^\infty(\partial\Omega))^3 \quad (3.24)$$

□

¹Detailed mathematical analyses of the drift-diffusion system can be found amongst others in in [74, 66, 65].

In principle there would be other assumptions, e.g. about the regularity of the domain Ω and the doping profile, which are not explicitly stated here [65].

Assumption (i) assures that the operators are uniformly elliptic. Assumption (iii) assures that the solution $(\varphi, \phi_n, \phi_p) \in (H^1)^3$ we are seeking is essentially bounded and therefore lies in $(H^1 \cap L^\infty)^3$, which is in fact necessary for a physically meaningful solution.

We rewrite \mathbf{A} as

$$\begin{aligned} \mathbf{A}(x, u, \nabla u) &= \alpha(x, u) \nabla u \\ &\text{with} \\ \alpha^\varphi(x, \varphi) &= \epsilon_r(x) \\ \alpha^{\phi_n}(x, \phi_n) &= \mu_n(x) n(x, \phi_n) \\ \alpha^{\phi_p}(x, \phi_p) &= \mu_p(x) p(x, \phi_p) \end{aligned} \tag{3.25}$$

From the above it is concluded that (i) $0 < \underline{\alpha} \leq \alpha(x, u) \leq \bar{\alpha} < \infty$ in Ω , assuming $u \in H^1 \cap L^\infty$ and $\mu_{n,p}, \epsilon_r > 0$, and (ii) $\mathbf{A}(x, u, \nabla u)$ and $B(x, u)$ are continuously (Fréchet) differentiable with respect to u and ∇u , which can be easily seen by examining their functional form. It follows from (i) that all Q are uniformly elliptic.

For the following description we decouple the system and consider every equation on its own. The case of the semilinear Poisson equation is rather simple. Noting that $B = \rho$ and using the usual expressions for the densities of free carriers and ionized dopants it follows that B is monotonically decreasing in φ , i.e. $\partial_\varphi \rho < 0$ for any given ϕ_n, ϕ_p . With this a comparison principle for the Poisson equation can be derived which assures existence and uniqueness of the solution (alternatively, upper and lower bounds for φ can be found which correspond to super- and subsolutions of the equation and for uniqueness use the maximum principle for $g = \varphi_1 - \varphi_2$, see [65]). In a similar way a comparison principle can be obtained also for the continuity equations, provided that the recombination rates satisfy certain hypotheses (see Appendix C for some more details). It should be noted, however, that although the single equations satisfy a comparison principle, this is by no means automatically the case also for the system of equations [21].

The fact that the continuous problem meets a comparison or a maximum principle should be reflected in the discretized system in order to obtain a numerically stable discrete problem. The discrete analogue of a comparison or maximum principle is the M -matrix property [69, 65]:

Definition 3.3. An M -matrix is a nonsingular matrix \mathbf{A} with $(\mathbf{A})_{ij} \leq 0, \forall i \neq j$ and $(\mathbf{A}^{-1})_{ij} \geq 0$.

A discretization which does not produce an M -matrix can lead to unphysical oscillations in the solution. This is for example the case for a standard finite difference or finite element discretization of the continuity equations when the densities are used as variables (see e.g. [72, 94]).

To study the solution approaches for the nonlinear system (3.20) we rewrite it in a different, more compact form as

$$\begin{aligned} \mathbf{F}(\varphi, \phi_n, \phi_p) &= \begin{pmatrix} F^\varphi(\varphi, \phi_n, \phi_p) \\ F^{\phi_n}(\varphi, \phi_n, \phi_p) \\ F^{\phi_p}(\varphi, \phi_n, \phi_p) \end{pmatrix} = \\ &\begin{pmatrix} a^\varphi(\varphi, v) - f(v) \\ a^{\phi_n}(\phi_n; \phi_n, v) - g(v) \\ a^{\phi_p}(\phi_p; \phi_p, v) + g(v) \end{pmatrix} = 0 \end{aligned} \quad (3.26)$$

i.e. we are seeking the root of a nonlinear system of equations. This system can formally be solved using a Newton method [15, 88]. Starting from an initial guess $\varphi^{(0)}$, $\phi_n^{(0)}$ and $\phi_p^{(0)}$, we try to find corrections δu , δv , and δw such that

$$\mathbf{F}(\varphi^{(0)} + \delta u, \phi_n^{(0)} + \delta v, \phi_p^{(0)} + \delta w) = 0 \quad (3.27)$$

We linearize \mathbf{F} around $(\varphi^{(0)}, \phi_n^{(0)}, \phi_p^{(0)})$ and write

$$\mathbf{F}^{(0)} + \mathbf{F}_\varphi \delta u + \mathbf{F}_{\phi_n} \delta v + \mathbf{F}_{\phi_p} \delta w = 0 \quad (3.28)$$

where $\mathbf{F}^{(0)}$ is short for $\mathbf{F}(\varphi^{(0)}, \phi_n^{(0)}, \phi_p^{(0)})$ and \mathbf{F}_φ , \mathbf{F}_{ϕ_n} and \mathbf{F}_{ϕ_p} are the Fréchet derivatives of \mathbf{F} at $(\varphi^{(0)}, \phi_n^{(0)}, \phi_p^{(0)})$, formally given by e.g. $\mathbf{F}_\varphi = \partial \mathbf{F} / \partial \varphi$ [40]. Note that the latter are linear operators acting on δu , δv and δw . The solution of eq. (3.28) gives a new guess, which hopefully is better than the old one. The method can formally be written as, using (3.26) and using the index k to enumerate the iterations

$$\begin{pmatrix} F_\varphi^\varphi & F_{\phi_n}^\varphi & F_{\phi_p}^\varphi \\ F_\varphi^{\phi_n} & F_{\phi_n}^{\phi_n} & F_{\phi_p}^{\phi_n} \\ F_\varphi^{\phi_p} & F_{\phi_n}^{\phi_p} & F_{\phi_p}^{\phi_p} \end{pmatrix}^{(k-1)} \begin{pmatrix} \delta u \\ \delta v \\ \delta w \end{pmatrix}^{(k)} = - \begin{pmatrix} F^\varphi \\ F^{\phi_n} \\ F^{\phi_p} \end{pmatrix}^{(k-1)} \quad (3.29)$$

with the update rule

$$\varphi^{(k)} = \varphi^{(k-1)} + \delta u^{(k)}, \quad \phi_n^{(k)} = \phi_n^{(k-1)} + \delta v^{(k)}, \quad \phi_p^{(k)} = \phi_p^{(k-1)} + \delta w^{(k)} \quad (3.30)$$

The matrix in (3.29) is the *Jacobian* of the system (3.26), and the vector F is called the *residual*. Note once again that (3.29) is only a formal notation, as the Jacobian is a matrix of linear operators acting on the elements of the update vector $(\delta u, \delta v, \delta w)^{(k)}$.

The Newton method is known to show quadratic convergence, provided the starting guess is sufficiently near to the true solution. This can be a problem, as in most cases one cannot guarantee to find such a guess. However there are modified approximate Newton methods that are globally convergent [15].

Principally there are two possibilities how to discretize eq. (3.27). Either the residual is discretized and then the Jacobian is calculated from the discretized system, or eq. (3.29), i.e. the Jacobian in operator form is discretized. The second

method is preferable, as it allows to see the operatorial form of the Jacobian and therefore the appropriate type of discretization can be chosen [16]. Moreover, it can suggest block-iterative solution approaches which cannot be easily seen in the other way. We therefore write down explicitly the Jacobian in operatorial form to be able to analyze its mathematical and numerical properties.

Using (3.20) and (3.26) we get for the Jacobian (omitting the boundary terms)

$$\mathbf{J} = \begin{pmatrix} a_\varphi^\varphi & 0 & 0 \\ a_\varphi^{\phi_n} & a_{\phi_n}^{\phi_n} & 0 \\ a_\varphi^{\phi_p} & 0 & a_{\phi_p}^{\phi_p} \end{pmatrix} + \begin{pmatrix} -f_\varphi & -f_{\phi_n} & -f_{\phi_p} \\ -g_\varphi & -g_{\phi_n} & -g_{\phi_p} \\ g_\varphi & g_{\phi_n} & g_{\phi_p} \end{pmatrix} \quad (3.31)$$

To simplify, we divide it into two parts, one originating from the right hand side and the other one from the left hand side. The Jacobian has an obvious block structure, and we will in the following explicitly write down every block.

Before doing so, we first remember that $\rho = p - n + N_d^+ - N_a^-$. Next, we make a few assumptions about the carrier densities and the ionized dopant densities. From section 2.2.1.1 we already know, that n and p are of the form $n = n(\varphi - \phi_n)$ and $p = p(\varphi - \phi_p)$ and therefore $\partial_{\phi_n} n = -\partial_\varphi n$ and $\partial_{\phi_p} p = -\partial_\varphi p$, respectively. Assuming for the densities of ionized dopants the following expressions [96, 56]

$$N_d^+ = \frac{N_d}{1 + g_d \exp\left(\frac{\varphi - \phi_n - E_c + \Delta E_d}{k_B T}\right)} \quad (3.32a)$$

$$N_a^- = \frac{N_a}{1 + g_a \exp\left(\frac{\phi_p - \varphi + E_v + \Delta E_a}{k_B T}\right)} \quad (3.32b)$$

we also get $\partial_{\phi_n} N_d^+ = -\partial_\varphi N_d^+$ and $\partial_{\phi_p} N_a^- = -\partial_\varphi N_a^-$.

Using the above assumptions, the first part of the Jacobian (3.31) gets

$$a_\varphi^\varphi \delta u = \lambda^2 \int_\Omega (\epsilon_r \nabla \delta u) \cdot \nabla v \, dx \quad (3.33a)$$

$$a_\varphi^{\phi_n} \delta u = \int_\Omega \left(\mu_n \frac{\partial n}{\partial \varphi} \delta u \nabla \phi_n \right) \cdot \nabla v \, dx \quad (3.33b)$$

$$a_{\phi_n}^{\phi_n} \delta v = \int_\Omega \left(\mu_n n \nabla \delta v - \mu_n \frac{\partial n}{\partial \varphi} \delta v \nabla \phi_n \right) \cdot \nabla v \, dx \quad (3.33c)$$

$$a_\varphi^{\phi_p} \delta u = \int_\Omega \left(\mu_p \frac{\partial p}{\partial \varphi} \delta u \nabla \phi_p \right) \cdot \nabla v \, dx \quad (3.33d)$$

$$a_{\phi_p}^{\phi_p} \delta w = \int_\Omega \left(\mu_p p \nabla \delta w - \mu_p \frac{\partial p}{\partial \varphi} \delta w \nabla \phi_p \right) \cdot \nabla v \, dx \quad (3.33e)$$

For the second part we obtain

$$f_\varphi \delta u = \int_{\Omega} \frac{\partial \rho}{\partial \varphi} \delta u v \, dx = \int_{\Omega} \left(\frac{\partial p}{\partial \varphi} - \frac{\partial n}{\partial \varphi} + \frac{\partial N_d^+}{\partial \varphi} - \frac{\partial N_a^-}{\partial \varphi} \right) \delta u v \, dx \quad (3.34a)$$

$$f_{\phi_n} \delta v = \int_{\Omega} \frac{\partial \rho}{\partial \phi_n} \delta v v \, dx = \int_{\Omega} \left(\frac{\partial n}{\partial \varphi} - \frac{\partial N_d^+}{\partial \varphi} \right) \delta u v \, dx \quad (3.34b)$$

$$f_{\phi_p} \delta w = \int_{\Omega} \frac{\partial \rho}{\partial \phi_p} \delta w v \, dx = \int_{\Omega} \left(-\frac{\partial p}{\partial \varphi} + \frac{\partial N_a^-}{\partial \varphi} \right) \delta w v \, dx \quad (3.34c)$$

$$g_\varphi \delta u = \int_{\Omega} \frac{\partial R}{\partial \varphi} \delta u v \, dx = \int_{\Omega} \left(\frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} + \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \right) \delta u v \, dx \quad (3.34d)$$

$$g_{\phi_n} \delta v = \int_{\Omega} \frac{\partial R}{\partial \phi_n} \delta v v \, dx = - \int_{\Omega} \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \delta v v \, dx \quad (3.34e)$$

$$g_{\phi_p} \delta w = \int_{\Omega} \frac{\partial R}{\partial \phi_p} \delta w v \, dx = - \int_{\Omega} \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \delta w v \, dx \quad (3.34f)$$

Note that $f_\varphi + f_{\phi_n} + f_{\phi_p} = 0$ and $g_\varphi + g_{\phi_n} + g_{\phi_p} = 0$.

We finally write down the first part of the Jacobian for the case of Boltzmann statistics, where $\partial_\varphi n = n$ (in scaled quantities), which allows to rewrite some of the expressions (3.33). We use the symbol \bullet as a placeholder for the arguments of the linear operators.

$$\begin{aligned} \mathbf{J}_a &= \begin{pmatrix} \lambda^2 \int_{\Omega} (\epsilon_r \nabla \bullet) \cdot \nabla v \, dx & 0 & 0 \\ \int_{\Omega} (\mu_n n \bullet \nabla \phi_n) \cdot \nabla v \, dx & \int_{\Omega} (\mu_n n \nabla \bullet - \mu_n n \bullet \nabla \phi_n) \cdot \nabla v \, dx & 0 \\ \int_{\Omega} (-\mu_p p \bullet \nabla \phi_p) \cdot \nabla v \, dx & 0 & \int_{\Omega} (\mu_p p \nabla \bullet + \mu_p p \bullet \nabla \phi_p) \cdot \nabla v \, dx \end{pmatrix} \\ &= \begin{pmatrix} \lambda^2 \int_{\Omega} (\epsilon_r \nabla \bullet) \cdot \nabla v \, dx & 0 & 0 \\ \int_{\Omega} (\mathbf{j}_n \cdot \nabla v) \bullet \, dx & \int_{\Omega} (\mu_n n \nabla \bullet - \bullet \mathbf{j}_n) \cdot \nabla v \, dx & 0 \\ \int_{\Omega} (\mathbf{j}_p \cdot \nabla v) \bullet \, dx & 0 & \int_{\Omega} (\mu_p p \nabla \bullet - \bullet \mathbf{j}_p) \cdot \nabla v \, dx \end{pmatrix} \end{aligned} \quad (3.35)$$

For comparison and future reference we write down the corresponding matrix for the case where the densities are used as dependent variables (see [4]):

$$\mathbf{J}'_a = \begin{pmatrix} \lambda^2 \int_{\Omega} (\epsilon_r \nabla \bullet) \cdot \nabla v \, dx & 0 & 0 \\ \int_{\Omega} (\mu_n n \nabla \bullet) \cdot \nabla v \, dx & \int_{\Omega} \mu_n (\nabla \bullet + \bullet E) \cdot \nabla v \, dx & 0 \\ \int_{\Omega} (\mu_p p \nabla \bullet) \cdot \nabla v \, dx & 0 & \int_{\Omega} \mu_p (\nabla \bullet - \bullet E) \cdot \nabla v \, dx \end{pmatrix} \quad (3.36)$$

We notice quite a similar structure with respect to the diagonal blocks in the sense that in both cases the blocks corresponding to the linearized continuity equations are of convection-diffusion form. This fact will be considered in the next section when discretizing the Jacobian.

For the sake of completeness, we write down the complete Jacobian on p. 61.

3.3.3 Application of FEM to the drift-diffusion equations

The application of the finite element method to the semiconductor equations dates back to the early nineteen-eighties [96, 71]. It was not widely used, however, as the standard Galerkin method applied to the continuity equations with the densities as dependent variables leads to an unstable system for not sufficiently refined meshes (as is the case for a standard finite difference discretization). This is due to the convection-diffusion nature of the equations, whose standard discretization fails to satisfy a discrete maximum principle, and is reflected in the fact, that the system matrix is not an M -matrix, unless the discretization mesh is very fine in regions of high electric field. The consequence of this are spurious oscillations, which can completely spoil the solution. This situation can be very easily demonstrated in a 1D test case [76, 96, 65, 94]. For this reason, FEM was usually considered as unsuccessful for semiconductor device simulation.

Subsequently, however, many methods have been proposed to get a stable FEM discretization by generalizing the Scharfetter-Gummel approach to a FEM setup [109, 23, 43, 102] or by using upwinding schemes [70, 103, 76]. In the mathematical community especially mixed finite element methods are considered, in which the equations of second order are written as a system of first order equations by introducing the particle fluxes as additional variables [22, 72].

We shall apply Galerkin's method to the linearized system (3.29). From the considerations made in section 3.2 we choose a standard Galerkin approach with H^1 -conforming Lagrange elements. The basis functions, denoted by $\psi_i(\mathbf{x})$, in this case are the piecewise linear functions

$$\left\{ \psi_i \in H^1 \mid \psi_i(\mathbf{x}_j) = \delta_{ij}, \text{ supp}\{\psi_i\} = \bigcup_{\substack{j \\ x_i \in K_j}} K_j \in \mathcal{T}_h \right\}, \quad (3.37)$$

$$\mathbf{J} = \begin{pmatrix} \int_{\Omega} \left[\lambda^2 (\epsilon_r \nabla \bullet) \cdot \nabla v - \left(\frac{\partial p}{\partial \varphi} - \frac{\partial n}{\partial \varphi} + \frac{\partial N_d^+}{\partial \varphi} - \frac{\partial N_a^-}{\partial \varphi} \right) \bullet v \right] dx & - \int_{\Omega} \left(\frac{\partial n}{\partial \varphi} - \frac{\partial N_d^+}{\partial \varphi} \right) \bullet v dx & - \int_{\Omega} \left(-\frac{\partial p}{\partial \varphi} + \frac{\partial N_a^-}{\partial \varphi} \right) \bullet v dx \\ \int_{\Omega} \left[(\mu_n n \bullet \nabla \phi_n) \cdot \nabla v - \left(\frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} + \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \right) \bullet v \right] dx & \int_{\Omega} \left[(\mu_n n \nabla \bullet - \mu_n n \bullet \nabla \phi_n) \cdot \nabla v + \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \bullet v \right] dx & \int_{\Omega} \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \bullet v dx \\ \int_{\Omega} \left[(-\mu_p p \bullet \nabla \phi_p) \cdot \nabla v + \left(\frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} + \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \right) \bullet v \right] dx & - \int_{\Omega} \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \bullet v dx & \int_{\Omega} \left[(\mu_p p \nabla \bullet + \mu_p p \bullet \nabla \phi_p) \cdot \nabla v - \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \bullet v \right] dx \end{pmatrix}$$

i.e. they are continuous across the element boundaries and they are one at the node they are associated with and zero on any other node. They clearly have “small” support in the sense of aspect (FEM 3) described in the last section. These functions, also called *hat functions*, are illustrated in Fig. 3.2. In the 1D-case, the

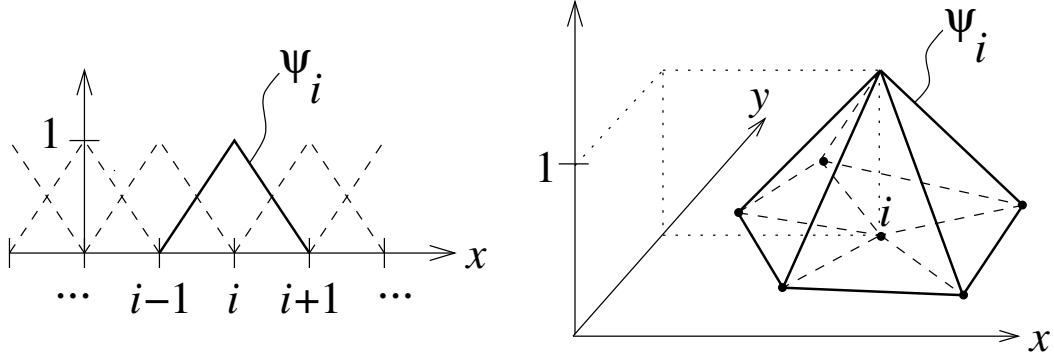


Figure 3.2: The basis functions of 1st-order Lagrange elements in 1D and 2D (hat functions).

finite elements are just the line segments connecting two neighbouring nodes. In the 2D-case in this example, they are formed by triangles, and the support for the i -th basis function is the union of the triangles containing the i -th node.

The expansions for the dependent variables in terms of the basis function is written as

$$\varphi_h(\mathbf{x}) = u^i \psi_i(\mathbf{x}) \quad (3.38a)$$

$$\phi_{n,h}(\mathbf{x}) = v^i \psi_i(\mathbf{x}) \quad (3.38b)$$

$$\phi_{p,h}(\mathbf{x}) = w^i \psi_i(\mathbf{x}) \quad (3.38c)$$

where we used the Einstein summing convention for ease of notation. With the basis (3.37) the degrees of freedom (i.e. the unknowns of the discretized system u^i , v^i and w^i) coincide with the values of the dependent variables at the mesh nodes.

Following Galerkin's method we choose for the test functions $v_h = \psi_i$. First, we calculate the residual of the discretized system:

$$F_h^\varphi = \int_{\Omega} (\lambda^2 \epsilon_r \nabla \varphi_h - \mathbf{P}) \cdot \nabla \psi_i \, dx - \int_{\partial\Omega} (\lambda^2 \epsilon_r \nabla \varphi - \mathbf{P}) \cdot \boldsymbol{\nu} \psi_i \, dx - \int_{\Omega} \rho_h \psi_i \, dx \quad (3.39a)$$

$$F_h^{\phi_n} = \int_{\Omega} (\mu_{n,h} n_h \nabla \phi_{n,h}) \cdot \nabla \psi_i \, dx - \int_{\partial\Omega} (\mu_n n \nabla \phi_n) \cdot \boldsymbol{\nu} \psi_i \, dx - \int_{\Omega} R_h \psi_i \, dx \quad (3.39b)$$

$$F_h^{\phi_p} = \int_{\Omega} (\mu_{p,h} p_h \nabla \phi_{p,h}) \cdot \nabla \psi_i \, dx - \int_{\partial\Omega} (\mu_p p \nabla \phi_p) \cdot \boldsymbol{\nu} \psi_i \, dx + \int_{\Omega} R_h \psi_i \, dx \quad (3.39c)$$

or, symbolically,

$$F_h^\varphi = a^\varphi(\varphi_h, \psi_i) - \int_{\partial\Omega} (\cdots) \psi_i \, dx - f_h(\psi_i) \quad (3.40a)$$

$$F_h^{\phi_n} = a_h^{\phi_n}(\phi_{n,h}; \phi_{n,h}, \psi_i) - \int_{\partial\Omega} (\cdots) \psi_i \, dx - g_h(\psi_i) \quad (3.40b)$$

$$F_h^{\phi_p} = a_h^{\phi_p}(\phi_{p,h}; \phi_{p,h}, \psi_i) - \int_{\partial\Omega} (\cdots) \psi_i \, dx + g_h(\psi_i) \quad (3.40c)$$

Note that all quantities depending on the dependent variables become approximations for the original expressions and get therefore labeled with the index h .

When substituting the expansions (3.38) into the above expressions, (3.40) can be written as

$$\mathbf{F} = \begin{pmatrix} F_h^\varphi \\ F_h^{\phi_n} \\ F_h^{\phi_p} \end{pmatrix} = \underbrace{\begin{pmatrix} K^u & 0 & 0 \\ 0 & K^v & 0 \\ 0 & 0 & K^w \end{pmatrix}}_{\mathbf{K}} \begin{pmatrix} u \\ v \\ w \end{pmatrix} - \begin{pmatrix} f_h \\ g_h \\ -g_h \end{pmatrix} \quad (3.41)$$

with

$$K_{ij}^u = a^\varphi(\psi_j, \psi_i), \quad K_{ij}^v = a_h^{\phi_n}(\phi_{n,h}; \psi_j, \psi_i), \quad K_{ij}^w = a_h^{\phi_p}(\phi_{p,h}; \psi_j, \psi_i) \quad (3.42)$$

We omitted the boundary terms for brevity in the above equations. The matrix \mathbf{K} in (3.41) is the *stiffness matrix* of the system and has block-diagonal form. We note that \mathbf{K} is an M -matrix whenever $\int \nabla \psi_i \nabla \psi_j$ is an M -matrix, or at least in the case (in 2D) when there is no triangular element with obtuse angle or rectangle with aspect ratio $> \sqrt{2}$. In other words, the M -property of \mathbf{K} depends on the mesh but not on the physical data or the solution. This is in contrast to a discretization of the continuity equation with the density as variable, where the system matrix can become non- M -matrix due to high electric fields. This will be illustrated below when discretizing the Jacobian.

As \mathbf{K} depends on the solution itself, the system (3.41) is a nonlinear system of order $3N$, where $N = \dim(V_h)$. This system can be solved using a Newton or approximate Newton method. We will follow what was mentioned in the last section and not calculate the Jacobian from (3.41) but rather discretize the Jacobian as given in (3.35).

In each Newton step we have to solve the linear system (cf. section 3.3.2)

$$\begin{aligned} \mathbf{J}^{(k-1)} \begin{pmatrix} \delta u \\ \delta v \\ \delta w \end{pmatrix}^{(k)} &= -\mathbf{F}^{(k-1)} \\ \begin{pmatrix} u \\ v \\ w \end{pmatrix}^{(k)} &= \begin{pmatrix} u \\ v \\ w \end{pmatrix}^{(k-1)} + t_k \begin{pmatrix} \delta u \\ \delta v \\ \delta w \end{pmatrix}^{(k)} \end{aligned} \quad (3.43)$$

We introduced a relaxation parameter t_k to underline the fact that usually a modified Newton method with adaptive t_k is used [16, 105]. For the discretization of \mathbf{J} we expand δu , δv and δw in the same way as u , v and w in (3.38). The contribution resulting from the stiffness matrix can immediately be written down, using (3.35)

$$\mathbf{J}_a = \begin{pmatrix} \lambda^2 \int_{\Omega} \nabla \psi_i \cdot (\epsilon_r \nabla \psi_j) \, dx & 0 & 0 \\ \int_{\Omega} (\mathbf{j}_n \cdot \nabla \psi_i) \psi_j \, dx & \int_{\Omega} \nabla \psi_i \cdot (\mu_n n \nabla \psi_j - \mathbf{j}_n \psi_j) \, dx & 0 \\ \int_{\Omega} (\mathbf{j}_p \cdot \nabla \psi_i) \psi_j \, dx & 0 & \int_{\Omega} \nabla \psi_i \cdot (\mu_p p \nabla \psi_j - \mathbf{j}_p \psi_j) \, dx \end{pmatrix} \quad (3.44)$$

Note that the above expression does not have any terms involving derivatives of the mobility as in the last section we considered only space-dependent mobility models. However, mobility models depending on the electric field are important in many cases and are therefore implemented in the software. It was found that omitting the necessary modifications of the Jacobian only slightly affects convergence behaviour. On p. 65, we once again state the complete Jacobian matrix to reveal some of its properties.

The Jacobian is usually not an M -matrix. This means that (3.43) is not monotone. Several techniques such as special integration rules or upwinding could be applied to increase the diagonal dominance or even transform the Jacobian into an M -matrix. Some of these techniques can be formulated in the framework of the finite element methods [37]. Such measures often change the Jacobian, leading explicitly to an approximate Newton scheme where the matrix \mathbf{J} in (3.43) is only an approximation of the true Jacobian. The gain in stability is therefore expected to implicate a loss of convergence speed. For this reason, and because in most practical cases no negative effect of the matrix properties on convergence has been observed, in TIBERCAD the Jacobian is implemented in TIBERCAD as given on p. 65 without further manipulation (apart from scaling).

The block-structure of the Jacobian could be utilized for a block-iterative method to solve eq. (3.43). For this purpose, the diagonal blocks of the Jacobian could be assembled in such a way that they have the M -property. Although this hasn't been done we shall make some comments about it.

We note that the operators on the diagonal of (3.44) (or (3.35)) are of convection-diffusion type, but with an important difference with respect to (3.36), where the densities are used as dependent variables. We illustrate this with the electrons, using the continuous equations:

$$\int_{\Omega} \mu_n n (\nabla \zeta - \zeta \nabla \phi_n) \nabla v \, dx \iff \int_{\Omega} \mu_n (\nabla \xi - \xi \nabla \varphi) \nabla v \, dx$$

$$\mathbf{J} = \begin{pmatrix} \int_{\Omega} \left[\lambda^2 \nabla \psi_i \cdot (\epsilon_r \nabla \psi_j) - \left(\frac{\partial p}{\partial \varphi} - \frac{\partial n}{\partial \varphi} + \frac{\partial N_d^+}{\partial \varphi} - \frac{\partial N_a^-}{\partial \varphi} \right) \psi_i \psi_j \right] dx & - \int_{\Omega} \left(\frac{\partial n}{\partial \varphi} - \frac{\partial N_d^+}{\partial \varphi} \right) \psi_i \psi_j dx & - \int_{\Omega} \left(-\frac{\partial p}{\partial \varphi} + \frac{\partial N_a^-}{\partial \varphi} \right) \psi_i \psi_j dx \\ \int_{\Omega} \left[(\mathbf{j}_n \cdot \nabla \psi_i) - \left(\frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} + \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \right) \psi_i \right] \psi_j dx & \int_{\Omega} \left[\nabla \psi_i \cdot (\mu_n n \nabla \psi_j - \mathbf{j}_n \psi_j) + \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \psi_i \psi_j \right] dx & \int_{\Omega} \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \psi_i \psi_j dx \\ \int_{\Omega} \left[(\mathbf{j}_p \cdot \nabla \psi_i) + \left(\frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} + \frac{\partial R}{\partial p} \frac{\partial p}{\partial \varphi} \right) \psi_i \right] \psi_j dx & - \int_{\Omega} \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \psi_i \psi_j dx & \int_{\Omega} \left[\nabla \psi_i \cdot (\mu_p p \nabla \psi_j - \mathbf{j}_p \psi_j) \cdot \nabla v - \frac{\partial R}{\partial n} \frac{\partial n}{\partial \varphi} \psi_i \psi_j \right] dx \end{pmatrix}$$

where on the left we have written the formulation in (φ, ϕ_n) and on the right in (φ, n) , and ζ and ξ denote the unknown functions. Apart from the possibly exponentially varying factor n , the main difference between the two formulations is given by the driving force of the convection. In the standard formulation using the densities as variables, the latter is given by the electric field whereas in our case it is given by the gradient of the electro-chemical potential. While in almost all practical cases the electric field can get very big (e.g. in space-charge regions), the electro-chemical potential is a much smoother quantity and only at heterointerfaces it can be expected to have big gradients (cf. section 3.2).

Consider a 1D case, assuming constant electron flux j and using finite differences in the point with index i (setting $\mu_n = 1$, using scaled quantities and letting h be the constant mesh spacing)

$$\begin{aligned} -\nabla(n\nabla\zeta - j\zeta) &= f \\ \Downarrow \\ -\frac{1}{h} \left(\frac{n_{i-\frac{1}{2}}}{h} + \frac{j}{2} \right) \zeta_{i-1} + \frac{n_{i-\frac{1}{2}} + n_{i+\frac{1}{2}}}{h^2} \zeta_i - \frac{1}{h} \left(\frac{n_{i+\frac{1}{2}}}{h} - \frac{j}{2} \right) \zeta_{i+1} &= f_i \end{aligned}$$

$n_{i\pm\frac{1}{2}}$ denote the densities in the center between the nodes i and $i+1$ and $i-1$ and i , respectively. From the above we find the condition for the discretization to be stable (for $j > 0$):

$$\frac{n_{i+\frac{1}{2}}}{h} > \frac{j}{2} \Rightarrow h < \frac{2n_{i+\frac{1}{2}}}{j} = \frac{2}{\nabla\phi_{n,i+\frac{1}{2}}}$$

where we have written the current as $j = n_{i+\frac{1}{2}} \nabla\phi_{n,i+\frac{1}{2}}$. This means that the electro-chemical potential should change less than one half of the thermal voltage in one element. This should be compared with the corresponding condition when the density is used as variable, where the electric potential should change less than one half of the thermal voltage in one element, which is usually much more restrictive on the mesh size.

The calculation of the matrix elements of the Jacobian involves integrations over potentially fast changing quantities. The integrals appearing in (3.44) could principally be evaluated analytically. The same is not true in general for the integrals resulting from the right-hand side. Currently the implementation in TIBERCAD uses a numerical Gauss integration for all integrals. Therefore the discretization as a whole can be regarded as non-conforming.

Up to now we did not address the conditioning of the jacobian. As has already been mentioned before (cf. section 3.2), the operators in the continuity equations are not well conditioned with our choice of variables. This fact leads to an ill-conditioned Jacobian which can result in convergence problems of the iterative solvers needed to solve the linear system (3.29). The next section will address conditioning in general and then propose a solution to produce a better scaled matrix.

3.3.4 Conditioning of the linearized system

The goal of this section is to examine the conditioning of the linearized semiconductor equations (3.43), which amounts to estimate the condition number $\kappa(\mathbf{J})$ of the Jacobian \mathbf{J} .

The conditioning of the drift-diffusion equations plays a crucial role for the performance of numerical device simulations. An ill-conditioned system can lead to very poor convergence properties or to instabilities or fail at all, because conditioning governs the accuracy of the correction steps in the Newton algorithm. Especially iterative solvers are sensitive to the conditioning of the system to be solved [17]. There is quite a lot of literature treating under several aspects the conditioning of the continuous semiconductor equations as well as of the discretized system [4, 66, 65, 10] and its effects on the numerics of the problem [32, 17, 44, 95].

The overall conditioning of the discretized system depends on several factors:

physical conditioning With this we mean the “intrinsic” conditioning of the continuous equations. It depends on the type of operators, the physical device properties (e.g. device diameter, intrinsic densities, device structure) and bias conditions (near breakdown e.g. the system becomes ill-conditioned as it moves towards a singularity).

numerical conditioning This includes effects of numerical calculations, i.e. the evaluation of the residual and the machine precision.

scaling The scaling of the system can be regarded as a preconditioning of the system. As such it influences the conditioning, and a proper scaling can reduce the condition number.

discretization The discretization, i.e. the mesh, plays an important role for the conditioning because shape and size of the elements together with device size affect directly the condition number of the Jacobian of the system.

Usually the user has no or little influence on the first three factors when using a device simulator, but the mesh is mostly user dependent. Finding a mesh of good quality with as little as possible nodes can be essential for a successful simulation [10].

Consider some function $f(x) \in V(\Omega) : \Omega \mapsto \Omega'$ where $\Omega \subset \mathbb{R}^n$, $\Omega' \subset \mathbb{R}^m$. We define the functional $C(x) : \Omega \mapsto \mathbb{R}$

$$C(x) = \sup_{x' \in B_\varepsilon(x)} \left(\frac{\|f(x') - f(x)\|_{\Omega'}}{\|f(x)\|_{\Omega'}} \right) / \frac{\|x' - x\|_\Omega}{\|x\|_\Omega}, \quad (3.45)$$

where $B_\varepsilon(x) \subset \Omega$ is a ball around x with (small) radius $\varepsilon > 0$. We can interpret this quantity as the (local) amplification factor of the relative error in the “input data” x , i.e. the relative error of the “output” divided by the relative error of the

“input”. The *condition number* C^f of an operator f is then defined as the upper bound of $C(x)$, i.e.

$$C^f = \sup_{x \in \Omega} C(x) \quad (3.46)$$

A well conditioned operator has a condition number near 1.

More specifically we consider the case of a differentiable function $f(x) \in C^1(\Omega) : \mathbb{R} \mapsto \mathbb{R}$. In this case (3.45) simplifies in the limit $\varepsilon \rightarrow 0$ to

$$C(x) = \left| \frac{f(x') - f(x)}{x' - x} \cdot \frac{x}{f(x)} \right| = \left| \frac{x}{f(x)} \cdot \frac{df(x)}{dx} \right| \quad (3.47)$$

We will assume now that $f(x)$ in (3.45) is a linear invertible operator $\mathbb{R}^n \mapsto \mathbb{R}^n$ and write it as $y := f(x) = \mathbf{A}x$, where \mathbf{A} is an invertible $n \times n$ Matrix. Now we can calculate $C(x)$ for the matrix multiplication as

$$\begin{aligned} C(x) &= \sup_{x' \in B_\varepsilon(x)} \left(\frac{\|\mathbf{A}x' - \mathbf{A}x\|}{\|\mathbf{A}x\|} \middle/ \frac{\|x' - x\|}{\|x\|} \right) \\ &= \sup_{x' \in B_\varepsilon(x)} \left(\frac{\|\mathbf{A}(x' - x)\|}{\|x' - x\|} \cdot \frac{\|x\|}{\|\mathbf{A}x\|} \right) \end{aligned} \quad (3.48)$$

Remembering that the norm of a matrix \mathbf{A} is defined as $\|\mathbf{A}\| = \sup(\|\mathbf{A}x\| / \|x\|)$ and considering the fact that the linearity of the matrix multiplication implies $\|\mathbf{A}(\alpha x)\| / \|\alpha x\| = \|\mathbf{A}x\| / \|x\|$, we are led to

$$C(x) = \sup_{\|x'\| < \varepsilon} \frac{\|\mathbf{A}x'\|}{\|x'\|} \cdot \frac{\|x\|}{\|\mathbf{A}x\|} = \|\mathbf{A}\| \frac{\|x\|}{\|\mathbf{A}x\|} \quad (3.49)$$

The condition number of the matrix multiplication is then according to (3.46) and using $y = \mathbf{A}^{-1}x$

$$C^{\mathbf{A}} = \sup_x \left(\|\mathbf{A}\| \frac{\|x\|}{\|\mathbf{A}x\|} \right) = \|\mathbf{A}\| \cdot \sup_y \frac{\|\mathbf{A}^{-1}y\|}{\|y\|} = \|\mathbf{A}\| \|\mathbf{A}^{-1}\| \quad (3.50)$$

We call $\kappa(\mathbf{A}) = C^{\mathbf{A}}$ the condition number of the matrix \mathbf{A} . For a symmetric matrix and using the l_2 norm the following relation holds

$$\kappa(\mathbf{A}) = \left| \frac{\lambda_{\max}(\mathbf{A})}{\lambda_{\min}(\mathbf{A})} \right| \quad (3.51)$$

In a nonsingular linear system $Ax = b$, the sensitivity of the solution x with respect to a perturbation in the data b or in the matrix A is closely related to the condition number $\kappa(\mathbf{A})$ of \mathbf{A} [30, 68].

Assume we are solving a nonlinear problem $g(u) = 0$ using a Newton-like method

$$g'(u_k)x = -g(u_k) \quad (3.52)$$

$$u_{k+1} = u_k + t_k x \quad (3.53)$$

where we identify $g_k = g(u_k)$ and $g'_k = g'(u_k)$ with the residual and the Jacobian at the k^{th} step, respectively. The correction x is given by $x = -(g'_k)^{-1}g_k$. The condition number of this operation is given by the product $C^J C^r$, where C^J and C^r are the condition numbers of the Jacobian and the residual, respectively. The numerical accuracy of the correction x is therefore $\varepsilon_x = C^J C^r \varepsilon_{CPU}$ where ε_{CPU} means the machine accuracy, which is about 10^{-16} for double precision. C^r measures how correctly we evaluate the residual and depends on implementation details of the models, on device type and biasing and is therefore difficult to estimate. C^J is given by (3.50), once the system is discretized. It can be estimated apriori for example in the sense that it will not be better than the condition number of a simpler operator, e.g. the Laplace operator, discretized on the same grid [10].

In the following we will examine the conditioning of the high-order diagonal terms of the jacobian (3.44). This part of the matrix corresponds to the system matrix \mathbf{K} of the discretized system:

$$\mathbf{K} = \begin{pmatrix} \lambda^2 \int_{\Omega} \nabla \psi_i \cdot (\epsilon_r \nabla \psi_j) \, dx & 0 & 0 \\ 0 & \int_{\Omega} (\mu_n n \nabla \psi_j) \cdot \nabla \psi_i \, dx & 0 \\ 0 & 0 & \int_{\Omega} (\mu_p p \nabla \psi_j) \cdot \nabla \psi_i \, dx \end{pmatrix} \quad (3.54)$$

The ill-conditioning of this matrix due to the dependence on the densities of the terms corresponding to the continuity equations was already mentioned before. To get a rough estimate of the condition number of \mathbf{K} we approximate the integrals in the following way. Consider a triangular mesh as given in Fig. 3.3. We now

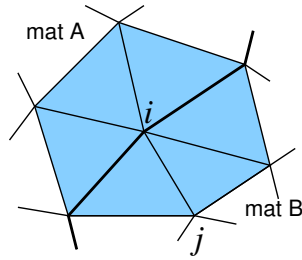


Figure 3.3: A triangular mesh around node i which lies on the boundary (indicated by the bold line) of two materials $mat A$ and $mat B$.

define mean values for the integrands ϵ_r , $\mu_n n$ and $\mu_p p$ over the supports of the basis functions $\text{supp}(\psi_i)$, indicated in the figure by the shaded region, and denote

them by $(\bar{\epsilon}_r)_i$, $(\bar{\mu}_n n)_i$ and $(\bar{\mu}_p p)_i$. With this, the system matrix (3.54) gets

$$\begin{aligned} \mathbf{K} &\approx \begin{pmatrix} \lambda^2(\bar{\epsilon}_r)_i \int_{\Omega} \nabla \psi_i \cdot \nabla \psi_j \, dx & 0 & 0 \\ 0 & (\bar{\mu}_n n)_i \int_{\Omega} \nabla \psi_i \cdot \nabla \psi_j \, dx & 0 \\ 0 & 0 & (\bar{\mu}_p p)_i \int_{\Omega} \nabla \psi_i \cdot \nabla \psi_j \, dx \end{pmatrix} \\ &= \underbrace{\begin{pmatrix} \lambda^2(\bar{\epsilon}_r)_i & 0 & 0 \\ 0 & (\bar{\mu}_n n)_i & 0 \\ 0 & 0 & (\bar{\mu}_p p)_i \end{pmatrix}}_{\mathbf{D}} \cdot \underbrace{\begin{pmatrix} \mathbf{K}_L & 0 & 0 \\ 0 & \mathbf{K}_L & 0 \\ 0 & 0 & \mathbf{K}_L \end{pmatrix}}_{\tilde{\mathbf{K}}} \end{aligned} \quad (3.55)$$

where we denoted by \mathbf{K}_L the matrix of the discretized Laplace operator in weak form, and \mathbf{D} is a positive diagonal matrix. The condition number of \mathbf{K} can therefore be estimated as

$$\kappa(\mathbf{K}) = \kappa(\mathbf{D})\kappa(\tilde{\mathbf{K}}) = \underbrace{\frac{\max(\lambda^2(\bar{\epsilon}_r)_i, (\bar{\mu}_n n)_i, (\bar{\mu}_p p)_i)}{\min(\lambda^2(\bar{\epsilon}_r)_i, (\bar{\mu}_n n)_i, (\bar{\mu}_p p)_i)}}_{\gamma} \kappa(\mathbf{K}_L) \quad (3.56)$$

This means the condition number of the system matrix \mathbf{K} increases by about a factor of γ with respect to the condition number of the Laplace operator discretized on the same mesh. It is especially critical for structures involving wide bandgap materials where the carrier densities can get very low. Generally the carrier densities can be estimated to lie between $N_{c,v} \exp(-E_g/k_B T)$ and $\sim 10^{20} \text{ cm}^{-3}$.

Based on the above observations, a diagonal scaling has been implemented in TIBERCAD. We denote the diagonal scaling matrix by \mathbf{D} . The linear system (3.43) that has to be solved in each step of the Newton algorithm is preconditioned using \mathbf{D} as follows:

$$\begin{aligned} \mathbf{J}\mathbf{x} &= -\mathbf{F} \\ \Downarrow \\ \mathbf{D}^{-1}\mathbf{J}\mathbf{x} &= -\mathbf{D}^{-1}\mathbf{F} \end{aligned} \quad (3.57)$$

In the practical implementation the preconditioning is not applied as matrix operation but rather included directly during the assembly of Jacobian and residual. For this purpose, the elements of \mathbf{D} have to be precalculated before assembly as they are not defined exclusively by element-local parameters. This is due to the fact that a node can lie on the boundary of two materials as shown in Fig. 3.3.

There is no unique or “best” choice for \mathbf{D} . We compare numerically a few possibilities, using as an example a GaN pn-junction as shown in Fig. 3.4.

The first possible implementation (1) uses as scaling matrix the main diagonal of the system matrix \mathbf{K} :

$$\mathbf{D} = \text{diag}(\mathbf{K}) \quad (3.58)$$

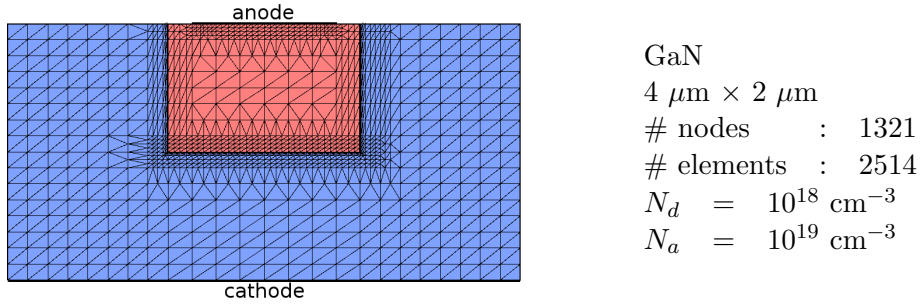


Figure 3.4: A simple GaN pn-diode.

where \mathbf{K} is the system matrix given in (3.54).

The second implementation (2) adds to the scaling of the Poisson equation the derivative of the total charge density as this term can dominate the corresponding part of the Jacobian. We write \mathbf{D} in symbolic form

$$\mathbf{D} = \text{diag}(\mathbf{K}) - \text{diag}\left(\int_{\Omega} \frac{\partial \rho}{\partial \varphi} \psi_i^2 dx, \mathbf{0}, \mathbf{0}\right) \quad (3.59)$$

This approach is currently implemented in TIBERCAD.

The numerical performance of the different scaling schemes are summarized in Table 3.2.

	Scheme		
	–	(1)	(2)
condest	$5.62 \cdot 10^{21}$	$2.56 \cdot 10^8$	$1.06 \cdot 10^7$

Table 3.2: Numerical performance of different diagonal scaling schemes. The scheme denoted by a dash (–) is the unconditioned Jacobian. *condest* means the estimated condition number obtained in Matlab using the `condest` command. The values are calculated for the Jacobian at 3.2 V.

Chapter 4

The TIBERCAD Software

4.1 Introduction

This chapter describes some implementation and usage details of TIBERCAD [85]. The TIBERCAD software is written in C++ for the following reasons:

- Its object-oriented features are well suited for the handling of complex data structures, for the implementation of hierarchies of simulation models and for their implementation independent handling (polymorphism).
- It can easily be intermixed with C to allow low level operations such as direct access to shared libraries, which allows for a modularization of the software package.
- It can interface to code written in other languages as Fortran without too much problems.
- The most important library used in TIBERCAD, libMesh [55], which implements and handles all finite element specific details and data structures (mesh, elements, integration rules etc.) is already written in C++.

4.2 Software structure

The structure of the TIBERCAD software is schematically shown in Fig. 4.1. The implementation uses heavily the object-oriented features of C++. Especially polymorphism is used in creating model hierarchies with a common interface.

The core of TIBERCAD is formed by modules that implement the physical models described in chapter 2. Each module is self-contained in the sense that it does not rely on implementation details of other modules but rather uses an abstract common interface to communicate with other modules. This allows for a highly modular software package that can be easily extended with new modules.

A control module manages the program flow. This task includes the creation of all data structures needed to describe the device and its properties, the creation

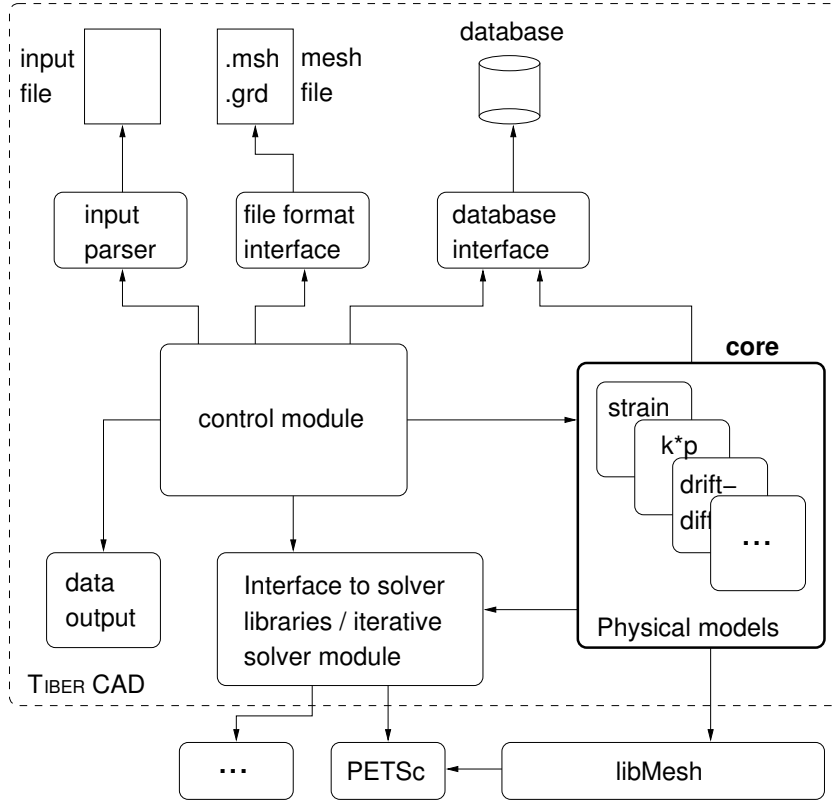


Figure 4.1: Structure of the TIBERCAD software.

of the objects representing the physical models and coordination of the different simulations.

The device description is based on two text files and a database (currently implemented as a collection of flat text files). A *mesh file* describes the geometry of the device in terms of physical (or geometrical) regions and the discretization mesh. An *input file* contains the definitions of the materials and models to be assigned to the different regions found in the mesh file, parameters for the solvers to be used, type of simulations to be done and general configuration options. The *database* finally contains all physical parameters needed for the different physical models. These parameters can generally be overridden from the input file.

An input parser is responsible for the parsing of the input file.

External libraries are used for the numerical solution of the equation systems. PETSc [13, 12, 14] is responsible for the solution of linear and nonlinear sparse systems, using iterative methods. SLEPc [45] solves standard and generalized eigenvalue problems and is based on PETSc.

4.2.1 Mesh handling

The geometrical description of a device to simulate is based on a finite element discretization. The corresponding mesh can be generated in different graphical tools (currently DEVISE of the former ISE TCAD toolchain [101] and Gmsh [39]). This mesh, called *parent mesh*, is common to all TIBERCAD modules to allow for easy interchange of data (cf. sec. 1.3). The mesh contains not only volumic elements (ND -elements where N is the dimension of the simulation domain), but also boundary elements used for definition of boundary conditions. The single elements are combined into non-overlapping geometrical regions, identified by unique strings or numbers, depending on the mesh generation tool and mesh file format. These regions are referred to in the input file, using their identifiers, to be reassembled into logical groups, called *physical regions* and *clusters*. The physical regions are device regions characterized by a single material with certain attributes such as constant doping and crystal orientation. Therefore the **Region** blocks in the input file (cf. Listing 4.1) contain the definition of the material including material parameters, crystal orientation and doping. Each region has a distinct user defined name that can be referred to in the input file. Clusters can be used to reassemble sets of geometrical regions into logical units irrespective of their membership to physical regions. They are mainly useful when doing quantum mechanical calculations on parts of the device. Fig. 4.2 illustrates these concepts schematically.

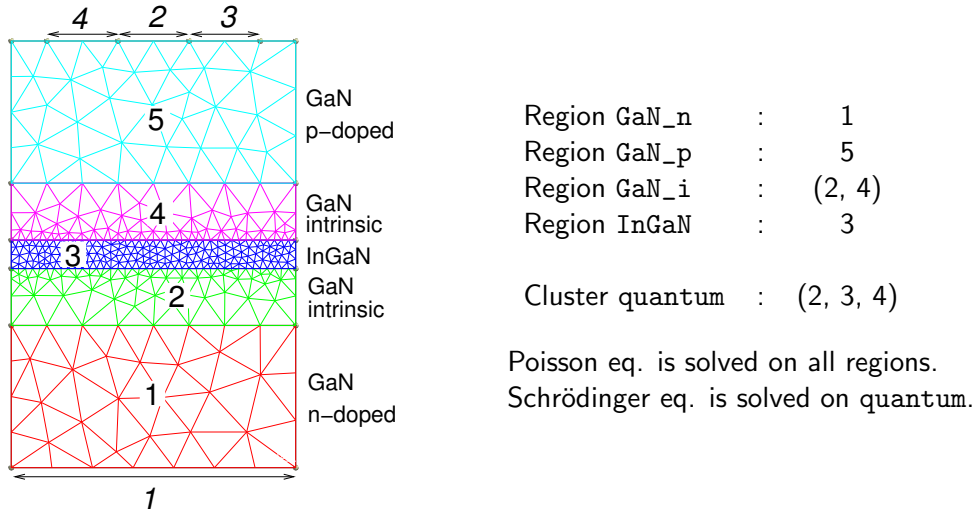


Figure 4.2: Illustration of the mesh with mesh regions. The device is a p-i-n structure with embedded quantum well. We assume a Schrödinger-Poisson calculation to be performed. On the mesh regions 2–4 a quantum mechanical simulation will be performed. Boundary regions are indicated by italic numbers.

4.2.2 Model hierarchy

The implementation of physical models is based on three class hierarchies, shown in Figs. 4.3–4.5. The core of each model is contained in a class inherited from `SimulationInterface`. Its task is the numerical implementation of the mathematical description of the physical model. This essentially means the discretization of operators, assembly of matrices and the interfacing to numerical libraries.

All physical parameters that depend on material (and thus on mesh regions) or on model details which formally do not affect discretization or matrix assembly are contained in classes derived from `PhysicalModel`. Objects of these classes are instantiated for each physical region defined in the input file (`Region` sections in the input file). To control model details they can make use of other classes derived directly from `PhysicalModelInterface`. Examples for this are mobility and recombination models which are encapsulated as members in the `PhysicalModel` of drift-diffusion. Usually the preparation for the numerical calculations involves loops over the mesh elements. The `PhysicalModel` objects can then easily be found as each element knows the physical region it belongs to.

Boundary conditions are handled in a similar way.

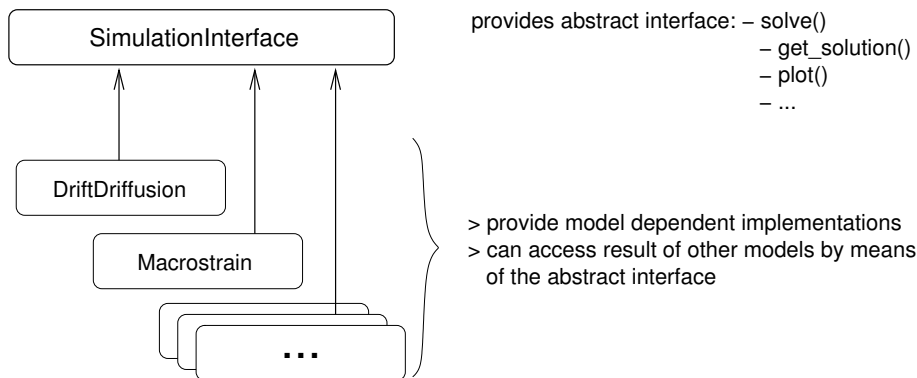


Figure 4.3: `SimulationInterface` hierarchy.

4.3 User interface

The user interface of TIBERCAD is given by a textual input file. This file is organized in different sections tagged by a section name preceded by the dollar sign, e.g. `$Device`. The content of the section is enclosed between curly brackets. Each section can contain any number of blocks, where each block is enclosed by curly brackets and preceded by a block name. Blocks can be nested. A block containing any number of key-value pairs is called a *parameter block*. Key-value pairs are assignments of the form

```
key = value
```

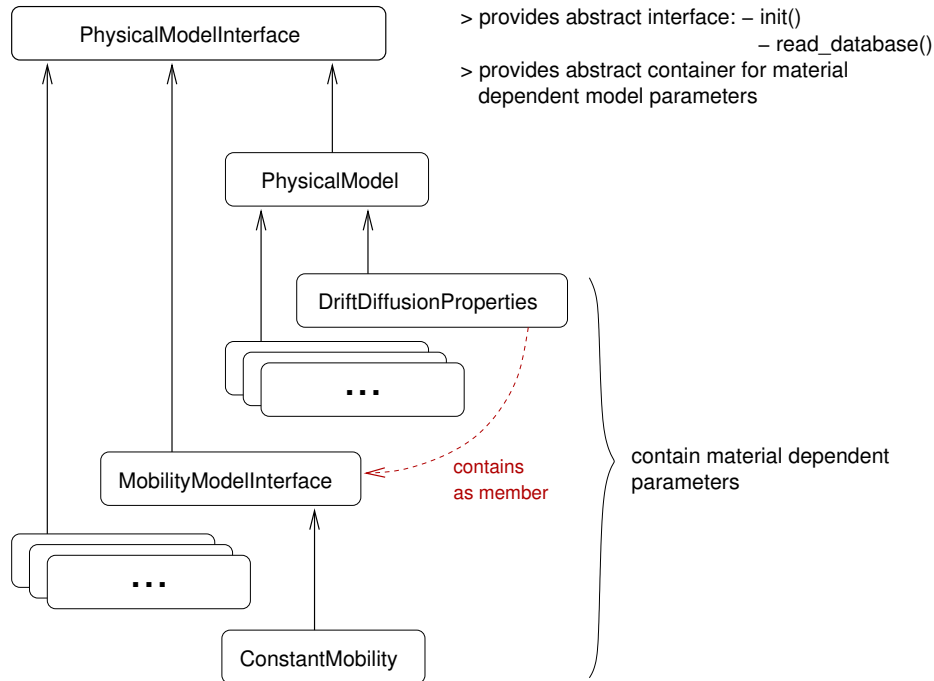


Figure 4.4: `PhysicalModelInterface` hierarchy. For each physical region (Device section) and each simulation (Models section) an instance of a class derived from `PhysicalModel` is created according to the options given in the Physics section.

where **key** is a string and **value** is a single numerical value or string or a list of values separated by commas and enclosed by parentheses. Different key-value pairs are separated by spaces, therefore values cannot contain white space.

A hash sign (#) initiates a comment extending until the end of line.

The following sections are defined:

Device Contains the description of the different geometrical regions of the device. It can contain two types of blocks.

The **Region** blocks are used to define the different regions that are present in the mesh. It defines the material, the crystal structure, growth direction, constant doping and optionally material properties that are needed for simulations based on continuous media descriptions.

The optional **Cluster** blocks can be used to logically unite different regions which can then be addressed by the cluster name.

Scale Contains the description of device regions or simulation entities that are not based on continuous media approaches handled in the **Device** section. This can be models for lumped circuit devices (not implemented yet) or regions that have to be simulated on an atomistic scale, e.g. tight-binding. The latter

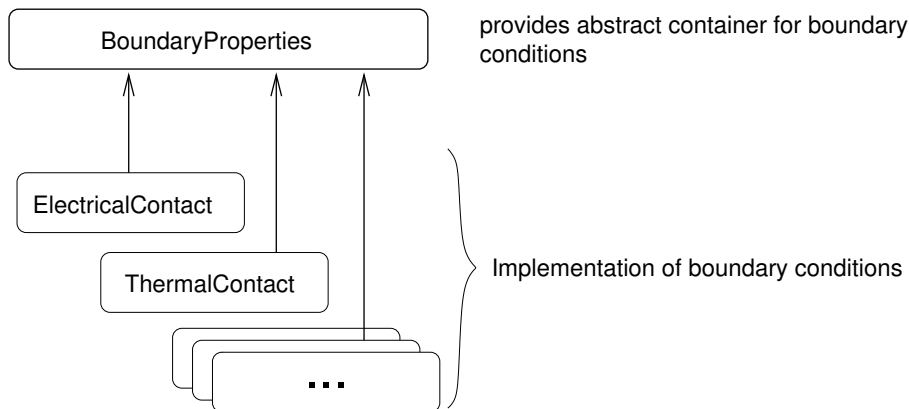


Figure 4.5: BoundaryProperties hierarchy.

is defined in an **Atomistic** block. The microscopic (atomic) structure inside the regions assigned in the **Atomistic** block will be generated automatically.

Models This section contains the definition of the models used to describe the entities defined in **Region** or **Atomistic** blocks. Each model is defined in its own block, tagged by the identifier of the model (e.g. **driftdiffusion**, **macrostrain**). These blocks can embody three types of nested blocks, namely **options**, **physical_model** and **BC_regions**.

The **options**-block contains the names of physical regions defined in **Region** or **Cluster** blocks and an optional user-defined identifier for the simulation.

The (optional) **physical_model**-block can contain parameters relevant to some aspect of the current model. Usually submodels such as carrier mobility or recombination models are defined in this way.

The **BC_regions**-block finally contains the description of the boundary conditions.

Solver Contains options and parameters relevant for the numerics of a model. For each model defined in the **Models** section, a block tagged with the model identifier or the user-provided model name can be defined.

There are two special blocks to define parameter sweeps and selfconsistent simulations.

Physics In this section one can define additional physical parameters relevant for each model or override material parameters from the database. For each model defined in the **Models** section, a block tagged with the model identifier or the user-provided model name can be defined. Material or model parameters given in **Regions** blocks will override parameters specified in block tagged with the user-provided simulation name which in turn overrides data provided in a block tagged with the model identifier.

Simulation This section defines global options for the simulation such as mesh file, temperature, path for output data and database, the variables to plot and the simulations to be performed.

Listing 4.1 shows an example of a simple input file for the pn-heterojunction as illustrated in Fig. 4.6. Certain values can be specified as variables using a special syntax. An example is the contact voltage for a voltage sweep (see Listing 4.1, line 56).

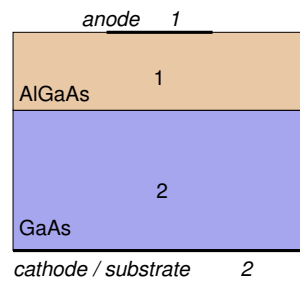


Figure 4.6: pn-heterojunction diode for Listing 4.1.

```

# diode example

# Description of the device physical regions
$Device
5 {
    Region n_side
    {
        reg_numb = 2
        mat = GaAs
10    doping = 5e18    doping_type = donor
    }

    Region p_side
    {
15    reg_numb = 1
        mat = AlGaAs
        x = 0.1
        doping = 1e19    doping_type = acceptor
    }
20 }

# Definition of Simulation Models and associated Boundary Conditions
$Models
25 {
    model driftdiffusion
    {

        options
30    {
        simulation_name = driftdiffusion
        physical_regions = all
    }
}

```

```

35     physical_model recombination
    {
        model = srh
    }

40     physical_model electron_mobility
    {
        model = doping_dependent
    }

45     physical_model hole_mobility
    {
        model = doping_dependent
    }

50     BC_Regions
    {
        BC_Region anode
        {
            BC_reg_numb = 1
55             type = ohmic
            voltage = @Vb
        }

        BC_Region cathode
60         {
            BC_reg_numb = 2
            type = ohmic
            voltage = 0.0
        }
65     }
}

model macrostrain
{
70     options
    {
        simulation_name = strain
        physical_regions = all
75     }

    BC_Regions
    {
        BC_Region substr
80         {
            BC_reg_numb = 2
            type = substrate
            material = GaAs
        }
85     }
}

90 # Definition of Model-dependent Solver parameters
$Solver
{
    driftdiffusion
    {
95         nonlin_step_tol = 1e-6
        nonlin_max_it = 30
    }
}

```

```

        ls_max_step = 1
    }

100    strain
    {
        substrate = substr
        max_iterations = 1000
    }

105    sweep
    {
        variable = Vb
        start = 0.0
110    stop = 3.2
        steps = 32
        plotvariable = current
        plot_data = true
    }
115 }

    # Definition of Model dependent physical parameters
    $Physics
120 {
        driftdiffusion
        {
            model = strained
            statistics = FD
125    strain_model = strain
        }
    }

130 # Definition of global simulation options
    $Simulation
    {
        searchpath = ../materials
        meshfile = diode_mdr.grd
135    dimension = 2
        temperature = 300
        solve = sweep
        resultpath = output
        output_format = vtk
140    plot = (Ec, Ev, QFermi_e, QFermi_h, eDensity, hDensity, eCurrent, hCurrent, strain)
    }

```

Listing 4.1: Input file for a simple pn-heterojunction diode.

Chapter 5

Simulation Examples

5.1 Piezoresistivity effects of HEMT structures

The piezoelectric effect plays a fundamental role for the electronic properties of devices especially based on wurtzite GaN/AlGaN heterostructures, and measurements of such structures have been reported in literature [33]. The experimental results were interpreted in [33] by the authors in the limits of analytical one-dimensional models assuming homogeneous distribution of strain induced by a mechanical force. A realistic strain pattern in AlGaN/GaN N-face and AlGaAs/InGaAs/GaAs A- and B-face FETs has been calculated and published in [6] using TIBERCAD with the aim of investigating the effect on the resistivity of the device.

The pseudomorphic heterojunction devices considered here possess strain deformation due to both lattice mismatch and an external pressure applied on top of the structure. We make use of the continuous media model as described in section 2.1 in order to compute the strain distribution.

In the simulation we assume that the device is grown on a thick substrate that remains unstrained. At the surface of the device we apply the following boundary condition for stress:

$$\sigma_{ij}(\mathbf{r})n_i = \begin{cases} \mathbf{f}_j, & \text{if the force } \mathbf{f} \text{ is applied at point } \mathbf{r} \\ 0, & \text{otherwise,} \end{cases} \quad (5.1)$$

where \mathbf{n} is the unit normal to the surface onto which the force \mathbf{f} acts. From the computed strain we get the piezoelectric polarization \mathbf{P}^{pz} as given in (2.16). The total polarization is then used in the Poisson equation for the transport calculation. In the latter electrons and holes are considered. The band parameters for both carriers are calculated according to section 2.4.1, including the effects of strain. The mobilities are assumed to be constant as the devices are simulated for small bias near equilibrium. The necessary material parameters were taken from [106, 107].

The first simulated structure, Fig. 5.1(a), is an AlGaN/GaN inverted HEMT grown along $[000\bar{1}]$ direction with N-face polarity without gate metallization, as considered in [33]. The second one, shown in Fig. 5.1(b), is an AlGaAs/InGaAs/GaAs

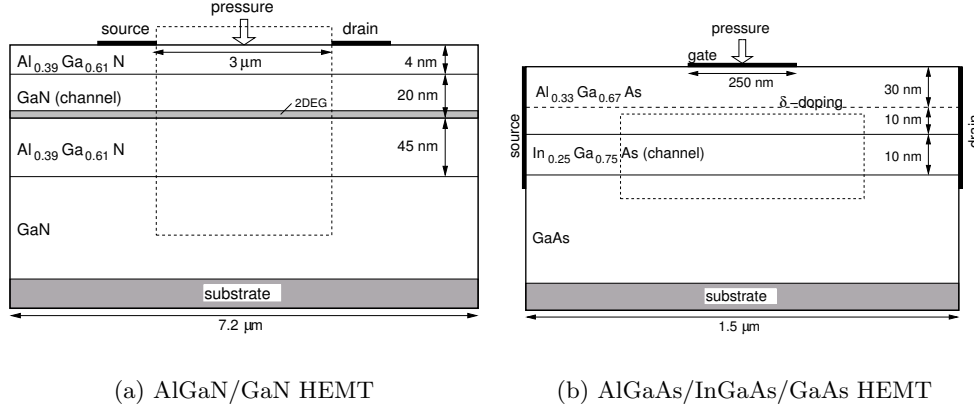


Figure 5.1: Schematic drawing of the simulated heterostructures.

HEMT grown along [311] direction with A- or B-face polarity [99]. The latter does not show spontaneous polarization as it is built of materials with cubic symmetry.

For both structures a well defined pressure was applied on a region of the surface at the center between the source and drain contacts. In the second structure this region coincides with the gate, whereas in the first one pressure is imposed on a line of 100 nm length. The resistance was calculated for a drain-source voltage of 0.1 V. The gate of the AlGaAs/InGaAs/GaAs structure was biased at 0 V, assuming the Schottky-barrier height to be 0.8 eV.

In Fig. 5.2 we show the pattern of the lateral strain component ε_{xx} of the AlGaIn/GaN structure without and with an external force of 0.04 Ncm^{-1} . Only the strain map of the AlGaIn layers are made visible in the picture. The GaIn, which is slightly compressed appears as white surface. The inhomogeneity in the strain pattern without pressure is due to the finite size of the device which causes a slight bending of the structure. The external pressure induces a highly nonhomogeneous strain distribution as can be seen in part (b) of the figure.

The simulated piezoresistivities of the two structures are shown in Fig. 5.3. For small pressures the dependence of the relative change of resistance on pressure is linear. In the AlGaIn/GaN structure the resistance increases with increasing pressure which is in agreement with the findings in [33]. The external pressure leads to a slight compression of the AlGaIn barrier as is clearly seen in Fig. 5.2. This causes a decrease of the piezoelectric polarization in the barrier which decreases the discontinuity in polarization and therefore the electron density in the channel [3].

For the AlGaAs/InGaAs/GaAs structure our simulation results show a strong dependence on substrate termination. For the case of A-face polarity the resistance increases for increasing stress on the gate. For B-face polarity and low pressure the resistance is decreasing but with the increasing of the pressure the effect becomes non-monotonic. Then the resistance continues to grow even faster than in the case of the A-face substrate.

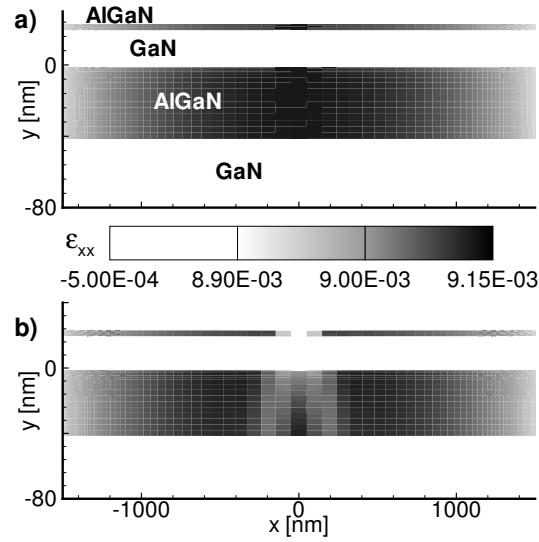


Figure 5.2: Contour plots of the lateral strain component ε_{xx} in a part (cf. dashed box in Fig. 5.1(a)) of the AlGaIn/GaN FET without (a) and with (b) pressure.

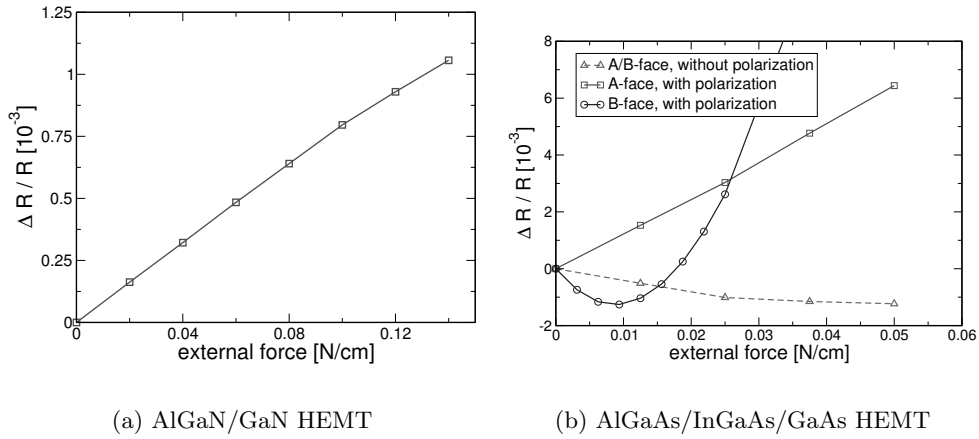


Figure 5.3: Relative change of resistance as a function of external force.

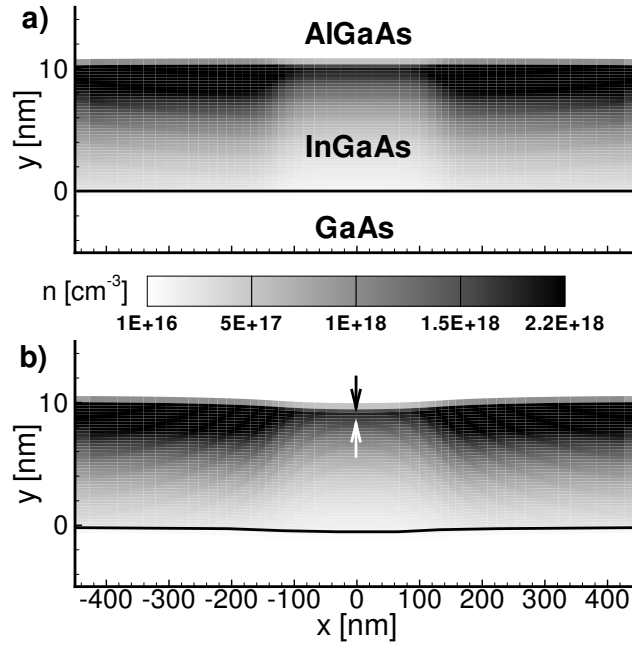


Figure 5.4: Electron density in a part (cf. dashed box in Fig. 5.1(b)) of the B-face GaAlAs/InGaAs/GaAs FET without (a) and with (b) an external pressure of 50 mN/cm.

In the B-face structure the 2DEG is more confined towards the upper InGaAs-AlGaAs interface, where both piezoelectric polarization and strain are significantly affected by the pressure. Due to that there are two effects that make the dependence non-monotonic: the effect of the closing of the channel by the piezoelectric field that dominates at high pressure and the effect of the deformation potential that dominates at lower pressure. For comparison, we also show the results of a simulation which does not consider the piezoelectric polarization where the deformation potential is the only effect that affects the device resistivity.

The piezoelectric effect appears to be stronger in the AlGaAs/InGaAs/GaAs structure. This can be explained by the fact that the external pressure affects the whole channel region which is the determinant part of the device. In the AlGaAs/GaN structure on the other hand the influence is limited to a small part of the channel.

Fig. 5.4 presents the electron density in the channel in the gate region of the B-face AlGaAs/InGaAs/GaAs HEMT. An external force of 50 mN/cm applied onto the gate leads to a convex deformation of the device and to a narrowing of the electron channel thus increasing the resistance.

5.2 The influence of gate tunneling in MOSFETs

The present scaling of Si-based technology [46] is leading to quantum mechanical tunneling which results in excessive gate leakage current in MOSFETs. In order to avoid this limitation, several alternative high- κ gate dielectrics have been studied. Among these, ZrO_2 and HfO_2 have attracted great interest and have been selected for MOSFET applications [89].

Quantum phenomena in MOS systems, both based on SiO_2 and high- κ oxides, are usually either neglected or studied within simplified schemes, such as the effective-mass approximation (EMA). However, a macroscopic description is often not fully satisfying for phenomena that occur on a characteristic length scale of a few nm such as tunneling through high- κ gate dielectric stacks [52]. Moreover, many of the effective parameters used in EMA are not known for this length scale and should be extracted by microscopic approaches.

Here we present simulation results of a high- κ MOSFET including gate tunneling calculated using a microscopic approach [8]. Fig. 5.5 shows a schematic drawing of the simulated device. The tunneling properties of the gate dielectric (SiO_2 , ZrO_2 and HfO_2) were calculated by applying quantum mechanical methods that include the full band structure of Si and the oxide materials [91].

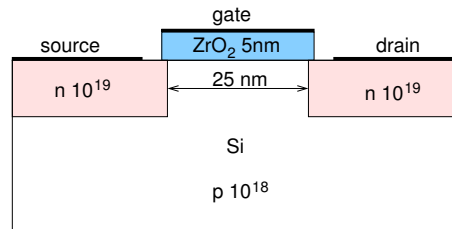


Figure 5.5: Schematic drawing of simulated device.

First, the electronic band structures of the oxide materials (in crystalline form) are calculated based on density functional theory (DFT). Then the parameters for a semiempirical $\text{sp}^3\text{s}^*\text{d}^5$ tight-binding (TB) parametrization are determined such as to reproduce the band dispersions obtained from the DFT calculations and the experimental band gap of approximately 5.7 eV. The atomic structure of $\text{Si}/\text{ZrO}_2/\text{Si}$ and $\text{Si}/\text{HfO}_2/\text{Si}$ is modeled as shown in Fig. 5.6.

The tunneling properties of these structures are calculated using the transfer matrix method based on the TB description of the structure [90, 26]. Although the oxides considered are amorphous materials, calculations based on a crystalline form are expected to provide reasonable results as has been found in SiO_2 -based MOS systems [90]. One of the main advantages of the TB approach is its ability to include any microscopic feature such as bond arrangement at the interfaces or the complex band structure, which is particularly important for modeling tunneling currents.

A $\text{n}^+\text{-Si}/\text{oxide}/\text{p-Si}$ MOS capacitor has been assumed for the tunneling calcu-

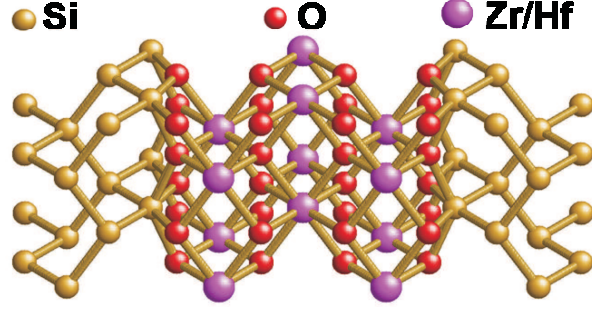


Figure 5.6: Microscopic model for the Si/high- κ oxide/Si heterostructure.

lations with n^+ -doped Si gate ($N_d = 3 \cdot 10^{20} \text{ cm}^{-3}$) and slightly p-doped Si bulk ($N_a = 10^{15} \text{ cm}^{-3}$), biased in accumulation region. Fig. 5.7 shows the tunneling current through different oxides in dependence on the equivalent oxide thickness (EOT). The latter is defined as

$$\text{EOT} = t_{\text{SiO}_2} = \frac{\epsilon_{\text{SiO}_2}}{\epsilon_{\text{high-}\kappa}} t_{\text{high-}\kappa}$$

where t_{SiO_2} and $t_{\text{high-}\kappa}$ are the oxide thicknesses for SiO_2 and a high- κ dielectric, respectively. Fig. 5.8 shows the tunneling through ZrO_2 against the applied voltage for several EOT.

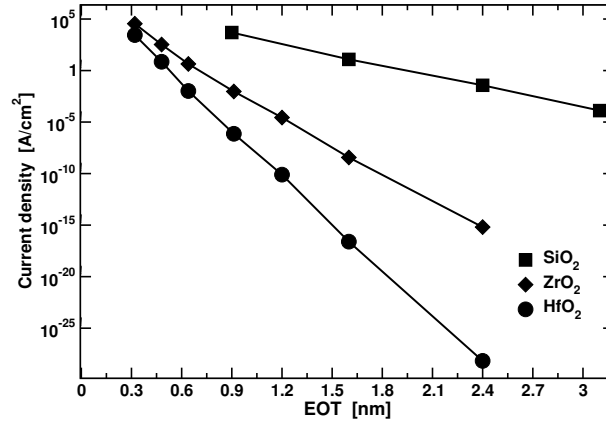


Figure 5.7: TB tunneling currents through ZrO_2 and HfO_2 against the EOT, compared to a SiO_2 -based MOS.

The multiscale simulation of the MOSFET shown in Fig. 5.5, coupling microscopic gate tunneling with semi-classical drift-diffusion, is done in the following way. First, the Poisson equation is solved on the whole device for a gate bias such that the transistor is in accumulation. Then the tunneling current density is calculated quantum-mechanically. For this we assume that the tunneling current density locally depends only on the barrier height and on the local electro-chemical potential

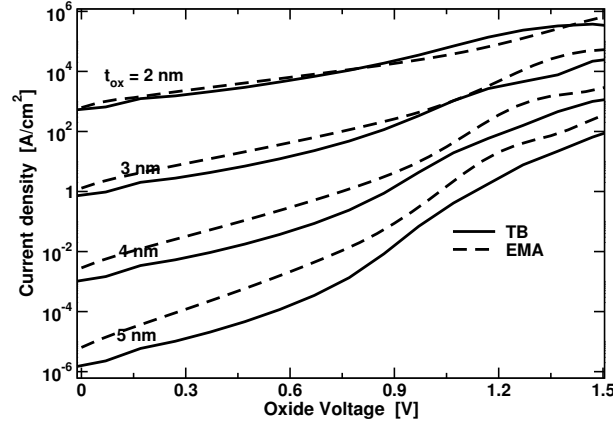


Figure 5.8: Tunneling currents through ZrO_2 against the applied potential.

at the semiconductor-oxide interface. This allows to calculate the tunneling current independently on slices perpendicular to the channel using the 1D approach as described before. The current density obtained from this is used as a boundary condition for the normal electron flux at the semiconductor-oxide interface during the drift-diffusion calculation. The procedure is iterated until convergence is reached. The charge density in the gate oxide resulting from the electrons tunneling through the gate is neglected in the calculation of the electric potential.

In Figs. 5.9 and 5.10 we show the subthreshold transfer characteristics for a source-drain voltage of 0.1 V and 1 V calculated with and without including the gate oxide tunneling current. Results for ZrO_2 are compared with HfO_2 and SiO_2 with the same EOT. As it is expected from Fig. 5.7, the SiO_2 gate shows significant tunneling with respect to the HfO_2 gate, and the ZrO_2 tunneling current is visible only at small source-drain voltage. From Fig. 5.9 follows, that in the case of a source drain voltage of 0.1 V the effect of tunneling becomes visible for SiO_2 at a gate voltage of -0.6 V and for ZrO_2 at -1.7 V. At such voltages the drain current for SiO_2 and ZrO_2 transistors becomes mostly controlled by the gate tunneling. For a higher drain-source bias of 1 V, which corresponds to the saturation of the channel current in the inversion regime, both ZrO_2 and HfO_2 characteristics are not influenced significantly by the tunneling effect, however the SiO_2 drain current begins to saturate at a gate voltage of -1 V due to tunneling.

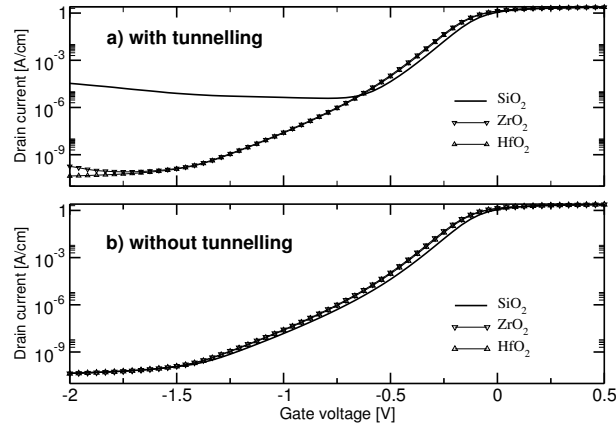


Figure 5.9: Drain current vs. gate voltage at 0.1 V drain-source bias.

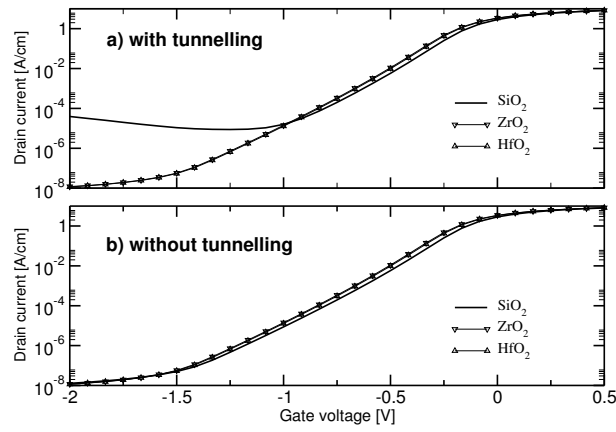


Figure 5.10: Drain current vs. gate voltage at 1 V drain-source bias.

5.3 GaAs-based pin-diodes for polariton LASER

Structures based on GaN and GaAs are currently investigated for use in polariton-Lasers. Polaritons are quasi-particles formed by the coupling of the electromagnetic field with excitons. As such they are composite bosons. To achieve polariton formation one has to create an exciton population inside an optical cavity such that they can couple to the electromagnetic field. Electrical injection of excitons into an $\text{In}_{0.05}\text{Ga}_{0.95}\text{As}$ quantum well embedded in the intrinsic region of a GaAs p-i-n photodiode has been demonstrated in 2007 [11]. This is an important step towards practical exploitation of the strong coupling for polariton-based optoelectronic devices.

Whereas in GaN strong coupling has been observed at room temperature [97], in GaAs based structures excitons are formed only at low temperatures. Excitonic emission was observed up to 70 K [11].

A GaAs based structure grown and characterised at the LPN-CNRS has been simulated at different temperatures using TIBERCAD. Fig. 5.11 shows the device structure.

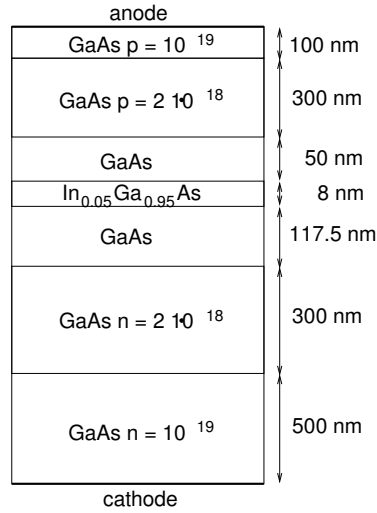


Figure 5.11: Schematic drawing of simulated device.

In the real device, square mesas of $300 \times 300 \mu\text{m}^2$ are formed by optical lithography and chemical etching and contact metallizations are applied on the substrate back side and on the top of the mesa. As the experiments do not show in-plane variations of emission, only 1D simulations along growth direction were performed.

The following equations are considered in the simulation:

$$-\nabla(\epsilon \nabla \varphi) = -e(n - p - N_d^+ + N_a^-) \quad (5.2a)$$

$$-\nabla(\mu_n n \nabla \phi_n) = R \quad (5.2b)$$

$$-\nabla(\mu_p p \nabla \phi_p) = -R \quad (5.2c)$$

$$-\nabla(\mu_x x \nabla \phi_x) = -R_x, \quad (5.2d)$$

where x and ϕ_x are the exciton density and an “effective” exciton potential defined in an analogous way as the electro-chemical potential for electrons and holes (cf. sections 2.2.1.2 and 2.2.1.3). The net recombination rates are modeled as follows, considering SRH recombination, direct (radiative) e-h recombination, exciton formation, radiative exciton recombination and exciton dissociation:

$$R = R_{SRH} + B(np - n_i^2) + \gamma np - x/\tau_{diss} \quad (5.3)$$

$$R_x = -\gamma np + x/\tau_{diss} + x/\tau_{rad} \quad (5.4)$$

In the bulk GaAs, excitons are not expected to be found, so eq. (5.2d) can be restricted to the quantum well. Furthermore, the electron-hole system is assumed to be in chemical equilibrium with the exciton gas. This latter assumption leads from statistical arguments to a law of mass action of the form

$$np = n^* x, \quad (5.5)$$

where n^* can be calculated from the statistics of the involved particles [113].

Assuming an approximately constant exciton density and exciton formation rate in the quantum well, we can neglect the exciton flux and (5.2d) leads to

$$x/\tau_{rad} = \gamma np - x/\tau_{diss}. \quad (5.6)$$

Together with (5.5) we can now eliminate the exciton density in the electron and hole continuity equations and the respective recombination rates reduce to

$$R = R_{SRH} + B(np - n_i^2) + \frac{np}{n^* \tau_{rad}} \quad (5.7)$$

For a 2D system, $n_{2D}^* = \frac{k_B T \mu}{2\pi \hbar^2} e^{-E_b/k_B T}$, where μ and E_b are the reduced exciton mass and the exciton binding energy, respectively. As the simulation does not consider 2D gases, $n^* = n_{3D}^*$ for the simulation is estimated to be $n_{3D}^* = n_{2D}^*/w$, where w is the quantum well width.

Figure 5.12 shows the simulated and measured IV characteristics. In fig. 5.13 the dependence of the estimated 2D electron, hole and exciton density on the current density at 50 K is shown. The 2D electron and hole densities are calculated by numerically integrating the 3D densities over the quantum well. The exciton density is then calculated from (5.5).

Two curves based on a different way of calculating the exciton density are also shown in the figure. It uses charge conservation and assumes that all injected

carriers recombine radiatively in the quantum well, which leads to $n_x = J\tau/e$. This formula gives very similar results at low temperatures and high current densities, but as it neglects the other recombination terms, it overestimates the exciton density especially at low currents. Subtracting the contribution from SRH-recombination, one gets a similar result to the one using (5.5).

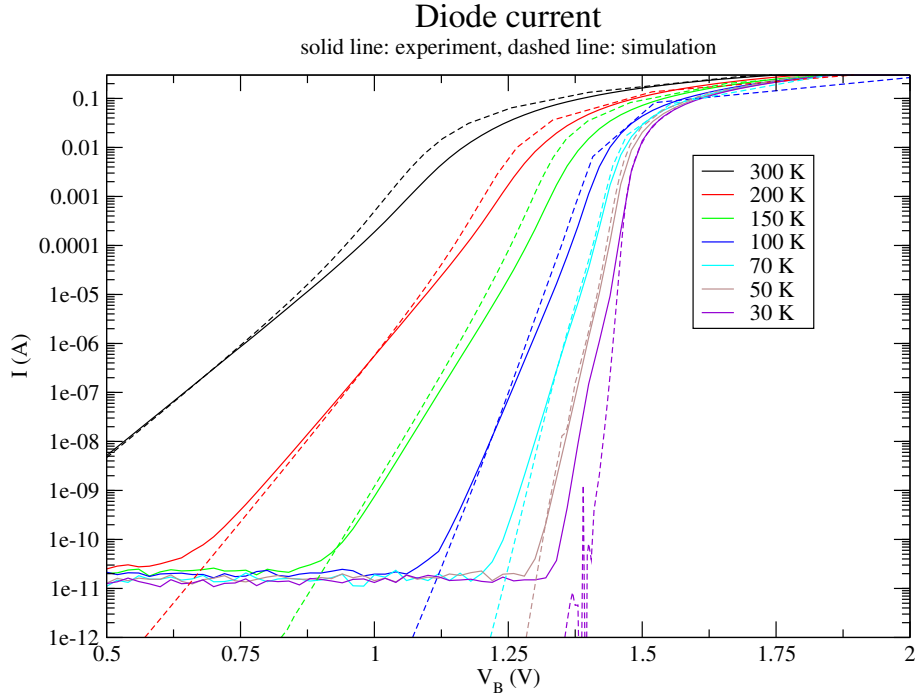


Figure 5.12: Simulated and measured IV characteristics at different temperatures.

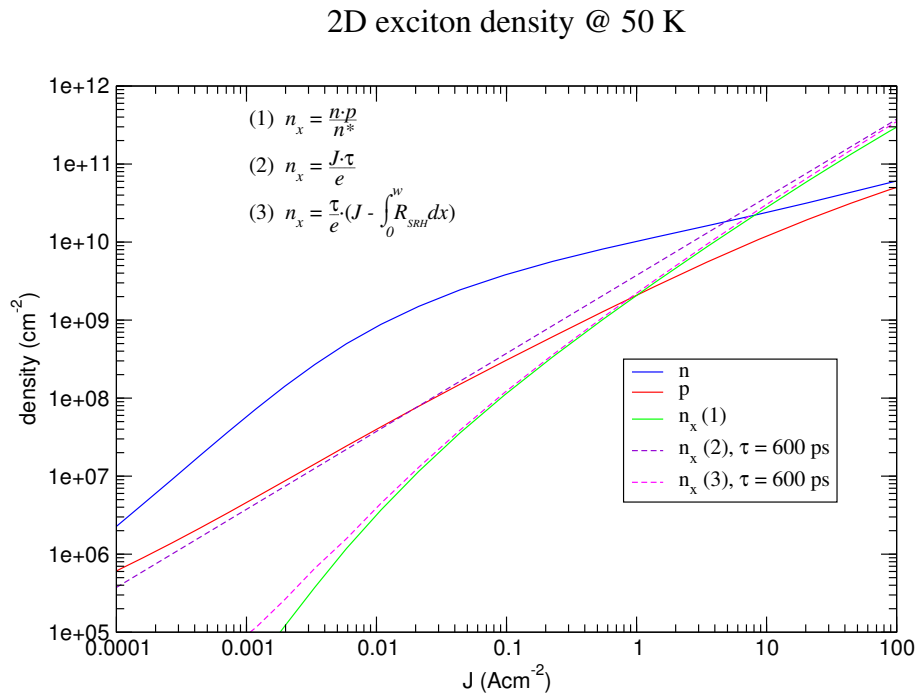


Figure 5.13: Electron, hole and exciton density versus current density at 50 K.

5.4 Structures for polariton LASERs and LEDs based on GaN

GaN is a very promising material for the realisation of exciton based optoelectronic devices and is therefore studied extensively [97, 28]. Due to a big exciton binding energy R of more than 20 meV exciton formation can be expected even at room temperature, in contrast to structures based on GaAs (see last section) which have to be cooled down to below 100 K.

Different structures proposed for GaN polariton lasers have been simulated with TIBERCAD, the simplest of which shall be presented here [81]. The simulated structure is schematically drawn in Fig. 5.14. An $\text{In}_{0.05}\text{Ga}_{0.95}\text{N}$ quantum well embed-

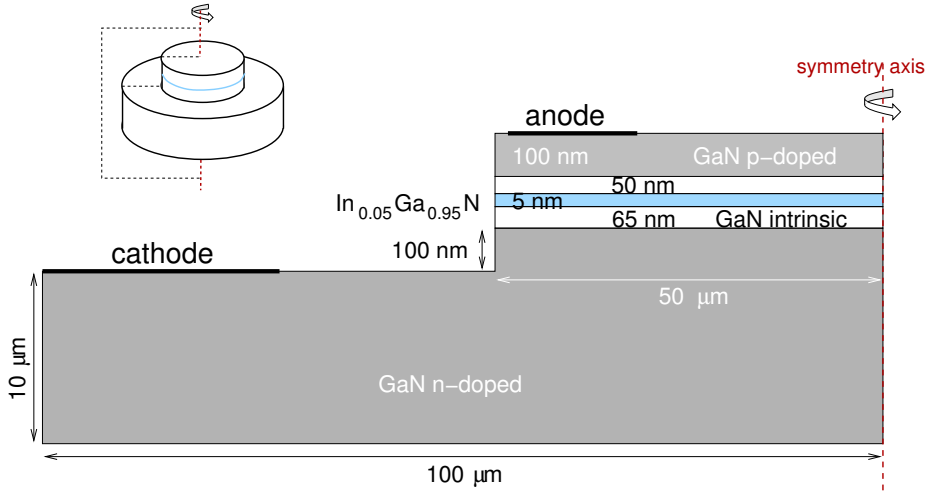


Figure 5.14: The structure of the simulated device.

ded in a 120 nm thick intrinsic GaN layer is grown on a n-doped GaN substrate ($N_d = 5 \cdot 10^{18} \text{ cm}^{-3}$), followed by a p-doped cap layer ($N_a = 5 \cdot 10^{18} \text{ cm}^{-3}$). We assume that a circular mesa with radius $R = 50 \mu\text{m}$ is defined by etching such as to build a VCSEL-like structure with anode on top and cathode on the side. The symmetry of the structure is considered in the simulation by formulating the problem in cylinder coordinates and assuming a solution without angular dependence (cf. Appendix B).

In this example we do not assume chemical equilibrium between the electron/hole and exciton gases, and we cannot neglect the exciton flux as the structure is inhomogeneous in x -direction. Exciton formation is expected to occur mainly in the InGa N quantum well under the anode, creating a laterally inhomogeneous exciton population and thus leading to an exciton drift towards the center of the mesa. For this reason the drift-diffusion equation for the excitons is solved explicitly, in contrast to the example in the last section, leading to a system of four equations (cf. eqns. (5.2)). The simulation has been performed self-consistently, using a mod-

ified Broyden algorithm [49]. For the excitons we assumed a binding energy of $R = 20.4$ meV, a mobility of $\mu_x = 1500$ cm²V⁻¹s⁻¹, an exciton generation rate parameter $\gamma = 2 \cdot 10^{-7}$ cm³s⁻¹, an exciton dissociation time of $\tau_{diss} = 6 \cdot 10^{-9}$ s, a non-radiative exciton recombination time of $\tau_{nr} = 1.2 \cdot 10^{-9}$ s and a radiative exciton recombination time of $\tau_{nr} = 1 \cdot 10^{-11}$ s in the undoped materials.

A voltage sweep has been done from 0 V to 3.5 V, in each step doing a self-consistent calculation of the coupled drift-diffusion/exciton system. Fig. 5.15 shows the resulting IV characteristics, comparing simulations with and without exciton formation.

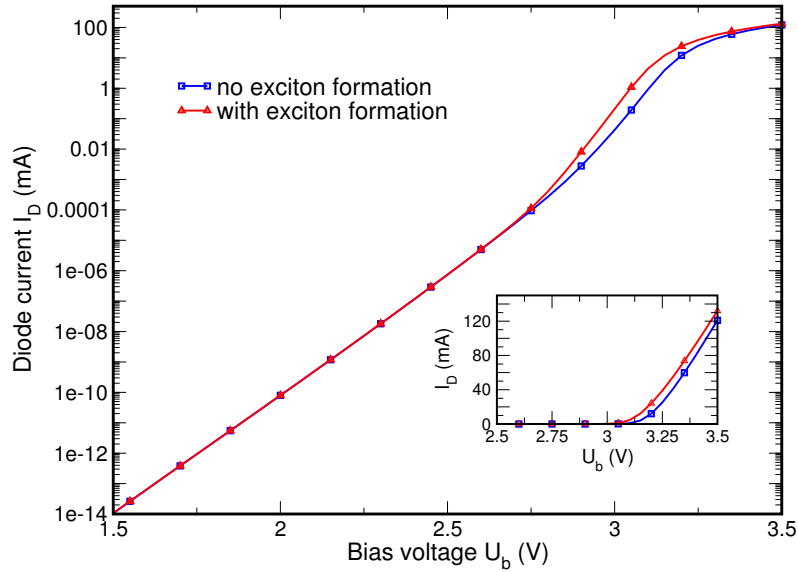


Figure 5.15: IV characteristics of the pin diode with and without considering exciton formation.

In Fig. 5.19 on p. 97 we show the band profiles and the electron and hole densities in the structure at a bias voltage of 3.35 V. We note that the highest carrier densities are found in the quantum well below the anode such that also the excitons will accumulate in the same region, which is an undesired effect as light generation due to radiative exciton recombination mostly occurs in the periphery of the mesa instead of in the center. The effect of diffusion is too weak to effectively populate the quantum well in the center of the device for the studied radius of the mesa.

Fig. 5.16 presents the electron, hole and exciton densities in the mesa for the same bias voltage. It can easily be seen that exciton generation happens most effectively below the anode, i.e. in regions of high current density. Therefore one would like to confine the current to the center of the device. In Fig. 5.17 we show a cut along the quantum well. We may notice that the excitons generated below the anode tend to diffuse towards the center of the mesa. A cut along the vertical direction y is shown in Fig. 5.18.

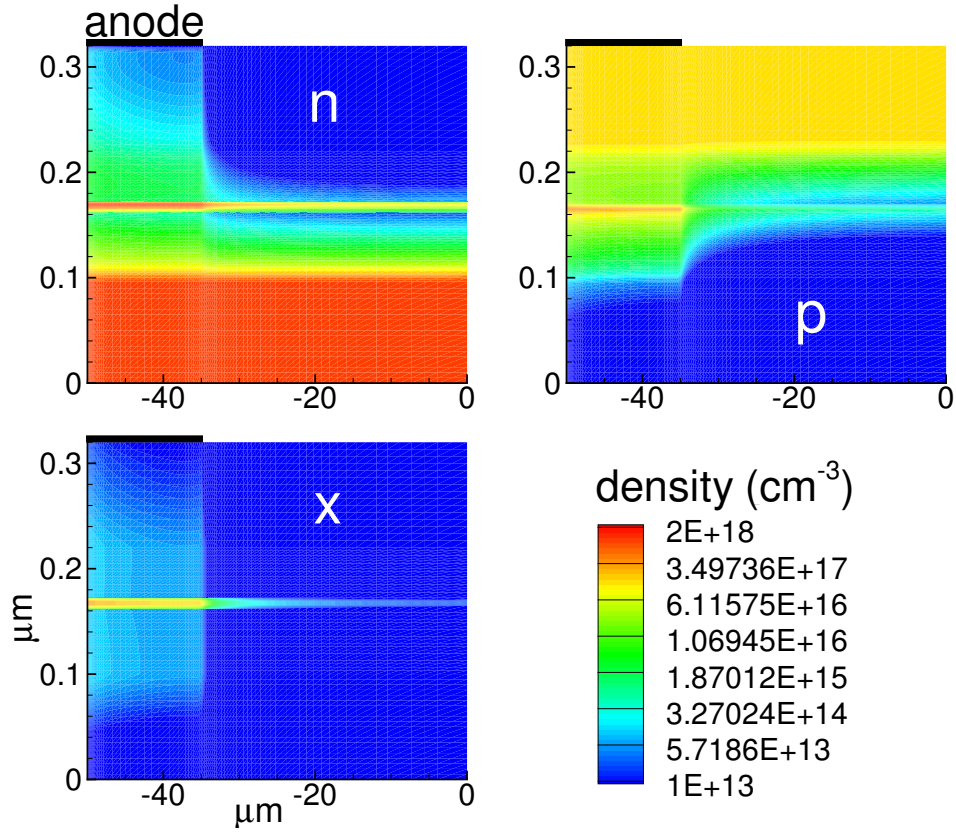


Figure 5.16: Electron, hole and exciton densities in the mesa.

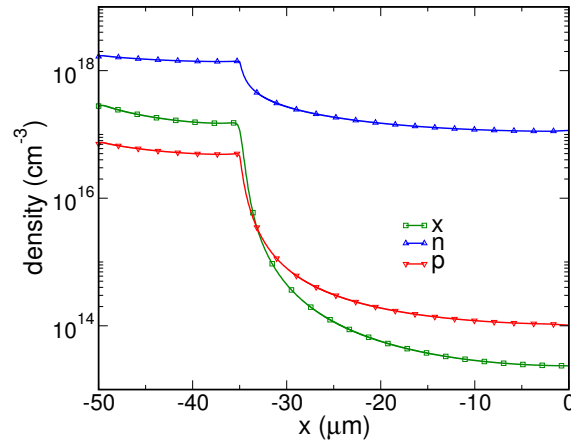
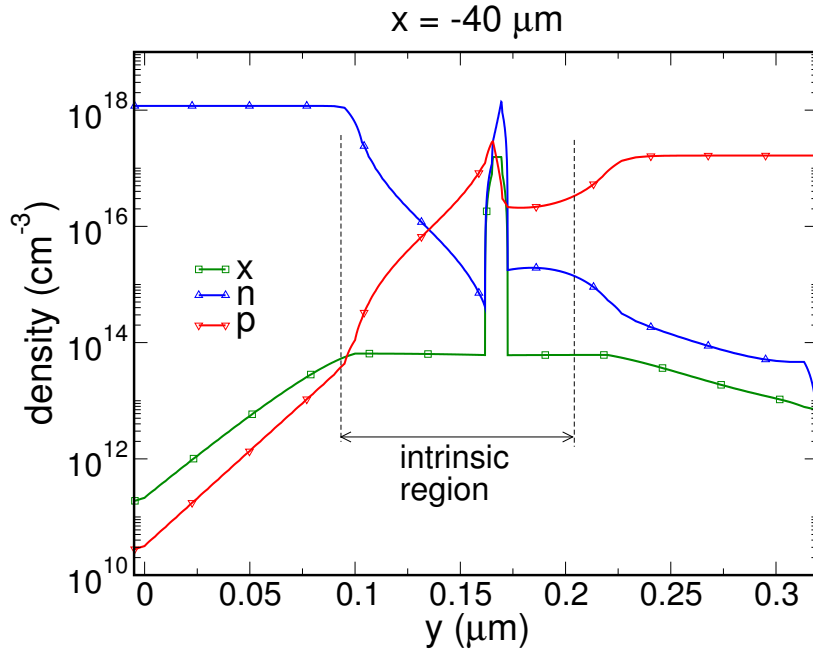
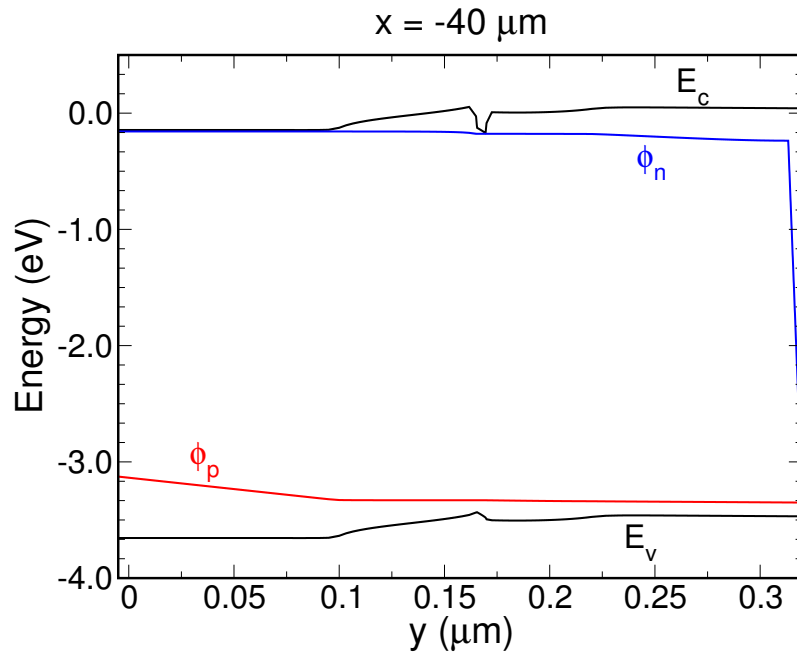


Figure 5.17: Electron, hole and exciton densities along a cutline in x -direction in the quantum well. Not that the data are extracted at the y -position of maximum electron density, i.e. towards the upper boundary of the quantum well, whereas the holes are confined towards the lower boundary (cf. Fig. 5.18).



(a) densities



(b) band diagram

Figure 5.18: Electron, hole and exciton densities along a cutline in y -direction at $x = -40 \mu\text{m}$.

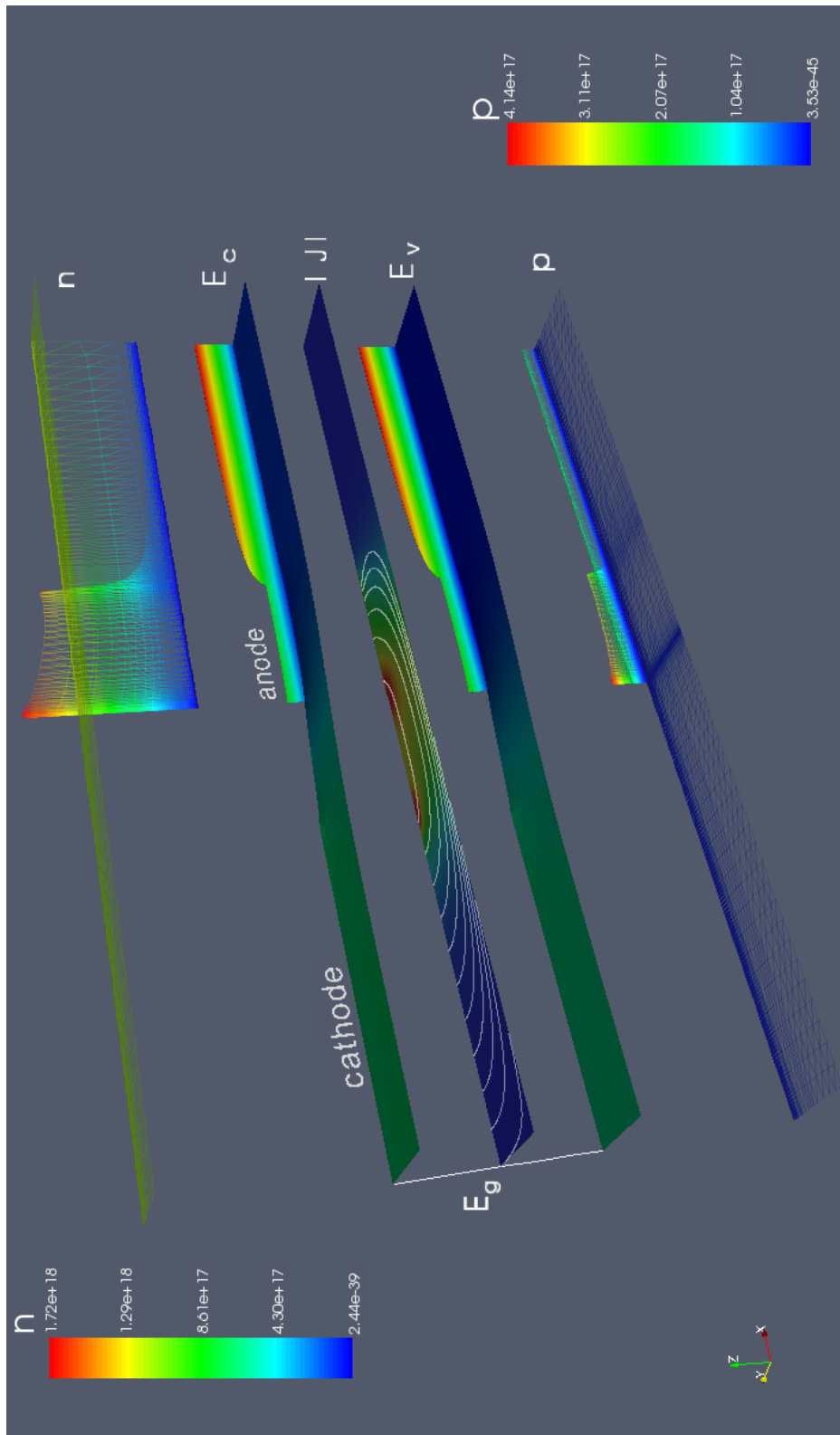


Figure 5.19: Plots of the band edges, electron and hole densities and of the total current density, including current flow lines. Simulation voltage is 3.35 V, corresponding to a diode current of approximately 60 mA.

5.5 AlGaAs/GaAs/AlGaAs Quantum well

In this last example we present a simple self-consistent calculation of a 20 nm GaAs quantum well in an $\text{Al}_{0.3}\text{Ga}_{0.7}\text{As}$ n-i-n structure (cf. Fig. 5.20). The quantum well

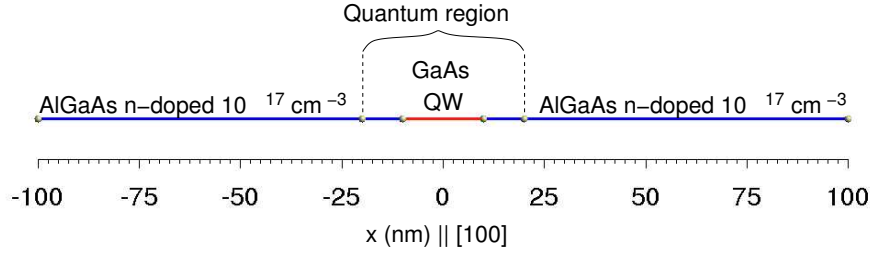


Figure 5.20: Schematical drawing of the simulated quantum well structure.

width is 20 nm, and the AlGaAs barriers are n-doped with $N_d = 10^{17} \text{ cm}^{-3}$. Only a part of the structure including the QW is considered for the quantum mechanical calculation. Outside a classical density is assumed. Strain is shown in Fig. 5.21. Due to the larger lattice constant of AlGaAs with respect to GaAs, the latter gets expansively strained in the plane perpendicular to the growth direction and thus compressively strained in growth direction. As the structure is grown along [100], all off-diagonal strain components vanish.

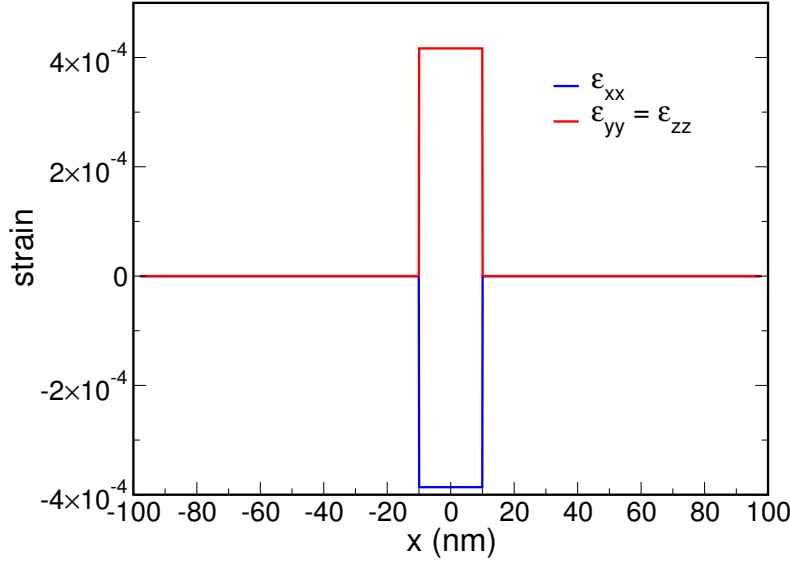
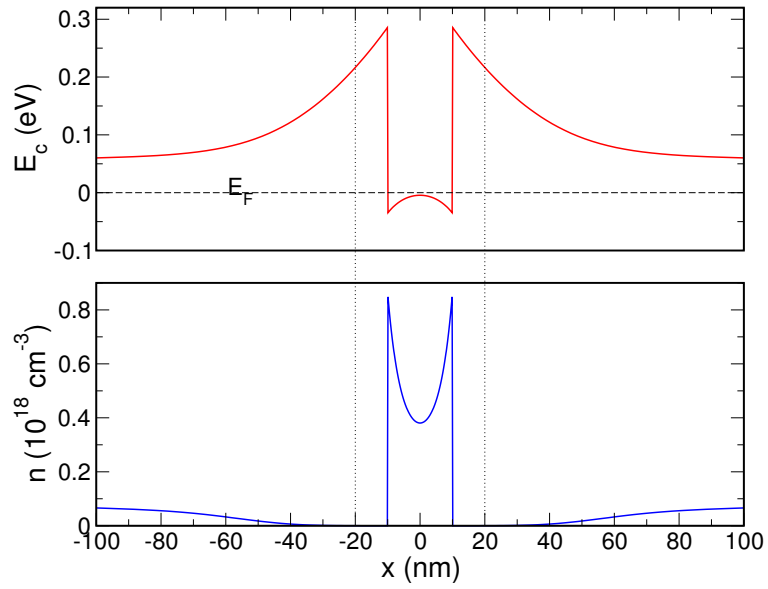
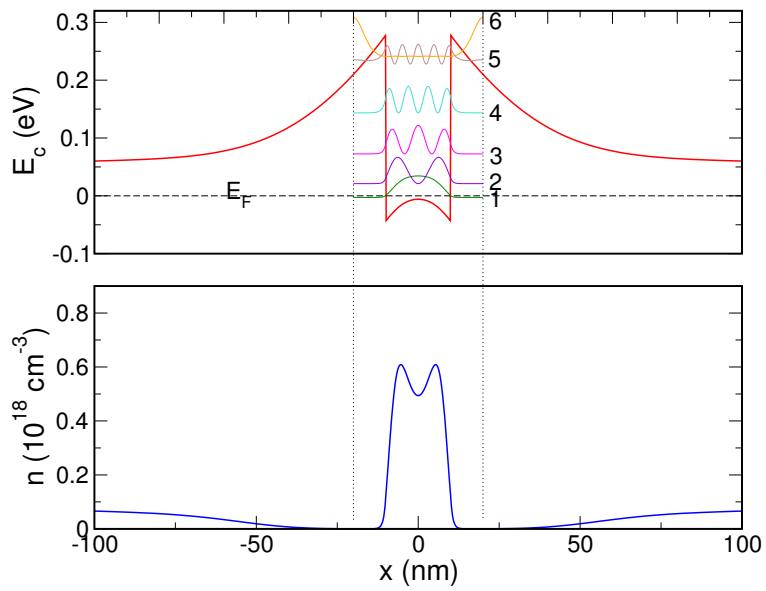


Figure 5.21: Non-zero strain components.

Fig. 5.22 shows a comparison between the classically calculated conduction band and electron density and the result of the self-consistent Schrödinger-Poisson calculus.



(a) Classical results



(b) Self-consistent results

Figure 5.22: Equilibrium results of classical and self-consistent Schrödinger-Poisson calculations.

Appendix A

Numerical Evaluation of Terminal Currents

The total current flowing out of the l -th contact of a device can be expressed (in the stationary case) as

$$I_l = \int_{\Gamma_l} (\mathbf{J}_n + \mathbf{J}_p) \cdot \mathbf{n} \, dA, \quad (\text{A.1})$$

where Γ_l , \mathbf{n} and dA are the boundary part, the outer normal vector and the surface element of the l -th contact. The direct evaluation of this integral however can lead to big numerical errors.

A more accurate evaluation of the terminal currents transforms the surface integral into a volume integral over the device volume Ω by means of so called *Ramo-Shockley test functions (RSTF)* [111, 112]. The RSTF h_l for the l -th contact has the following property

$$h_l|_{\Gamma_m} = \delta_{lm}, \quad h_l \in H^1(\Omega), \quad (\text{A.2})$$

i.e. it is 1 on the l -th contact and 0 on all the other contacts and belongs to the Hilbert space $H^1(\Omega)$ [111].

Using h_l as test functions in the continuity equation for the total current density

$$\nabla \cdot (\mathbf{J}_n + \mathbf{J}_p) = 0, \quad (\text{A.3})$$

one can write equation (A.3) in weak form

$$\int_{\Omega} h_l \nabla \cdot (\mathbf{J}_n + \mathbf{J}_p) \, dV = 0 \quad (\text{A.4})$$

Integrating by parts leads to

$$0 = \underbrace{\int_{\partial\Omega} h_l (\mathbf{J}_n + \mathbf{J}_p) \cdot \mathbf{n} \, dA}_{=I_l} - \int_{\Omega} \nabla h_l \cdot (\mathbf{J}_n + \mathbf{J}_p) \, dV \quad (\text{A.5})$$

Due to the special properties of the RSTF, the boundary integral in equation (A.5) is exactly the current flowing out from the l -th contact, so one gets

$$I_l = \int_{\Omega} \nabla h_l (\mathbf{J}_n + \mathbf{J}_p) \, dV \quad (\text{A.6})$$

The RSTFs can be chosen arbitrarily as long as they have the properties given in (A.2), but their choice will have an impact on the precision of the calculation. When the Poisson and continuity equations are discretized using the finite element method (FEM) it is natural to construct them from the finite element basis:

$$h_l = \sum_i \alpha_i \psi_i, \quad \psi_i \in V_N \quad (\text{A.7})$$

where $V_N \subset H^1(\Omega)$ is a space of continuous, piecewise linear functions. The coefficients are chosen in TIBERCAD to be

$$\alpha_i = \begin{cases} 1 & \text{if } x_i \in \Gamma_l \\ 0 & \text{else} \end{cases}$$

With this choice and using $\mathbf{J} = -e(\mu_n n \nabla \phi_n + \mu_p p \nabla \phi_p)$, equation (A.6) can be written as

$$\begin{aligned} I_l &= -e \sum_{i \in \Gamma_l} \int_{\Omega} (\mu_n n \nabla \phi_n + \mu_p p \nabla \phi_p) \, dV \\ &= -e \sum_{i \in \Gamma_l} \sum_k \int_{\Omega} (\mu_n n v_k + \mu_p p w_k) \nabla \psi_i \nabla \psi_k \, dV \end{aligned} \quad (\text{A.8})$$

For the last expression we used the expansion of the electro-chemical potentials in terms of the finite element basis functions. The numerical evaluation can be implemented quite easily using an element wise approach.

For device simulations one usually asks for overall current conservation, i.e. $\sum_l I_l = 0$. To prove that in our method current is conserved we observe from equation (A.5) that

$$\int_{\Omega} \nabla s (\mathbf{J}_n + \mathbf{J}_p) \, dV = 0, \quad \forall s \in U_N, \quad U_N = \left\{ s \mid s \in V_N, s|_{\Gamma_{\text{Dirichlet}}} = 0 \right\}.$$

We can define a function $p := 1 - \sum_l h_l$ which is clearly an element of U_N , as it vanishes on all Dirichlet boundaries. So we can write (defining the total current density $\mathbf{J} = \mathbf{J}_n + \mathbf{J}_p$)

$$\begin{aligned} 0 &= \int_{\Omega} \nabla p \mathbf{J} \, dV = \int_{\Omega} \nabla \left(1 - \sum_l h_l \right) \mathbf{J} \, dV \\ &= - \sum_l \int_{\Omega} \nabla h_l \mathbf{J} \, dV = - \sum_l I_l, \end{aligned} \quad (\text{A.9})$$

which proves that total current is conserved.

A proper choice of the test functions can lead to improved accuracy. In particular, the RSTF can be chosen to satisfy the poisson equation $\nabla(\epsilon \nabla h_l) = 0$, which would follow quite intuitively from the original formulation of Ramo and Shockley [111]. It remains to note that this method can be extended also for time dependent calculations.

Appendix B

Implementation of Cylindrical Symmetry

The computing time needed for a electronic device simulation depends very much on the dimension of the problem. 2D and 3D simulations need much more nodes than calculations in 1D and 2D, respectively, and they also produce matrices with much bigger bandwidth. Therefore one tries to reduce the dimensionality for the simulation by using symmetries of the physical device.

Usually the symmetry is not exact and does neglect e.g. boundary effects. MOSFETs for example are often simulated in only 2D assuming an infinitely wide gate. In other cases the fabrication process induces a spatial symmetry as is the case for VCSELs and other vertical structures with cylindrically shaped mesa. Provided that the physical properties of the constituent materials have the same symmetry, one can assume that the quantities describing the device behaviour have no angular dependence.

Whereas the implementation of specular symmetry is straightforward (it is mainly a question of boundary conditions as no coordinate transformation is involved), the case of cylinder symmetry is slightly more complicated. The transformation between cartesian and cylinder coordinates can be stated as

$$\begin{aligned}x &= \rho \cos(\varphi) \\y &= \rho \sin(\varphi) \\z &= \zeta\end{aligned}\tag{B.1}$$

with jacobian

$$\mathbf{J} = \begin{pmatrix} \frac{\partial x}{\partial \rho} & \frac{\partial x}{\partial \varphi} & \frac{\partial x}{\partial \zeta} \\ \frac{\partial y}{\partial \rho} & \frac{\partial y}{\partial \varphi} & \frac{\partial y}{\partial \zeta} \\ \frac{\partial z}{\partial \rho} & \frac{\partial z}{\partial \varphi} & \frac{\partial z}{\partial \zeta} \end{pmatrix} = \begin{pmatrix} \cos(\varphi) & -\rho \sin(\varphi) & 0 \\ \sin(\varphi) & \rho \cos(\varphi) & 0 \\ 0 & 0 & 1 \end{pmatrix}, |\mathbf{J}| = \rho \tag{B.2}$$

The Laplace operator $\Delta = \sum_i \frac{\partial^2}{\partial x_i^2}$ in cylinder coordinates reads

$$\Delta = \frac{1}{\rho} \frac{\partial}{\partial \rho} \left(\rho \frac{\partial}{\partial \rho} \right) + \frac{1}{\rho^2} \frac{\partial^2}{\partial \varphi^2} + \frac{\partial^2}{\partial z^2} \tag{B.3}$$

The implementation is simplified a lot because in the finite element method (and also in the box integration method) only first order partial derivatives appear in the equations (after applying the Gauss theorem).

Writing $x_i := x, y, z$ $\underline{x} = (x, y, z)$ and $\xi_i := \xi, \eta, \zeta$ $\underline{\xi} = (\xi, \eta, \zeta)$, the first order partial derivatives can be written as

$$\frac{\partial}{\partial x_i} = \frac{\partial \xi_k}{\partial x_i} \frac{\partial}{\partial \xi_k} \quad (\text{B.4})$$

or

$$\nabla_{\underline{x}} = (\mathbf{J}^{-1})^T \nabla_{\underline{\xi}} \quad (\text{B.5})$$

Assuming no angular dependence, i.e. $f(\rho, \varphi, z) = f(\rho, z)$, one can choose $\varphi = 0$ such that ρ and x coincide. In this way the gradient is formally invariant under the coordinate transformation:

$$\nabla_{\underline{x}} = \left(\frac{\partial}{\partial \rho}, 0, \frac{\partial}{\partial z} \right) \quad (\text{B.6})$$

and one can map (ρ, z) to (x, y) .

As a consequence, for the implementation of cylindrical symmetry only the proper volume elements dV need to be considered. The poisson equation $\nabla \epsilon \nabla \psi = f$ in weak form for example transforms in the following way:

$$- \iiint \epsilon \nabla \psi \nabla \nu \, dx \, dy \, dz = \iiint f \nu \, dx \, dy \, dz \quad (\text{B.7})$$

\Downarrow

$$- \iiint \epsilon \nabla \psi \nabla \nu |\mathbf{J}| \, d\rho \, d\varphi \, dz = \iiint f \nu |\mathbf{J}| \, d\rho \, d\varphi \, dz \quad (\text{B.8})$$

\Downarrow

$$-2\pi \iint \epsilon \nabla \psi \nabla \nu \rho \, d\rho \, dz = 2\pi \iint f \nu \rho \, d\rho \, dz \quad (\text{B.9})$$

Using the mapping $(\rho, z) \rightarrow (x, y)$ the last equation can be implemented very easily using an existent 2D implementation in cartesian coordinates.

Appendix C

A comparison principle for quasi-linear elliptic equations

We consider the quasi-linear equation in divergence form

$$Q(u, \eta) \doteq \int_{\Omega} [\mathbf{A}(x, u, \nabla u) \cdot \nabla \eta - B(x, u)\eta] dx = 0 \quad (\text{C.1})$$

where $\mathbf{A}(x, u, q)$ and $B(x, u)$ (writing $q \equiv \nabla u$) are measurable and continuously differentiable with respect to u, q in $\bar{\Omega} \times \mathbb{R} \times \mathbb{R}^n$, Q is uniformly elliptic and $u \in H^1(\bar{\Omega})$, $\eta \in H_0^1(\bar{\Omega})$, i.e.

$$\begin{aligned} &\exists \lambda(u), \Lambda(u) > 0 \text{ such that} \\ &\lambda(u) |\xi|^2 \leq \sum_{i,j} \frac{\partial A_i}{\partial q_j} \xi_i \xi_j \leq \Lambda(u) |\xi|^2 \quad \forall u, \xi \end{aligned} \quad (\text{C.2})$$

Eq. (C.1) corresponds to a boundary value problem

$$\begin{aligned} Qu &= -\nabla \mathbf{A}(x, u, \nabla u) - B(x, u) = 0 \text{ in } \Omega \\ u &= g \text{ on } \partial\Omega \end{aligned} \quad (\text{C.3})$$

We can formulate the following *comparison principle* for the problem (C.3) or its weak formulation (C.1) (see [40, 21])

Theorem C.1 (comparison principle). *Let $u, v \in H^1(\bar{\Omega})$ satisfy $Qu \leq Qv$ in Ω and $u \leq v$ on $\partial\Omega$, then if B is non-increasing in u for fixed x it follows that $u \leq v$ in Ω .*

For the proof of the theorem we use the following definition:

Definition C.1 (sub- and super-solution). *A function u is called a (weak) sub-solution of (C.3) if $Qu \leq 0$ ($Q(u, \eta) \leq 0 \forall \eta$). A function v is called a (weak) super-solution of (C.3) if $Qv \geq 0$ ($Q(v, \eta) \geq 0 \forall \eta$).*

Proof. Let $u, v \in H^1(\bar{\Omega})$ be two weak solutions of the above boundary value problem such that $Qu \leq Qv$ in a weak sense, i.e. $Q(u, \eta) \leq Q(v, \eta)$. This is especially the case if u, v are sub- and super-solutions. We define

$$\begin{aligned} w &= u - v \\ u_t &= v + t(u - v), \quad 0 \leq t \leq 1 \end{aligned} \quad (\text{C.4})$$

We will consider in the following the difference $Q(u, \nabla u) - Q(v, \nabla v) \leq 0$ and write it in terms of the difference $w = u - v$.

First we remember that for a function $g(u)$ that is continuous in its argument we can write

$$g(u) - g(v) = \int_v^u \frac{dg(\zeta)}{d\zeta} d\zeta = \int_0^1 \frac{dg(u_t)}{du_t} \frac{du_t}{dt} dt = \int_0^1 \frac{dg(u_t)}{du_t} dt \cdot w \quad (\text{C.5})$$

Then we calculate $Q(u, \nabla u) - Q(v, \nabla v)$:

$$\begin{aligned} Q(u, \nabla u) - Q(v, \nabla v) &= \int_{\Omega} \{ [A(x, u, \nabla u) - A(x, v, \nabla v)] \cdot \nabla \eta \\ &\quad - [B(x, u) - B(x, v)] \eta \} dx \end{aligned} \quad (\text{C.6})$$

As A and B are differentiable in u and ∇u we can apply (C.5) so that we get (in index notation, writing $\partial_i u = (\nabla u)_i$)

$$\begin{aligned} A_i(x, u, \nabla u) - A_i(x, v, \nabla v) &= \\ &= \underbrace{\int_0^1 \frac{\partial A_i(x, u_t, \nabla u_t)}{\partial u_t} dt}_{b_i} \cdot w + \underbrace{\int_0^1 \frac{\partial A_i(x, u_t, \nabla u_t)}{\partial (\partial_j u_t)} dt}_{a_{ij}} \cdot \partial_j w \end{aligned} \quad (\text{C.7})$$

$$\begin{aligned} B(x, u) - B(x, v) &= \\ &= \underbrace{\int_0^1 \frac{\partial B(x, u_t, \nabla u_t)}{\partial u_t} dt}_{d} \cdot w \end{aligned}$$

Using these expressions we can finally write eq. (C.6) in terms of w :

$$L(w, \eta) = \int_{\Omega} [(a_{ij} \partial_j w + b_i w) \partial_i \eta - d w \eta] dx \leq 0 \quad (\text{C.8})$$

In other words, if u, v such that $Qu \leq Qv$ and $u \leq v$ on $\partial\Omega$, then $w = u - v$ is a (weak) sub-solution of

$$Lw = -\partial_i (a_{ij} \partial_j w + b_i w) - dw = 0 \text{ in } \Omega, \quad w \leq 0 \text{ on } \partial\Omega$$

Due to (C.2) the above equation is uniformly elliptic and $d \leq 0$ as B is supposed to be non-increasing in u . Hence the linear equation for w can be treated by the theory for linear elliptic equations [40], leading to a maximum principle. Therefore, if w is a solution of (C.8) with $w \leq 0$ on $\partial\Omega$, then $w \leq 0$ in Ω by virtue of the maximum principle and consequently $u \leq v$ in Ω . \blacksquare

To apply Theorem C.1 to the drift-diffusion equations it has to be shown that $\mathbf{A}(x, u, \nabla u)$ and $B(x, u)$ are continuously differentiable with respect to $u, \nabla u$ and that B is non-increasing in u . The former can be seen easily using the explicit formulas for \mathbf{A} as stated in (3.25) and noting that n and p depend continuously on their arguments. In the following we show the monotonicity of B with respect to u , assuming for the recombination models a form $R = R(x, n, p) = g(x, n, p)(np - n_i^2)$ with $g \geq 0$. We make the following

Proposition C.1. *All the functions $B(x, u)$ are strictly decreasing in u , i.e.*

$$\begin{aligned} \text{(i)} \quad & \frac{\partial B^\varphi}{\partial \varphi} = \frac{\partial \rho}{\partial \varphi} < 0 \\ \text{(ii)} \quad & \frac{\partial B^{\phi_n}}{\partial \phi_n} = \frac{\partial R}{\partial \phi_n} < 0 \\ \text{(iii)} \quad & \frac{\partial B^{\phi_p}}{\partial \phi_p} = -\frac{\partial R}{\partial \phi_p} < 0 \end{aligned}$$

- (i) Using the expressions for the densities (2.59) we can immediately deduce that n is strictly increasing and p strictly decreasing in φ , i.e. $\partial_\varphi n > 0$ and $\partial_\varphi p < 0$. From (3.32) we get (with $N_d = g_d = N_a = g_a = k_B T = 1$ for convenience)

$$\begin{aligned} \partial_\varphi N_d^+ &= -(N_d^+)^2 \exp(\varphi - \phi_n - E_d) < 0 \\ \partial_\varphi N_a^- &= (N_a^-)^2 \exp(\phi_p - \varphi + E_a) > 0 \end{aligned}$$

Thus,

$$\partial_\varphi \rho = \partial_\varphi B^\varphi = \partial_\varphi (p - n + N_d^+ - N_a^-) < 0$$

- (ii) We use $\partial_{\phi_n} R = \partial_n R \partial_{\phi_n} n$ and we allow for the following recombination mechanisms:

- $\gamma(np - n_i^2)$, e.g. radiative recombination
- $(np - n_i^2)/[\tau_p(n + n_i) + \tau_n(p + n_i)]$, Shockley-Read-Hall (SRH)
- $(C_n n + C_p p)(np - n_i^2)$, Auger

Furthermore we use $\partial_{\phi_n} n < 0$.

For the first case it is easily obtained that $\partial_{\phi_n} R = p \partial_{\phi_n} n < 0$.

For the SRH model we get, abbreviating $D = \tau_p(n + n_i) + \tau_n(p + n_i)$,

$$\begin{aligned} \partial_n R &= \frac{\partial}{\partial n} \frac{np - n_i^2}{\tau_p(n + n_i) + \tau_n(p + n_i)} = \frac{p}{D} + \frac{-\tau_p(np - n_i^2)}{D^2} \\ &= \frac{p[\tau_p(n + n_i) + \tau_n(p + n_i)] - \tau_p(np - n_i^2)}{D^2} \\ &= \frac{pn_i(\tau_n + \tau_p) + p^2\tau_n + \tau_p n_i^2}{D^2} > 0 \end{aligned}$$

and thus $\partial_{\phi_n} R < 0$.

In the case of Auger recombination we get

$$\begin{aligned}\partial_n R &= \frac{\partial}{\partial n} [(C_n n + C_p p)(np - n_i^2)] = C_n(np - n_i^2) + (C_n n + C_p p)p \\ &= 2C_n np + C_p p^2 - C_n n_i^2 > -C_n n_i^2\end{aligned}$$

This quantity is therefore not always positive and strictly speaking does not satisfy the the Proposition [C.1](#). Auger recombination however is only important in the high-injection regime, and the coefficients $C_{n,p}$ are very small ($\approx 10^{-31}$ in silicon [\[94\]](#)). Therefore in almost all cases the contribution from SRH will compensate the above small negative term and ensure the monotonicity of the total recombination rate.

(iii) follows in the same way as (ii)

We repeat once again that even if the single equations of the drift-diffusion system satisfy a comparison principle, this does by no means imply such a principle to hold for the system of equations.

Bibliography

- [1] R. A. Adams and J. J.F. Fournier, *Sobolev Spaces*, 2nd ed., Pure and Applied Mathematics, vol. 140, Academic Press, 2003.
- [2] F. Amato, C. Cosentino, S. Pricl, M Ferrone, M. Fermeglia, M. M.-C. Cheng, R. Walczak, and M. Ferrari, *Multiscale modeling of protein transport in silicon membrane nanochannels. Part 2. From molecular parameters to a predictive continuum diffusion model*, Biomed Microdevices **8** (2006), 291–298.
- [3] O. Ambacher, J. Smart, J. R. Shealy, N. G. Weimann, K. Chu, M. Murphy, W. J. Schaff, L. F. Eastman, R. Dimitrov, L. Wittmer, M. Stutzmann, W. Rieger, and J. Hilsenbeck, *Two-dimensional electron gases induced by spontaneous and piezoelectric polarization charges in N- and Ga-face Al-GaN/GaN hetrostructures*, Journal of Applied Physics **85** (1999), 3222.
- [4] U. Ascher, P. A. Markowich, C. Schmeiser, H. Steinrück, and R. Weiss, *Conditioning of the steady state semiconductor problem*, SIAM J. Appl. Math. **49** (1989), no. 1, 165–185.
- [5] Neil W. Ashcroft and N. David Mermin, *Solid State Physics*, Thomson Learning, 1976.
- [6] M. Auf der Maur, M. Povolotskyi, F. Sacconi, and A. Di Carlo, *Simulation of piezoresistivity effect in FETs*, J. Comput. Electron. **5** (2006), 323–326.
- [7] ———, *TIBERCAD: A new multiscale simulator for electronic and optoelectronic devices*, Superlattices and Microstructures **41** (2007), 381–385.
- [8] M. Auf der Maur, M. Povolotskyi, F. Sacconi, A. Pecchia, and A. Di Carlo, *Multiscale Simulation of MOS Systems based on High- κ Oxides*, J. Comput. Electron. (2007).
- [9] M. Auf der Maur, M. Povolotskyi, F. Sacconi, G. Romano, E. Petrolati, and A. Di Carlo, *Multiscale Simulation of Electronic and Optoelectronic Devices with TIBERCAD*, Simulation of Semiconductor Processes and Devices, 2007, pp. 245–248.

- [10] Valery Axelrad, *Grid quality and its influence on accuracy and convergence in device simulation.*, IEEE Transactions on Computer Aided Design of Integrated Circuits and Systems **17** (1998), no. 2, 149–157.
- [11] D. Bajoni, A. Miard, A. Lemaître, S. Bouchoule, J. Bloch, and J. Tignon, *Nonresonant Electrical Injection of Excitons in an InGaAs Quantum Well*, Applied Physics Letters **90** (2007), no. 12, 211114.
- [12] Satish Balay, Kris Buschelman, Victor Eijkhout, William D. Gropp, Dinesh Kaushik, Matthew G. Knepley, Lois Curfman McInnes, Barry F. Smith, and Hong Zhang, *PETSc users manual*, Tech. Report ANL-95/11 - Revision 2.1.5, Argonne National Laboratory, 2004.
- [13] Satish Balay, Kris Buschelman, William D. Gropp, Dinesh Kaushik, Matthew G. Knepley, Lois Curfman McInnes, Barry F. Smith, and Hong Zhang, *PETSc Web page*, 2001, <http://www.mcs.anl.gov/petsc>.
- [14] Satish Balay, William D. Gropp, Lois Curfman McInnes, and Barry F. Smith, *Efficient management of parallelism in object oriented numerical software libraries*, Modern Software Tools in Scientific Computing (E. Arge, A. M. Bruaset, and H. P. Langtangen, eds.), Birkhäuser Press, 1997, pp. 163–202.
- [15] R. E. Bank and D. J. Rose, *Global approximate Newton methods*, Numerische Mathematik **37** (1981), no. 2, 279–295.
- [16] R. E. Bank, D. J. Rose, and W. Fichtner, *Numerical methods for semiconductor device simulation*, IEEE Transactions on Electron Devices **30** (1983), no. 9, 1031–1041.
- [17] R. E. Bank and L. R. Scott, *On the conditioning of finite element equations with highly refined meshes*, SIAM J. Numer. Anal. **26** (1989), no. 6, 1383–1394.
- [18] C. W. J. Beenakker and H. van Houten, *Semiconductor heterostructures and nanostructures*, Solid State Physics: Advances in Research and Applications, vol. 44, ch. Quantum Transport in Semiconductor Nanostructures, Academic Press, 1991.
- [19] J. S. Blakemore, *Semiconductor Statistics*, Dover Publications, New York, 1987.
- [20] K. Bløtekjær, *Transport Equations for Electrons in Two-Valley Semiconductors*, IEEE Transactions on Electron Devices **17** (1970), no. 1, 38–47.
- [21] G. Boyadzhiev, *Comparison principle for non - cooperative elliptic systems*, ArXiv Mathematics e-prints (2007).

- [22] F. Brezzi, L. D. Marini, S. Micheletti, P. Pietra, R. Sacco, and S. Wang, *Finite element and finite volume discretizations of Drift-Diffusion type fluid models for semiconductors*, Handbook of Numerical Analysis (W. H. A. Schilders and E. J. W. Maten, eds.), vol. XIII: Numerical Methods in Electromagnetics, 2005.
- [23] F. Brezzi, L. D. Marini, and P. Pietra, *Two-dimensional exponential fitting and applications to drift-diffusion models*, SIAM J. Numer. Anal. **26** (1989), no. 6, 1342–1355.
- [24] M. G. Burt, *The justification for applying the effective-mass approximation to microstructures*, Journal of Physics: Condensed Matter **4** (1992), no. 32, 6651–6690.
- [25] ———, *Direct derivation of effective-mass equations for microstructures with atomically abrupt boundaries*, Physical Review B **50** (1994), no. 11.
- [26] A. Di Carlo, *Microscopic theory of nanostructured semiconductor devices: beyond the envelope-function approximation*, Semiconductor Science and Technology **18** (2003), 1.
- [27] Z. Chen, *Finite Element Methods and Their Applications*, Scientific Computation, Springer, 2005.
- [28] T.-Y. Chung and K. J. Chang, *Exciton binding energies in GaN/AlGaIn quantum-well structures*, Semicond. Sci. Technol. **13** (1998), 876–881.
- [29] P. G. Ciarlet, *The Finite Element Method for Elliptic Problems*, Classics in Applied Mathematics, vol. 40, SIAM, 1978.
- [30] A. K. Cline, C. B. Moler, G. W. Stewart, and J. H. Wilkinson, *An estimate for the condition number of a matrix*, SIAM J. Numer. Anal. **16** (1979), no. 2, 368–375.
- [31] S. Datta, *Electronic Transport in Mesoscopic Systems*, Cambridge University Press, 1995.
- [32] S. P. Edwards and K. De Meyer, *A charge damping algorithm applied to a newton solver for solving SOI devices and other ill-conditioned problems*, Numerical Analysis of Semiconductor Devices and Integrated Circuits, 1987. NASCODE V. Proceedings of the Fifth International Conference on the (1987), 174–186.
- [33] T. Eickhoff, O. Ambacher, G. Krötz, and M. Stutzmann, *Piezoresistivity of $\text{Al}_x\text{Ga}_{1-x}\text{N}$ layers and $\text{Al}_x\text{Ga}_{1-x}\text{N}/\text{GaN}$ heterostructures*, Journal of Applied Physics **90** (2001), 3383.

- [34] P. Enders, A. Bärwolf, M. Woerner, and D. Suisky, *k·p theory of energy bands, wave functions, and optical selection rules in strained tetrahedral semiconductors*, Physical Review B **51** (1995), no. 23.
- [35] A. Ern and J.-L. Guermond, *Theory and Practice of Finite Elements*, Applied Mathematical Sciences, vol. 159, Springer, 2004.
- [36] W. Fichtner, D. J. Rose, and R. E. Bank, *Semiconductor device simulation*, IEEE Transactions on Electron Devices **30** (1983), no. 9, 1018–1030.
- [37] L. P. Franca and A. Russo, *Deriving upwinding, mass lumping and selective reduced integration by residual-free bubbles*, Applied Mathematics Letters **9** (1996), no. 5, 83–88.
- [38] Antonio J García-Loureiro and Juan M López-González, *Electron quasi-Fermi level splitting at the base-emitter junction of HBTs and DHBTs*, Semiconductor Science and Technology **19** (2004), no. 3, 552–557.
- [39] C. Geuzaine, *Gmsh 3d finite element mesh generator*, <http://www.geuz.org/gmsh>.
- [40] D. Gilbarg and N. S. Trudinger, *Elliptic Partial Differential Equations of Second Order*, Classics Mathematics, Springer, 1998.
- [41] Z. Xiao Guo (ed.), *Multiscale materials modelling: Fundamentals and applications*, Woodhead Publishing, 2007.
- [42] Stefan Hackenbuchner, *Elektronische Struktur von Halbleiter-Nanobauelementen im thermodynamischen Nichtgleichgewicht*, Ph.D. thesis, Walter Schottky Institut, Technische Universität München, 2002.
- [43] Y. He and G. Cao, *A Generalized Scharfetter-Gummel Method to Eliminate Crosswind Effects*, IEEE Transactions on Computer-Aided Design **10** (1991), no. 12, 1579–1582.
- [44] G. Heiser, C. Pommerell, J. Weis, and W. Fichtner, *Three-Dimensional Numerical Semiconductor Device Simulation: Algorithms, Architectures, Results*, IEEE Transactions on Computer-Aided Design **10** (1991), 1218.
- [45] V. Hernandez, J. E. Roman, and V. Vidal, *SLEPc Web page*, <http://www.grycap.upv.es/slepc>.
- [46] International Technology Roadmap for Semiconductors (ITRS), 2005, <http://www.itrs.net>.
- [47] C. Jacoboni and P. Lugli, *The Monte Carlo Method for Semiconductor Device Simulation*, 1st ed., Springer-Verlag Wien New York, 1989.

- [48] J.-M. Jancu, R. Scholz, F. Beltram, and F. Bassani, *Empirical spds* tight-binding calculation for cubic semiconductors; General method and material parameters*, Physical Review B **57** (1998), no. 11.
- [49] D. D. Johnson, *Modified Broyden's method for accelerating convergence in self-consistent calculations*, Physical Review B **38** (1988), no. 18.
- [50] Leo P. Kadanoff and Gordon Baym, *Quantum Statistical Mechanics*, Perseus Books, Cambridge, 1962.
- [51] D. Kahng and M. M. Atalla, *Silicon-Silicon Dioxide Field Induced Surface Devices*, IRE Solid-State Device Research Conference (1960).
- [52] A. Kawamoto, J. Jameson, K. Cho, and R. W. Dutton, *Challenges for atomic scale modeling in alternative gate stack engineering*, IEEE Transactions on Electron Devices **47** (2000), no. 10, 1787–1794.
- [53] P. N. Keating, *Effect of invariance requirements on the elastic strain energy of crystals with application to the diamond structure*, Phys. Rev. **145** (1966), no. 2, 637–645.
- [54] T. Kerkhoven, *A proof of convergence of Gummel's algorithm for realistic device geometries*, SIAM J. Numer. Anal. **23** (1986), no. 6, 1121–1137.
- [55] B. S. Kirk and J. W. Peterson, *the libMesh library*, <http://libmesh.sourceforge.net>.
- [56] Charles Kittel, *Introduction to Solid State Physics*, 7th ed., John Wiley & Sons, New York, 1953.
- [57] L. D. Landau and E. M. Lifschitz, *Lehrbuch der theoretischen Physik: Statistische Physik, Teil 1*, 8th ed., Akademie Verlag, Berlin, 1978.
- [58] ———, *Lehrbuch der theoretischen Physik: Quantenmechanik*, 9th ed., Akademie Verlag, Berlin, 1979.
- [59] ———, *Lehrbuch der theoretischen Physik: Elektrodynamik der Kontinua*, 5th ed., Akademie Verlag, Berlin, 1985.
- [60] ———, *Lehrbuch der theoretischen Physik: Elastizitätstheorie*, 7th ed., Akademie Verlag, Berlin, 1987.
- [61] ———, *Lehrbuch der theoretischen Physik: Statistische Physik, Teil 2*, 4th ed., Akademie Verlag, Berlin, 1992.
- [62] Juin J. Liou and Frank Schwierz, *RF MOSFETS: recent advances, current status and future trends*, Solid State Electronics **47** (2003), 1881–1895.

- [63] P. O. Löwdin, *On the Non-Orthogonality Problem Connected with the Use of Atomic Wave Functions in the Theory of Molecules and Crystals*, J. Chem. Phys. **18** (1950), 365.
- [64] J. M. Luttinger and W. Kohn, *Motion of Electrons and Holes in Perturbed Periodic Fields*, Physical Review **97** (1955), no. 4, 869–883.
- [65] Peter A. Markowich, *The Stationary Semiconductor Device Equations*, 1st ed., Springer-Verlag Wien New York, 1986.
- [66] Peter A. Markowich, Christian A. Ringhofer, and Christian Schmeiser, *Semiconductor Equations*, 1st ed., Springer-Verlag Wien New York, 1990.
- [67] Richard M. Martin, *Electronic Structure: Basic Theory and Practical Methods*, Cambridge University Press, 2004.
- [68] C. McCarthy and G. Strang, *Optimal conditioning of matrices*, SIAM J. Numer. Anal. **10** (1973), no. 2, 370–388.
- [69] C. D. Meyer, *Matrix Analysis and Applied Linear Algebra*, SIAM, 2000.
- [70] S. Micheletti, *Stabilized finite elements for semiconductor device simulation*, Comput. Visual Sci. **3** (2001), 177–183.
- [71] J. J. H. Miller (ed.), *Finite Element Programming with Special Emphasis on Semiconductor Device and Process Modelling*, Lecture Notes of a Short Course held in Galway on 13th and 14th June, 1983 in association with the NASECODEII Conference, Boole Press, Dublin, 1983.
- [72] J. J. H. Miller, W. H. A. Schilders, and S. Wang, *Application of finite element methods to the simulation of semiconductor devices*, Rep. Prog. Phys. **62** (1999), 277–353.
- [73] K. D. Mish, L. R. Herrmann, and L. Haws, *Finite Element Procedures in Applied Mechanics, lecture notes*, 2000.
- [74] M. S. Mock, *Analysis of Mathematical Models of Semiconductor Devices*, Boole Press, 1983.
- [75] Gordon E. Moore, *Cramming more components onto integrated circuits*, Electronics **38** (1965), no. 8, 114–117.
- [76] K. W. Morton and D. F. Mayers, *Numerical Solution of Partial Differential Equations*, 2nd ed., Cambridge University Press, 1985.
- [77] L. Onsager, *Reciprocal relations in irreversible processes*, Physical Review **37** (1931), 405–426.
- [78] A. Pecchia, *personal communication*.

- [79] A. Pecchia and A. Di Carlo, *Atomistic theory of transport in organic and inorganic nanostructures*, Rep. Prog. Phys. **67** (2004), 1497–1561.
- [80] S. S. Perlman and D. L. Feucht, *n-p heterojunctions*, Solid-State Electronics **7** (1964), 911–923.
- [81] E. Petrolati, M. Auf der Maur, M. Povolotskyi, and A. Di Carlo, *Simulation of Exciton Formation and Transport in Electrically Driven Polariton Laser Structures*, Superlattices and Microstructures **41** (2007), 364–367.
- [82] M. Povolotskyi, *Theoretical Study of Electronic and Optical Properties of Low-Dimensional Semiconductor Nanostructures*, Ph.D. thesis, University of Rome "Tor Vergata", 2004.
- [83] M. Povolotskyi and A. Di Carlo, *Elasticity theory of pseudomorphic heterostructures grown on substrates of arbitrary thickness*, J. Appl. Phys. **100** (2006), 063514.
- [84] M. Povolotskyi, M. Auf der Maur, and A. Di Carlo, *Strain effects in freestanding three-dimensional nitride nanostructures*, Phys. Stat. Sol. (c) **2** (2005), no. 11, 3891–3894.
- [85] The TIBERCAD project, <http://www.tibercad.org>.
- [86] D. Querlioz, J. Saint-Martin, K. Huet, A. Bournel, V. Aubry-Fortuna, C. Chassat, S. Galdin-Retailleau, and P. Dollfus, *On the ability of the particle monte carlo technique to include quantum effects in nano-mosfet simulation*, Electron Devices, IEEE Transactions on **54** (Sept. 2007), no. 9, 2232–2242.
- [87] K. Rim, S. Koester, M. Hargrove, J. Chu, P.M. Mooney, J. Ott, T. Kanarsky, P. Ronsheim, M. Jeong, A. Grill, and H.-S.P. Wong, *Strained Si NMOSFETs for high performance CMOS technology*, 2001 Symposium on VLSI Technology. Digest of Technical Papers (IEEE Cat. No.01 CH37184) (Kyoto, Japan), 2001, pp. 59 – 60.
- [88] C. Ringhofer and C. Schmeiser, *An approximate newton method for the solution of the basic semiconductor device equations*, SIAM J. Numer. Anal. **26** (1989), no. 3, 507–516.
- [89] J. Robertson, *High dielectric constant gate oxides for metal oxide Si transistors*, Rep. Prog. Phys. **69** (2006), 327–396.
- [90] F. Sacconi, A. Di Carlo, P. Lugli, M. Stadele, and Jancu J. M., *Full-band approach to tunneling in MOS structures*, IEEE Transactions on Electron Devices **51** (2004), no. 5, 741–748.
- [91] F. Sacconi, J. M. Jancu, M. Povolotskyi, and A. Di Carlo, *Full-band tunneling in high- κ oxide MOS structures*, IEEE Transactions on Electron Devices **54** (2007), no. 12, 3168–3176.

- [92] J. J. Sakurai, *Modern Quantum Mechanics*, Addison Wesley, 1994.
- [93] D. L. Scharfetter and H. K. Gummel, *Large-Signal Analysis of a Silicon Read Diode Oscillator*, IEEE Transactions on Electron Devices **16** (1969), 64–77.
- [94] A. Schenk, *Halbleiter-Bauelemente: Physikalische Grundlagen und simulation, lecture notes, ETH Zurich*, 2002.
- [95] O. Schenk, S. Rölli, and A. Gupta, *The Effects of Unsymmetric Matrix Permutations and Scalings in Semiconductor Device and Circuit Simulation*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **23** (2004), no. 3, 400–410.
- [96] Siegfried Selberherr, *Analysis and Simulation of Semiconductor Devices*, 1st ed., Springer-Verlag Wien New York, 1984.
- [97] F. Semon, I. R. Sellers, F. Natali, D. Byrne, M. Leroux, J. Massies, N. Ollier, J. Leymarie, P. Disseix, and A. Vasson, *Strong light-matter coupling at room temperature in simple geometry GaN microcavities grown on silicon*, Applied Physics Letters **87** (2005), 1102–+.
- [98] Silvaco TCAD Solutions, <http://www.silvaco.com>.
- [99] M. Y. Simmons, A. C. Churchill, G. H. Kim, A. R. Hamilton, A. Kurobe, D. R. Mace, D. A. Ritchie, and M. Pepper, *Growth of high mobility heterostructures on (311)B GaAs*, Microelectronics Journal **26** (1995), 897.
- [100] D. L. Smith and C. Mailhot, *Theory of semiconductor superlattice electronic structure*, Reviews of Modern Physics **62** (1990), no. 1, 173–234.
- [101] Synopsis Inc., <http://www.synopsis.com>.
- [102] G.-L. Tan, X.-L. Yuan, Q.-M. Zhang, W. H. Ku, and A.-N. Shey, *Two-Dimensional Semiconductor Device Analysis Based on New Finite-Element Discretization Employing the S-G Scheme*, IEEE Transactions on Computer-Aided Design **8** (1989), no. 5.
- [103] T.-W. Tang and M.-K. Jeong, *A Generalized Scharfetter-Gummel Method to Eliminate Crosswind Effects*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **14** (1995), no. 11, 1301–1315.
- [104] A. Trellakis, A. T. Galick, A. Pancelli, and U. Ravaioli, *Iteration scheme for the solution of the two-dimensional Schrödinger-Poisson equations in quantum structures*, J. Appl. Phys. **81** (1997), 7880.
- [105] H. Vande Sande, H. De Gersem, F. Henrotte, and K. Hameyer, *Solving nonlinear magnetic problems using newton trust region methods*, IEEE Transactions on Magnetics **39** (2003), no. 3, 1709 – 1712.

-
- [106] I. Vurgaftman, J. R. Meyer, and L. R. Ram-Mohan, *Band parameters for III-V compound semiconductors and their alloys*, Journal of Applied Physics **89** (2001), 5815.
 - [107] I. Vurgaftman and J.R. Meyer, *Band parameters for nitrogen-containing semiconductors*, Journal of Applied Physics **94** (2003), 3675.
 - [108] G. K. Wachutka, *Rigorous Thermodynamic Treatment of Heat Generation and Conduction in Semiconductor Device Modeling*, IEEE Transactions on Computer-Aided Design **11** (1990), 1141–1149.
 - [109] S. Wang, *A new exponentially fitted triangular finite element method for the continuity equations in the drift-diffusion model of semiconductor devices*, Mathematical Modelling and Numerical Analysis **33** (1999), no. 1, 99–112.
 - [110] B. M. Wolbert, G. K. Wachutka, B. H. Krabbenborg, and T. J. Mouthaan, *Nonisothermal Device Simulation Using the 2-D Numerical Process/Device Simulator TRENDY and Application to SOI-devices*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **13** (1994), no. 3, 293–302.
 - [111] P. D. Yoder, K. Gärtner, and W. Fichtner, *A generalized Ramo-Shockley theorem for classical to quantum transport at arbitrary frequencies*, Journal of Applied Physics **79** (1996), 1951–1954.
 - [112] P. D. Yoder, K. Gärtner, U. Krumbein, and W. Fichtner, *Optimized Terminal Current Calculation for Monte Carlo Device Simulation*, IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **16** (1997), 1082–1087.
 - [113] H. W. Yoon, D. R. Wake, and J. P. Wolfe, *Effect of exciton-carrier thermodynamics on the GaAs quantum well photoluminescence*, Physical Review B **54** (1996), 2763–2774.
 - [114] O. C. Zienkiewicz and R. L. Taylor, *The Finite Element Method*, 4th ed., vol. 1, McGraw-Hill, 1994.

List of Figures

1.1	Moore's law	1
1.2	The scales in device simulation	4
1.3	Schematic flow chart for a multiscale simulation.	5
1.4	Data interpolation between different meshes.	6
2.1	Iterative procedure for the calculation of the deformed equilibrium shape.	12
2.2	Deformed free standing GaN/AlGa _N heterostructure.	12
2.3	Pyropolarization in a hexagonal lattice.	13
2.4	Electron, hole and exciton dispersion in a semiconductor.	23
2.5	Band structure of GaAs	37
2.6	A three-terminal device	38
2.7	A two-lead device	39
3.1	The reference element for a triangular finite element.	54
3.2	Lagrange element basis functions	62
3.3	A triangular mesh around node i	69
3.4	A simple GaN pn-diode.	71
4.1	Structure of the TIBERCAD software.	73
4.2	Mesh and mesh regions	74
4.3	<code>SimulationInterface</code> hierarchy.	75
4.4	<code>PhysicalModelInterface</code> hierarchy.	76
4.5	<code>BoundaryProperties</code> hierarchy.	77
4.6	pn-heterojunction diode for Listing 4.1.	78
5.1	Schematic drawing of the simulated heterostructures.	82
5.2	Contour plots of the lateral strain component ε_{xx} in a part (cf. dashed box in Fig. 5.1(a)) of the AlGa _N /Ga _N FET without (a) and with (b) pressure.	83
5.3	Relative change of resistance as a function of external force.	83
5.4	Electron density in a part (cf. dashed box in Fig. 5.1(b)) of the B-face GaAlAs/InGaAs/GaAs FET without (a) and with (b) an external pressure of 50 mN/cm.	84

5.5	Schematic drawing of simulated device.	85
5.6	Microscopic model for the Si/high- κ oxide/Si heterostructure.	86
5.7	TB tunneling currents through ZrO_2 and HfO_2 against the EOT, compared to a SiO_2 -based MOS.	86
5.8	Tunneling currents through ZrO_2 against the applied potential.	87
5.9	Drain current vs. gate voltage at 0.1 V drain-source bias.	88
5.10	Drain current vs. gate voltage at 1 V drain-source bias.	88
5.11	Schematic drawing of simulated device.	89
5.12	Simulated and measured IV characteristics at different temperatures.	91
5.13	Electron, hole and exciton density versus current density at 50 K.	92
5.14	The structure of the simulated device.	93
5.15	IV characteristics of the pin diode with and without considering ex- citon formation.	94
5.16	Electron, hole and exciton densities in the mesa.	95
5.17	Electron, hole and exciton densities along a cutline in x -direction in the quantum well.	95
5.18	Electron, hole and exciton densities (1) and band diagram (b) along a cutline in x -direction in the quantum well.	96
5.19	Plots of the band edges, electron and hole densities and of the total current density.	97
5.20	Schematical drawing of the simulated quantum well structure.	98
5.21	Non-zero strain components.	98
5.22	Equilibrium results of classical and self-consistent Schrödinger-Poisson calculations.	99

List of Tables

3.1	The scaling factors	45
3.2	Numerical performance of different diagonal scaling schemes.	71